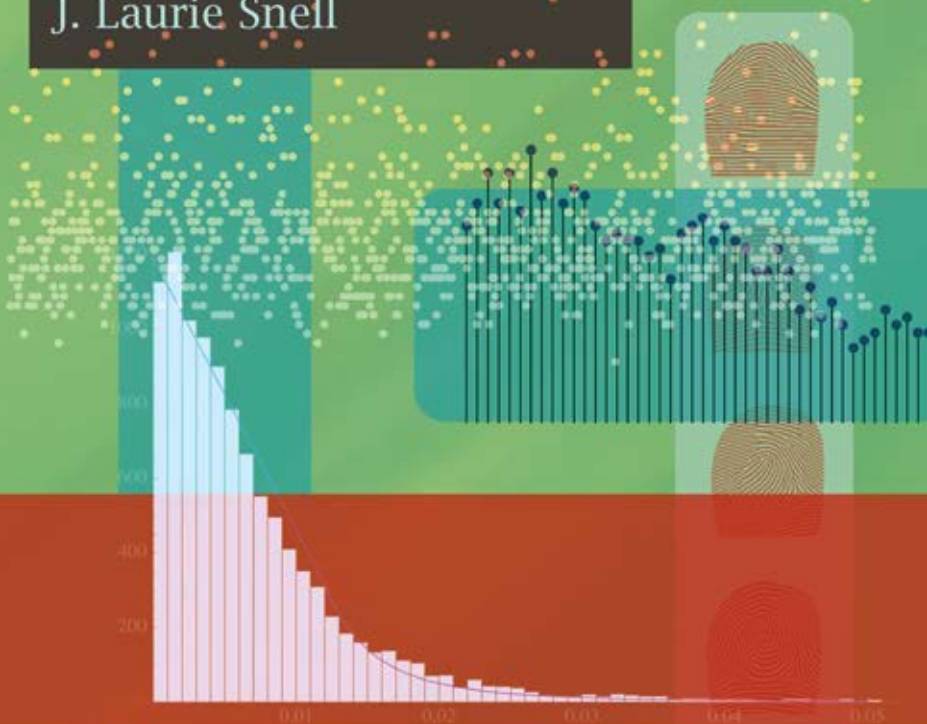


STUDENT MATHEMATICAL LIBRARY
Volume 57

Probability Tales

Charles M. Grinstead
William P. Peterson
J. Laurie Snell



STUDENT MATHEMATICAL LIBRARY
Volume 57

Probability Tales

Charles M. Grinstead
William P. Peterson
J. Laurie Snell



American Mathematical Society
Providence, Rhode Island

Editorial Board

Gerald B. Folland
Robin Forman

Brad G. Osgood (Chair)
John Stillwell

2010 *Mathematics Subject Classification*. Primary 60–01, 62–01.

For additional information and updates on this book, visit
www.ams.org/bookpages/stml-57

Library of Congress Cataloging-in-Publication Data

Grinstead, Charles M. (Charles Miller), 1952–

Probability tales / Charles M. Grinstead, William P. Peterson, J. Laurie Snell.

p. cm. — (Student mathematical library ; v. 57)

Includes bibliographical references and index.

ISBN 978-0-8218-5261-3 (alk. paper)

1. Probabilities. 2. Stochastic processes. I. Peterson, William Paul. II. Snell, J. Laurie (James Laurie), 1925– III. Title.

QA274.G765 2011

519.2—dc22

2010038517

Copying and reprinting. Individual readers of this publication, and nonprofit libraries acting for them, are permitted to make fair use of the material, such as to copy a chapter for use in teaching or research. Permission is granted to quote brief passages from this publication in reviews, provided the customary acknowledgment of the source is given.

Republication, systematic copying, or multiple reproduction of any material in this publication is permitted only under license from the American Mathematical Society. Requests for such permission should be addressed to the Acquisitions Department, American Mathematical Society, 201 Charles Street, Providence, Rhode Island 02904-2294 USA. Requests can also be made by e-mail to reprint-permission@ams.org.

© 2011 by the authors. All rights reserved.

Printed in the United States of America.

∞ The paper used in this book is acid-free and falls within the guidelines established to ensure permanence and durability.

Visit the AMS home page at <http://www.ams.org/>

10 9 8 7 6 5 4 3 2 1 16 15 14 13 12 11

To our wives

Contents

Preface	vii
Chapter 1. Streaks	1
§1. Introduction	1
§2. Models for Repeated Trials	4
§3. Runs in Repeated Trials	6
§4. Statistical Tests on Repeated Trials	9
§5. Data from Various Sports	24
§6. Runs in the Stock Market	71
§7. Appendix	79
Chapter 2. Modeling the Stock Market	97
§1. Stock Prices	97
§2. Variations in the Price of a Stock	103
§3. The Normal Distribution and Power Laws	106
§4. Distribution of Returns	114
§5. Independence of Returns	128
§6. Is the Power Law Exponent Intrinsic?	136
§7. Appendix	139

Chapter 3. Lotteries	149
§1. Rules of the Powerball Lottery	149
§2. Calculating the Probabilities of Winning	152
§3. What is Your Expected Winning for a \$1 Ticket?	157
§4. Does a Ticket's Expected Value Ever Exceed \$1?	165
§5. What Kind of Numbers Do Lottery Buyers Choose?	168
§6. Finding Patterns	173
§7. How Often is the Jackpot Won?	179
§8. Other Lotteries Pose New Questions	181
§9. Using Lottery Stories to Discuss Coincidences	182
§10. Lottery Systems	184
§11. Lottery Stories from Chance News	185
§12. Lottery Questions from John Haigh	189
Chapter 4. Fingerprints	191
§1. Introduction	191
§2. History of Fingerprinting	191
§3. Models of Fingerprints	198
§4. Latent Fingerprints	203
§5. The 50K Study	217
Answers to John Haigh's Lottery Questions	225
Bibliography	231
Index	235

Preface

This book was begun when two of the authors decided that their collaboration on their introductory probability book was both productive and fun. We cast about to find some other project on which we could continue this collaboration. At the same time, all three of the authors were working on Chance News, which is an ongoing website containing reviews of news and journal articles pertaining to probability and statistics in the real world. It was decided that we should try to write a book in which we would take some of our favorite articles in Chance News and revise and expand them.

The result is not exactly what was originally planned. Instead of having many short chapters on a variety of subjects, the present book consists of four chapters in which we go into some depth on the covered topics. In fact, it is this depth that we think makes this book different than other books that cover applications of probability and statistics to the real world.

In the first two chapters, we use some ideas from calculus. Much of the more technical mathematics has been placed in appendices so as not to break the flow of the chapters. We think that even if a reader has not studied calculus, he or she will be able to read and understand most of the material in these chapters. In particular, one

can gain much insight in these chapters from the numerous graphs found therein.

Our first chapter concerns the idea of streaks. Many participants in and observers of sports believe that individuals and teams can be “hot” or “cold.” For example, a basketball player who has made many consecutive field goals is frequently described at the time as having a “hot hand.” It turns out that in many sports, the observed streaks can be shown to fit a very simple coin-tossing model. In other words, coins (fair or otherwise) that are flipped repeatedly exhibit the same kinds of streaks, with the same distributions of lengths of these streaks as those observed in sports. The reader will note that we are not saying that such streaks do not exist, but rather that it is not necessary to posit a model that is any more sophisticated than a simple coin-tossing model to explain these streaks.

The second chapter introduces the reader to some aspects of the U.S. stock market. This is an area in which a vast amount of research has been conducted. We explain how the important class of probability distributions known as power laws can help in understanding the movements of stock prices.

The third chapter is concerned with lotteries. We take the reader through the calculations that are necessary to understand the probabilities of winning various prizes in a typical lottery. We use the Powerball lottery as an example. We also consider the effects that income tax, present value, and possible sharing of the prize have on the value of a lottery jackpot.

The last chapter contains a short history of fingerprinting and discusses some of the problems with the use of fingerprints in fighting crime. These problems are still extant and defy easy solutions.

We think that almost all of the material in this book is accessible to those who have had one semester of calculus (some of the material in the appendices requires some knowledge of power series) and much of it is accessible to all interested readers. We hope that the

material in this book is used to supplement the material in a standard probability or statistics course at the undergraduate level.

We thank Dartmouth, Middlebury, and Swarthmore Colleges for their financial support of this endeavor. We also thank Mike Saitas and Peter Sykes, of the American Mathematical Society, for their editing and design help. Finally, we especially thank Sergei Gelfand, of the American Mathematical Society, for his support and patience.

Chapter 1

Streaks

1. Introduction

Most people who watch or participate in sports think that hot and cold streaks occur. Such streaks may be at the individual or the team level and may occur in the course of one contest or over many consecutive contests. As we will see in the next section, there are different probability models that might explain such observations. Statistics can be used to help us decide which model does the best job of describing (and therefore predicting) the observations.

As an example of a streak, suppose that a professional basketball player has a lifetime free throw percentage of 85%. We assume that over the course of her career, this probability has remained roughly the same. Now suppose that over the course of several games, she makes 20 free throws in a row. Even though she shoots free throws quite well, most sports fans would say that she is “hot.” Most fans would also say that because she is hot, the probability of making her next free throw is higher than her career free throw percentage. Some fans might say that she is “due” for a miss. The most basic question we look at in this chapter is whether the data show that in such situations, a significant number of players make the next shot (or get a hit in baseball, etc.) with a higher probability than might

be expected, given the players' lifetime (or season) percentages. In other words, is the player streaky? One can also ask whether the opposite is true, namely that many players make the next shot with a lower probability than might be expected. We might call such behavior "anti-streaky." Both types of behavior are examples of non-independence between consecutive trials.

An argument in favor of dependence might run something like this. Suppose that the player's shot attempts are modeled by a sequence of Bernoulli trials, i.e. on each shot, she has an 85% chance of making the shot, and this percentage is not affected by the outcomes of her previous shot attempts. In this model, the probability that she makes 20 shots in a row, starting at a particular shot attempt, is $(.85)^{20}$, which is approximately .0388. This is a highly improbable event under the assumption of independence, so this model does not do a good job of explaining the observation.

An argument in favor of the Bernoulli trials model might run as follows. In a sequence of outcomes, a run is a set of consecutive outcomes of the same type that is not contained in any larger such set. So, for example, in the sequence

$$TTFFF TTTTF,$$

the three F 's in positions three through five form a run of length three, but two F 's in positions three and four do not form a run. It can be shown that in a sequence of 200 independent trials, where the probability of a success on a given trial is .85, the average length of the longest run of successes is about 24.0. (We will discuss this calculation a little later in this chapter.) Since many players shoot 200 or more free throws in a given season, it is not surprising that this player has a success run of 20. We will say more below about the length of the longest run in this model.

There are several probability models that might be used to detect streaky behavior. We will consider two of these models in this chapter. The first uses Markov chains. For a thorough introduction to the idea of a Markov chain, the reader should consult [19]. A Markov chain

consists of a set of states, and for each pair of states (including the pairs where the states are the same) there is a certain probability of moving from the first state to the second state.

In the case at hand, there are only two states, success and failure. We define p_1 to be the probability that a success follows a success (i.e. that the chain moves from the success state to the success state) and we define p_2 to be the probability that a success follows a failure (i.e. that the chain moves from the failure state to the success state). If $p_1 > p_2$, then one might expect to see streaky behavior in the model. If $p_1 = p_2$, the model is the same as the Bernoulli model. If $p_1 < p_2$, then one might expect to see “anti-streaky” behavior in the model.

For example, suppose that in the basketball example given above, the player has a 95% chance of making a free throw if she has made the previous free throw. It is possible to show that in order for her overall success rate to be 85%, the probability that she makes a free throw after missing a free throw is 28.3%. It should be clear that in this model, once she makes a free throw, she will usually make quite a few more before she misses, and once she misses, she may miss several in a row. In fact, in this model, once she has made a free throw, the number of successful free throws until her first miss is a geometric random variable, and has expected value 19, so including the first free throw that she made, she will have streaks of made free throws of average length 20. In the Bernoulli trials model, the average length of her streaks will be only 6.7. Since these two average lengths are so different, the data should allow us to say which model is closer to the truth.

A second possible meaning of streakiness, which we call block-Bernoulli, refers to the possibility that in a long sequence of independent trials, the probability of success varies in some way. For example, there may be blocks of consecutive trials, perhaps of varying lengths, such that in each block, the success probability is constant, but the success probabilities in different blocks might be unequal. As an example, suppose the basketball player has a season free throw percentage of 85%, and assume, in addition, that during the season,

there are some blocks of consecutive free throws of length between 20 and 40 in which her probability of success is 95%. This means there must be other blocks in which her success probability is less than 95%. If we compute the observed proportion of successes over all blocks of length 30, say, we should see under these assumptions a wide variation in these proportions. The question we need to answer is how much wider this variation will be than in the Bernoulli model with constant probability of success. The greater the difference between the two variation sizes, the easier it will be to say which model fits the data better.

It is natural to look at success and failure runs under the various models described above, since the ideas of runs and streakiness are closely related. We will describe some statistical tests that have been used in attempts to decide which model most accurately reflects the data. We will then look at data from various sports and the financial markets.

Exercise.

1. Would you be more likely to say that the basketball player in the example given above was “hot” if she made 20 free throws in a row during one game, rather than over a stretch of several games?

2. Models for Repeated Trials

The simplest probability model for a sequence of repeated trials is the Bernoulli trials model. In this model, each trial is assumed to be independent of all of the others, and the probability of a success on any given trial is a constant, usually denoted by p . This means, in particular, that the probability of a success following one success, or even ten successes, is unchanged; it equals p . It is this last statement that causes many people to doubt that this model can explain what actually happens in sports. However, as we shall see in the next section, even in this model, there are hot and cold streaks. The

question is whether the numbers and durations of hot and cold streaks observed in real data exceed the predicted numbers under this model.

This model is a simple one, so one can certainly give reasons why it should not be expected to apply very well in certain situations. For example, in a set of consecutive at-bats in baseball, a batter will face a variety of different pitchers, in a variety of game situations, and at different times of the day. It is reasonable to assert that some of these variables will affect the probability that the batter gets a hit. Some baseball models have been proposed that have many such situational, or explanatory, variables. In other situations, such as free throw shooting in basketball, or horseshoes, the conditions that prevail in any given trial probably do not change very much.

Another relatively simple model is the Markov chain model, described above. It is, of course, possible to define similar, but more complicated, models where the probability of success on a given trial depends upon the outcomes in the preceding k trials, where k is some integer greater than 1. In the case of baseball, some statisticians have considered such models for various values of k , and have also assumed that the dependence weakens with the distance between the trials under consideration.

Another model, which has been used to study tennis, is called the odds model. Under this model, the probability $p_{(0,0)}$ that player A wins the first set in a match against player B might depend upon their rankings, if they are professionals, or on their past head-to-head history. If the set score is (i, j) (meaning that player A has won i sets and player B has won j sets), the probability that player A wins the next set is denoted by $p_{(i,j)}$. In the particular model we will consider, odds, instead of probabilities, are used. The odds O_{ij} that player A wins a set, if the set score is (i, j) , is defined by the equation

$$O_{ij} = k^{i-j} O_{00} ,$$

where k is a parameter that is estimated from the data. If $k > 1$, then this means that a player does better as the set score becomes more and more favorable to him; in other words, he has “momentum.” The

relatively simple form of the above equation is the reason that odds, and not probabilities, are used in this model. The corresponding equation involving probabilities is more complicated.

Finally, models have been proposed that add in “random effects” to one of the above models, i.e. at each trial, or possibly after a set of trials of a given length, a random real number is added to the probability of success.

3. Runs in Repeated Trials

Suppose we have an experiment that has several possible outcomes, and we repeat the experiment many times. For example, we might roll a die and record the face that comes up. Suppose the sequence of rolls is

$$1, 4, 4, 3, 5, 6, 2, 2, 2, 3, 3, 5, 6, 5, 6, 1, 1, 2, 2.$$

We define a run to be a consecutive set of trials with the same outcome that is not contained in any longer such set. So, in the sequence above, there is one run of length 3, namely the subsequence of three consecutive 2’s, and there are four runs of length 2.

If we wish to compare various models of repeated trials, with each trial having two possible outcomes, we might look at the length of the longest success run (or failure run), or the number of runs. Here it makes little difference whether one looks at the number of success runs or the total number of runs, since the second is within one of being twice the first in any sequence. One might also look at the average length of the success runs. When considering whether a process is Markovian, one might look at the observed success probabilities following successes and failures (or perhaps following sequences of consecutive successes or failures). When considering the block-Bernoulli model, one might compute observed values of the probability of success over blocks of consecutive trials.

In order to use statistical tests on any of the above parameters, one needs to compute or simulate the sampling distributions of the statistics under the models being considered. For example, suppose

that we have a set of data that consists of many strings of 0's and 1's, with each string being of length around 500. For each string, we can determine the length of the longest run of 1's. At this point, we could compare the observations with the theoretical distribution of the longest success run in a Bernoulli trial, where the parameters are $n = 500$ and p equaling the observed probability of a 1 in the strings. This comparison between the data and the theoretical distribution would yield, for each string, a p-value. The reader will recall that if we are testing a hypothesis, the p-value of an observation is the probability, assuming the hypothesis is true, that we would observe a result that is at least as extreme as the actual observation.

For example, suppose that in a sequence of 500 0's and 1's, we observe 245 0's and 255 1's, and we observe that the longest run of 1's is of length 11. One can show that if $n = 500$ and $p = .51$, then about 86% of the time, the longest success run of 1's in a Bernoulli trials sequence with these parameters is between 5 and 10, inclusive. Thus, the p-value of this observation is about .14, which means that we might be somewhat skeptical that the model does a good job of explaining the data.

Exercises.

1. (Classroom exercise). Split the class into two groups. Each student in the first group should flip a coin 200 times, recording the sequence of results. Each student in the other half should write down sequences of length 200 that they think look like typical sequences of coin flips. There are two related questions that one can consider here. First, can a person “make up” a sequence of flips that looks as if it came from actual coin flips? Second, can we use probability to distinguish, in many cases, between the sequences that come from actual experiments and those that are made up? The next two exercises give some insights on the latter question; we will have more to say about this later in the chapter.

2. Suppose that we have a Bernoulli trials process in which the probability of a success equals p . If there are n trials, what is the expected number of runs? Hint: In any outcome sequence, the number of runs is equal to one more than the number of two consecutive unequal trials. For example, if the outcome sequence is $SFFFSSFS$, then there are four pairs of consecutive unequal trials (these pairs correspond to the trials numbered $(1, 2)$, $(4, 5)$, $(6, 7)$, and $(7, 8)$). There are five runs in this sequence. If we let X_i be a random variable which equals 1 if the outcomes of trials i and $i + 1$ are different, then the number of runs R is given by

$$R = 1 + \sum_{i=1}^{n-1} X_i .$$

Thus, to find the expected value of R , it suffices to find the expected value of the right-hand sum; this equals

$$1 + \sum_{i=1}^{n-1} E(X_i) .$$

This can be used to help distinguish between actual and made-up coin toss sequences.

3. Suppose that we have a Bernoulli trials process in which the probability of success equals p . Let us call a sequence of k consecutive successes a k -string of successes. Note that this differs from a run in that a k -string of successes might be part of a longer string of successes, while a run of successes of length k is not part of any longer run. If there are n trials, and k is a positive integer, what is the expected number of k -strings of successes? Hint: For each i between 1 and $n - k + 1$, let Y_i be the random variable which is 1 if the tosses numbered between i and $i + k - 1$ are all successes, and is 0 otherwise. Then the expected number of k -strings

of successes equals

$$\sum_{i=1}^{n-k+1} E(Y_i).$$

It is easy to compute the value of $E(Y_i)$.

4. Statistical Tests on Repeated Trials

We now proceed to describe the distributions for the statistics that were mentioned above. We begin by recalling that if we let $p_1 = p_2$ in the Markov chain model, we obtain the Bernoulli model. If we are trying to test whether there is evidence of streakiness in a set of data, we might set the null hypothesis to be the statement that $p_1 = p_2$, and the alternative hypothesis to be the statement that $p_1 > p_2$, since this statement is one definition of streakiness. If we are testing whether the Bernoulli trials model fits the data, without any prejudice towards streakiness as an alternative, we might set the alternative hypothesis to $p_1 \neq p_2$.

The distribution for the number of runs in the Markov chain model is derived in the appendix of this chapter. Plots of both distributions for the case $n = 50$, $p = .2$, $p_1 = .3$, and $p_2 = .175$ are shown in Figures 1 and 2. The reason for the choices of p_1 and p_2 is that if $p_1 = .3$, then in order to make the long-range percentage of successes equal .2, as it is in the first model, we must choose $p_2 = .175$.

The most obvious characteristic of these two plots is the up and down fluctuation in probabilities between an odd and an even number of runs. This difference is easily explained. If $p = .2$, as it is in the first plot, then in a sequence of 50 trials, it is much more likely that the first and last trials are both failures than it is that any of the other three possibilities of success and failure occur. If the first and last trials are the same, there must be an odd number of streaks in the sequence.

By comparing these two figures, one can see that the distribution for the Bernoulli model is centered slightly to the right of the center

of the Markov chain model. One can see why this is the case. In the second model, a success is more likely to be followed by a success than in the first model. Similarly, a failure is more likely to be followed by failure in the second model than in the first. These statements make it reasonable to infer that in the second model, there will, on average, be fewer changes from success to failure, or vice versa, in successive trials, than there will be in the first model. This leads to the conclusion that on average, there will be fewer streaks in the second model than in the first model (and, as we will discuss below, the lengths of the streaks in the second model will be, on average, longer than in the first model).

This observation means that one might count the runs in a data set and use an interval of the form $[1, a]$ as the rejection region and the interval $[a+1, \infty)$ as the acceptance region. One can also compute a p-value of a data set under the Bernoulli model.

Another parameter that we will consider is the length of the longest success run. The distribution of this parameter in the case of Bernoulli trials was computed by Schilling [37]. This paper contains many interesting results concerning runs. We will content ourselves here with the following result from this paper. (See Exercise 2 for more on the distribution of the longest run of successes.) In a Bernoulli trials process with parameters n and p , we let X denote the length of the longest observed run of successes. Also, let $q = 1 - p$ and let γ denote Euler's constant (the value of this constant is about .577). Then we have

$$E(X) = \frac{(\log nq) + \gamma}{\log(1/p)} - \frac{1}{2} + r_1(n) + \epsilon_1(n),$$

where for $p \in [.1, 1]$ and all n , $|r_1(n)| < .005$, and $\epsilon_1(n) \rightarrow 0$ as $n \rightarrow \infty$. Thus, for example, if a baseball player's batting average is .300 in a given season, and he has 600 at-bats during that season, the above equation says that the expected value of the longest sequence of consecutive hits by this player is about

$$\frac{\log(600 \cdot .7) + \gamma}{\log(1/.3)} - \frac{1}{2} \approx 4.996.$$

If the player's batting average is .400, then the expected length of the longest run of consecutive hits is 6.554.

In the appendix, we derive the corresponding distribution for the Markov chain model. In Figure 3 we show the distribution of the longest run in both the Bernoulli trials model and the Markov chain model, using the parameters $n = 100$, $p = .5$, $p_1 = .7$, and $p_2 = .3$. The Markov chain distribution is the one shown with the larger dots. The means of the two distributions for these values of the parameters are 5.99 and 9.57. Again, one can test the hypothesis that $p_1 = p_2$ against the hypothesis that $p_1 > p_2$, using this parameter.

In [2], Albright used a χ^2 -test to test for dependence of consecutive at-bats for baseball players; we now give a general description of this test. Given a sequence of 0's and 1's, define n to be the length of the sequence, n_0 to be the number of 0's in the sequence, n_1 the number of 1's, and n_{ij} to be the number of pairs of consecutive of terms in the sequence with the first equal to i and the second equal to j . Then the statistic

$$\chi^2 = \frac{n(n_{00}n_{11} - n_{10}n_{01})^2}{n_0^2 n_1^2}$$

is, under the assumption that the sequence has been generated by a Bernoulli trials process, approximately χ^2 -distributed with 1 degree of freedom. So, to apply this test, we compute the value of χ^2 and see if the value exceeds the critical value of the χ^2 -distribution at a given level of significance, or we report the p-value corresponding to the observed value.

When trying to decide whether $p \neq p_1$ another parameter that is useful is the number of success doubletons (hereafter called doubletons), i.e. pairs of consecutive successes. If we let \hat{d} denote the number of doubletons in a sequence of length n , then Exercise 1 shows that \hat{d} is within one of $n\hat{p}\hat{p}_1$, where \hat{p} and \hat{p}_1 are the observed values of p and p_1 . Thus, using the number of doubletons to study a sequence is very similar to using the value of the parameter p_1 .

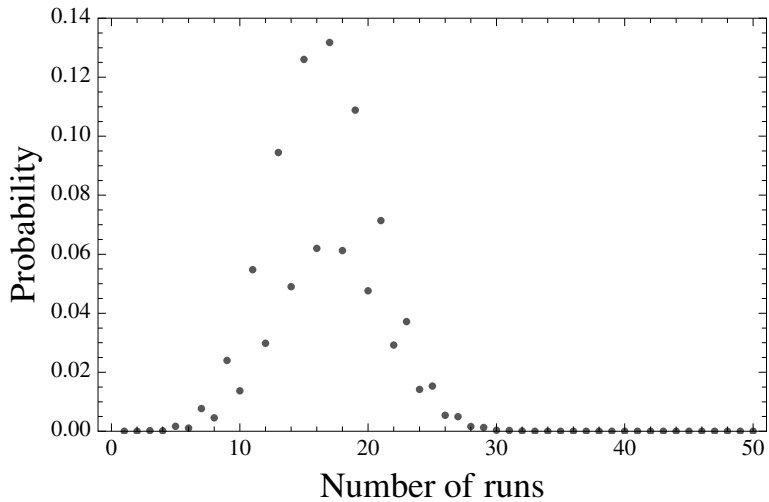


Figure 1. Distribution of number of runs in the Bernoulli model

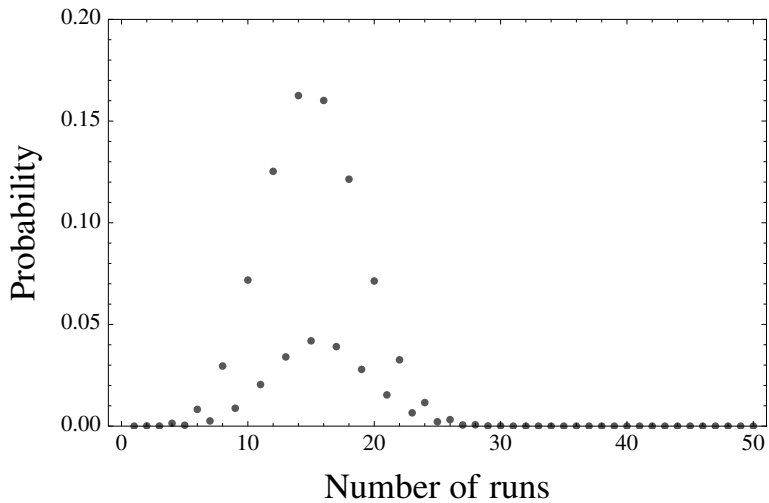


Figure 2. Distribution of the number of runs in the Markov model

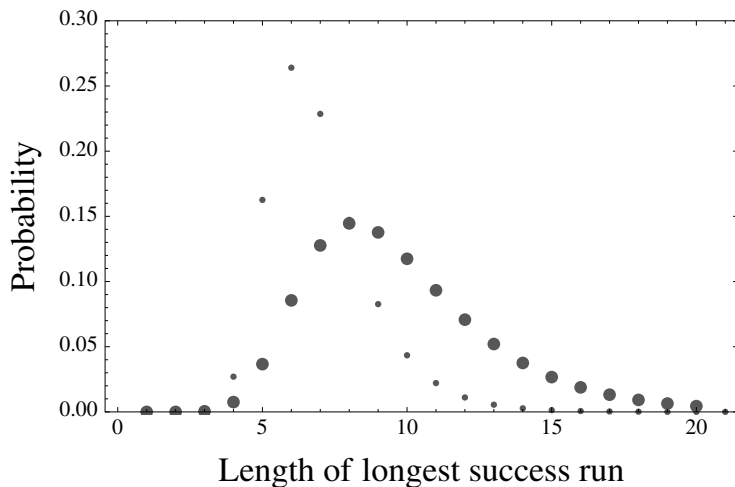


Figure 3. Distributions of longest success run for Markov and Bernoulli models

In the appendix, we explain why the number of doubletons is asymptotically normally distributed in the Bernoulli trials model and give expressions for the asymptotic mean and standard deviation. We now go through an example to show how we can use this distribution. We assume that the sequence is of length n , and that the observed value of p (which we will take as the value of p) is $.3$. In Figure 4, we show the normal distributions corresponding to values of $p_1 = .3$ (the Bernoulli case) and $p_1 = .4$; the first of these is on the left. The vertical line in the graph marks the right-hand endpoint of a 95% confidence interval. The horizontal coordinate of this line is 57.7. Although this distribution and the next one are discrete, we have drawn them as continuous to make it easier to see the vertical line.

In this case, we are testing the hypothesis that $p_1 = .3$ (the null hypothesis) against the alternative hypothesis $p_1 = .4$. If we have an observed sequence of length 500, we count the number of doubletons. If the null hypothesis were true, the number of doubletons would be greater than 57.7 only 5% of the time. Thus, if the observed number

of doubletons is greater than 57.7, we reject the null hypothesis at the 5% level of significance.

When one is carrying out a test of a hypothesis, there are two types of errors that can occur. A type I error occurs when the null hypothesis is rejected, even though it is true. The probability of a type I error is typically denoted by α . We see that in the present case, $\alpha = .05$. A type II error occurs if the null hypothesis is not rejected, even though it is false. The probability of such an error is denoted by β . In order to estimate β , one needs an alternative value of the parameter being estimated. In the present case, if we take this alternative value to be $p_1 = .4$, then one can calculate that $\beta = .404$. It is possible to visualize both α and β . In Figure 4, α is the area under the left-hand curve to the right of the vertical line and β is the area under the right-hand curve to the left of the vertical line.

The power of a hypothesis test is defined to be $1 - \beta$; it is the probability that the null hypothesis will be rejected at the α level of significance when the alternative hypothesis is true. In this case, the power is .596. The higher the power of a test, the more confident we are that we would be able to detect a real departure from the null hypothesis. If we change the value of n to 2000 and graph the same distributions again, we obtain Figure 5. In this case, the power of the test is .964.

The reader will recall that in Section 1 we defined a block-Bernoulli process as a way to model sequences of successes and failures having success probabilities that change over time. For example, over the course of a baseball season, a batter might have periods in which he has different probabilities of getting a hit.

We will take as our model one that consists of blocks (intervals) of consecutive trials, such that in each block, the individual trials are mutually independent and the success probability is constant. Of course, for such a model to be interesting, the lengths of the blocks in which the success probability is constant must be fairly long. Otherwise one could not hope to differentiate this process from a Bernoulli

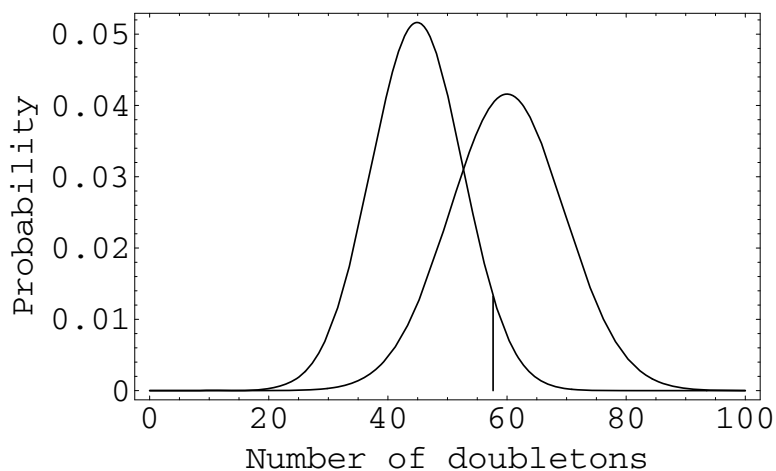


Figure 4. Number of Doubletons for $n = 500$, $p = .3$ and either $p_1 = .3$ or $p_1 = .4$

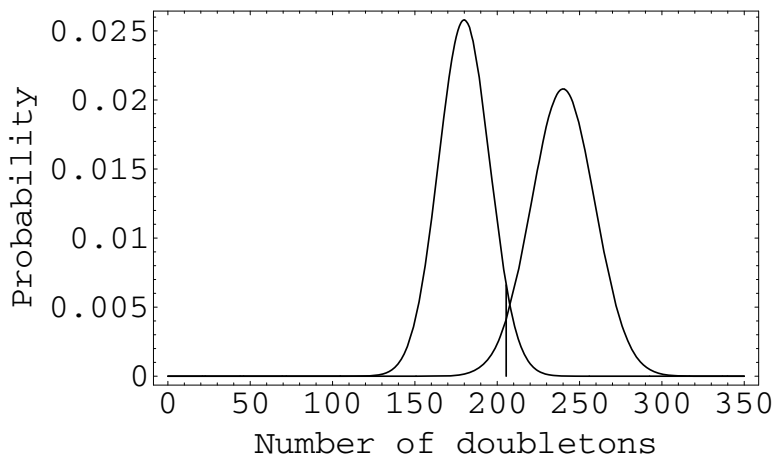


Figure 5. Number of Doubletons for $n = 2000$, $p = .3$ and either $p_1 = .3$ or $p_1 = .4$

process. For example, suppose that we define a process that, for each block of size 1, has a success probability that is picked uniformly at random from the interval $[0, 1]$. This process will be completely indistinguishable from a sequence of coin flips of a fair coin. The reason for this is that on each trial, the probability that the coin will come up heads is exactly $1/2$, by symmetry considerations.

We will assume that the lengths of the blocks vary uniformly between a and b , and the success probabilities on these blocks vary uniformly between p_{min} and p_{max} . It is possible to imagine other assumptions one might make; for example, perhaps the success probabilities are more likely to be near the middle of the interval $[p_{min}, p_{max}]$ than near the end.

What statistic might we use to distinguish block-Bernoulli trials processes from Bernoulli ones? One fairly obvious parameter can be obtained in the following way. Suppose we think that the blocks in a block-Bernoulli trials process are of average length 40. In a sequence of length n , generated by this process, there will be $n - 39$ blocks of length 40. (These blocks are allowed to overlap; they must start between position 1 and position $n - 39$.) For each of these blocks, we compute the observed probability of success by dividing the number of successes in the block by 40, the length of the block. We then take the difference between the maximum and minimum observed probabilities. We will call this statistic the windowed difference of the sequence.

In the Bernoulli trials model, if the success probability is p , then the observed success probability in a block of length m is, for moderate-sized m , approximately normally distributed with mean p and variance $p(1 - p)/m$. If we have a sequence of length n , we wish to know the distribution of the difference between the maximum and minimum observed success probabilities over all blocks of length m .

We do not know whether an exact expression for this distribution has been calculated by other authors. However, it can certainly be simulated. In Figures 6 and 7, we show simulations of the success

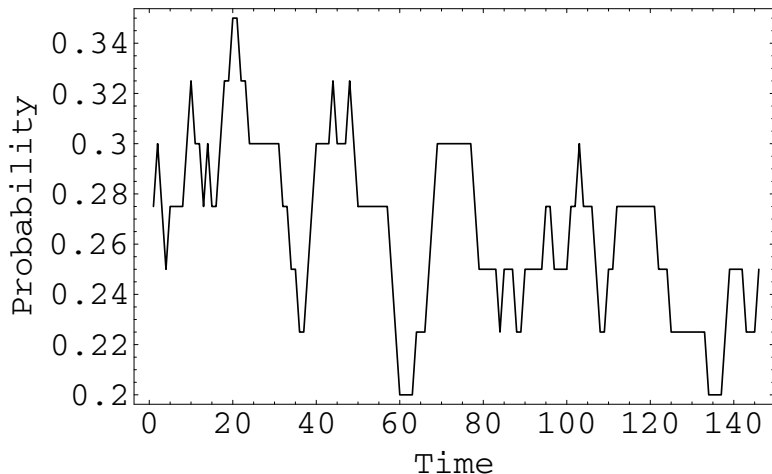


Figure 6. Running success probabilities for block-Bernoulli model

probabilities for blocks of length 40 in both a block-Bernoulli and a Bernoulli trials sequence. In both simulations, $n = 200$. In the Bernoulli case, $p = .3$, and in the block-Bernoulli case, we let $a = 30$, $b = 50$, $p_{min} = .25$, and $p_{max} = .35$. In the sequence corresponding to Figure 6, the block sizes are 38, 40, 38, 39, and 30, and the success probabilities in these blocks are .250, .293, .279, .319, and .288.

The reader will notice that the two figures look similar in the amount of variation in the success probabilities. This similarity should make the reader wonder whether the statistic described above, the difference between the maximum and minimum success probabilities over the blocks, does a very good job of distinguishing between the two models.

One way to answer this question is to compute the power of a test. In the present case, we let the null hypothesis be the statement that the success probability is constant over the entire sequence, i.e. that our process is a Bernoulli process. We will take $n = 500$ and $p = .3$. The alternative hypothesis deals with the block-Bernoulli process. Here we assume that $a = 30$ and $b = 50$, although one

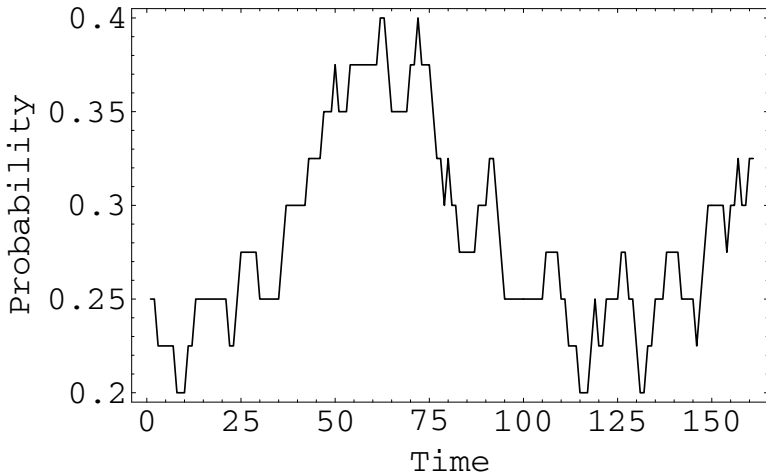


Figure 7. Running success probabilities for Bernoulli model

can certainly imagine varying these values. The specific parameter whose value we will let vary is p_{max} , the maximum allowable success probability. We then let $p_{min} = 2p - p_{max}$, so that p is the midpoint of the interval $[p_{min}, p_{max}]$. Note that if we let $p_{max} = .3$, then the block-Bernoulli process reduces to the Bernoulli process.

We now wish to find a critical region, at the 5% level of significance, for our hypothesis test. To do this, we simulate the Bernoulli process 1000 times, and determine a value, called the critical value, below which we find 95% of the values of the windowed difference. In Figure 8, we show a histogram of the values of the windowed difference for the Bernoulli process. The critical value is .425, so the region $[\text{.425}, 1]$ is the critical region, meaning that if we obtain a value of the parameter that exceeds .425, we reject the null hypothesis.

In Figure 9, we show the corresponding histogram for the block-Bernoulli process for the parameter value $p_{max} = .4$. It can be seen that the values in this case are slightly larger than in the Bernoulli case, but the two distributions overlap quite a bit. The power of this test for the specific value of $p_{max} = .4$ is one minus the probability

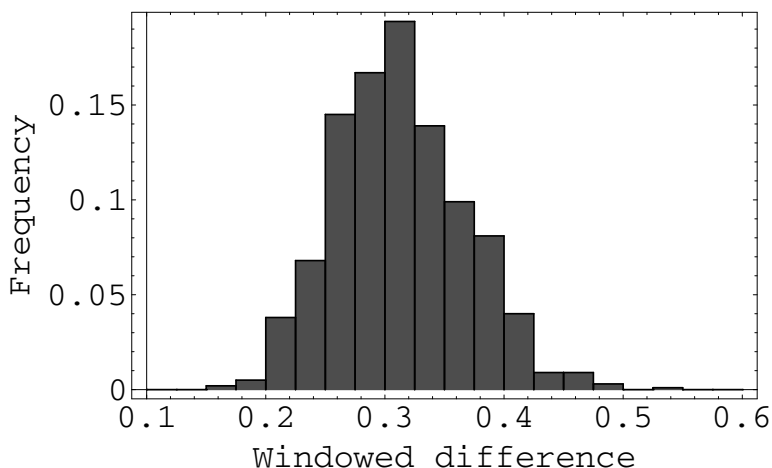


Figure 8. Simulated values of windowed difference for Bernoulli process

that we will fail to reject the null hypothesis when, in fact, the alternative hypothesis is true. In this case, this means that we want to estimate the probability that the windowed difference falls in the interval $[0, .425]$ for the block-Bernoulli process. Our simulation gives us a value of .84 for this probability, so the power of this test is .16, which isn't very good. One way to increase the power is to increase the size of n , but this may or may not be feasible, depending upon how we come upon the data.

Suppose the null hypothesis is of the form " $p = p_0$ " for some parameter p , and the alternative hypothesis is of the form " $p > p_0$." If we observe a test statistic with value \hat{p} , then the p-value of the test is the probability, given the null hypothesis is true, that we would obtain a value of the test statistic at least as large as the \hat{p} that we observed.

On the other hand, suppose that the alternative hypothesis is of the form " $p \neq p_0$." In this case the p-value is the probability, given the null hypothesis is true, that we would obtain a value of the test

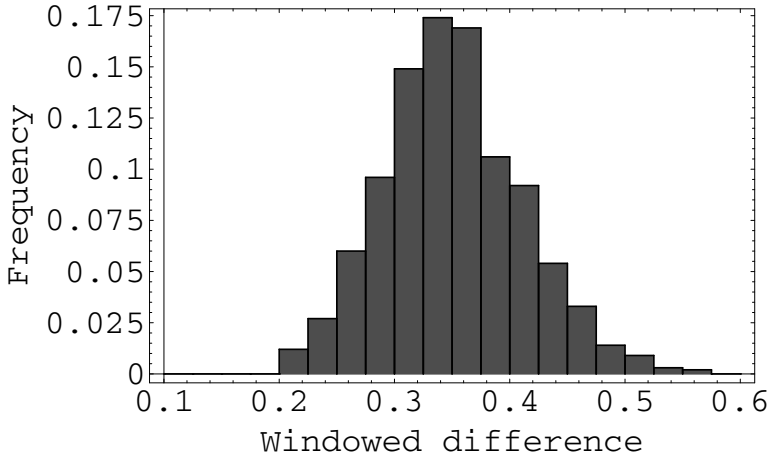


Figure 9. Simulated values of windowed difference for block-Bernoulli process

statistic whose distance from the hypothesized value is at least as large as the distance $|\hat{p} - p_0|$ that we observed.

For a given hypothesis, the smaller the p-value, the less faith we have in the model as an explanation of the observed data. On the other hand, if we have many data sets that are supposedly explained by a certain model, then some of the data sets will probably have small p-values, even if the model is correct. To see why this is true, consider the following simple experiment. We have a coin, and we are testing the hypothesis that it is a fair coin, i.e. that the probability of a head coming up when we flip the coin is .5. Suppose the coin is fair. Suppose we perform 100 tests in each of which the coin is flipped 500 times. Then about 5% of the tests will report a p-value of less than .05. The reason for this is that the observed value \hat{p} has a certain known distribution, because the null hypothesis is true, and we are using this distribution to calculate the p-values of the observations. Since the observations are distributed according to this known distribution, about 5% of them will have p-values less than .05 (just as about 45% of them, say, will have p-values less than .45).

Another way of looking at this situation is as follows: If we have a large number of observations under a single model, and the p-values of these observations are fairly uniformly distributed across the interval $[0, 1]$, then the model is doing a good job of fitting the data. Of course this does not mean that the model is the only one that fits the data well or that it the “correct” one. One should consider other factors, such as the simplicity (or lack thereof) of the various models under consideration, before deciding which model one likes best.

4.1. Selection Bias. If one calculates a statistic from observed data, and the value is seen to be closer to its expected value under one model than under a second model, one is tempted to state that the first model does a better job of explaining the observation than does the second model. However, if one selects only those observations for which the parameter is closer to the expected value in the first model than in the second model, then one is guilty of selection bias.

It is also, in general, the case that if one considers more parameters in creating a model, the model may fit the observed data better than the original, simpler model. For example, since the Markov model subsumes the Bernoulli model, if one includes p_1 and p_2 as parameters, thereby considering the Markov model, this model will do a better job of fitting a sequence of 0's and 1's than will the Bernoulli model. The trade-off is that the Markov model is more complicated than the Bernoulli model. If the Bernoulli model does a good job of describing the observations, then its simplicity should be an argument against its rejection.

Exercises.

1. Suppose that we have a sequence of successes and failures of length n . Define \hat{p} to be the observed proportion of successes and define \hat{p}_1 to be the observed proportion of successes (excluding the last entry in the sequence, if it is a success) that are followed by another success. Define \hat{d} to be the observed number of consecutive pairs of entries that are both

successes (i.e. the number of doubletons). Show that if the last entry in the sequence is a failure, then

$$\hat{d} = n\hat{p}\hat{p}_1 ,$$

while if the last entry is a success, then

$$\hat{d} = (n\hat{p} - 1)\hat{p}_1 .$$

Thus, in either case, \hat{d} is within one of $n\hat{p}\hat{p}_1$.

- 2.** In Schilling [37], an approximation is given for the distribution of the length of the longest success run in a Bernoulli trials process with parameters n and p (as usual, n denotes the number of trials and p is the probability of success in any one trial). A more precise description of the behavior of this quantity as n gets large is described in a paper by Gordon, Schilling, and Waterman [17]. Here is the approximation. We assume that $0 < p < 1$, and we set $q = 1 - p$. We define

$$F_p(x) = e^{-(p^x)} .$$

If we let R_n denote the random variable whose value is the length of the longest success run in a sequence of n coin tosses, then we have

$$P(R_n = x) \approx F_p\left(x + 1 - \frac{\log(nq)}{\log(1/p)}\right) - F_p\left(x - \frac{\log(nq)}{\log(1/p)}\right) .$$

Even for moderate-sized values of n , this approximation is fairly accurate. In the table below, we show the exact values and the approximation, for $n = 100$ and $p = 1/3$.

x	Exact Probability	Approximation
0	0	0
1	0.0001	0.0006
2	0.0692	0.0841
3	0.3654	0.3544
4	0.3297	0.3210
5	0.1515	0.1525
6	0.0556	0.0574
7	0.0190	0.0199
8	0.0063	0.0067
9	0.0021	0.0023
10	0.0002	0.0003

- (a) Show that the above approximation can be written as

$$P(R_n = x) \approx e^{-p^{x+1}nq} - e^{-p^x nq}.$$

- (b) Use the expression in part a) to show the following.

Given n , p , and x , the value of $P(R_n = x)$ is approximately the same as the value of $P(R_{n/p} = x+1)$. This statement can be interpreted as follows, using $p = 1/2$ to make the interpretation easier to understand. In a sequence of n flips of a fair coin, the probability that the longest run of heads equals a certain value x is approximately the same as the probability, in a sequence of $2n$ flips of a fair coin, that the longest run of heads equals $x+1$. Thus, every time we double the number of coin flips, for a fair coin, the distribution for the longest run of heads moves one to the right on the x -axis.

3. In the text, we asserted that if a block Bernoulli process is created in which the blocks are of length 1 and the success probabilities for the blocks are chosen uniformly in $[0, 1]$, then this process is identical to the Bernoulli trials process with success probability $1/2$. Explain why this assertion is true.

5. Data from Various Sports

5.1. Baseball. In baseball, the sequence of hits and outs recorded by a batter in consecutive at-bats is a candidate for study. A plate appearance that results in a walk (or the batter being hit by a pitch, or a sacrifice bunt) is not counted as an at-bat, because in those cases, it is thought that the batter has not been given a chance to hit a well-pitched ball. We remind the reader that even if one uses a Bernoulli trials approach to model such sequences, one should admit that this is only a first approximation to what is actually occurring on the field; presumably the batter's chance of success in a given at-bat is affected by who is pitching, whether the game is being played during the day or at night, how many runners are on base, and many other factors. However, it is still of interest to test the Bernoulli trials model to see how well it fits the data.

We will now apply the various tests described above to our data set. This data set consists of all players who played in the major leagues between 1978 and 1992. We obtained this data from the website www.retrosheet.org. We greatly appreciate their labor in creating the files of raw data that we used.

The first test we will apply is the chi-squared test of Albright, described above. For each player and for each season in which that player had at least 150 at-bats, we computed the chi-squared value of his sequence of hits and outs, and then computed the p-value of this chi-squared value. If there is no overall departure from the Bernoulli trials model, we would expect to see the p-values uniformly distributed across the interval $[0, 1]$. The point here is that even though one might find some seasons for some players for which the p-value is less than some prescribed level of significance, say the 5% level, we should not reject the hypothesis that the Bernoulli trials model fits the data well unless we observe many more than 5% of the p-values below .05.

There were 4786 player-seasons in our data set with the at-bat cutoff of 150. Figure 10 shows a histogram of the p-values associated

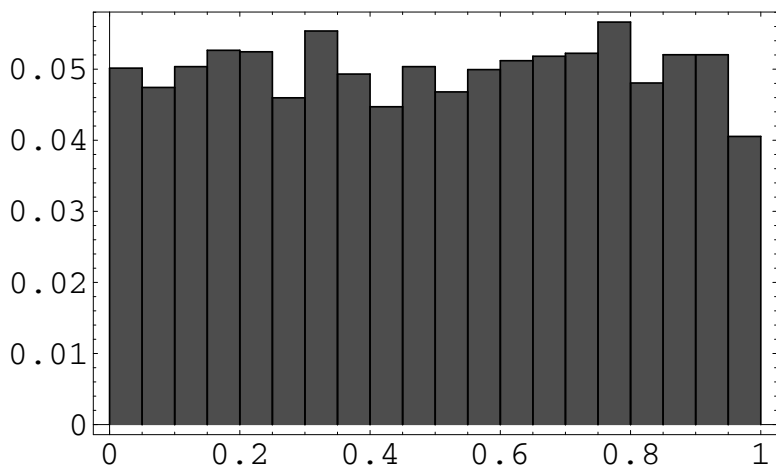


Figure 10. p-Values for Albright's chi-squared test in baseball

with the chi-squared values for these player-seasons. Since the interval $[0, 1]$ has been partitioned into 20 subintervals, uniformity would imply that about 5% of the values occur in each sub-interval. It can be seen that the p-values are essentially uniform across the unit interval. In particular, there is no apparent excess of small p-values that would correspond to large values of the chi-squared parameter (and therefore to either streaky or anti-streaky behavior).

Another pair of statistics to look at is the pair (\hat{p}_1, \hat{p}_2) ; these are the observed probabilities of a hit after a hit and after an out, respectively. Under the assumption of Bernoulli trials, it can be shown that $\hat{p}_1 - \hat{p}_2$ is approximately normally distributed with a mean and a variance that depend on n and p . (We have shown this in the appendix.) In fact, the mean is 0 and the variance is

$$\frac{n + 3np - 1 - 5p}{qn^2}.$$

We can test to see if the hits-outs sequences of baseball players fit the Bernoulli model by looking at the values of $\hat{p}_1 - \hat{p}_2$. It should be noted that the n and p values change as we move from player to player, so

we will aggregate the values in the following way. We partition the interval $[0, .750]$ into 15 subintervals of length .050, and we partition the interval $[0, .400]$ into 16 subintervals of length .025. Given a pair of subintervals, one of each type, we take all player-seasons in which the number of at-bats is in the first subinterval and the batting average (the value of p) is in the second interval. So, for example, the pair $[400, 450]$ and $[\text{.300}, \text{.325}]$ consists of all player-seasons in which the player had between 400 and 450 at-bats and a batting average of between .300 and .325. To avoid duplications, the subintervals do not contain their left-hand endpoint. The reason for this partitioning scheme is so that we have enough data for certain values of n and p to be able to say something meaningful.

Of the 240 possible pairs of subintervals, 22 of them have data sets of size at least 100. There were 3188 player-seasons in these data sets (which is more than two-thirds of the data). For each of these data sets, we calculated the average value of $\hat{p}_1 - \hat{p}_2$. (Recall that in the Bernoulli model, the mean of $\hat{p}_1 - \hat{p}_2$ is 0.) The variance of the average value of $\hat{p}_1 - \hat{p}_2$ is the variance of $\hat{p}_1 - \hat{p}_2$ divided by the square of the size of the data set. Since $\hat{p}_1 - \hat{p}_2$ is approximately normal, so is the average value of $\hat{p}_1 - \hat{p}_2$. So one way to see how well the data fits the Bernoulli model is to compute, for each data set, the z -value of the average value of $\hat{p}_1 - \hat{p}_2$; a z -value that is less than -2 or larger than 2 is significant at the 5% level. The z -value is obtained by dividing the observed value of the average of $\hat{p}_1 - \hat{p}_2$ in the data set by the standard deviation of the average value of $\hat{p}_1 - \hat{p}_2$.

The results are shown in Figure 11. Of the 22 z -values, one is of absolute value greater than 2 (which is about what one would expect if one were to choose 22 random values from a normal distribution).

We turn next to the length of the longest success run (i.e. the longest run of hits). For each player-season in our data, we can compute the length of the longest success run. However, in the Bernoulli model, the distribution of this length depends upon both n , the length of the sequence (in this case, the number of at-bats), and p , the probability of a success in an individual trial (in this case, the batting

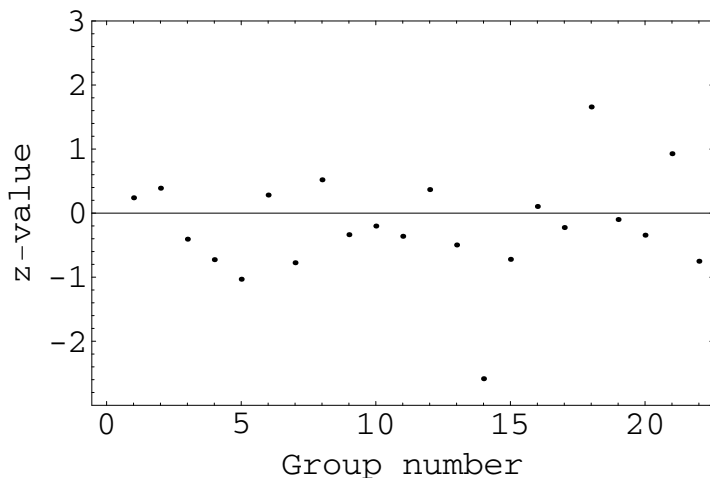


Figure 11. z-Values for the Average of $\hat{p}_1 - \hat{p}_2$

average). The Markov model depends upon n , p , and also on p_1 , the probability of a success following a success. We would like to be able to display the aggregate data against the predictions made by both models.

In the case of the Bernoulli model, we can proceed in several ways. Our first method of comparison will be carried out as follows. For each of the 4786 player-seasons with at least 150 at-bats, we will simulate a season using Bernoulli trials with the same number of at-bats and batting average as the actual season. Then we will observe the longest run of successes in each of these simulated seasons. Finally, we will compare the distribution of the lengths of the longest run of successes in the simulated seasons with the corresponding distribution in the actual data.

When we carry this procedure out, we obtain the results shown in Figure 12. The horizontal coordinate represents the length of the longest success run in each player-season, and the vertical coordinate represents the observed frequencies. The dots correspond to the simulated data and the horizontal lines correspond to the actual data. The

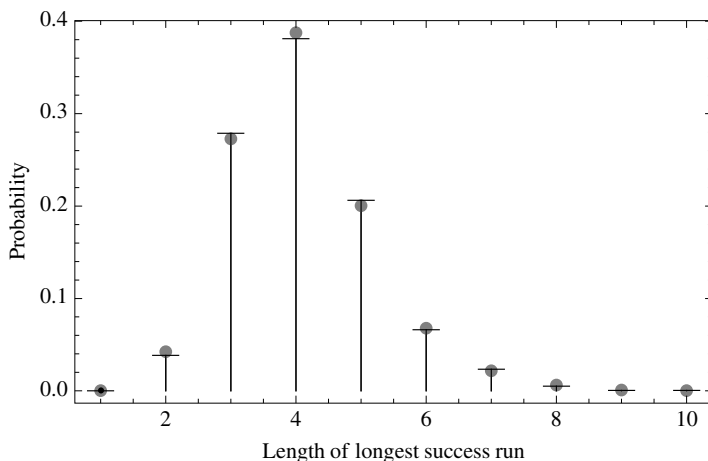


Figure 12. Distributions of simulated and actual longest success runs

fit between the simulated and actual data is quite good. In addition, the reader will recall Figure 3, which shows that in the Markov model, if $p_1 > p$, then the distribution of the longest success run is shifted to the right from the corresponding distribution in the Bernoulli model. The aggregate data does not show any such shift.

Another way to compare the Bernoulli model with the data is to proceed as we did above with the distribution of $\hat{p}_1 - \hat{p}_2$. We will group the data using pairs of subintervals, so that the numbers of at-bats and the batting averages are close to equal across the group. Then, for each group, we will compute the distribution of the longest success run and compare it with the theoretical distribution. This comparison will be carried out using a chi-square test. For each of the 22 groups we used above (those containing at least 100 player-seasons) we will report a p-value for the chi-square value.

In Figure 13, we show the results of the above calculations. For each of the 22 groups of data containing at least 100 player-seasons, the observed and theoretical average length of the longest success run; the theoretical average is obtained by using, for each group, the

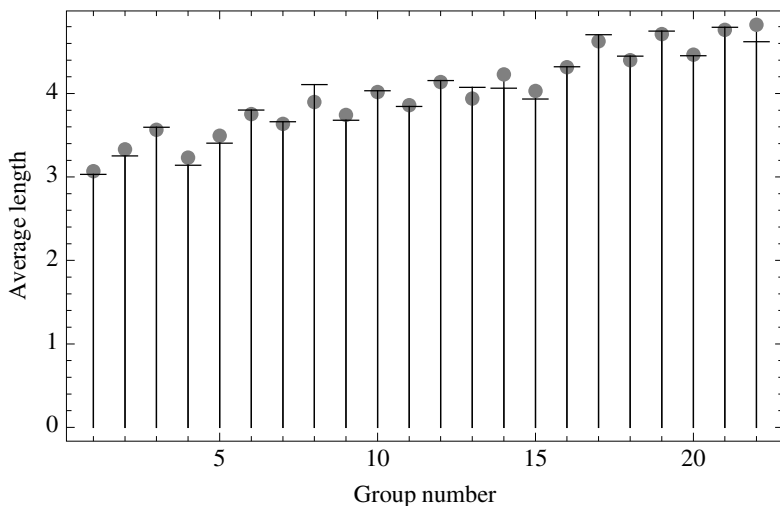


Figure 13. Observed vs. theoretical longest success run lengths

average number of at-bats and the average success probability. As can be seen, the fit is very good; the relative error is never greater than 6%, and only twice is it as large as 4%. For each of the 22 data sets, we also compute the chi-squared value of the comparison between the theoretical and observed distributions and then compute the p-values corresponding to these chi-squared values. Figure 14 shows the 22 p-values, sorted by size. If there were a perfect fit between the theoretical and observed distributions, one would see the points lying close to a diagonal line from the bottom left-hand corner to the top right-hand corner. This figure shows that the Bernoulli model does a fairly good job of fitting the data.

It has been suggested that there might be a qualitative difference between success runs and failure runs. The thought was that while success runs are limited in length by the skill of the batter, failure runs might, in some cases, continue longer than predicted because the hitter will spiral downwards psychologically during a slump. We

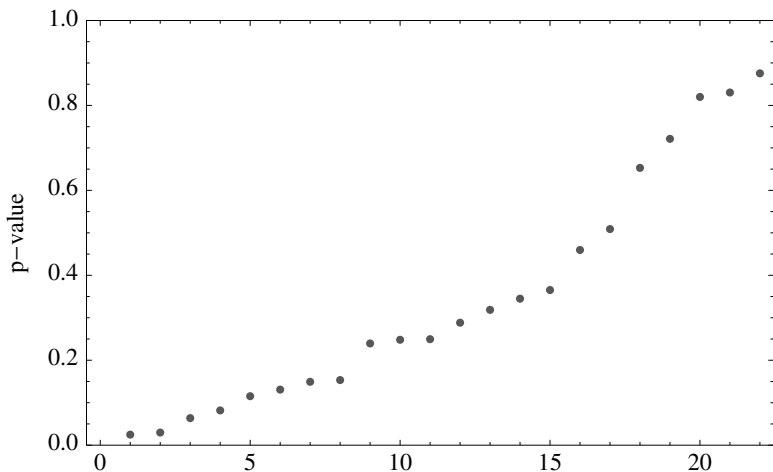


Figure 14. p-Values for chi-squared comparison of longest success run data

ran the above tests again on the same 22 groups, this time looking at the longest failure runs. The p-values, corresponding to those in Figure 14, are shown in Figure 15. It can be seen that once again, the fit is very good.

Another aspect of baseball that has been taken to provide evidence in favor of streakiness is the idea of a hitting streak. A hitting streak by a player is a set of consecutive games played by the player in which he gets at least one hit in each game. In order to qualify as a hitting streak, the games must all occur in a single season. The longest hitting streak that has ever occurred in the major leagues lasted 56 games; it happened in 1941, and the player who accomplished this streak was Joe DiMaggio. Much has been written about this streak (see, for example, [4], [18], and [27]). The question that we wish to consider is whether the fact that this streak occurred provides evidence against the Bernoulli trials model.

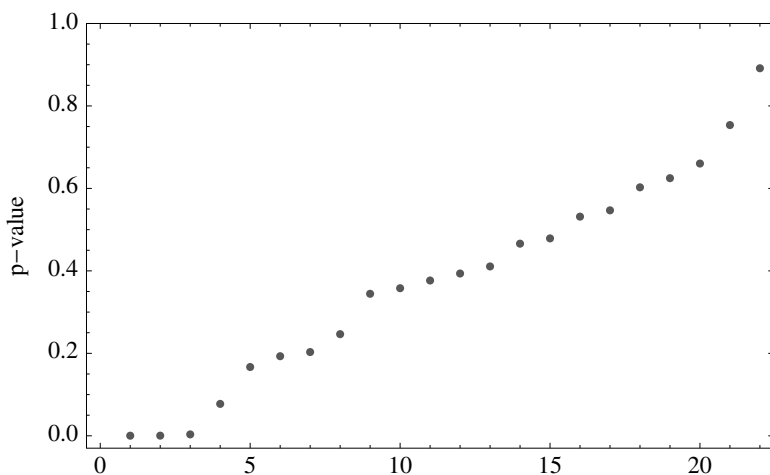


Figure 15. p-Values for chi-squared comparison of longest failure run data.

If a hitter has a fixed number of at-bats in each game, then it is fairly easy to calculate the probability that he will hit safely in n consecutive games, assuming the Bernoulli model. For example, if a hitter has a .340 batting average, and has four at-bats in each game, then, under the Bernoulli model, the probability that he gets at least one hit in a given game is

$$1 - (.660)^4 = .8103 .$$

Thus, in a given set of 56 consecutive games, the probability that the hitter gets at least one hit in each game is

$$(.8103)^{56} = .00000766 .$$

This sounds like a very unlikely event, but we're not really asking the right question. A slightly better question might be the following: Given that a hitter plays in 154 games in a season, and has a season batting average of .340, what is the probability that he bats safely for at least 56 consecutive games? (We use 154 for the length of a

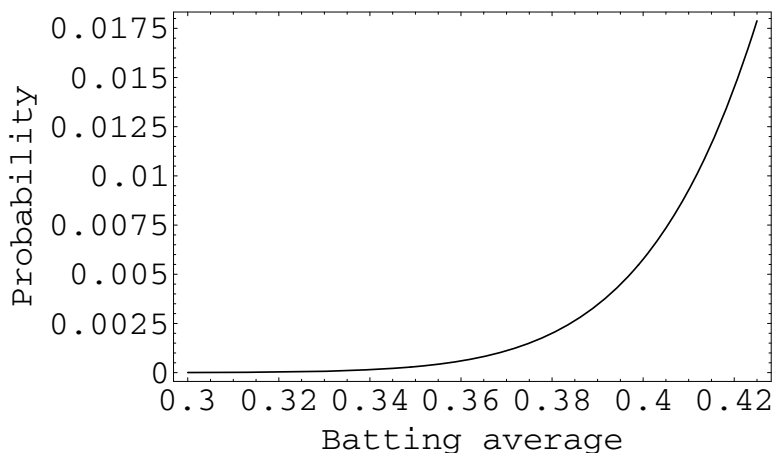


Figure 16. Probability of a hitting streak of at least 56 games

season, since almost all of the player seasons with batting averages of at least .340 occurred before the seasons were lengthened to 162 games.) Once again, under the Bernoulli model, we can answer this question rather easily. Putting this question in terms of ideas that we have already introduced, this question is equivalent to asking what is the probability, in a Bernoulli trials sequence of length 154, and with success probability .8103, that the longest success run is at least of length 56. The question, in a form similar to this one, was considered in [38]. The answer, in this case, is .0003.

Of course, there have been many batters who have had season batting averages much higher than .340, and it is clear that as the batting average climbs, so does the probability of having a long hitting streak. Figure 16 shows the probability of a hitting streak of at least 56 games in a 154-game season as a function of the season batting average. This figure shows that the season batting average has a large effect on the probability of a long hitting streak. Figure 17 shows a histogram of player seasons since 1901 for which the batting average was at least .340.

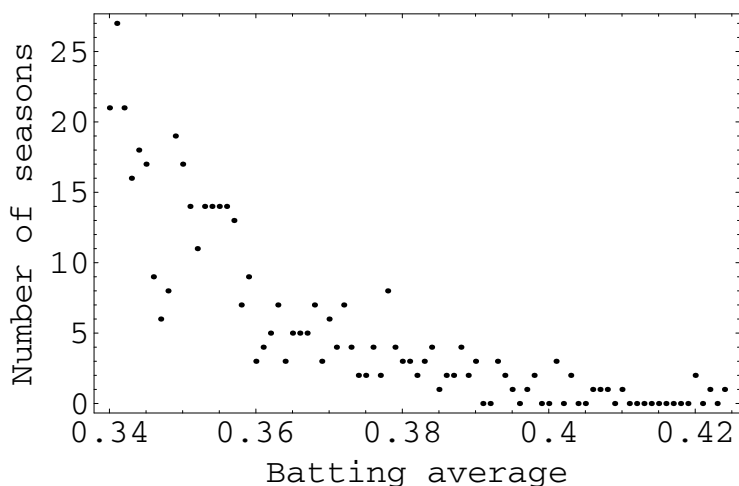


Figure 17. Number of player seasons with a given batting average

We still are not asking quite the right question. There are two changes that we should make. First, we do not want to know the probability that a given hitter will have a long hitting streak; rather, we want to know the probability that at least one entry in a set of player seasons will contain a long hitting streak. We will take this set of hitters to correspond to all hitters since 1901 who batted at least .340 in a season and had at least 300 at-bats. (There were 430 such seasons.) We are arbitrarily disregarding hitters whose batting averages are below .340. (By doing so, we will underestimate the probability we are considering.)

The second change that should be made concerns the variability of the number of at-bats belonging to an individual player over a set of games. It should be clear that if a batter averages four at-bats per game, but the number of at-bats varies widely over a set of games, then it is less likely that he will have a long hitting streak. As a simple example to help show this, consider a batter with a .340 batting average under the following two scenarios: first, he gets exactly four

at-bats in each of 10 consecutive games, and second, he gets, alternatively, two and six at-bats in 10 consecutive games. In the first case, the probability that he gets at least one hit in each game is

$$\left(1 - (.660)^4\right)^{10} = .12196 .$$

In the second case, the probability is

$$\left(1 - (.660)^2\right)^5 \left(1 - (.660)^6\right)^5 = .03721 .$$

So, here is the question that we want to try to answer: Suppose we restrict our attention to those player-seasons in which the player batted at least .340 (we will call this set our good-hitter data set), and suppose we take into account the typical variation between games of the number of at-bats of a player in a game. What is the probability, under the Bernoulli trials model, that at least one of the players would have a hitting streak of at least 56 games?

This is too complicated a question to hope for a theoretical answer, but we can simulate the seasons of these hitters, using the observed batting averages and numbers of games and at-bats. Here is how the simulation is carried out. We begin by finding an estimate for the standard deviation in the sequence of numbers of at-bats in a game, over a typical player season. We can use our original data set from the years 1978 to 1992 to obtain this estimate. We restrict our attention to those players who batted at least .300 for the season and who had at least 300 at-bats. There were 369 player seasons that satisfied these requirements. For each such season, we found the sequence of numbers of at-bats in all games *started* by the player. The reason for this restriction is that the players in our good-hitter data set started almost every game in which they played. For each of these 369 seasons, we compute the standard deviation of the sequence of at-bats. Then we compute the average of these standard deviations. When we do this, we obtain a value of .8523.

For each player-season in our good-hitter data set, we will produce a sequence of at-bats per game for the number of games in which the player participated that season. The terms of this sequence will

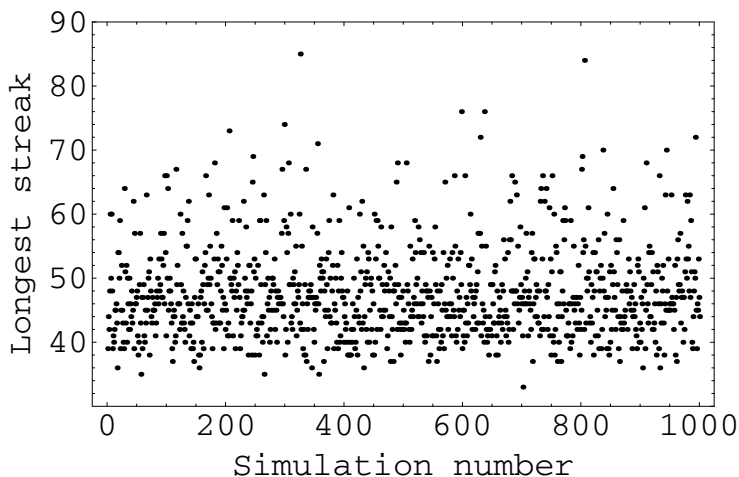


Figure 18. Simulated longest hitting streak data

be drawn from a normal distribution with mean equal to the average number of at-bats by the player during that season, and with standard deviation equal to .8523. The terms will be rounded to the nearest integer, so that they represent numbers of at-bats (which must be integers). If an integer in the sequence is less than or equal to 0, we throw it out. The reason for this is that a hitting streak is not considered to be interrupted if a player appears in a game but does not record any official at-bats.

Since we are operating under the assumption of Bernoulli trials, we use the player's season batting average to simulate his season with the above simulated at-bat sequence. We then record the longest hitting streak in the season. We do this for each player, and record the longest hitting streak among all of the player seasons. To estimate the probability that someone would have a hitting streak of at least 56 games, we carry out the above procedure many times and observe the fraction of those trials that lead to at least one such streak.

The results of the simulation are shown in Figure 18. The above procedure was carried out 1000 times. In the simulated data, the

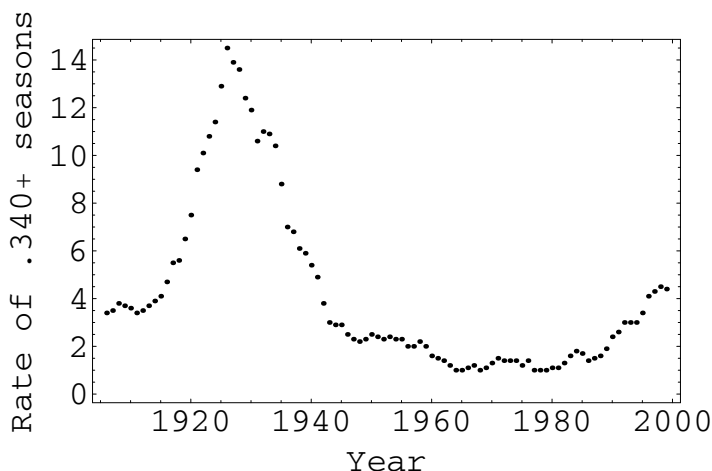


Figure 19. Moving average number of .340+ player seasons

length of the longest hitting streak ranged from 33 to 85, and 126 times the length was at least 56 games. This simulation shows that the probability, under the Bernoulli trials model, that we would observe a hitting streak of at least 56 games in length is about $1/8$ ($\approx 126/1000$) and thus a streak like Joe DiMaggio's would not be unusual under this model. The simulation also shows that viewed in this way, DiMaggio's feat is probably not the most amazing feat in all of sports, or even in baseball.

One can consider another question that is usually raised when DiMaggio's streak is discussed; that is the question of whether anyone will ever break his record. Predicting the future is a dangerous business, but one can use the above simulations, together with known data, to say something of interest. Of the 430 player seasons in which the player batted at least .340, about three-quarters of them occurred in the period 1901-1950. Thus, the rate at which .340-plus seasons occur has changed dramatically over the years. Figure 19 shows the moving average of the number of .340-plus player seasons, averaged over ten-year windows. Currently, the average number of such player

seasons is about 4 per year. Thus we assume that there will be about 4 such player seasons in each year in the future. We have seen that in 430 such player seasons, the probability is about $1/8$ that someone will have a 56-game hitting streak. Thus, we should expect to see one in the next 3440 ($= 8 \times 430$) such player seasons, which translates to 860 years. Putting it another way, it seems very unlikely that we will see such a streak in any of our lifetimes. Viewed from this perspective, DiMaggio's streak seems much more impressive than before.

There is another way to estimate how long it might be until the record is broken. It was stated above that the probability is .0003 that a .350 hitter will have a hitting streak of at least 56 games in a season. If there are about 4 such seasons per year in the future, we would expect to see such a streak, assuming the Bernoulli trials model, every $1/ (.0003 \times 4) = 833$ years.

The above argument does not take into account the incredible pressure that will surely be felt by any player who approaches the record. Joe DiMaggio must have felt pressure as well during the streak, but once he had broken the existing record of 44 consecutive games, the pressure probably abated quite a bit. Thus, any player who threatens DiMaggio's record will feel pressure for quite a bit longer than DiMaggio did. Of course, those who believe that Bernoulli trials are a good model for hitting in baseball might argue that the pressure under which a player finds himself in such a situation is irrelevant.

We performed the above simulation in 2005. In 2008, a similar simulation was carried out by S. Arbesman and S. Strogatz, and a summary of their simulation was published in the *New York Times* [3]. Quite a few readers wrote to the *Times* on this subject after the article was published. We will discuss a few of their responses below. In the Arbesman-Strogatz simulation, the authors found that in 42% of their simulated baseball histories, the longest hitting streak was at least 56 games in length. We note that this estimate is quite different than ours, but the difference is easily explained. The first difference is that in their simulation, they considered the period from 1871 to the

present. We used data from 1901 to the present. More than half of the longest streaks in their simulation occurred before 1901, and it is likely that an even higher percentage of those streaks that were longer than 56 games occurred in that era. The second difference is that in their simulation, they did not take into account, as we did, the variation in the number of at-bats that occur in a given player's record. In 2009, D. Rockoff and P. Yates [36] carried out a similar simulation in which they took into account this variation (so their simulation was very much like ours). They reported the probability of seeing at least one hitting streak of at least 56 games as 2.5%.

Many readers responded to the Arbesman-Strogatz article. One reader said "Numbers, no matter how statistically significant, do not measure DNA or individual eye/hand coordination." This was an argument against using statistics to model hitting in baseball. However, one could also argue that eye/hand coordination is expressed by an individual batter's batting average, so in fact statistics do give some information about an individual's abilities. The question that we are considering is how well a simple Bernoulli trials model fits the data in the real world.

Another reader asked about the design of the experiment. For example, were the abilities of the pitchers whom the batters faced varied over time? Were injuries and fatigue simulated as well? These are interesting questions. Our response is that one can always think of ways to make probability models more realistic. However, making a model more realistic almost always makes it more complicated as well. If a simple coin-tossing model fits a data set very well, and a more complicated model fits the data better, then the statistician must give up some simplicity to obtain greater accuracy. Each of these two approaches has some virtue. In this chapter, we show that even a very simple coin-tossing model fits the data very well in many instances. The fact that in such a model streaks can occur should help convince the reader that it may not be necessary to turn to more complicated models to explain the streaks observed in real data.

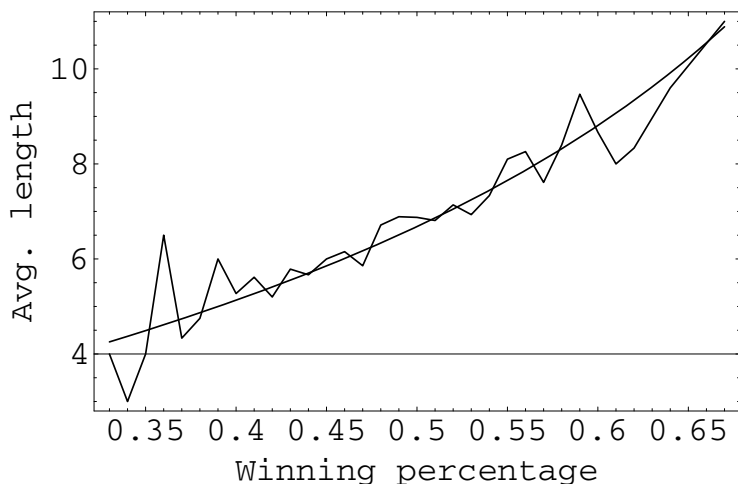


Figure 20. Actual and theoretical longest winning streak lengths

We now turn to team statistics; we seek to answer the question of whether teams exhibit streaky behavior. One way to check this is to recall that if a process has a positive autocorrelation (meaning that the probability of a success in a given trial increases if the previous trial is a success), then the lengths of the longest streaks will, on average, be longer than in the Bernoulli trials model. This will be true of both winning and losing streaks.

There are 390 team seasons in our data set. For each of these seasons, we computed the lengths of the longest winning and losing streaks. Then we grouped the team seasons by rounding the winning percentages to the nearest multiple of .01. (So, for example, a winning percentage of .564 is rounded to .56.) For each group, we calculated the average lengths of the longest winning and losing streaks. Figures 20 and 21 show these average lengths, as well as the theoretical average lengths in the Bernoulli trials model. (The theoretical lengths are shown as a solid curve.) We note that the fit between the actual and the theoretical longest streak lengths is very good, except perhaps at the ends of the graphs. But even at the ends, there is

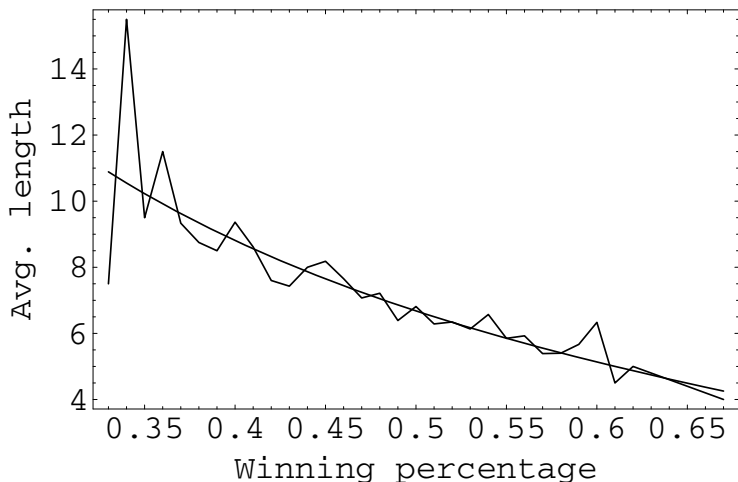


Figure 21. Actual and theoretical longest losing streak lengths

no bias. In addition, the large deviations at the ends should not be too surprising when one notes that the group sizes are so small. Of the 390 team seasons in the data set, only 9 have rounded winning percentages that lie outside of the interval $[\cdot36, \cdot64]$.

In the book *Curve Ball* [1], the authors, Jim Albert and Jay Bennett, study the question of team streakiness, in the sense of the block-Bernoulli process. Specifically, their model of streaky behavior is as follows: A hypothetical season is split into nine non-overlapping blocks of 18 games each (so in our notation, $a = b = 18$). Three winning percentages are computed, denoted by p_C , p_{av} , and p_H (for cold, average, and hot). The percentage p_{av} is the team's season winning percentage. The percentages p_C and p_H are defined to be .1 less than and more than p_{av} , respectively. So this is somewhat like our p_{min} and p_{max} , except that in their model, the block winning percentages are not uniformly chosen in the interval $[p_C, p_H]$; rather, there are only three possible block winning percentages. Finally, each of these three block winning percentages is assigned randomly to three

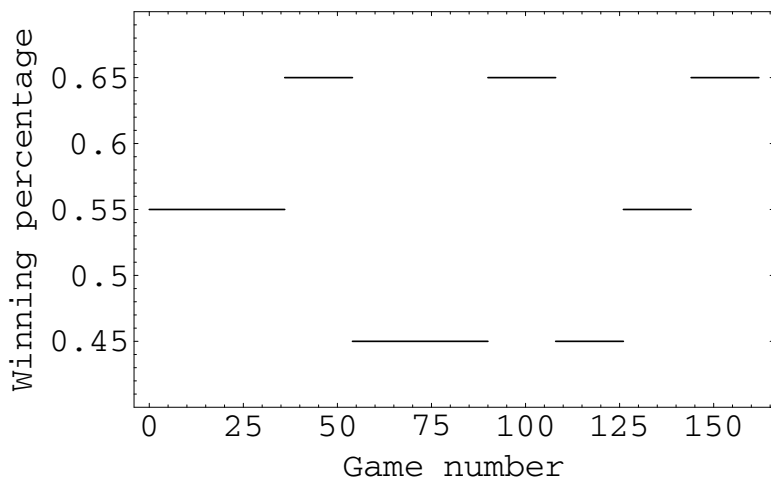


Figure 22. Block winning percentages for a streaky team

of the nine 18-game blocks. Figure 22 shows an example of the block winning percentages for a team whose season percentage is .55.

A team whose win-loss sequence arises from the above block-Bernoulli process will be said to be *streaky*, while one whose win-loss sequence arises from a Bernoulli process will be called *consistent*. Of course, we do not expect all teams to behave in one of these two ways; we need a parameter that can be estimated and that does a good job of distinguishing between these two models. As we said above, our windowed difference statistic does not distinguish very well between the two models we posited.

Albert and Bennett define their parameter, called *Black*, as follows. Given a team's win-loss sequence for a season, they compute the windowed winning percentages for all blocks of 12 games. (So in a 162-game season, there are 151 such blocks, starting at games 1 to 151.) They plot these percentages, along with a horizontal line whose y -value is the season winning percentage. An example of this plot is shown in Figure 23; the team is the 1984 Baltimore Orioles. The

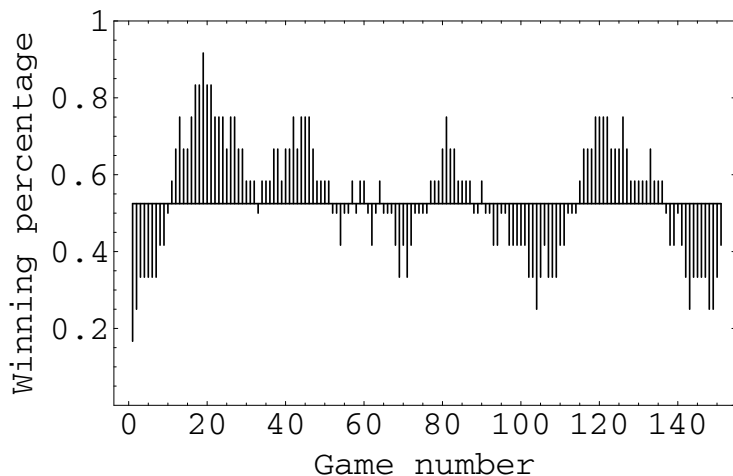


Figure 23. Windowed winning percentage for 1984 Baltimore Orioles

parameter Black is easy to describe in terms of this figure; it is the area of the black region.

It is clear that if a team is consistent, then Black will be very small, while if a team is streaky, Black will probably be large. Thus, this parameter might be able to distinguish between the block-Bernoulli and the Bernoulli models. To decide which model fits a given team better, Albert and Bennett compute the winning percentage of the team. Then they obtain, by simulation, the distribution of the parameter Black under both models. The actual value of Black for the team is computed and compared with the two distributions. For each distribution, a p-value is reported. If the first p-value of the observation is larger than the second p-value, then they claim that the first model fits the team's performance better than the second model.

In fact, they go further and use the ratio of the p-values as the odds that the team is consistent (or streaky). For example, if the p-values for a given team, under the consistent and streaky models, are .08 and .30, respectively, then they say that the probability that the

team is streaky is $.79 (= .30/ (.08 + .30))$. This language corresponds to a Bayesian approach for comparing the two models.

We will proceed in a different direction with the parameter Black. We are trying to determine whether the Bernoulli model does a good job of fitting the win-loss sequences in our data set. For each of the 390 teams in this data set, we can calculate a p-value for the observed value of Black, under the null hypothesis of the Bernoulli model. For each team, we use the winning percentage of that team, rounded to the nearest .01, to compute a simulated distribution for Black. If there are more streaky teams than can be explained well by the Bernoulli model, we should see more small p-values than expected, i.e. the set of p-values should not be uniformly distributed across the interval $[0, 1]$. (The reason that the p-values of streaky teams are small is because the parameter Black is very large for such teams, and large values of Black correspond to small p-values.) Figure 24 shows this set of p-values. The set has been sorted. If the set were perfectly uniformly distributed across $[0, 1]$, the graph would be the straight line shown in the figure. One sees that the set of p-values is very close to uniform, and in addition, there are certainly no more small p-values than expected under the Bernoulli model. In fact, we see that there are slightly more teams with large p-values than might be expected. Since large p-values correspond to small values of Black, this indicates that the number of very consistent teams is slightly more than expected.

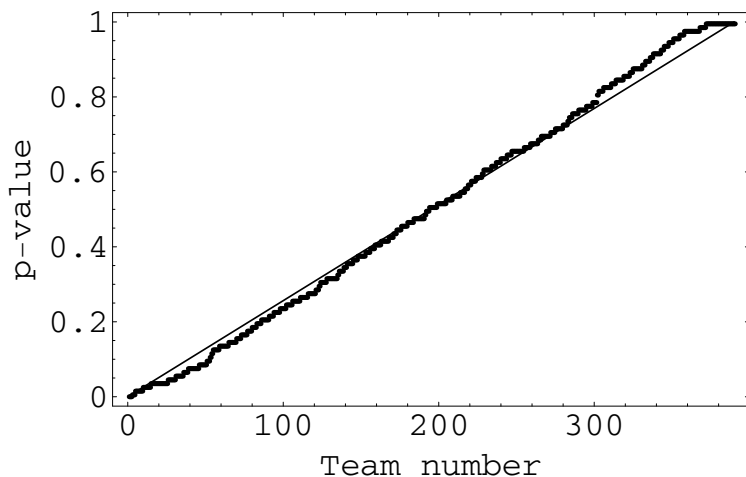


Figure 24. p-values for the Black parameter

Exercises.

1. The major league record for the number of consecutive hits by a batter (over several games) is 12, by Walt Dropo in 1952. How can we decide how surprising this record is? One way is to proceed as follows. First, we should understand how a hitter's batting average affects the probability that his longest streak of hits is at least 12.
 - (a) Using the approximation in Exercise 4.2, find the probability of a batter with a lifetime batting average of .300 achieving a streak of at least 12 hits in a row. Assume the number of at-bats in the hitter's career is 5000. Do the same calculation for batters with lifetime batting averages of .325 and .350. (You should obtain the values .0019, .0047, and .0109.) This shows that a .350 hitter is five times as likely to have such a streak as a .300 hitter. Of course, there are many more of the latter than the former in major league history.

- (b) Now suppose we restrict our attention to the 430 player seasons, mentioned above (and shown in Figure 17) in which the player's batting average was at least .340 and he had at least 300 at-bats. The average of all of the batters during these seasons was .358. If we assume the average number of at-bats was 450, then we can think of this set of seasons as one sequence of Bernoulli trials, with $n = 450 \cdot 430 = 193500$ and $p = .350$. Using the approximation in Exercise 4.2, show that in this sequence of Bernoulli trials, the probability that the longest success run is at least 12 is about 0.423. This shows that this record is not very surprising. (What is more surprising is that Dropo's batting average in 1952, when he set this record, was .279.)
2. On five different occasions, there have been home runs hit by four consecutive batters. Is this predicted by a runs calculation? To answer this question, use the following approximate data. Since 1900, there have been approximately 170000 games, with an average of about 80 plate appearances per game. In 2006, the probability of a home run in a plate appearance was .02857. These data come from the article by Cleary [8]. This last probability is probably greater than the corresponding probability for the period from 1900 to 2010; we will consider various alternative values below.

Using the approximation in Exercise 4.2, find the probabilities of a run of 3, 4, 5, and 6 home runs in

$$170000 \times 80 = 13600000$$

plate appearances. Use values for p of .02857, .025, and .020. After computing these approximations, are you surprised that there have been four home runs hit in succession on at least one occasion? Are you surprised there have never been five home runs hit in succession?

3. Can you describe a win-loss sequence for a baseball team that you would say exhibits streakiness, yet for which the parameter Black is very small?
4. Suppose that a .300 hitter averages four at-bats over a six-game stretch. Using the arithmetic-geometric inequality, prove that the probability he has at least one hit in each of the six games is maximized in the case that he has exactly four at-bats in each game.

5.2. Basketball. The most widely referenced article on streaks in sports is undoubtedly the paper by Gilovich, Vallone, and Tversky [16]. This article examines whether or not there is auto-correlation in the shooting sequences of basketball players. The authors surveyed fans, collected data from games at both the professional and college level, and carried out a set of experiments with amateur subjects. We will give some of their most intriguing findings here.

In a survey of more than 100 basketball fans, recruited from the student bodies of Cornell and Stanford, 91% believed that a player was more likely to make his next shot if he had made the last two or three shots than if he had missed the last two or three shots. In addition, 68% said the same thing when asked about free throws. The fans were then asked to consider a hypothetical basketball player who has a field goal percentage of 50% (i.e. he makes about one-half of his shots). They were asked to estimate the player's field goal percentage for those shots that occur after he has made a shot, and for those shots that occur after he has missed a shot. (These numbers are our p_1 and p_2 .) The average value of their estimate for p_1 was .61, and for p_2 it was .42. In addition, the former estimate was greater than the latter estimate for all of the fans in the survey.

The authors compared these estimates with data obtained from 48 games involving the Philadelphia 76ers in the 1980-81 season. The data consisted of the sequences of hits and misses for all of the field goal attempts for all of the players on the 76ers. Among the more interesting observations is the fact that for eight of the nine major

players on the team, the probability that they made a shot was higher after a miss than after a hit. For these players, the weighted mean of their field goal percentages (the average of their probabilities of making a field goal, weighted by the number of attempts) was .52. The weighted mean of their field goal percentages after a miss and a hit were .54 and .51, respectively.

In their paper, the authors performed other statistical tests, including a runs test, a test of stationarity, and a test for stability of shooting percentage across games. None of these tests provided any strong evidence for rejecting the Bernoulli model.

They also looked at free throw data from the Boston Celtics teams in the 1980-81 and 1981-82 seasons and again found no evidence of streakiness. Finally, they recruited members of the Cornell basketball teams (both men and women) to take part in a controlled shooting experiment. Each participant was asked to take 100 shots from a distance at which their accuracy was roughly 50%. They were paid money, where the amount was based both on how accurately they shot and how accurately they predicted their next shot. Once again, there was no evidence suggesting the existence of streakiness.

We now turn to another aspect of basketball in which there is some evidence that fans think there is streakiness. In much of the world, money is bet on almost all aspects of sports. In the United States, for example, one can bet on professional basketball in several ways. One of the most popular ways proceeds as follows. Professional odds-makers (called bookies) set a point spread before the beginning of a game. For example, if the Lakers are playing the Spurs, the bookies may say that the Spurs are favored by 5.5 points. A bettor can bet on either team; if he bets on the Lakers, he is given 5.5 points, meaning that if the Lakers win the game, or if they lose by 5 points or fewer, then the bettor wins. Conversely, if he bets on the Spurs, then the bettor wins only if the Spurs win by at least 6 points.

The bookies take 10% (called a vigorish, or vig) of the winning bets and all of the losing bets. In many cases, the bookies do not

set the point spread at their prediction for the game. The bookies' objective is to have an equal amount of money bet on both teams. To see why they want this, suppose first that \$10,000 is bet on each team. Then the bookies make \$1,000, regardless of which team beats the spread. However, if \$15,000 is bet on the Lakers and \$5,000 is bet on the Spurs and the Lakers beat the spread, then the bookies must pay out \$13,500 and only take in \$5,000, so they lose \$8,500.

In a paper that appeared in 1989 [5], Colin Camerer showed that the bettors believed in streaky behavior in the NBA (the professional basketball league in the United States). Camerer collected data on all games played in three seasons of the NBA (from 1983 to 1986). His data set consists of the scores of all of the games and the point spreads set by a popular bookie in Las Vegas.

At the beginning of every game (except for the first game played by a team in a season), each team has a winning streak or a losing streak. Each game is put into a subset depending upon the lengths of these two streaks. For example, if the first team has a three-game winning streak and the second team has a two-game losing streak, then the game is put into the $(+3, -2)$ subset. The longer of the two streaks determines which of the two teams is the "first" team; if the streaks are of equal length, then a coin flip determines which team is the first team. Thus, each game appears in exactly one subset. For each subset, the fraction of those games in which the first team beat the spread was calculated.

This fraction is very important for the following reason. Camerer's data show, for example, that the subsets of the form $(+4, k)$, with $1 \leq k \leq 4$, have a combined fraction of .46. This means that 46% of the teams with four-game winning streaks who played teams with winning streaks of length at most four managed to beat the spread. This can be taken as evidence that the bettors (and hence the bookies) overvalued the teams with four-game winning streaks. The question of whether such a fraction is significant can be answered using standard statistical methods. We will now show how this is done using the above data.

There were 159 games in the subsets that we are dealing with. Suppose that we assume the probability that the first team will beat the spread is .5. What is the probability that in 159 trials, the first team will actually beat the spread 46% of the time or less, or equivalently, in at most 73 games? The answer is about .17, which a statistician would not regard as significant.

The data can be pooled by considering all of the subsets (j, k) with j positive (and $j \geq |k|$). There were 1208 games of this type, and the first team beat the spread 47.9% of the time. We can once again ask the question: If a fair coin is flipped 1208 times, what is the probability that we would see no more than 579 ($= .479 \times 1208$) heads? We can calculate this exactly, using a computer, or we can find an accurate approximation, by recalling that the number of heads N_H is a binomially distributed random variable with mean equal to 604 ($= 1208 \times .5$) and standard deviation equal to 17.38 ($= \sqrt{1208 \times .5 \times .5}$). Thus, the expression

$$\frac{N_H - 604}{17.38}$$

has, approximately, a standard normal distribution. The value of this expression in the present case is -1.44. The probability that a standard normal distribution takes on a value of -1.44 or less is about .0749, which is thus the p-value of the observation. This observation is therefore significant at the 10% level, but not at the 5% level.

If one instead takes all of the subsets of the form (j, k) with $j \geq 3$ and $j \geq |k|$, there are 698 games, and the first team beat the spread in 318 of these games. This represents 45.6% of the games. What is the p-value of this observation? It turns out to be about .0095, which is a very small p-value.

One can also look at the corresponding pooled data for teams with losing streaks. There were 1140 games in the subsets of the form (j, k) , with $j \leq -1$ and $|j| \geq |k|$, and the first team beat the spread in 597 of these games, or 52.4% of the time. The p-value of this observation is .055. If we restrict our attention to those subsets for which $j \leq -3$ and $|j| \geq |k|$, there were 643 games, and the first

team beat the spread in 340 of these games, or 52.9% of the time. The p-value of this observation is .0721.

There are two more collections of subsets that are worth mentioning. The first is the collection of subsets of the form (j, k) , with j negative and k positive (and $|j| \geq |k|$). If one does not believe in the predictive value of streaks, then one would think that in this case bettors would undervalue the first team's chances of beating the spread, since the first team is on a losing streak and the second team is on a winning streak. Thus, the first team should beat the spread more than half the time. There were 670 games in this collection of subsets, and the first team beat the spread in 356 of these games, or 53.1% of the time, leading to a p-value of .0526. Interestingly, in the corresponding collection, in which the first team is on a winning streak and the second team is on a losing streak of equal or smaller length, the first team beat the spread in 358 games, which is 50.1% of the time. The p-value of this observation is, of course, .5.

Although only one of the six p-values reported above is less than .05, five of them are rather small, and thus there is some evidence that the bettors are overvaluing teams with winning streaks and undervaluing those with losing streaks.

If the astute bettor realizes that the average bettor is overvaluing teams with long winning streaks and undervaluing teams with long losing streaks, can he or she make money on this information? In Exercise 1, the reader is asked to show that if there is a 10% vig, then the bettor needs to win 52.4% of the time (assuming he is betting constant amounts) to break even. A few of the pooled subsets mentioned above had winning percentages that were at least 2.4% away from 50% in one direction or the other, meaning that had one made bets, in the correct direction, on all of the games in that pooled subset, one could have made money. Unfortunately, if one tests the hypothesis that the winning percentage in any of the pooled subsets is at least 52.4%, one finds that none of the results are significant at the 5% level.

The results for the 2003-04 season are qualitatively different than those found by Camerer. Through the All-Star break (February 15, 2004), teams having a winning streak, playing teams with equal or shorter streaks, beat the spread in 193 out of 370 games, or 52.2% of the games. Teams having a losing streak, playing teams with equal or shorter streaks, beat the spread in 179 out of 375 games, or 47.7% of the games. Thus, betting on teams with winning streaks, and betting against teams with losing streaks, would have resulted in a winning percentage of 52.2%. Note that this is the opposite of what happened in Camerer's data. Does the reader think that this change will persist, and if so, is it because the bettors have gotten smarter (i.e. have they incorporated the fact that streaks are over-rated into their betting practices)?

Exercise.

1. Show that if there is a 10% vig, and a bettor makes bets of a constant size, then the bettor needs to win 52.6% of the time to break even.

5.3. Horseshoes. The game of horseshoes differs in many ways from the games of baseball and basketball. In studying streakiness, some of these differences make it easier to decide whether the results in horseshoes diverge from those that would be predicted under the assumptions of the Bernoulli trials model.

In the game of horseshoes, two contestants face each other in a match. A match consists of an indefinite number of innings. In each inning, one player pitches two shoes, and then the other player pitches two shoes. The shoes are pitched at a stake that is 37 feet from the pitching area. If the shoe encircles the stake, it is called a ringer. A nonringer that is within 6 inches of the stake is called a "shoe in count." The former is worth three points and the latter is worth one point. If one player throws j ringers, and the other player throws k ringers, where $j \geq k$ and $j \geq 1$, then the first player gets $3(j - k)$ points, and in this case, shoes in count do not count. If neither player

throws a ringer, then the closest shoe in count is worth one point for the player who threw it. If that player threw two shoes in count that are closer than either of his opponent's shoes, that player gets two points. The first player to score 40 points or more is the winner of the match.

Unlike the game situations of baseball and basketball, the game situation of horseshoes does not affect a player's strategy. In addition, horseshoe players typically throw many times in a relatively short time period, and are attempting to do the same thing every time they throw.

Gary Smith has analyzed the data from the 2000 and 2001 Horse-shoe World Championships (see [39]). He was particularly interested in whether the data exhibit streakiness. In the cited article, Smith concentrated on doubles (i.e. two ringers thrown by one player in one inning) and non-doubles (i.e. everything else). At the championship level, players throw doubles about half the time. We show some of the data from this paper in Table 1.

Group	After 1 Nondouble	After 1 Double
Men 2000	.480	.514
Women 2000	.501	.544
Men 2001	.505	.587
Women 2001	.516	.573

Table 1. Frequency of doubles following non-doubles or doubles

Each line of the table represents 16 players. It can be seen that the players were more likely to throw a double following a double than following a non-double. The table gives the average frequencies over sets of players. A breakdown of the data shows that of the 64 players, 25 of the men and 26 of the women were more likely to throw a double following a double than following a non-double. (We will refer to this situation as positive auto-correlation, as before.) That this is evidence of streakiness can be seen as follows. If players

were not affected by their throws in the preceding inning, then about half of them would have positive auto-correlation. Continuing under the assumption of independence between innings, we see that the probability of observing as many as 51 of the 64 players with positive auto-correlation is the same as the probability that if a fair coin is flipped 64 times, it comes up heads at least 51 times. This probability is approximately .0000009.

At the championship level, the players' probabilities of throwing a double in any given inning is so high that using either the length of the longest run of doubles or the number of runs of doubles and non-doubles cannot be used to reject the null hypothesis of Bernoulli trials at the 5% level. For example, one of the players in the 2000 championships pitched 13 doubles in 14 innings. This means that there were either two or three runs of either type. But under the null hypothesis of Bernoulli trials, there is a probability of $2/14 \approx .143$ of two runs, since this happens if and only if the failure occurs in the first or last inning. So in this case, even if there are two runs, the p-value is much larger than .05.

Smith gets around this problem by calculating, under the null hypothesis, the expected number of runs by a given player in a given game and then tabulating the number of games in which the actual number of runs was above or below the expected number. Fewer runs than expected means that the data exhibits streakiness. Table 2 shows the number of games with fewer or more runs than expected for each championship, together with the p-values for the observations, under the null hypothesis.

Group	Fewer Runs	More Runs	p-value
Men 2000	129	107	0.0857
Women 2000	136	103	0.0191
Men 2001	137	98	0.0065
Women 2001	138	99	0.0067

Table 2. Number of games with fewer or more runs than expected

For each group in the data, the p-values are calculated by using the null hypothesis to compute the distribution of the number of runs and then computing the probability, using this distribution, of observing an outcome as extreme as the actual outcome. The small sizes of the p-values show that the null hypothesis should be rejected, i.e. there is streakiness in championship-level horseshoes.

5.4. Tennis. The game of tennis is interesting in probability theory because it provides an example of a nested set of Markov chains. The reader will recall that, roughly speaking, a Markov chain is a process in which there is a set of states and a transition matrix whose entries give the probabilities of moving from any state to any other state in one step. The chain can either be started in a specific state or it can be started with a certain initial distribution among the states. We will describe the various Markov chains that make up a tennis match and then give some results about tennis that follow from elementary Markov chain theory. We will then look at whether or not tennis is streaky at the professional level.

A tennis match is divided into sets; in most cases, the first person to win two sets is the winner of the match. (There are a few professional tournaments in which the winner is the first person to win three sets.) The set scores in an on-going tennis match can be thought of as labels in a Markov chain. The possible scores, from the point of view of one player, are 0-0 (at the beginning of the match), 1-0, 0-1, 1-1, 2-0, 2-1, 1-2, and 0-2. The last four of these states are said to be absorbing states, because once the match enters one of these states, it never leaves the state. The other four states are called transient states, because the match does not end in any of those states. (A Markov chain is said to be an absorbing chain if it contains at least one absorbing state and it is possible to go, in one or more steps, from every state to some absorbing state. In an absorbing Markov chain, a state is called transient if it is not an absorbing state.)

Suppose that player A has probability p of winning a set against player B . Then, for example, the probability that the Markov chain will move from state 0-1 to state 1-1 is p , while the probability that it will move from state 1-1 to state 1-2 is $1 - p$. The reader can check that if the above states are numbered from 1 to 8 and we denote by p_{ij} the probability of moving from state i to state j in one step, then the transition matrix $\mathbf{P} = (p_{ij})$ is given by

$$\mathbf{P} = \begin{matrix} & \begin{matrix} 0-0 & 1-0 & 0-1 & 1-1 & 2-0 & 2-1 & 1-2 & 0-2 \end{matrix} \\ \begin{matrix} 0-0 \\ 1-0 \\ 0-1 \\ 1-1 \\ 2-0 \\ 2-1 \\ 1-2 \\ 0-2 \end{matrix} & \begin{pmatrix} 0 & p & 1-p & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1-p & p & 0 & 0 & 0 \\ 0 & 0 & 0 & p & 0 & 0 & 0 & 1-p \\ 0 & 0 & 0 & 0 & 0 & p & 1-p & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix} \end{matrix}.$$

It turns out that we will want to split most of the above states into two states in our final model, because it is thought that the serve in tennis has a large effect on who wins a given game; thus, it will be important for us to record who is serving at the beginning of each set. We will discuss this in more detail below.

Each set in tennis consists of a number of games. The first player to win six games, if his or her opponent has won at most four games, wins the set. If the score is 6-5, the set continues until it is either 7-5, in which case the player with seven games wins, or it is 6-6, in

which case a tie-breaker is played. Whoever wins the tiebreaker wins the set. Thus a set can be modeled by an absorbing Markov chain.

Both a regular game and a tie-breaker are themselves absorbing Markov chains. In a game, the first person to win at least four points and be at least two points ahead wins the game. In a tie-breaker, the first person to win at least seven points and be at least two points ahead wins the tie-breaker.

At the professional level, the person who is serving a point has a much greater than even chance of winning the point. Thus, we need to take account of the serve. In each non-tie-break game, one person serves all of the points. The service rotation in a tie-break is more complicated; one player starts the tie-break by serving one point, and then the players take turns serving two consecutive points.

We show, in Figure 25, the probability that a player wins a game that he is serving, if he has a probability p_{point} of winning any given point. This graph shows that if $p_{\text{point}} = .6$, then the player will win the game with probability .74, and if $p_{\text{point}} = .8$, then the player will win the game with probability .98.

Now suppose that two players are playing a best-of-three set match in which the first player has a probability of p_{point} of winning any given point (so in this case, we are assuming that the serve does not affect the outcome of the point). In Figure 26, we show the probability of the first player winning the match, as a function of p_{point} . Note that even if $p_{\text{point}} = .51$, the probability that the first player wins the match is .748.

How are such probabilities determined? We will briefly describe the calculations that are needed, and the theorems on which these calculations are based. For more examples, and for proofs of the theorems, the reader is referred to [19].

If \mathbf{P} is the transition matrix of an absorbing Markov chain, then we can relabel the states so that the first set of states are the transient ones and the last set of states are the absorbing ones. If we do so, the matrix \mathbf{P} assumes the following canonical form:

$$\mathbf{P} = \begin{array}{c} \text{TR.} \\ \text{ABS.} \end{array} \left(\begin{array}{c|c} \text{TR.} & \text{ABS.} \\ \hline \mathbf{Q} & \mathbf{R} \\ \hline \mathbf{0} & \mathbf{I} \end{array} \right)$$

The four expressions \mathbf{Q} , \mathbf{R} , $\mathbf{0}$, and \mathbf{I} are rectangular submatrices of \mathbf{P} . If there are t transient states and r absorbing states, then \mathbf{I} , for example, is an r -by- r matrix. The reason that it is denoted by \mathbf{I} is that it is an identity matrix, since the probability of moving from one absorbing state to a different absorbing state is 0, while the probability of remaining in an absorbing state is 1.

If \mathbf{P} denotes the transition matrix for the set scores in tennis (which was given above), then the matrices \mathbf{Q} and \mathbf{R} are as follows:

$$\mathbf{Q} = \begin{array}{c} \begin{array}{cc} & \begin{array}{cccc} 0-0 & 1-0 & 0-1 & 1-1 \end{array} \end{array} \\ \begin{array}{c} 0-1 \\ 1-0 \\ 0-1 \\ 1-1 \end{array} \left(\begin{array}{cccc} 0 & p & 1-p & 0 \\ 0 & 0 & 0 & 1-p \\ 0 & 0 & 0 & p \\ 0 & 0 & 0 & 0 \end{array} \right), \end{array}$$

$$\mathbf{R} = \begin{array}{c} \begin{array}{cccc} 2-0 & 2-1 & 1-2 & 0-2 \end{array} \\ \begin{array}{c} 0-0 \\ 1-0 \\ 0-1 \\ 1-1 \end{array} \left(\begin{array}{cccc} 0 & 0 & 0 & 0 \\ p & 0 & 0 & 0 \\ 0 & 0 & 0 & 1-p \\ 0 & p & 1-p & 0 \end{array} \right).$$

Let \mathbf{P} be the transition matrix, in canonical form, for an absorbing Markov chain. The matrix \mathbf{N} , defined by the equation

$$\mathbf{N} = (\mathbf{I} - \mathbf{Q})^{-1},$$

is called the fundamental matrix for the chain. The following theorem shows one reason that \mathbf{N} is useful.

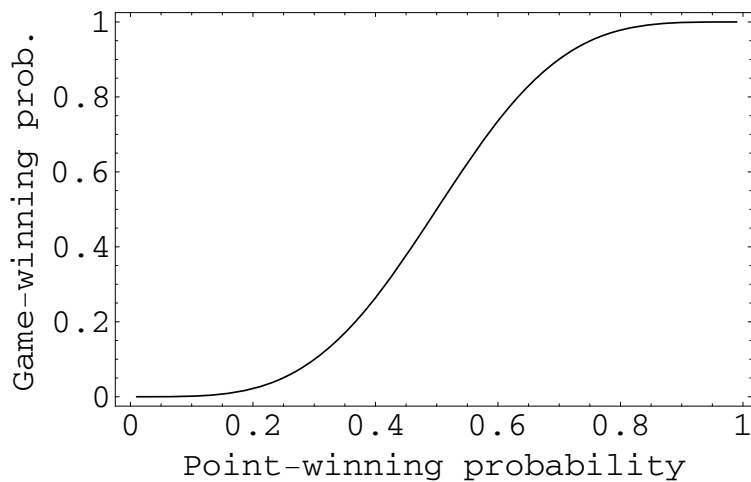


Figure 25. Game-winning vs. point-winning probabilities in tennis

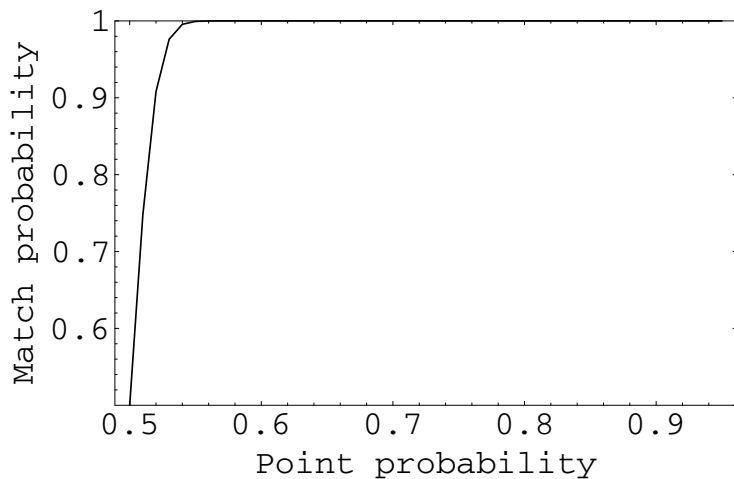


Figure 26. Match-winning vs. point-winning probabilities in tennis

Theorem 1. *Let b_{ij} be the probability that an absorbing chain will be absorbed in the j 'th absorbing state if it starts in the i 'th transient state. Let \mathbf{B} be the matrix with entries b_{ij} . Then \mathbf{B} is an t -by- r matrix, and*

$$\mathbf{B} = \mathbf{NR} ,$$

where \mathbf{N} is the fundamental matrix and \mathbf{R} is as in the canonical form.

As an example of how this theorem is used, suppose that the first player has probability $p = .6$ of winning a given set against his opponent. Then

$$\mathbf{Q} = \begin{pmatrix} 0 & .6 & .4 & 0 \\ 0 & 0 & 0 & .4 \\ 0 & 0 & 0 & .6 \\ 0 & 0 & 0 & 0 \end{pmatrix} ,$$

so one can calculate that

$$\mathbf{N} = \begin{pmatrix} 1 & .6 & .4 & .48 \\ 0 & 1 & 0 & .4 \\ 0 & 0 & 1 & .6 \\ 0 & 0 & 0 & 1 \end{pmatrix} .$$

Thus, the matrix $\mathbf{B} = \mathbf{NR}$ is given by

$$\mathbf{B} = \begin{matrix} & \begin{matrix} 2-0 & 2-1 & 1-2 & 0-2 \end{matrix} \\ \begin{matrix} 0-0 \\ 1-0 \\ 0-1 \\ 1-1 \end{matrix} & \begin{pmatrix} .36 & .288 & .192 & .16 \\ .6 & .24 & .16 & 0 \\ 0 & .36 & .24 & .4 \\ 0 & .6 & .4 & 0 \end{pmatrix} \end{matrix} .$$

The first row of this matrix is of particular interest, since it contains the probabilities of ending in each of the absorbing states, if the chain starts in the state 0-0 (i.e. the match starts with no score). We see that with the given value of p , the probability that the first player wins both of the first two sets is .36 and the probability that he wins the match is $.36 + .288 = .648$.

There are 43 states in the Markov chain representing a set of tennis. There are four absorbing states, with two corresponding to a win by the first player and two to a win by the second player. The

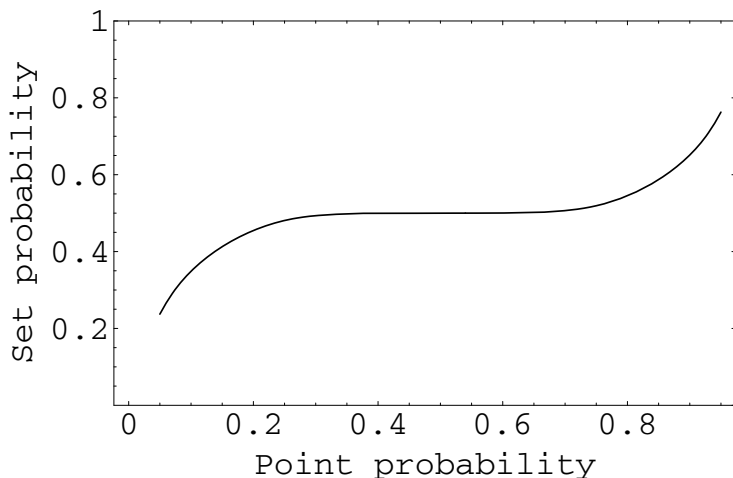


Figure 27. Set-winning probabilities for players of equal abilities

reason that we need two winning states for each player is that we need to keep track of who serves first in the subsequent set (if there is one). The transition probabilities are found by using the graph shown in Figure 25; for example, if the first player is serving, the game score is 2-1, and he has a probability of .6 of winning a given point, then the game score will become 3-1 with probability .74.

Once again, we are interested in the probability that the player who serves first wins the set. This time, there are two parameters, namely the probabilities that each of the players wins a given point when they are serving. We denote these two parameters by p_1 and p_2 . If we let $p_1 = p_2$, then Figure 27 shows the probability that the first player wins the set as a function of p_1 .

There is nothing very remarkable about this graph. Suppose instead that $p_1 = p_2 + .1$, i.e. the first player is somewhat better than the second player. In this case, Figure 28 shows the probability that the first player wins the set as a function of p_1 (given that the first player serves the first game). If, for example, $p_1 = .55$ and $p_2 = .45$, then the probability that the first player wins the set is .81.

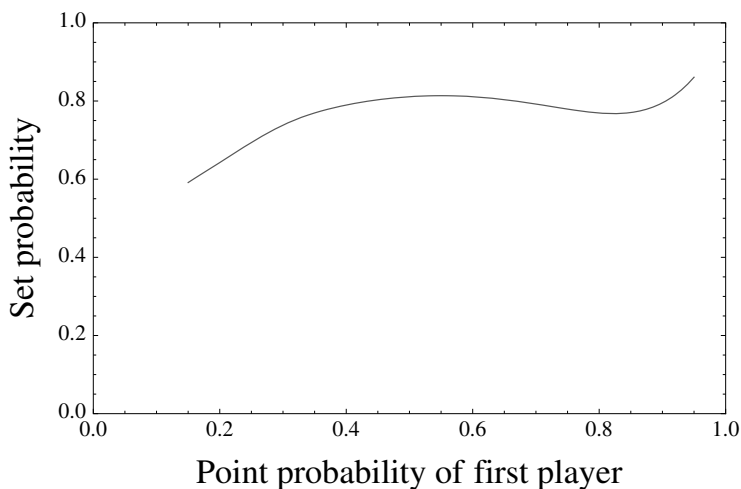


Figure 28. Set-winning probabilities for players of unequal abilities

Finally, suppose that two players play a best-of-three set match. The Markov chain corresponding to this situation has 11 states, because we must keep track of who begins serving each set. Suppose that the first player has probability p_{point} of winning a given point on his serve, and his opponent has probability $p_{\text{point}} - .1$ of winning a given point on his serve. Figure 29 shows the probability that the first player will win the match as a function of p_{point} . Note that even though there is a dip in the graph, the probability that the first player wins the match is always at least .9, even though the probability that he wins a given point is only slightly greater than his opponent's probability.

What kind of streakiness, if any, is evident in professional tennis? The above models show that even with slight differences in the players' abilities, many of the matches are likely to be one-sided. One can turn this around and say that the fact that there are so many close matches among the top professional tennis players means that these players must be very close in ability.

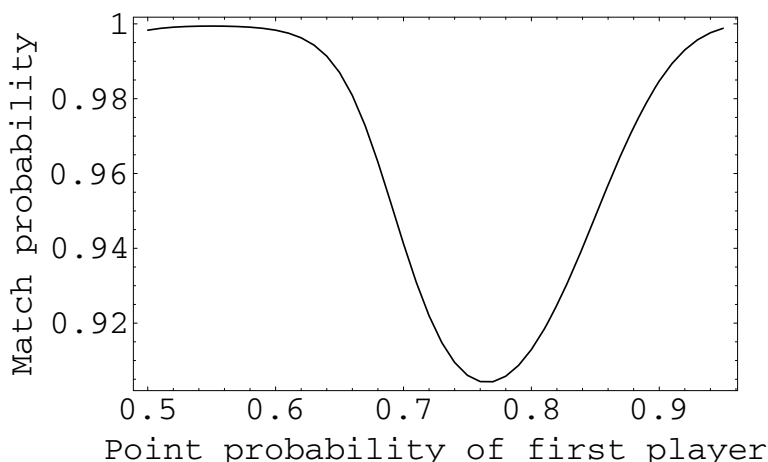


Figure 29. Match-winning probabilities for players of unequal abilities

In [22], Jackson and Mosurski describe two models, and some variations on these models, that deal with match scores in tennis. They were interested in the apparent overabundance of “heavy defeats” in tennis matches, i.e. best-of-three matches in which the loser does not win a set, or best-of-five matches in which the loser wins 0 or 1 sets.

One can model the sequence of sets in a tennis match in several ways. The simplest model is a Bernoulli trials model, in which a given player has the same probability p of winning each set. In Exercise 1, we ask the reader to determine the distribution of match scores under this assumption. In order to see if this model does a good job of explaining the data, we need to have some way of estimating p , since we certainly do not assume, even among professional tennis players, that $p = .5$.

A reasonable way to proceed is to use the rankings of the players to estimate p . Tennis rankings of professional players are determined by the players’ performance in the previous twelve months. These

rankings are updated each week. They are not perfect because, for example, if a very good player is injured and doesn't play for a while, his ranking can change by quite a bit, even though once he has recovered, he is presumably as good as he was before the injury.

We have obtained data from the website www.stevegetennis.com. We very much appreciate the work that was done to create the raw data files that we used. Our data set consists of all men's singles matches on the professional tour between 2000 and 2002.

There are at least two ways to use player rankings to estimate p . One idea is to fit the observed probabilities to the difference of the two players' ranks. A plot of the observed probabilities for the first sets in men's singles matches in 2002 versus the difference in the ranks of the two players is shown in Figure 30. We show only the results of matches where the rank difference was at most 300, since there are comparatively few matches for each difference larger than this. There are 3925 matches in the resulting data set. This figure also shows a best fit line. The reader can see that the fit is not very good; in fact, the correlation coefficient is .0298.

Jackson and Mosurski adopt a different strategy for estimating p . Denote by r and s the two players' ranks and assume that $r \leq s$. They then define $O(r, s)$ to be the odds (not the probability) that the better player wins the first set. The reason that they define these odds in terms of the first set, rather than the match, is that they are concerned about the outcomes of sets affecting the outcomes of later sets in the same match. If we are considering a model in which the sets are considered to be independent events, then we can use $O(r, s)$ as the odds that the better player wins any given set in the match. However, we must be careful with our use of $O(r, s)$ if we are not assuming the sets are independent.

Jackson and Mosurski next state the following assumed form for $O(r, s)$:

$$O(r, s) = (\text{ratio of ranks})^\alpha = \left(\frac{s}{r}\right)^\alpha.$$

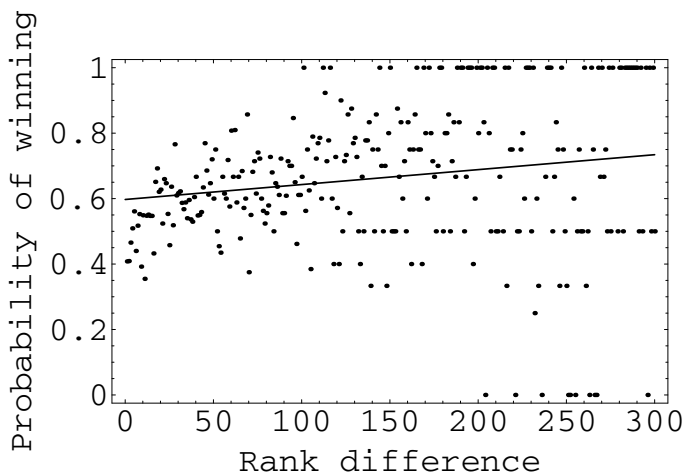


Figure 30. Observed winning probabilities vs. rank difference in 2002

As a function of $\log(\text{ratio})$, this is a line through the origin with slope α . The parameter α is to be determined by the data. The above relation is equivalent to the relation

$$\log(O(r, s)) = \log(\text{odds of success}) = \alpha \log(s/r) .$$

This is a line through the origin with slope α .

We carry out this regression for all matches in the years 2000 through 2002 in which the ratio of the ranks of the players is at most 7.4. (This is a completely arbitrary cut-off; it is about e^2 .) There were 12608 matches in this data set. We obtain a value of $\alpha = .480$ and a correlation coefficient of .363. Figure 31 shows the observed probability of winning the match versus the ratio of the ranks, together with the graph of the function $(\text{ratio})^\alpha$.

Now that we have a way to estimate the probability that a given player will win the first set against a certain opponent, we can proceed in several directions. First, we can explain and test the model of Jackson and Mosurski, which they call the odds model. Second,

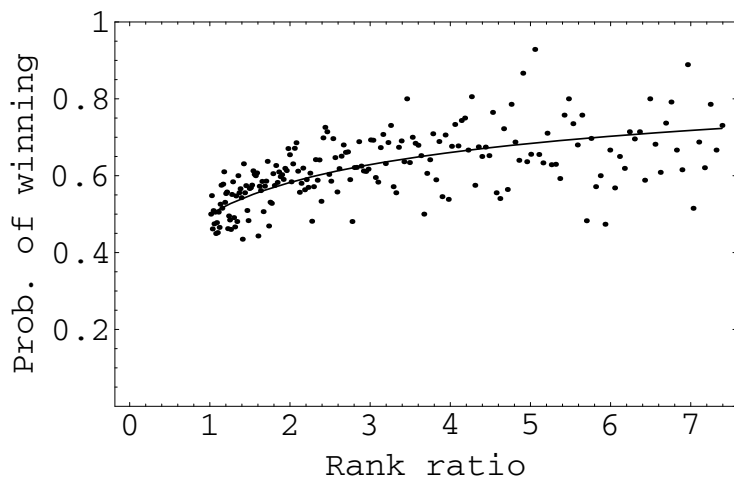


Figure 31. Observed winning probabilities vs. rank ratio in 2000-2002

we can test for independence among the sets in professional tennis matches.

In the odds model, it is assumed that the odds the better player will win a given set depend upon the set score of the match. We let O_{ij} denote the odds that the better player will win the next set if the set score is i (for the better player) to j . The discussion above gives us an estimate for O_{00} . (We are now suppressing the dependence on r and s .) The model asserts that each time a player wins a set, the odds for the next set change (in one direction or the other) by a factor of k , where k is to be determined by the data. So, for example, if the better player wins the first set, then the odds that he wins the second set equal kO_{00} . In general, we have

$$O_{ij} = k^{i-j} O_{00} .$$

Using our estimate for O_{00} , and taking logarithms, we have

$$\log(O_{ij}) = \alpha \log(s/r) + (i - j) \log(k) .$$

Note that if we take $k = 1$, then we have a model in which the sets are assumed to be independent.

The above equation can be used, along with the data, to estimate α and k . The procedure we now describe results in the maximum likelihood estimates for α and k . In a nutshell, a maximum likelihood estimate for a set of parameters for a family of distributions is the set of values for those parameters that leads to the highest probability, among all distributions in the family, that the actual data set would occur. As an example of this idea, suppose that we flip a coin 20 times and observe 12 heads. We wish to find the maximum likelihood estimate for p , the probability of a head on a single toss. We certainly hope that this estimate is 12/20. Let us see if this is the case. For each p between 0 and 1, we compute the probability that if a coin has probability p of coming up heads, it will come up heads 12 times in 20 tosses. This probability is

$$\binom{20}{12} p^{12} (1-p)^8 .$$

We wish to maximize this expression over all p . If we denote this expression by $f(p)$, then we have

$$f'(p) = \binom{20}{12} \left(12p^{11}(1-p)^8 - 8p^{12}(1-p)^7 \right) .$$

Setting this equal to 0 and solving for p , we obtain

$$p = \frac{12}{20} ,$$

as we hoped. (It is easy to check that this value of p corresponds to a maximum value of $f(p)$.) To reiterate, this means that the probability that we would actually obtain 12 heads in 20 tosses is largest if $p = 12/20$.

In the case of set scores in tennis, we treat each match as an independent event and use the relationship given above among the odds that the better player wins a given set, α , and k , to compute the probability that we would obtain the actual data set. Since the matches are assumed to be independent, this probability is the product of the probability that each match in the data set occurs, given

values of α and k . This calculation is easy to carry out with a computer. The number of completed best-of-three set matches in our data set is 12608. The maximum likelihood estimates for the parameters are

$$\alpha = .37$$

and

$$k = 1.78 .$$

Note that this value of α is quite a bit different than the one we obtained using just the first sets of the matches.

Now that we have estimated α and k , we can proceed, as in Jackson and Mosurski [22], to see how well the model fits the observed distribution of set scores in our data. We do this by simulating the matches many times on a computer, using the actual pairs of rankings in each match in our data set. When we did this 100 times, the distribution of match results was

$$\{5271.77, 1196.99, 881.01, 941.8, 1317.93, 2998.5\} ,$$

where the i 'th number in the above list is the average number of matches in which the outcomes of the sets are given by the i 'th element in the following list:

$$\{(1, 1), (1, 0, 1), (1, 0, 0), (0, 1, 0), (0, 1, 1), (0, 0)\} .$$

(In these outcomes, a 1 in the j 'th entry means the better player won the j 'th set.) This simulation should give us a reasonable approximation to the theoretical distribution.

We also carry out the above calculation under the assumption of independence of sets. To do this, we simply assume that the odds, O_{00} , of the better player winning the first set continue to hold for all subsequent sets in this match. This corresponds to setting $k = 1$ in the odds model. However, we need to re-estimate α , since our earlier maximum likelihood estimate of α was found simultaneously with our maximum likelihood estimate of k . When we do this, we obtain a new estimate: $\alpha = .41$. Using this value, a simulation gives the following

distribution for match results:

$$\{4459.56, 1742.79, 1235.17, 1241.94, 1746.56, 2181.98\}.$$

The graph in Figure 32 compares these two simulated distributions with the actual distribution of the matches in our data:

$$\{5453, 1157, 993, 813, 1338, 2854\}.$$

The simulation from the odds model gives a good fit to the actual

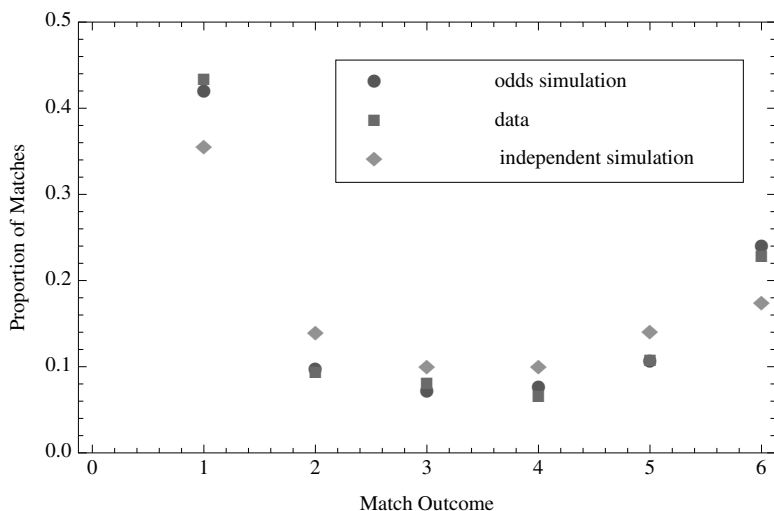


Figure 32. Simulated and actual distribution of match outcomes

data and is clearly superior to the independent sets model. We see that, consistent with the findings of the article [22], the independent set model underestimates the proportion of “heavy defeats”; these correspond to the set outcomes $(1, 1)$ and $(0, 0)$.

Before turning to the best-of-five set data, we digress slightly to suggest another way to estimate α and k , which involves minimizing the chi-squared goodness-of-fit statistic for the model. Recall that if we let Obs_i denote the i ’th count in the observed distribution of

match results, and Exp_i be the expected count under the model, then the chi-squared statistic is given by

$$\sum_{i=1}^6 \frac{(\text{Obs}_i - \text{Exp}_i)^2}{\text{Exp}_i}.$$

Under the null hypothesis that the model is correct, this statistic should have an approximate chi-squared distribution with $6 - 1 = 5$ degrees of freedom. See [19] for examples of goodness-of-fit tests. We note that it would not be advisable to use this test in the present context because the number of observations is so large. In such cases, the chi-squared test will detect even very small departures from a hypothesized model. (In fact, it would reject both the odds model and the independence model here.)

However, the chi-square statistic can still be used as a basis for estimation. Regarding the statistic as a function of α and k , we seek values for these parameters that minimize the statistic. This method is due to Karl Pearson, who was the originator of chi-squared procedures. Carrying out the minimization with our data leads to the values $\alpha = .41$ and $k = 1.75$, which are very close to our earlier maximum likelihood estimates.

Moving on, we repeated the maximum likelihood analysis for all of the best-of-five set matches played between 2000 and 2002. There were 1571 completed best-of-five set matches in the data set. The maximum likelihood estimates for α and k are .31 and 1.51. Noting that there are 20 possible outcomes in a best-of-five set match, we found that the chi-squared statistic for the odds model was 17.54. When we attempt to fit the independent set model to the data, we found a maximum likelihood estimate for α of .36 and a chi-squared statistic of 235.2. Thus we see again that the odds model does a much better job fitting the data than the independent set model.

There is another way to help us decide whether professional tennis matches are streaky. Consider a best-of-three set match. If each player has won one set in the match, then we might think that each of them is equally likely to win the third, and deciding, set. Of

the 12608 completed best-of-three set matches in our data set, 4301 required three sets to complete. Of these 4301 matches, 2331 matches were won by the player who won the second set. Thus, among the matches that took three sets to complete, if a player won the second set, the probability that he won the third set is .542. If the sets were independent, we might think that this probability should be close to .5. Since the actual probability is greater than .5, one could say that the player who won the second set has “momentum” or is “on a streak.” Before concluding that this is evidence of streaky behavior, we should consider the relative abilities of the players involved in the matches. It is possible that the reason that the winner of the second set does so well in the third set is because he is the better player.

If we simulate a distribution of set scores, using the value of $\alpha = .480$ (obtained earlier by fitting the odds model to just the first sets in the matches in the data), we obtain the following (the average of 20 simulations, using the same sets of opponents as in the actual data set):

$$\{4676.65, 1740.05, 1187.75, 1180.4, 1746.85, 2076.3\} .$$

In this simulated distribution, there are 5855 matches that took three sets to complete. In these matches, there were 2935 (which is almost exactly one-half of 5855) in which the player who won the second set won the third. Thus, we can discount any effect due to the relative ranks of the players.

It is also the case that this percentage does not change very much if we vary α . For $\alpha = .3, .31, \dots, .5$, the percentage stays between .497 and .505. Thus, we may assume that in this model, the player who won the second set has about a 50% chance of winning the third set.

How likely is it that in as many as 2331 matches out of 4301, the player who wins the second set wins the third set as well, given our assumption about independence of sets? This is essentially the same question as asking how likely a fair coin, if tossed 4301 times, will come up heads at least 2331 times. The number of heads in a long sequence of coin flips is approximately normal. If the coin is fair, and

the number of tosses is n , then the mean is $n/2$ and the standard deviation is $\sqrt{n}/2$. The standardized value corresponding to 2331 is

$$\frac{2331 - n/2}{\sqrt{n}/2} = 5.505 .$$

The probability that a standard normal random variable takes on a value greater than 5.505 is less than 2×10^{-8} . Thus, we can claim that the data exhibits streakiness.

Exercise.

1. Assume that two players are playing a best-of-three set tennis match and the first player has probability p of winning each set (i.e. the sets are independent trials). Find the distribution of the four possible match outcomes: 2-0, 2-1, 1-2, and 0-2.

6. Runs in the Stock Market

It is accepted that over the long run, most stocks go up in price. Thus, one way to model the price of a stock is to imagine that there is a line, with positive slope, that represents the long-term trend of the stock and then consider the stock's variation about this line. It is typically the case that instead of using the daily prices of the stock, one uses the logarithms of these prices. When logarithms are used, a straight trend line corresponds to the stock price changing by a constant percentage each day. (This is explained more fully in Chapter 2.) The slope of the trend line can be found by either fitting a straight line to the log price data or simply by using the starting and ending points of the data as anchor points for the line. In either case, the slope cannot be found unless one knows the data.

As might be imagined, an incredible amount of effort has been directed to the problem of modeling stock prices. In Chapter 2, we will focus on what is known about the distribution of the residual movements (those that remain after the trend line has been subtracted from the data) of stock prices. In this section, we will consider whether

the stock market exhibits streaky behavior (or perhaps other forms of non-randomness).

One obvious way in which stock prices could be streaky concerns their daily up-and-down motions. Most stocks have more up days than down days, so one might model them with coins whose success probabilities are greater than .5. Using the data, one can compute the observed success probability and then count the number of success (or failure) streaks of a given length, or of any length. Then one can compare these observed values with the theoretical values.

Our data set consists of daily prices, over a period of 14 years, of 439 of the stocks that make up the S&P 500 list. These prices have been adjusted to take into account dividends and stock splits. For this reason, many of the stock prices are very low near the beginning of the data set and hence are subject to rounding errors. We typically get around this problem by using only the part of the data set in which a stock's price is above some cutoff value. Also, we throw out all days for a given stock on which the stock price was unchanged.

Suppose that we want to compare the expected and the observed number of pairs of consecutive days in which a given stock's price went up on both days. If we have n data points, then we can define the random variable X_i to equal 1 if the i 'th and $(i + 1)$ 'st price changes are both positive, and 0 otherwise. Under the assumption that the signs of the daily price changes are independent events, the probability that $X_i = 1$ is just p^2 , where p is the observed long-range probability of success for that stock. Thus, it is easy to calculate the expected number of up-up pairs. There are $n - 1$ daily changes (since there are n data points), so there are $n - 2$ random variables X_i , each with the same distribution. Thus, the expected number of up-up pairs is just

$$(n - 2)p^2 .$$

Unfortunately, the X_i 's are not mutually independent. For example, if $X_6 = 1$ and $X_8 = 1$, then it is not possible for X_7 to equal 0. Nonetheless, the X_i 's are m -independent, for $m = 2$. The sequence

$\{X_i\}$ is m -independent if it is possible to partition the sequence into m subsets such that the random variables in each subset are mutually independent. In this case, we can use the partition $\{X_1, X_3, \dots\}$ and $\{X_2, X_4, \dots\}$. If a sequence of random variables is m -independent, then it satisfies a modified version of the Central Limit Theorem (see [7]). (The modification comes in how the variance of the sum of the random variables is calculated.)

Using this modified version of the Central Limit Theorem, we can transform the sum $X_1 + X_2 + \dots + X_{n-2}$ into a random variable that is approximately standard normal. (This is done in the usual way; we subtract the mean and divide by the standard deviation.) This gives us, for each stock, a z -value, i.e. a value of a standard normal distribution that represents how far above or below the mean the observed number of up-up pairs is. This approach works for any other pattern, such as down-down, or up-down-up, etc.

If stocks are streaky, then the z -values for up-up pairs (and for down-down pairs) should be significantly greater than 0. The same thing should be true for up-up-up triples. For each of several different patterns, we calculated the set of z -values for all 439 stocks in our data set. We used 1 as our cutoff value, meaning that for each stock, we used only those log prices after the last time the log price failed to exceed 1. The average number of log prices per stock that this gives us is 3861, or almost 15 years' worth of data.

At this point, we have 439 z -values. If these were drawn from a standard normal distribution, the distribution of their average would have mean 0 and standard deviation $1/\sqrt{439} \approx .0477$. Table 3 shows, for various patterns (up-up is denoted by UU, for example), the average z -value over our set of stocks. For each pattern, the distance of the average from 0, in units of standard deviation, and the percentage of stocks whose z -value is positive, are given. If the price movements are independent, then one would expect about half of the z -values to be positive.

We see that these stocks, in general, are anti-streaky. The numbers of UU and DD streaks are not significantly less than what would be

Pattern	Average z -Value	Standard Deviation Units	Percent Positive
UU	-0.017	-0.35	49.7
DD	-0.016	-0.34	50.1
UD	0.066	1.37	50.8
DU	0.079	1.65	51.0
UUU	-0.192	-4.03	43.3
UUUU	-0.335	-7.02	37.6
UUUUU	-0.461	-9.67	33.0
DDD	-0.280	-5.86	39.6
DDDD	-0.485	-10.16	35.3
DDDDD	-0.621	-13.02	27.6

Table 3. z -Values for various patterns in daily price changes

Pattern	Average z -Value	Standard Deviation Units	Percent Positive
UU	-0.452	-9.46	18.0
DD	-0.520	-10.90	17.5
UD	1.070	22.42	81.1
DU	1.064	22.30	81.5
UUU	-0.587	-12.29	19.6
UUUU	-0.577	-12.09	21.9
DDD	-0.606	-12.70	21.0
DDDD	-0.579	-12.14	21.9

Table 4. z -Values for various patterns in weekly price changes

expected, but the numbers of streaks of length three, four, and five of both types are significantly smaller than expected.

Table 4 shows the results of similar calculations involving weekly prices. Specifically, the closing price at the end of the first day of each week of trading was used. Once again, our set of stocks show strong evidence of being anti-streaky.

Pattern	Average z -Value	Standard Deviation	Units	Percent Positive
UU	-0.010	-0.21		50.1
DD	-0.074	1.56		51.7
UD	0.005	0.10		51.0
DU	0.002	0.04		51.0
UUU	-0.035	-0.74		44.9
UUUU	-0.020	-0.41		46.2
DDD	0.064	1.34		50.8
DDDD	-0.060	-1.26		43.7

Table 5. Simulated z -values for various patterns in weekly price changes

The results in these tables should be compared with simulated data from the model in which weekly changes for a given stock are mutually independent events. For each stock in our set, we used the observed probabilities of a positive and a negative weekly change in the stock price to create a simulated set of stock prices. Then, using the same algorithms as were used above, we calculated the average z -values for various patterns over our set of stocks. The results are shown in Table 5.

The simulated data from this model is much different than the actual data. This supports the observation that both daily and weekly stock prices are anti-streaky.

In their book [29], Andrew Lo and Craig MacKinlay discuss a parameter they call the variance ratio. Here is a description of this concept. Suppose that $\{X_i\}_{i=0}^n$ is a sequence of n logarithms of a stock's price. The time increment between successive values might be days, or weeks, or even something as small as minutes. The log price increments are the values $\{X_{i+1} - X_i\}_{i=0}^{n-1}$. A central question that concerns this sequence of log price increments is whether it can be

modeled well by a process in which the increments are assumed to be mutually independent.

Lo and MacKinlay begin by fixing a stock and a sequence of log prices $\{X_i\}_{i=0}^n$ for that stock. All prices have been adjusted for splits and dividends, and we assume in what follows that the time increment is a week. Next, they calculate the estimator

$$\hat{\mu} = \frac{1}{n} \sum_{i=0}^{n-1} (X_{i+1} - X_i) ,$$

which is an estimate of the average change per week in the logarithm of the stock price. This expression can be simplified to

$$\hat{\mu} = \frac{1}{n} (X_n - X_0) .$$

Next, they calculate the estimator

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=0}^{n-1} (X_{i+1} - X_i - \hat{\mu})^2 ,$$

which is an estimator of the variance of the weekly change in the logarithm of the stock price. If we assume for the moment that n is even, we could instead look at the estimator

$$\frac{1}{n/2} \sum_{i=0}^{n/2-1} (X_{2i+2} - X_{2i} - 2\hat{\mu})^2 ,$$

which is an estimator of the variance of the increments in even-numbered observations. Under the assumption that the increments are independent, the variance of the differences

$$\{X_{2i+2} - X_{2i}\}$$

is twice the variance σ^2 of the differences

$$\{X_{i+1} - X_i\} .$$

This provides a way to test whether the assumption of independent increments in stock prices is a reasonable one. One can compute both estimators for a given stock and see how close the ratio is to two.

Lo and MacKinlay slightly change the second variance estimator above, by dividing by two. So let us define

$$\hat{\sigma}_2^2 = \frac{1}{n} \sum_{i=0}^{n/2-1} (X_{2i+2} - X_{2i} - 2\hat{\mu})^2.$$

Under the assumption of independent increments, the theoretical value of σ_2^2 (for which $\hat{\sigma}_2^2$ is an estimator) equals the value of σ^2 . Thus, still under this assumption, the ratio of the estimators should be close to 1. Of course, one can, for any integer $q > 1$, define the estimator $\hat{\sigma}_q^2$ in the same way. Once again, under the assumption of independent increments, $\sigma_q^2 = \sigma^2$.

Lo and MacKinlay modify the above set-up in one additional way. Instead of using non-overlapping time increments in the definition of $\hat{\sigma}_q^2$, they use all of the time increments of length q , obtaining the following definition (note that we no longer need to assume that q divides n):

$$\hat{\sigma}_q^2 = \frac{1}{q(n-q+1)} \sum_{i=0}^{(n-q+1)} (X_{i+q-1} - X_i - q\hat{\mu})^2.$$

The test statistic for the variance ratio is denoted by $M_r(q)$ (the subscript refers to the fact that we are dealing with a ratio), and is defined by

$$M_r(q) = \frac{\hat{\sigma}_q^2}{\hat{\sigma}^2} - 1.$$

We have seen that under the assumption of independent increments, this statistic should typically be close to 0. In order for it to be a useful statistic in testing this assumption, we need to know something about how this statistic is distributed. Lo and MacKinlay show that for large n , the statistic $\sqrt{nq}M_r(q)$ is approximately normally distributed with mean 0 and variance

$$\frac{2(2q-1)(q-1)}{3q}.$$

(In fact, to be strictly accurate, they make one further modification to ensure that the individual variance estimators are unbiased. We

will ignore this modification, as the above asymptotic statement does not change when this modification is made.)

Lo and MacKinlay calculated the variance ratio for both individual stocks and for various sets of stocks on the New York and American stock exchanges. Their data consisted of weekly prices for 625 stocks from 1962 to 1985 (1216 weeks). For the values $q = 2, 4, 8$, and 16, their observed values of $M_r(q)$ for the set of all stocks in their data set were .30, .64, .94, and 1.05 respectively. These were all statistically different from 0 at the 5% level of significance.

They also computed the average variance ratio for the individual stocks. For the above values of q , these average variance ratios were $-.03, -.06, -.08$, and $-.11$. These observed values were not statistically different from 0 at the 5% level of significance. It is interesting that these observed values are all of opposite sign from those corresponding to the set of all stocks.

It is possible to describe a situation in which the variance ratio statistic of a stock (for $q = 2$, say) would be negative. Suppose that over blocks of two consecutive weeks, there were more up-down and down-up pairs than expected. This might mean that the average net change in the stock's price over two-week periods might be somewhat less than twice the average net change over one-week periods, and the same might be true of the average variance.

The above is admittedly only speculation. We carried out two sets of calculations with our stocks to see if any of this speculation is correct. First, we calculated the variance ratios for our stocks, for the values of q listed above. The average variance ratios were $-.060, -.103, -.135$, and $-.159$.

Our second set of calculations consisted of simulating our stocks by using a two-state Markov chain, in which the transition matrix entries are obtained from the observed probabilities for the four possible pairs up-up, up-down, down-up, and down-down. (These were discussed above.) For each of the above q , we simulated each of our stocks ten times, and computed the average variance ratio. Then we

computed the average, over all of the stocks, of these averages. The values obtained were -.088, -.125, -.142, and -.153. Furthermore, of the 439 stocks, the numbers whose average simulated variance ratio had the same sign as the actual variance ratio, for $q = 2, 4, 8$, and 16, were 386, 384, 373, and 358. So the Markov chain model behaves similarly, in terms of the variance ratio, to the actual stock prices.

7. Appendix

In this section, we will show how one can use generating functions to derive the theoretical distributions of several parameters of interest in this chapter. We will begin by deriving the distribution of the total number of runs of successes and failures in a sequence of n Bernoulli trials with probability of success p . (Although we do not claim any credit for the derivation of this distribution, we do not know of a reference in the literature. In Mood [31], the author derives a similar distribution in the case that the number of successes is known.) An example of this distribution, with $n = 50$ and $p = .2$, is shown in Figure 1.

One sees two roughly normal-shaped distributions. The reason for this is that the probability that the number of runs is even does not equal the probability that the number of runs is odd (except if $p = 1/2$). In fact, the number of runs is even if and only if the first trial and the last trial disagree, which happens with probability $2p(1 - p)$. Thus, the sum of the distribution values corresponding to even-valued outcomes equals this value, and the odd distribution values sum to 1 minus this value. We will derive this distribution using a method that will generalize to the distribution of the number of success runs (i.e. the number of runs in the first state) of a Markov chain with two states. The distribution for Markov chains was first derived by Zaharov and Sarmanov [45].

If we denote by $r_{n,k}$ the probability of k runs in n trials (this number also depends upon p), then we will show that the generating

function

$$r(x, y, p) = \sum_{n=1}^{\infty} \sum_{k=1}^n r_{n,k} x^n y^k$$

equals

$$\frac{xy(1 + 2px(-1 + y) - 2p^2x(-1 + y))}{1 - x + (-1 + p)px^2(-1 + y^2)} .$$

To derive this expression, we start by defining $f(n, p, k, S)$ to be the probability that in a sequence of n Bernoulli trials with success probability p , there are exactly k runs and the sequence ends in a success. The quantity $f(n, p, k, F)$ is defined similarly. Next, we define

$$G(x, y) = \sum_{i=1}^{\infty} \sum_{j=1}^i f(i, p, j, S) x^i y^j$$

and

$$H(x, y) = \sum_{i=1}^{\infty} \sum_{j=1}^i f(i, p, j, F) x^i y^j .$$

How are the coefficients related to one another? If a sequence of length at least two ends in a success, and it has k runs, and we chop off the last term in the sequence, we obtain a sequence of length $n - 1$ that still has k runs and ends in a success or it has $k - 1$ runs and ends in a failure. Thus, for $n \geq 2$, we have

$$f(n, p, k, S) = f(n - 1, p, k, S)p + f(n - 1, p, k - 1, F)q .$$

Similarly, we have, for $n \geq 2$,

$$f(n, p, k, F) = f(n - 1, p, k, F)q + f(n - 1, p, k - 1, S)p .$$

These recursions do not hold when $n = 1$, because f is not defined if the first parameter is 0. When $n = 1$, we have

$$f(1, p, 1, S) = p$$

and

$$f(1, p, 1, F) = q .$$

Using the recursions and initial conditions given above, we can write

$$G(x, y) = pxyH(x, y) + pxG(x, y) + pxy$$

and

$$H(x, y) = qxyG(x, y) + qxH(x, y) + qxy .$$

One can solve these equations for $G(x, y)$ and $H(x, y)$, obtaining

$$G(x, y) = \frac{pxy(1 - qx + qxy)}{1 - x + pqx^2(1 - y^2)}$$

and

$$H(x, y) = \frac{qxy(1 - px + pxy)}{1 - x + pqx^2(1 - y^2)} .$$

Finally, note that

$$r(n, k) = f(n, p, k, S) + f(n, p, k, F) ,$$

so

$$r(x, y, p) = G(x, y) + H(x, y) = \frac{xy(1 - 2pqx(1 - y))}{1 - x + pqx^2(1 - y^2)} ,$$

giving us the required expression for $r(x, y, p)$.

The above expression for the generating function $r(x, y, p)$ allows one to write an exact expression for $r_{n,k}$. The expression for $r(x, y, p)$ is of the form

$$xy \frac{A + Bx}{1 + Dx + Ex^2}$$

where the capital letters represent expressions that do not involve x . Using the algebraic method of partial fractions (usually first learned in calculus to help integrate rational functions), one can write this expression as

$$xy \left(\frac{F}{1 - Ix} + \frac{J}{1 - Kx} \right) .$$

The two summands can be rewritten, using geometric series, to obtain

$$xy(F + FIx + FI^2x^2 + \dots + J + JKx + JK^2x^2 + \dots) .$$

We want the coefficient of $x^n y^k$ in this expression. The coefficient of x^n is

$$y(FI^{n-1} + JK^{n-1}) .$$

This is a series involving y but not x . We want the coefficient of y^k in this series. The answer, which is obtained after some gruesome

algebra (best performed by a computer algebra package) is as follows. If $p \neq 1/2$ and k is odd, then

$$\begin{aligned}
 r_{n,k} = & \frac{1}{2^{n-1}} \left[(1-4pq) \sum_{v=0}^{(k-1)/2} \left[\binom{-1/2}{v} (1-4pq)^{-1/2} \left(\frac{4pq}{1-4pq} \right)^v \right. \right. \\
 & * \left(\sum_{\substack{u=1 \\ u \text{ odd}}}^{n-1} \binom{n-1}{u} \binom{u/2}{(k-1-2v)/2} (1-4pq)^{u/2} \right. \\
 & * \left. \left. \left(\frac{4pq}{1-4pq} \right)^{(k-1-2v)/2} \right) \right] + \sum_{\substack{u=0 \\ u \text{ even}}}^{n-1} \binom{n-1}{u} \binom{u/2}{(k-1)/2} \\
 & * (1-4pq)^{u/2} \left(\frac{4pq}{1-4pq} \right)^{(k-1)/2} \Bigg] ,
 \end{aligned}$$

while if $p \neq 1/2$ and k is even, then

$$\begin{aligned}
 r_{n,k} = & \frac{1}{2^{n-1}} (4pq) \sum_{v=0}^{(k-2)/2} \left[\binom{-1/2}{v} (1-4pq)^{-1/2} \left(\frac{4pq}{1-4pq} \right)^v * \right. \\
 & \left(\sum_{\substack{u=1 \\ u \text{ odd}}}^{n-1} \binom{n-1}{u} \binom{u/2}{(k-2-2v)/2} (1-4pq)^{u/2} \right. \\
 & * \left. \left. \left(\frac{4pq}{1-4pq} \right)^{(k-2-2v)/2} \right) \right] .
 \end{aligned}$$

If $p = 1/2$, the expression for $r_{n,k}$ is much simpler (see Exercise 3). These expressions were used to generate Figure 1.

We can use the expression for $r(x, y, p)$ to calculate the mean (and variance) of the distribution. We recall that for fixed n , the mean of the distribution $\{r_{n,k}\}$ equals

$$\sum_{k=1}^n k r_{n,k} .$$

The value of this sum can be obtained from the generating function $r(x, y, p)$ by using calculus. If we compute the partial of $r(x, y, p)$

with respect to y , and then set $y = 1$, we obtain the expression

$$\sum_{n=1}^{\infty} \sum_{k=1}^n k r_{n,k} x^n ,$$

which can be written as

$$\sum_{n=1}^{\infty} x^n \sum_{k=1}^n k r_{n,k} .$$

Thus the mean of the distribution for sequences of length n is just the coefficient of x^n in the above expression. The point is that we do not need to use the formulas for $r_{n,k}$ to calculate the mean. Rather, we use the closed-form expression for $r(x, y, p)$, and apply the ideas above to this expression.

If we perform these calculations, we obtain the expression

$$\frac{x}{1-x} + \frac{2pqx^3}{(1-x)^2} + \frac{2pqx^2}{1-x} .$$

Using the facts that

$$\frac{1}{1-x} = 1 + x + x^2 + \dots$$

and

$$\frac{1}{(1-x)^2} = 1 + 2x + 3x^2 + \dots ,$$

in a suitable interval containing the origin, we can expand each of the three summands above as series; they are, respectively,

$$\begin{aligned} x + x^2 + x^3 + \dots , \\ 2pq(x^3 + 2x^4 + 3x^5 + \dots) , \end{aligned}$$

and

$$2pq(x^2 + x^3 + x^4 + \dots) .$$

Now we can easily write down the coefficient of x^n ; it is

$$1 + 2pq(n-2) + 2pq = 1 + 2pq(n-1) ,$$

if $n \geq 2$.

There is an easy way to check this. In fact, the calculation below is an easier way to *find* the mean in this case, but the above method can be used to find other moments (including the variance) and the

calculation below does not generalize. Let X_i , for $1 \leq i \leq n-1$, denote the random variable that is 1 if the i 'th and $(i+1)$ 'st outcomes disagree. Then the average number of runs is just

$$\sum_{i=1}^{n-1} X_i ,$$

so

$$\mu = \sum_{i=1}^{n-1} E(X_i) .$$

But for each i , the probability that $X_i = 1$ is just $2pq$, so for each i ,

$$E(X_i) = 2pq ,$$

so the average number of runs is just

$$= 1 + 2p(1-p)(n-1) .$$

In the Markov chain model, the situation is more complicated(!). We will not go into the details here, but rather give an outline of how to proceed. We write

$$f_{S,S}(n, k)$$

for the probability that a sequence of n trials begins and ends with a success and has k runs. The quantities $f_{S,F}$, $f_{F,S}$ and $f_{F,F}$ are defined similarly. We define the generating function

$$G_{S,S}(x, y) = \sum_{n=1}^{\infty} \sum_{k=1}^{\infty} f_{S,S}(n, k) x^n y^k ;$$

the functions $G_{S,F}$, $G_{F,S}$, and $G_{F,F}$ are defined in a similar manner. One can show that

$$\begin{aligned} G_{S,S}(x, y) &= p_2 xy G_{S,F}(x, y) + p_1 x G_{S,S}(x, y) + xy , \\ G_{S,F}(x, y) &= (1 - p_1) xt G_{S,S}(x, y) + (1 - p_2) x G_{S,F}(x, y) , \\ G_{F,S}(x, y) &= p_2 xy G_{F,F}(x, y) + p_1 x G_{F,S}(x, y) , \\ G_{F,F}(x, y) &= (1 - p_1) xy G_{F,S}(x, y) + (1 - p_2) x G_{F,F}(x, y) + xy . \end{aligned}$$

Note that two of these equations are homogeneous (there are no summands on the right-hand sides that do not involve the functions $G_{*,*}$). Thus, for example, using the second equation above, it is easy

to find $G_{S,F}$ once we have found $G_{S,S}$. In addition, by switching the roles of success and failure, one sees that it is easy to find $G_{F,F}$ once we know $G_{S,S}$ (and $G_{F,S}$ once we know $G_{S,F}$).

Thus, we need to find only the coefficients $f_{S,S}(n, k)$, in a manner similar to the one used to obtain the results for Bernoulli trials. One can show that

$$G_{S,S}(x, y) = \frac{xy}{1 - (p_2(1 - p_1)x^2y^2)/(1 - (1 - p_2)x) - p_1x}.$$

Note that since the sequences corresponding to $G_{S,S}(x, y)$ begin and end with a success, the only positive probabilities are those corresponding to odd k . If k is odd, then

$$\begin{aligned} f_{S,S}(n, k) = & \frac{1}{2^{n-1}} \left[-(1-p_1-p_2) \sum_{j=0}^{\lfloor \frac{(n-2)}{2} \rfloor} \binom{n-1}{2j+1} (1+p_1-p_2)^{n-1-(2j+1)} \right. \\ & * \left(\sum_{i=\frac{k-1}{2}}^j \binom{j}{i} (1+p_1-p_2)^{2(j-i)} (-4)^i (-(1-p_1)p_2)^{(k-1)/2} \right. \\ & * (p_1(1-p_2))^{i-(k-1)/2} \binom{i}{(k-1)/2} \Bigg) + \sum_{j=0}^{\lfloor \frac{n-1}{2} \rfloor} \binom{n-1}{2j} (1+p_1-p_2)^{n-1-2j} \\ & * \left(\sum_{i=\frac{k-1}{2}}^j \binom{j}{i} (1+p_1-p_2)^{2(j-i)} (-4)^i (-(1-p_1)p_2)^{(k-1)/2} \right. \\ & \left. \left. * (p_1(1-p_2))^{i-(k-1)/2} \binom{i}{(k-1)/2} \right) \right]. \end{aligned}$$

One can use the recursions relating $G_{S,S}(x, y)$ and $G_{S,F}(x, y)$ given above to obtain the following formula for $f_{S,F}(n, k)$:

$$f_{S,F}(n, k) = (1 - p_1) \sum_{l=0}^{n-1} (1 - p_2)^l f_{S,S}(n - l - 1, k - 1).$$

These formulas are useful in plotting the distributions (depending upon whether the first trial is a success or a failure) of the number of runs. In Figure 2 a plot of the distribution for $n = 50$, $p_1 = .3$, and $p_2 = .1$, where the first trial is a success.

We note that the fixed probability vector of the transition matrix in the Markov chain model equals

$$\left(\frac{p_2}{1 - p_1 + p_2}, \frac{1 - p_1}{1 - p_1 + p_2} \right).$$

This means that the long-term proportion of successes equals

$$\frac{p_2}{1 - p_1 + p_2}.$$

Note that in this model, we have three quantities that can be estimated from the data, namely p_1 , p_2 , and the long-term proportion of successes (which we will call p), but there is a relation among them. We have no idea how a statistician would deal with this, since it is unlikely that the observed values satisfy the relation.

In order to create a distribution for the Markov chain model, we will weight the distributions corresponding to starting with a success or a failure by p and $(1 - p)$. The resulting generating function is

$$p \left(G_{S,S}(x, y) + G_{S,F}(x, y) \right) + (1 - p) \left(G_{F,S}(x, y) + G_{F,F}(x, y) \right).$$

If we do this, the expected number of runs equals

$$n\Omega_1 + p\Omega_2 + (1 - p)\Omega_4 + (p\Omega_3 + (1 - p)\Omega_5)(p_1 - p_2)^{n-1},$$

where

$$\begin{aligned} \Omega_1 &= \frac{2p_2(1 - p_1)}{1 - p_1 + p_2}, \\ \Omega_2 &= \frac{2 - 4p_1 + 2p_1^2 - p_2 + 3p_1p_2 - 2p_1^2p_2 - p_2^2 + 2p_1p_2^2}{(1 - p_1 + p_2)^2}, \\ \Omega_3 &= \frac{-1 + 2p_1 - p_1^2 + p_2 - p_1p_2}{(1 - p_1 + p_2)^2}, \\ \Omega_4 &= \frac{1 - 2p_1 + p_1^2 - p_2 + 3p_1p_2 - 2p_1^2p_2 + 2p_1p_2^2}{(1 - p_1 + p_2)^2}, \\ \Omega_5 &= \frac{p_2(1 - p_1 - p_2)}{(1 - p_1 + p_2)^2}. \end{aligned}$$

Thus, the expected number of runs is asymptotic to

$$n\Omega_1 + \Omega_4.$$

If one lets $p_1 = p_2$, so that the Bernoulli model is obtained, one can check that the result agrees with the one already obtained.

The two graphs show what one might expect, namely that the expected number of runs in the Markov model is less than in the Bernoulli trials model. It is easy to show this; one need only show that if $p_1 > p_2$, then

$$\Omega_1 < 2p(1 - p) .$$

Next, it might be nice to show that the runs distribution corresponding to $G_{S,S}(x, y)$ is asymptotically normal. One might need to do this to discuss the power of the test that compares the two models.

We also might try writing approximation algorithms for the calculation of these probabilities, since the probabilities fall off rapidly and it is time-consuming to calculate the sums using the exact expressions.

The distribution of the length of the longest success run in the Markov chain model can be calculated using recursions. We define $A(n, x, k, p_1, p_2)$ to be the probability that a Markov sequence, beginning with a success, has no success run exceeding x in length, and ends with k successes, for $0 \leq k \leq x$. This function satisfies the following equations. First, if $n = 1$, then the function is 1 if $k = 1$ and 0 otherwise. If $n > 1$, then if $k > 1$,

$$A(n, x, k, p_1, p_2) = A(n - 1, x, k - 1, p_1, p_2)p_1 ,$$

since a sequence of length n that ends in k successes is obtained from one of length $n - 1$ that ends in $k - 1$ successes by adding one success to the end, and this happens with probability p_1 , since $k > 1$. If $n > 1$ and $k = 1$, then

$$A(n, x, 1, p_1, p_2) = A(n - 1, x, 0, p_1, p_2)p_2 ,$$

since in this case we are adding a success to the end of a sequence whose last state was the failure state. Finally, if $n > 1$ and $k = 0$,

then

$$A(n, x, 0, p_1, p_2) = \left(\sum_{j=1}^x A(n-1, x, j, p_1, p_2)(1-p_1) \right) + A(n-1, x, 0, p_1, p_2)(1-p_2),$$

since in this case either there were j successes immediately preceding the last trial, which was a failure, for some j between 1 and x , or else the penultimate trial was also a failure.

These equations allow one to compute the values of the function A . To obtain the desired distribution, namely the set of probabilities that a Markov sequence, beginning with a success, has longest success run exactly equal to x , we compute the quantity

$$\sum_{k=0}^x A(n, x, k, p_1, p_2) - \sum_{k=0}^{x-1} A(n, x-1, k, p_1, p_2).$$

This adds the weights of all of the sequences with longest success run at most x and subtracts from this the weights of all of the sequences with longest success run at most $x-1$.

7.1. Doubletons. We now turn to the question of the distribution of d , the number of doubletons (i.e. consecutive pairs of successes in a sequence of trials). The reason for our interest in this quantity is because, as was stated in Exercise 4.1, the number d is closely related to p_1 , the probability of a success following a success, and p_1 is of interest when studying autocorrelation. We will first consider asymptotic behavior of the distribution of the number of doubletons, in both the Bernoulli trials model and the Markov model, and then we will derive a recursion for the exact distribution in the Bernoulli case.

Since the Bernoulli trials model is a special case of the Markov model, we will work with the Markov model. In this model, the probabilities that a success follows a success or a failure are defined to be, respectively, p_1 and p_2 .

We are interested in the distribution of d for large values of n , the length of the sequence. It turns out that if n is large, it does

not matter very much in which state the Markov chain started (i.e. whether the first trial resulted in a success or a failure). We have stated above that the long-term fraction of the time that the chain is in state 1 is equal to

$$\frac{p_2}{1 - p_1 + p_2}.$$

Note that if $p_1 = p_2$ (i.e. we are in the Bernoulli model) then this expression reduces to p_2 ($= p_1$).

We can find the distribution of the number of doubletons in a sequence of length n generated by the Markov process by considering the related Markov chain that consists of four states: SS, SF, FS, and FF. The first state, SS, means that the first two trials in the original sequence are both successes. It is straightforward to determine the transition probabilities for this new Markov chain. For example, if the chain is in state SS, it stays in state SS or moves to state SF depending upon whether the next trial in the original chain is a success or a failure. Thus, these two transitions occur with probability p_1 and moves to state SF with probability $1 - p_1$.

In the new Markov chain, we wish to know the distribution of the number of times $Y_{SS}^{(n)}$ that the chain is in state SS in the first n trials. Of course, this distribution depends upon the starting state (or starting distribution), but it turns out that the limiting distribution is independent of the starting state. The Central Limit Theorem for Markov Chains (see [24], p. 89) states that $Y_{SS}^{(n)}$ is asymptotically normally distributed, with a mean and standard deviation that are straightforward to calculate. Examples of how the mean and standard deviation are calculated are given in [24]. In the present case, the asymptotic value of the mean and standard deviation are

$$\frac{np_1p_2}{1 - p_1 + p_2}$$

and

$$\frac{np_1p_2(p_1^2(-1 + p_2) + (1 + p_2)^2 - p_1p_2(3 + p_2))}{(1 - p_1 + p_2)^3}.$$

Now suppose that we have a process that presents us with a sequence of successes and failures and suppose that the observed probability of a success is .3. We assume that $p = .3$, and we wish to test the hypothesis that $p_1 = .3$ against the alternative hypothesis that $p_1 > .3$. To carry out this test, using the number of doubletons as a parameter, we first choose an acceptance region around the value np^2 in which 95% of the values will fall. (Remember, in the case that the null hypothesis is true, then $p_1 = p_2 = p$, and we are dealing with Bernoulli trials.) Since $Y_{SS}^{(n)}$ is asymptotically normal, it is easy to pick this acceptance region. It is an interval of the form $[0, c]$, because the form of the alternative hypothesis is a one-way inequality. The number c is np^2 plus 1.65 times the standard deviation of $Y_{SS}^{(n)}$. We obtain the value of

$$c = n(.3)^2 + 1.65\sqrt{.1197n} .$$

This leads us to the distributions shown in Figures 4 and 5.

We now give an outline of the method used to find the exact distribution of the number of doubletons in the Bernoulli model, with parameters n and p . We define $r(n, k, p)$ and $s(n, k, p)$ to be the probabilities of exactly k doubleton successes in n trials, with success probability p , and with the sequence ending in a failure or a success, respectively. Then the following recursions hold:

$$r(n, k, p) = q * r(n - 1, k, p) + q * s(n - 1, k, p) ,$$

$$s(n, k, p) = p * r(n - 1, k, p) + p * s(n - 1, k - 1, p) .$$

The first of these equations says that sequences of length n with exactly k doubletons and which end in a failure arise from sequences of length $n - 1$ with exactly k doubletons by adding a failure to the end. The second equation says that sequences of length n with exactly k doubletons and which end in a success arise by adding a success to the end of either a sequence of length $n - 1$ with exactly k doubletons, ending in a failure, or to the end of a sequence of length $n - 1$ with exactly $k - 1$ doubletons, ending in a success.

Next, define

$$R(p, x, y) = \sum_{n=1}^{\infty} \sum_{k=0}^{n-1} r(n, k, p) x^n y^k$$

and

$$S(p, x, y) = \sum_{n=1}^{\infty} \sum_{k=0}^{n-1} s(n, k, p) x^n y^k .$$

The distribution of the number of doubletons in sequences of length n is given by the set of values

$$\{r(n, k, p) + s(n, k, p)\} .$$

The above recursions, together with attention to initial conditions, give rise to the following functional equations:

$$R(p, x, y) = qx \left(R(p, x, y) + S(p, x, y) \right) + qx$$

and

$$S(p, x, y) = px \left(R(p, x, y) + S(p, x, y) \right) + px .$$

These equations can be solved for $R(p, x, y)$ and $S(p, x, y)$:

$$R(p, x, y) = -1 + \frac{1 - pxy}{1 - qx - pqx^2 - pxy + pqx^2y}$$

and

$$S(p, x, y) = \frac{px}{1 - qx - pqx^2 - pxy + pqx^2y} ,$$

where $q = 1 - p$.

We can write

$$1 + R(p, x, y) = \frac{A}{1 - \gamma_1 x} + \frac{B}{1 - \gamma_2 x} ,$$

for suitable choices of constants A , B , γ_1 , and γ_2 . The right-hand side can be expanded to yield

$$(A + B) + (A\gamma_1 + B\gamma_2)x + (A\gamma_1^2 + B\gamma_2^2)x^2 + \dots ,$$

allowing us to determine the coefficient of x^n . A similar method allows us to deal with $S(p, x, y)$.

If we let

$$\Delta = q^2 - 2pq(-2 + y) + p^2y^2 ,$$

then one can show, after some work, that the coefficient of x^n in the power series for $R(p, x, y)$ equals

$$\frac{1}{2^n} \left((q-py) \sum_{j=0}^{\lfloor n/2 \rfloor} \binom{n}{2j+1} (q+py)^{n-2j+1} \Delta^j + \sum_{t=0}^{\lfloor n/2 \rfloor} (q+py)^{n-2t} \Delta^t \right),$$

and the coefficient of x^n in the power series for $S(p, x, y)$ equals

$$\frac{p}{2^{n-1}} \sum_{j=0}^{\lfloor n/2 \rfloor} \binom{n}{2j+1} (q+py)^{n-2j-1} \Delta^j.$$

To obtain $r(n, k, p)$ and $s(n, k, p)$ from these expressions, we need to find the coefficients of y^k in the above expressions. Writing Δ_l^j for the coefficient of y^l in Δ^j , one can show that

$$\Delta_l^j = \sum_{h=0}^l \binom{j}{h} p^l (-1)^l (q - 2\sqrt{-pq})^{j-h} (q + 2\sqrt{-pq})^{j-l+h}.$$

Using this abbreviation, one can then show that

$$\begin{aligned} r(n, k, p) = & \frac{1}{2^n} \left(q \sum_{j=0}^{\lfloor n/2 \rfloor} \binom{n}{2j+1} \sum_{l=0}^k \Delta_l^j \binom{n-2j-1}{k-l} p^{k-l} q^{n-2j-1-(k-l)} \right. \\ & - p \sum_{j=0}^{\lfloor n/2 \rfloor} \binom{n}{2j+1} \sum_{l=0}^{k-1} \Delta_l^j \binom{n-2j-1}{k-1-l} p^{k-1-l} q^{n-2j-1-(k-1-l)} \\ & \left. + \sum_{j=0}^{\lfloor n/2 \rfloor} \binom{n}{2j} \sum_{l=0}^k \Delta_l^j \binom{n-2j}{k-l} p^{k-l} q^{n-2j-(k-l)} \right), \end{aligned}$$

and

$$\begin{aligned} s(n, k, p) = & \frac{p}{2^{n-1}} \sum_{j=0}^{\lfloor n/2 \rfloor} \binom{n}{2j+1} \sum_{l=0}^k \Delta_l^j \binom{n-2j-1}{k-l} p^{k-l} q^{n-2j-1-(k-l)}. \end{aligned}$$

As before, we can use the closed-form expressions for $R(p, x, y)$ and $S(p, x, y)$ to find the mean and variance of the distribution of the number of doubletons. If we compare these expressions with the

asymptotic ones obtained above for the Markov case, we can partially check the accuracy of our calculations. If we write

$$F(p, x, y) = R(p, x, y) + S(p, x, y) ,$$

then the mean of the distribution of doubletons is obtained by differentiating $F(p, x, y)$ with respect to y , setting $y = 1$, and asking for the coefficient of x^n . To see why this works, note first that since

$$F(p, x, y) = \sum_{n=1}^{\infty} \sum_{k=0}^n (r(n, k, p) + s(n, k, p)) x^n y^k ,$$

if we differentiate $F(p, x, y)$ with respect to y , we obtain

$$\frac{\partial}{\partial y} F(p, x, y) = \sum_{n=1}^{\infty} \sum_{k=1}^n k(r(n, k, p) + s(n, k, p)) x^n y^{k-1} .$$

If we set $y = 1$ and consider the coefficient of x^n , we find that it equals

$$\sum_{k=1}^n k(r(n, k, p) + s(n, k, p)) ,$$

which is clearly the mean of the distribution. In the present case, we find the value of the mean to equal $(n-1)p^2$, which can be checked as being asymptotically equal to the expression we obtained for the asymptotic value of the mean in the Markov case (with $p_1 = p_2 = p$).

A similar, but more complicated, calculation leads to the variance of the distribution of doubletons; we obtain the expression

$$p^2 q(n + 3np - 1 - 5p) .$$

One can check that this is asymptotic to the expression obtained for the asymptotic value of the variance in the Markov case.

We would like to show that the random variable that counts the number of doubletons in the Bernoulli trials case is asymptotically normal. If this were true, then we could use the above values for the mean and standard deviation to give a precise asymptotic description of the distribution of the number of doubletons. It is typically the case that one tries to use the Central Limit Theorem to show that a given sequence of distributions is asymptotically normal. In order

to use the Central Limit Theorem, one must write the terms of the sequence as sums of mutually independent random variables.

In the present situation, if we consider sequences of n Bernoulli trials, and we let X_i denote the random variable that is 1 if the i 'th and $i + 1$ 'st trials are successes, then the number of doubletons in the first n trials is

$$X_1 + X_2 + \dots + X_{n-1} .$$

Unfortunately, the X_i 's are not mutually independent. For example, if $X_{i-1} = X_{i+1} = 1$, then X_i must be 1 as well. Nevertheless, it is possible to salvage the situation, because the X_i 's are "independent enough." More precisely, the sequence $\{X_i\}$ is m -independent, i.e. it is possible to find an m (in this case, $m = 2$) such that the sequence can be partitioned into m subsets such that the random variables in each subset are mutually independent. In this case, we can take the sets $\{X_1, X_3, X_5, \dots\}$ and $\{X_2, X_4, X_6, \dots\}$. If some other conditions (which we will not state here) are satisfied, then the sequence satisfies the Central Limit Theorem, i.e. the sum $S_n = X_1 + X_2 + \dots + X_n$ is asymptotically normal. In the present case, all of the necessary conditions are satisfied, so the distribution of the number of doubletons is asymptotically normal.

Exercises.

1. Let $h_S(n, k)$ denote the probability that in the Markov model, a sequence of n trials begins with a success and has exactly k success runs. Explain how one can write $h_S(n, k)$ in terms of $f_{S,S}(n, i)$ and $f_{S,F}(n, j)$ for appropriate choices of i and j .
2. Define the function $F(r_1, r_2)$ by

$$F(r_1, r_2) = \begin{cases} 2 & \text{if } r_1 = r_2 \\ 1 & \text{if } |r_1 - r_2| = 1 \\ 0 & \text{otherwise.} \end{cases}$$

Suppose that we have a sequence of n_1 successes and $n - n_1$ failures in a Bernoulli trials process with parameters n and

p. Suppose, in addition, that there are r_1 success runs and r_2 failure runs.

- (a) Show that if $r_1 \geq 1$, then there are $\binom{n_1-1}{r_1-1}$ ways to split the n_1 successes into runs.
- (b) Using the fact that runs of successes and failures must alternate in a sequence, show that if $r_1 \geq 1$ and $r_2 \geq 1$, then the number of sequences of the given type equals

$$\binom{n_1-1}{r_1-1} \binom{n-n_1-1}{r_2-1} F(r_1, r_2) .$$

- (c) Show that any two sequences of length n with n_1 successes are equally likely.
- (d) Show that the probability that a Bernoulli trials sequence of length n has exactly $r_1 \geq 1$ success runs and $r_2 \geq 1$ failure runs, given that it has n_1 successes and $n - n_1$ failures, equals

$$P_{n_1}(r_1, r_2) = \frac{\binom{n_1-1}{r_1-1} \binom{n-n_1-1}{r_2-1} F(r_1, r_2)}{\binom{n}{n_1}} .$$

- (e) Determine the corresponding formula if either $r_1 = 0$ or $r_2 = 0$.
- (f) If E and F are events in a sample space, then

$$P(E \cap F) = P(E | F)P(F) .$$

In the space of all Bernoulli sequences of length n , let E denote the event that the number of success runs equals r_1 and the number of failure runs equals r_2 , and let F denote the event that the number of successes equals n_1 and the number of failures equals $n - n_1$. Parts (d) and (e) calculate $P(E | F)$. Find $P(F)$ and use this to find a formula for $P(E \cap F)$.

- (g) By summing over the appropriate set of n_1 's, find a summation that gives the probability that a Bernoulli sequence of length n has exactly r_1 success runs. (This is the form of the expression for this distribution in [31].)

- (h) Write a summation that gives the probability that a Bernoulli sequence of length n has exactly r runs of both types (i.e. the total number of runs is r).

3. Prove that if $p = 1/2$, then

$$r_{n,k} = \binom{n-1}{k-1} \frac{1}{2^{n-1}},$$

where, as above, $r_{n,k}$ denotes the probability of exactly k runs in a sequence of n coin tosses, where the probability of success on one toss equals p .

Chapter 2

Modeling the Stock Market

1. Stock Prices

Stock markets occupy a central position in modern finance. Hundreds of millions of people worldwide have money invested in them. Many economists, mathematicians, and statisticians find them to be a fascinating object of study. In this chapter, we will introduce the reader to some interesting and much-studied questions concerning stock markets.

Stock in a company represents partial ownership of that company. Companies issue shares of stock to raise money to run the company. Many companies are privately held, meaning that all of the stock in that company is owned by a small group of people (perhaps the managers of the company or the founding family of the company). However, most of the large companies in the world are publicly held, meaning that shares of their stock are traded in stock markets.

In a stock market, certain people are willing to sell shares in a given company, while other people are willing to buy shares in that company. The entities that buy and sell stock can also be mutual funds and pension funds. If a seller and a buyer agree on the price of

a share, that share changes hands, and a small commission is paid to the broker who effects the trade.

The central question of importance to most investors when selling or buying stock is whether the price of a share of that stock will move up or down (in the next day, the next month, or the next year). The time period (or horizon) in this question depends upon both the buyer and the seller. For example, if the buyer is investing for his or her retirement, the horizon may be as long as thirty years. At the other extreme, there are certain types of investors, known as day traders, whose horizons can, for some transactions, be as short as a minute.

Clearly, answering this central question is equivalent to predicting the future, which for most human endeavors, is difficult or impossible. But a smart investor knows this, and instead attempts to make predictions that, it is hoped, will usually be close to the actual outcome. The words “usually” and “close to” in the preceding sentence signal the entrance of probability and statistics into the prediction process.

Over a given time period, the discrete return of a stock is the ratio of the stock prices at the end and the beginning of the time period, decreased by one. So, for example, if the price of a stock at the beginning of a given year is \$10, and its price at the beginning of the next year is \$12, then its discrete return for that year is .2, or 20%. Predicting the future of a stock’s price is equivalent to predicting the future discrete returns of that stock.

Discrete returns can clearly be positive or negative. If one wants to model the possible discrete returns using a probability distribution, then the range of this distribution must therefore include both positive and negative numbers. Clearly, it makes no sense to have a discrete return of less than -1 , since this would correspond to a negative ending price. The value -1 for a discrete return means that the stock price fell 100%, i.e. its ending price is 0.

For reasons explained shortly, geometric returns, rather than discrete returns, are typically used. The geometric return of a stock over a given time period is the (natural) logarithm of the ratio of the

stock's prices at the end and the beginning of the time period. So, if a stock's price is \$10 at the beginning of a year and \$12 at the beginning of the next year, then its geometric return over that year is

$$\log\left(\frac{12}{10}\right) \approx .182,$$

or 18.2%. Note that this is close to, but not equal to, the discrete return of 20%. For most stocks and short durations, the two are usually close. The reason for this is that if we let d and g denote the discrete and geometric return and we assume that d is small (say between -20% and 20%), then

$$g = \log(1 + d) = d - \frac{d^2}{2} + \frac{d^3}{3} - \dots,$$

and the summands after the first one are very small in comparison with the first one.

If one uses a certain probability distribution (such as the normal distribution) to model the geometric returns of a stock, then large negative values of the distribution no longer lead to meaningless results. For example, suppose that the geometric return over a certain time period is predicted to be -3 . If the prices of the stock at the beginning and the end of the time period are P_0 and P_1 , then we have

$$-3 = \log\left(\frac{P_1}{P_0}\right),$$

which is equivalent to the statement

$$P_1 = e^{-3}P_0,$$

or

$$P_1 \approx .0498P_0.$$

Thus, under this prediction, the stock would lose about 95% of its value in the given time period.

Another reason geometric returns are used is that they behave better mathematically than do discrete returns. Suppose, for example, that we are modeling the price of a stock over two consecutive time periods. If the stock prices at the three endpoints of these two

time periods are P_0 , P_1 , and P_2 , then the geometric returns over the two time periods are

$$g_{0,1} = \log\left(\frac{P_1}{P_0}\right)$$

and

$$g_{1,2} = \log\left(\frac{P_2}{P_1}\right).$$

Note that the return over the combined time period is

$$g_{0,2} = \log\left(\frac{P_2}{P_0}\right) = g_{0,1} + g_{1,2},$$

i.e. the geometric returns simply add over consecutive time periods. The corresponding statement does not hold for discrete returns. For example, a 10% discrete return over each of two consecutive time periods does not produce a 20% discrete return over the combination of the time periods. In this example, if a stock were initially priced at \$10, then after one year it would be worth \$11 and after two years it would be worth \$12.10, leading to a two-year discrete return of .21, or 21%.

Next, suppose that we are trying to model the exchange rate between two currencies, say between the dollar and the euro. Let the number of euros per dollar at the beginning a time period i be denoted by a_i . Then the geometric changes in this quantity are

$$\log\left(\frac{a_2}{a_1}\right), \log\left(\frac{a_3}{a_2}\right), \dots$$

Viewed from the other perspective, the geometric changes in the number of dollars per euro are

$$\log\left(\frac{1/a_2}{1/a_1}\right), \log\left(\frac{1/a_3}{1/a_2}\right), \dots,$$

which are seen to be the negatives of the geometric changes of the a_i 's. Thus, the distributions of these two series are simply related.

What happens if we use discrete changes instead? Suppose, for example, that $a_1 = .7$ and $a_2 = .75$ (so one dollar is worth .7 euros at

the beginning of the first time period and .75 euros at the beginning of the next time period). Then the discrete change is

$$\frac{.75}{.7} - 1 \approx .0714.$$

The corresponding numbers of dollars per euro are 1.429 and 1.333, so the discrete change is

$$\frac{1.333}{1.429} - 1 \approx -.0672.$$

Although geometric returns are thus nicer mathematically than discrete returns, it is easy to translate each to the other. In other words, there is nothing that precludes using either quantity in a model. In the rest of this chapter, we will use the word “return” to mean “geometric return,” unless expressly stated otherwise.

We have not yet mentioned stock dividends. Many companies issue dividends, which are payments made every so often to the owners of the shares of the companies’ stocks. Typically, these dividends can be taken as cash or can be reinvested in additional shares of stock. For example, suppose that a certain stock is worth \$10 at the beginning of the year. Suppose that at the end of the year, the company issues a 43-cent dividend for each share. Finally, suppose that the stock is worth \$11 at the beginning of the next year. If the share owner takes the dividend in cash, then he has made \$0.43 in cash and has also made \$1.00 in gain on the price of a share. This last gain is sometimes said to be *unrealized*, in that the share owner does not have this \$1.00 in cash unless he sells the share (realizes the gain). If instead the share owner reinvests the dividend, he would receive an additional .039 (= 0.43/11.00) shares of stock. In either case, at the beginning of the next year, he would have some combination of stock and cash worth \$11.43, so his geometric return would be

$$\log\left(\frac{11.43}{10.00}\right) \approx 0.1337.$$

In what follows, we assume that all dividends are reinvested.

In this chapter, we are primarily interested in changes in the prices of stocks over short time intervals, such as days or weeks. Suppose that a stock's closing price on a certain day is \$10.00 and its closing price at the end of the next day is \$10.20. Then the return on that day is 2.0%. However, if the company issued a dividend on that day of 30 cents per share, then one might argue that the return on that day is 5.0% (since this share has increased the wealth of its owner by 50 cents on that day). However, the dividend is not the result of any trading in a market; rather, it is set by the company's management. Thus, if we are interested in the distribution of daily returns for this stock, it makes more sense (to us, anyway) to ignore all dividends and simply use the closing prices to compute the returns. This is what we shall do in what follows.

Finally, it is occasionally the case that a company will split its shares. For example, it might split each share into three shares (each worth one-third of the original share's value). This has no effect on the return. To see why, assume for example that a share in a certain company is worth \$30 at the beginning of the day and suppose that the value of the share increases by 1% on that day. If the share splits into three shares, then each are worth \$10.10 at the end of the day, and since each of the three shares was worth \$10 at the beginning of the day, there was a 1% return that day for each share. Because of this, we disregard stock splits.

There are other types of monetary quantities, such as stock index prices and mutual fund prices, that are of interest in finance. A stock index is a number that is a weighted average of the stock prices in a certain set of stocks.

For example, the S&P 500 Index uses the stock prices for a set of 500 stocks. The weights are affected by dividends and stock splits, among other things. For example, if we start with equal weights for all of the stocks, we can imagine the average represents $1/500$ of the total value of a portfolio consisting of one stock in each of the 500 represented companies. However, over the years, the different stocks pay out different dividends, and if, as we are assuming, these

dividends are reinvested, then in order for the average to continue to represent $1/500$ of the total value of the portfolio, the weights must change. Similarly, if a stock splits 2 shares for 1, then the price is halved, but there are now two shares of this stock in our portfolio, so the weight for this stock should double.

A mutual fund consists of a portfolio of stocks that are bought by an investment company, using money from investors. Over time, the portfolio represented by a given mutual fund changes. In addition, the various stocks in the fund sometimes pay dividends, which must, by law, be distributed to the investors. (These investors may choose to reinvest these dividends with the mutual fund, but taxes must typically be paid on these dividends whether or not they are reinvested.) Since so much money is invested in mutual funds today, much attention is paid to the prices of these funds.

2. Variations in the Price of a Stock

The first attempt to model stock market returns is generally considered to have occurred in a doctoral thesis of a French mathematics student named Louis Bachelier. His work concerned bond prices in the Paris exchange in the late 1800's. It had been noticed earlier that the price variations of a stock (or a bond) are, on average, larger over long time intervals than over short time intervals. Bachelier noted that the same idea occurs in certain types of random walks. We will illustrate this with the simplest random walk, a sequence of flips of a fair coin.

Suppose we flip a fair coin repeatedly, assigning a value of $+1$ for a head and -1 for a tail, and we keep track of the total value (which equals, at any time, the number of heads minus the number of tails). The name "random walk" refers to a physical model of this sequence of coin flips. We imagine a person on a road, who walks one step to the right or left, corresponding to a flip of heads or tails. If the person starts at a mark of 0 on the road, his position after n steps is the same as the total value of the first n coin flips.

The distribution of ending positions of the random walker after n steps is well known (and is undoubtedly known to the reader). It is closely related to the binomial distribution with parameters n and $p = 1/2$. An example when $n = 3$ will make this relationship clear. The binomial distribution gives the probability of obtaining exactly h heads in n tosses, if the probability of a given toss coming up heads is p . This probability equals

$$b(h; n, p) = \binom{n}{h} p^h (1 - p)^{n-h}.$$

If there are h heads in n tosses, then there are $n - h$ tails, so the value of the coin tossing sequence as originally defined equals $h - (n - h) = 2h - n$. Thus, the value of the coin tossing sequence is a linear function of the number of heads.

The reader will recall that the variance for the binomial distribution is $np(1-p)$. The variance of the value of the coin tossing sequence is simply four times this amount, since the value equals $2h - n$, doubling the quantity h multiplies its variance by four and subtracting n from $2h$ has no effect on the variance. Thus, since the variance is a constant multiple of n , we see that if one coin toss game is twice as long as another, the variance of the first is twice the variance of the second (since the value of n for the first game is twice the value of n for the second game). This last statement sounds much like the observed behavior of stock prices (at least in France in the 1800's).

Of course, the price of a given stock or bond does not move up or down by one amount only. So in order to model stock prices, it is necessary to make the coin-tossing game more complicated. This is a central issue in mathematical modeling. On the one hand, one wants to create a model that is simple enough so that there is some chance one can mathematically analyze it. On the other hand, one wants the model to be of sufficient accuracy to be of some use in understanding the process that is being modeled.

Bachelier assumed that the returns in the price of a given bond, in a given time interval, are normally distributed. In other words, if

the price of a bond at the beginning of day 1 is p_0 , then the price of the bond at the end of day 1 is $p_0 + a_1$, where a_1 is an *arithmetic* return that is normally distributed with a certain mean and variance. Although Bachelier used arithmetic returns, modern treatments of stock prices use geometric returns, for the reasons discussed above.

Bachelier also made the assumption that the returns in the price of a given bond in non-overlapping time intervals were independent. This corresponds to the situation in the coin-flipping game, in that the flips are independent events.

Next, Bachelier posited, as in the coin-flipping game (and as had been observed in the stock and bond markets), that the variance of the returns increased as the length of the time intervals increased. One can accomplish this in such a model by assuming that the variance is proportional to the length of the time interval. Thus, for example, the variances for the returns for time intervals of lengths one and two days are in the ratio of one to two.

Finally, one can assume that the means of the distributions that represent the returns are zero or non-zero. In the former case, the prices of the bonds will tend to fluctuate around a fixed number (as does the value of the coin-flip game). In the latter case, if the mean of the distribution that represents one-day returns is μ , then the mean of the distribution that represents k -day returns is $k\mu$. In this case, the prices of the bonds will tend to fluctuate about the line $y = \mu t$, where t measures the number of days from the beginning of the process. In this latter case, the number μ is sometimes called the drift, because it causes the prices to drift away from their starting value.

The model described above is quite complicated, and it is not even clear whether one can rigorously devise a process that satisfies the properties given. In fact, this process is called Brownian motion, and it was rediscovered slightly after Bachelier by Albert Einstein. Neither Bachelier nor Einstein gave a rigorous demonstration of the existence of such a process; such a demonstration was given later by Norbert Wiener. The name of this process comes from the botanist

Robert Brown, who observed erratic motion of a pollen grain in a drop of water under a microscope. This motion is due to water molecules hitting the grain. This explanation was first given by Einstein. Figure 1 shows an approximation of a sample Brownian motion path on the interval $[0, 1]$.

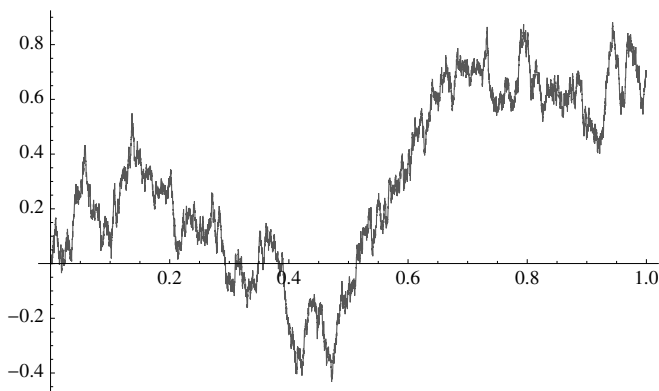


Figure 1. A sample Brownian motion path

In the next three sections, we will turn our attention to the two major assumptions made in the above model, namely the normality of the returns in a given time interval and the independence of returns corresponding to non-overlapping time intervals.

3. The Normal Distribution and Power Laws

In many probabilistic models of the real world, the normal distribution is used. In some cases, use of this distribution can be justified because of the Central Limit Theorem. Roughly speaking, this theorem says that if a sequence of independent experiments is carried out, and all of the experiments have the same numerical distribution, then the average of the resulting values will be approximately normally distributed. (We should add that for this theorem to apply, it is necessary that the numerical distribution of the individual experiments has a finite mean and a finite variance.) For example, if one

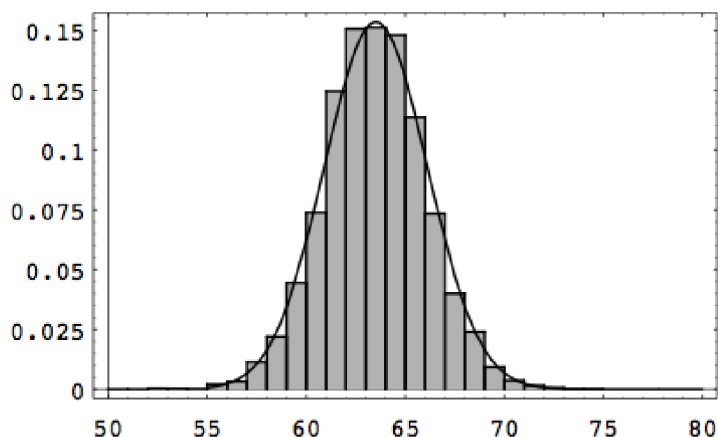


Figure 2. Heights of adult women in the U.S.

considers a sequence of flips of a fair coin, then the average number of heads in n flips is approximately normally distributed with mean $n/2$ and variance $n/4$.

In other cases, the normal distribution is assumed because an empirical perusal of the data suggests that the normal fits the data fairly well. For example, in Figure 2, we show a histogram of the heights of adult females in the United States, together with a fitted normal density. One can see that the fit is fairly good.

In still other cases, one might use the normal distribution because it has nice mathematical properties. For example, if the outcomes of two independent experiments are each normally distributed (perhaps with different means and variances), then the sum of the outcomes is also normally distributed. The mean and variance of the sum can be obtained by adding the means and variances of the distributions of the individual experiments.

One can see that the above property of the normal distribution is reminiscent of the observed property concerning stocks that over two non-overlapping time intervals, the variances add (and the means

add as well). Thus, the normal distribution is a natural distribution to use when attempting to model returns of stock prices.

The reader will recall that in any attempt to model real-world phenomena, there is usually a trade-off between simplicity and accuracy. The Brownian motion process is already quite complicated, and if one uses a set of distributions that are not as mathematically nice as the normal distribution in this process, the resulting process is still more difficult to handle. Nevertheless, after positing this model, one should check to see how well the model agrees with the data and whether there might be another process that fits the data better. One of the uses of statistics is exactly this checking procedure; there are many ways to quantify how well a given theoretical distribution fits a set of data.

In the early 1960's Benoit Mandelbrot, a mathematician working at IBM, was studying the distribution of income in a society. A Harvard economist, Hendrik S. Houthakker, had been studying the time series of cotton prices for several years, trying to make it fit the Bachelier model. Mandelbrot was invited to give a talk at Harvard on his work. While in Houthakker's office, he noticed a diagram on the blackboard that reminded him of the topic of his talk. Of course, the diagram referred to cotton prices, not income distribution. Mandelbrot was intrigued by the apparent similarity of the two distributions, which led to his study of these prices (and many other price sets as well).

Mandelbrot found a distribution that fit the geometric returns, derived from the cotton price data, much more closely than did the normal distribution. This distribution is called a power law. The density function of a power law is approximately of the form

$$f(t) = C_1 |t|^{-k},$$

for $t \neq 0$, where $k > 0$ and C_1 is a positive normalizing constant. One can see from this expression that as $|t|$ gets large, the density function gets small. This is also true of the normal density functions. The most important difference between the two classes of densities is

the rate at which they go to 0 as $|t|$ gets large. This difference will be studied in more depth below.

Elementary calculus shows that it is not possible for this function to be suitably normalized so that the area underneath its graph and above the real axis will equal one. If $0 < k \leq 1$, then the area below the graph of $f(t)$ and above the interval $[1, \infty)$ will be infinite, and if $k \geq 1$, then the area below the graph of $f(t)$ and above the interval $(0, 1]$ will be infinite.

The way around this problem is to realize that such distributions are useful because they fit the part of the data set that is not close to 0, i.e. the tails of the data set. So, when using a power law to describe a data set, it is to be understood that in a certain interval containing 0, the formula for $f(t)$ will be modified. We will typically not care very much about such distributions near the origin; as a result, we will assume in what follows that $k \geq 1$. In fact, in applications in finance, the values of k that fit the data are all much greater than 1. In the appendix, we have defined a set of power laws which we will use in what follows.

To understand the most salient difference between the normal distribution and a power law distribution, we consider the probability, for each distribution, that a value greater than a certain fixed positive quantity x is obtained. More precisely, suppose that N and L are values of experiments that are distributed according to the normal distribution and a power law distribution, respectively. We will use the standard normal density function ϕ to represent N and the above density function f to represent L . We want to compare the quantities

$$P(N \geq x)$$

and

$$P(L \geq x).$$

These are right-tail probabilities. To obtain such quantities from a density function, one integrates the density function over the appropriate interval. Thus, these two quantities equal

$$\int_x^\infty \phi(t) dt$$

and

$$\int_x^\infty f(t) dt.$$

To those readers for whom this idea is unfamiliar, we will say that we will shortly return to the non-calculus realm.

The expression for $\phi(t)$ is given by the formula

$$\phi(t) = \frac{1}{\sqrt{2\pi}} e^{-t^2/2}.$$

One cannot write an antiderivative of this function in terms of elementary functions, so one cannot use the Fundamental Theorem of Calculus directly to find the required probability. Instead, we proceed as follows. Using integration by parts, with the settings

$$u = \frac{1}{t}$$

and

$$dv = te^{-t^2/2} dt,$$

we obtain the equation

$$\int_x^\infty e^{-t^2/2} dt = \left[-\frac{1}{t} e^{-t^2/2} \right]_x^\infty - \int_x^\infty \frac{1}{t^2} e^{-t^2/2} dt.$$

The first summand on the right-hand side equals

$$\frac{1}{x} e^{-x^2/2},$$

and if we compare the integral on the right-hand side with the one on the left-hand side, we see that the integrand of the right-hand integral is always less than $1/x^2$ times the integrand of the left-hand integral. Thus, the value of the right-hand integral is less than $1/x^2$

times the value of the left-hand integral. Since we are interested in the left-hand integral for large values of x , we can say that

$$\int_x^\infty e^{-t^2/2} dt \sim \frac{1}{x} e^{-x^2/2},$$

where the \sim sign means that the two quantities are asymptotically equal, as x gets large. (Two functions of x are said to be asymptotically equal as x gets large if their ratio approaches 1 as x goes to ∞ .) So, we have

$$(1) \quad P(N \geq x) \sim \frac{C_2}{x} e^{-x^2/2},$$

where $C_2 = 1/\sqrt{2\pi}$.

The tail probability for the power law is easier to calculate. We have

$$\begin{aligned} P(L \geq x) &= \int_x^\infty C_1 |t|^{-k} dt \\ &= \left[-\frac{C_1}{k-1} t^{-k+1} \right]_x^\infty \\ &= \frac{C_1}{k-1} x^{-(k-1)}. \end{aligned}$$

If we replace $C_1/(k-1)$ by a constant C_3 , and we let $\alpha = k-1$, then we see that we have

$$(2) \quad P(L \geq x) = C_3 x^{-\alpha}.$$

We are now finished with the calculus, so various readers can exhale. The salient difference between the normal distribution and the power law distribution is described in Equations 1 and 2. We will spend some time explaining this difference, because it is at the heart of some of Mandelbrot's work (and much subsequent research as well).

Suppose, for example, we are looking at daily returns of a certain stock over a period of ten years. This set will contain about 2500 observations. Now let's suppose that, using certain units, these observations are described well by a standard normal distribution.

What fractions of these observations will be larger than $x = 1, 2, 3, 4$ and 5? Using Equation 1, we see that the answers are

$$.24, .027, .0015, .000033, .00000030.$$

(To show that this equation provides a good approximation to the actual values, we note that the actual values are

$$.15, .023, .0013, .000032, .00000029.)$$

Since the standard deviation of the standard normal distribution is one, we can interpret the above calculations as saying, for example, that the probability of observing a value at least 3 standard deviations above 0 (i.e. $x \geq 3$) is about .0015. In a data set of 2500 observations, this corresponds to about four observations.

Now let's assume that the data set is described well by a power law with $\alpha = 2$. For power laws, we don't know the value of C_3 , but this doesn't affect the main idea being discussed here, as we shall see. Suppose, for example, that $C_3 = .1$. Then the fractions of observations that will be larger than $x = 1, 2, 3, 4$ and 5 are

$$.10, .025, .011, .006, .0040.$$

Thus, in a data set of 2500 observations, we would expect to see about 28 ($= .011 * 2500$) observations that are greater than 3. The reader should object here that this might not be accurate, since we arbitrarily picked a value for C_3 . This is indeed correct; a different value of C_3 would lead to a different expected number of observations.

However, the value of C_3 does not affect the truth of the following statement. In a power law distribution, with $\alpha = 2$, if one doubles the value of x , one multiplies the upper tail probability by $1/4$. This can be seen by considering Equation 2. We have

$$P(L \geq x) = C_3 x^{-2}$$

and

$$P(L \geq 2x) = C_3 (2x)^{-2},$$

and it is easily seen that the second probability is one-quarter of the first probability, irrespective of the value of C_3 .

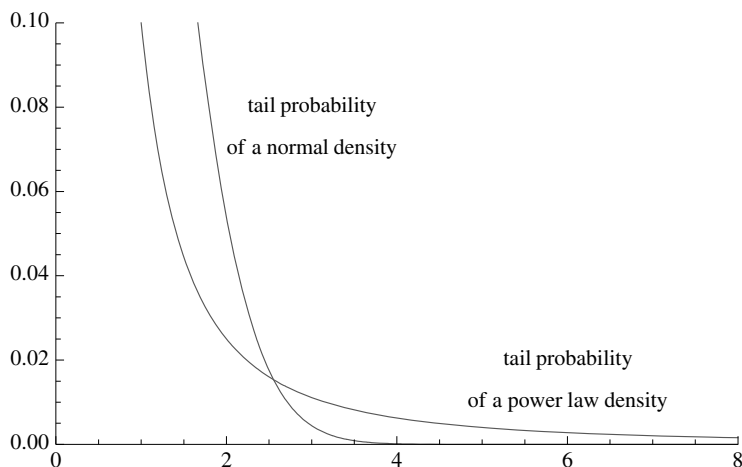


Figure 3. Tail probabilities of the normal and power law distributions

The situation is completely different for tails of the normal distribution. For example, if one compares the normal tail probabilities for $x = 2$ and $x = 4$ (computed above), one sees that the values are

$$P(N \geq 2) = .023 \quad \text{and} \quad P(N \geq 4) = .000032.$$

Thus, the second probability is about 1/70 of the first probability. If one compares the values for $x = 3$ and $x = 6$, one finds that the second probability is about 1/1,400,000 of the first probability.

Another way to understand the difference between the sizes of the tails of the normal distribution and a power law is by considering a graph. Figure 3 shows the two expressions $P(N \geq x)$ and $P(L \geq x)$ as functions of x , where we have chosen $\alpha = 2$ and $C_3 = 1/10$ for the power law. One sees that the normal distribution tail goes to 0 much more rapidly than does the power law tail. Analytically, we are claiming that the ratio of the expression in Equation 1 to the expression in Equation 2 goes to 0 as x goes to ∞ , for any positive constants C_2 , C_3 , and α . This can be shown using elementary calculus.

The expression “fat tails,” which occurs frequently in the study of stock prices, refers to the tail of a distribution where the tail falls

off to 0 as a constant power of x (as in the case of a power law) as opposed to an exponential rate in x (as is the case of the normal distribution).

We note that although we have used the standard normal distribution and a specific power law in the above exposition, nothing of consequence changes if we instead take a general normal distribution or a power law with a different value of α .

4. Distribution of Returns

We now turn to the question of how to determine if a normal distribution or a power law distribution describes the sequence of geometric returns of a stock (or another financial instrument). We note that at this stage, we are not concerned with the order in which the returns occur in time; rather, we are only interested in their overall distribution.

Before considering any actual data sets, we consider how such data sets would look if they were governed by a normal distribution or a power law. In the last section, we saw that if A denotes a random observation in the data set, then the quantity $P(A \geq x)$, as a function of x , looks much different under the two classes of distributions. Suppose, for example, that the distribution is fit well by a power law with constants C_3 and α . We do not assume that these constants are known to us yet. If we graph the values of $P(A \geq x)$ for different values of x , they should look like one of the graphs in Figure 3. But both of the graphs in this figure are curves, and it might be difficult to determine which type of distribution does a better job of fitting the tail probabilities.

A way around this problem is to proceed as follows. Let $g(x) = P(A \geq x)$ be the theoretical tail probability. Then we know from Equation 2 that

$$g(x) = C_3 x^{-\alpha}.$$

If we take logarithms of both sides of this equation, we obtain

$$(3) \quad \log g(x) = \log C_3 - \alpha \log x.$$

This says that the logarithm of tail probability function is a linear function of the logarithm of the input x . Thus, if we graph the values of the tail probabilities versus x on log-log paper, or graph $\log g(x)$ versus $\log x$ on a standard pair of axes, the result should look like a straight line with slope $-\alpha$.

Next, suppose that the geometric returns, represented by the random variable A , obey a normal distribution with mean μ and standard deviation σ . Once again the values of these parameters are unknown to us. If we let

$$N = \frac{A - \mu}{\sigma},$$

then N obeys the standard normal distribution. Hence,

$$\begin{aligned} g(x) &= P(A \geq x) \\ &= P\left(N \geq \frac{x - \mu}{\sigma}\right), \end{aligned}$$

and using Equation 1, we find that

$$g(x) \sim \frac{\sigma}{(x - \mu)\sqrt{2\pi}} e^{-(x - \mu)^2 / (2\sigma^2)}.$$

Although this is a somewhat ghastly equation, we merely want to take the logarithms of both sides to see how much the resulting equation differs from Equation 3. We obtain the equation

$$\log g(x) \sim \log C_4 - \log(x - \mu) - \frac{(x - \mu)^2}{2\sigma^2}.$$

In this case, one sees that the logarithm of the tail probability is definitely not a linear function of the logarithm of x .

We can illustrate the above idea of distinguishing between normal distributions and power law distributions by simulating some data from these two types of distributions. In the appendix, we have given one method of generating data according to a power law. The two key parameters are α , described above, and a , which is the value of x above which Equation 2 for the upper tail probability is supposed to hold. It also holds below $-a$. Actually, we construct a symmetric density, so a similar equation holds for x below $-a$. A graph of a

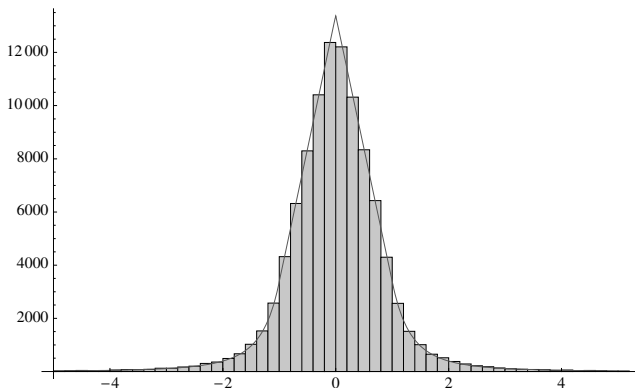


Figure 4. Simulation of a power law with $\alpha = 2.1$ and $a = 1$

power law density, with $\alpha = 2$ and $a = 1$, is shown in the appendix in Figure 27.

Figure 4 shows the result of simulating 100,000 values from a power law distribution with $\alpha = 2.1$ and $a = 1$, together with the appropriately scaled density function. As expected, the simulated results closely match the density function.

Next, we count, for various positive values of x (we are dealing with the right tail for the moment), the fraction of simulated values that equal or exceed x . These fractions should be close to the theoretical tail probability $g(x) = C_3 x^{-\alpha}$ discussed above. Following the discussion concerning Equation 3, we have plotted the logarithm of these fractions versus the logarithm of x in Figure 5.

It can be seen that in the range $\log x \geq 0$, which corresponds to $x \geq 1$, the region in which the power law part of the density function is operational, the graph is virtually straight. (The strange shape of the graph to the right of $\log x = 2.8$ is due to the fact that only 34 values of the 100,000 in the data set lie in this region.)

We will now repeat the above exposition for an experiment whose values, denoted by N , are distributed according to the standard normal distribution. Figure 6 shows a histogram of 100,000 simulated

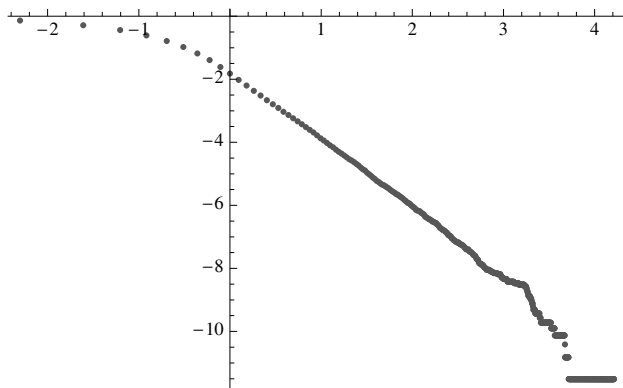


Figure 5. Log of tail probability versus log of input for power law

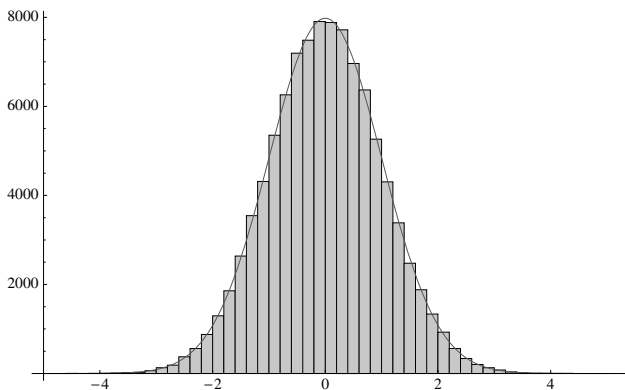


Figure 6. Simulation of the standard normal distribution

data points, together with a scaled version of the standard normal density function. As before, we count, for various positive values of x , the fraction of values of the experiment that exceed x . We have plotted the logarithm of these fractions versus the logarithm of x in Figure 7.

A comparison of Figures 5 and 7 shows a significant difference in the shapes of the curves, and corroborates what was said above about

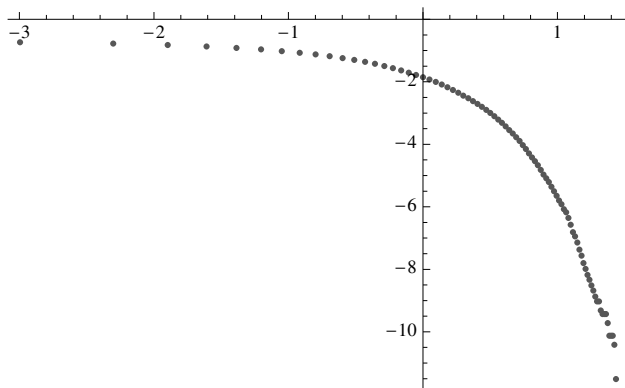


Figure 7. Log of tail probability versus log of input for normal distribution

the relationship between the logarithm of the tail probability and the logarithm of x for each of the two types of distributions.

In Figure 8, we show the graph of the logarithm of the tail probability versus the logarithm of x for the absolute values of the daily geometric returns for General Electric stock, between January 1, 1962 and May 22, 2009. (WE are fitting a symmetric density to the data, so the effect of taking absolute values is to combine the data from both the lower and upper tails to fit a value for the power law exponent α .) The number of returns in this data set is 11930. To cement the reader's understanding of this graph, we consider the point $(-3.96, -1.77)$, which lies on the graph. These two numbers are logarithms of .019 and .170. Thus, the fact that the point is on the graph means that the fraction of absolute returns exceeding .019 is .170. Next, we note that on the interval $[-4, -2.2]$, the graph is almost a straight line, as is the case with power laws. Both of the endpoint values -4 and -2.2 are arbitrary, and if we choose different endpoints, we will obtain different estimates for α . In such graphs, the right-hand edge of the plot will be somewhat ragged. There is a good reason for this; in the present case, for example, the point $(-2.20, -6.90)$ is about where we would claim that the graph becomes ragged. But

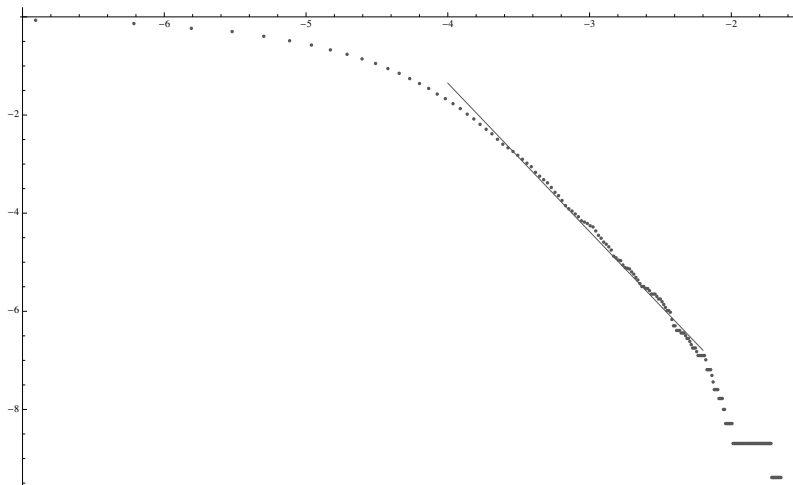


Figure 8. Logarithm of tail probabilities for returns of General Electric stock

$e^{-6.90}$ is about .00101, which means that the points on this graph with first coordinates exceeding -2.20 correspond to only .101% of the points in the data set (in fact, exactly 12 returns). These 12 returns are certainly important (they are the largest 12 in absolute value), but there are so few of them that it shouldn't be surprising that they don't fit a power law distribution very well. To put it another way, if we added one more large return to our data set, the right-hand edge of the above plot might look quite different.

We can fit a straight line to this part of the graph corresponding to the horizontal interval $[-4, -2.2]$; the negative of the slope of this line is our estimate for α . This best-fit line is shown in Figure 8. In this case, we obtain a value of $\alpha = 3.02$. In Section 6, we consider the question of whether or not the value of α changes over time for a given stock.

Carrying out the same calculations for IBM stock results in the graph shown in Figure 9. Again we use the horizontal interval $[-4, -2.2]$. We obtain a value of $\alpha = 2.97$. The best-fit line is shown in the figure.

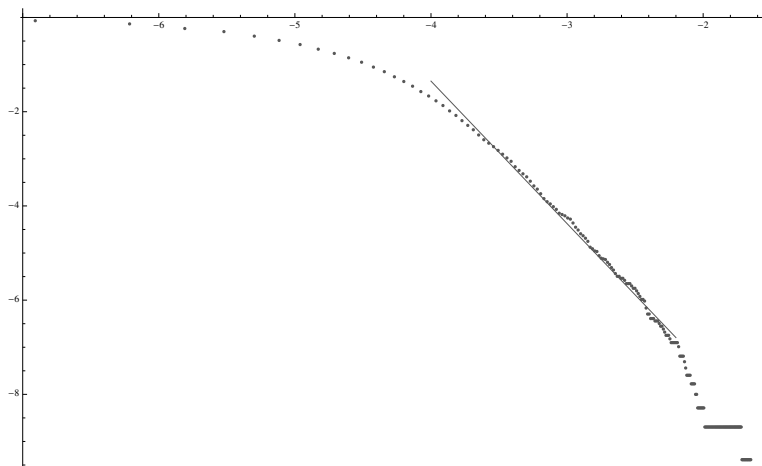


Figure 9. Logarithm of tail probabilities for returns of IBM stock

We will now attempt to fit the two types of distributions under discussion to the returns of the Dow Jones Industrial Average for the period between October 1, 1928, and March 17, 2009. A graph of the value of this average at the end of each day in this period is shown in Figure 10. The set of daily returns for this period is of size 20202. The largest positive return is .1427, and the largest negative return is $-.2563$. These returns occurred on March 3, 1933 and October 16, 1987. Of these returns, all but 12 are between $-.10$ and $.10$, and all but 139 are between $-.05$ and $.05$. A histogram of this data set, for the interval $[-.05, .05]$, is shown in Figure 11. We next plot the logarithms of the observed tail probabilities versus the logarithm of the input x , for the set of positive returns and the set of negative returns. The graphs are shown in Figures 12 and 13. The reader will certainly agree that both of these graphs look much more like Figure 5 than Figure 7.

In Figure 14, we have graphed the logarithms of the tail probabilities versus the logarithm of the input x for the set of absolute values of the returns that do not exceed $.10$. We have also shown a

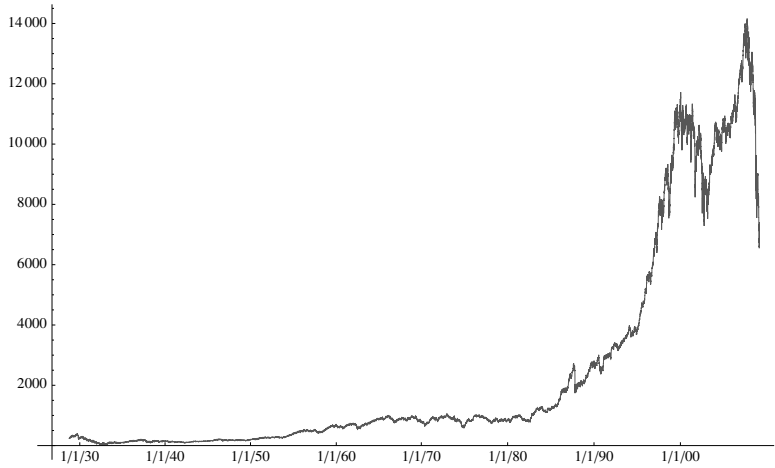


Figure 10. The Dow Jones from October 1, 1928 to March 17, 2009

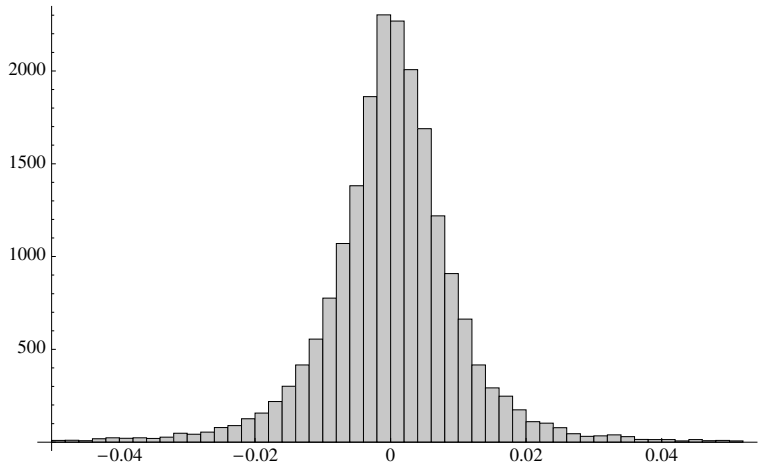


Figure 11. Daily returns of the Dow Jones in the interval $[-.05, .05]$

best-fit line to this set of points. The slope of the line, which is the negative of the estimate for α , is -2.75552 .

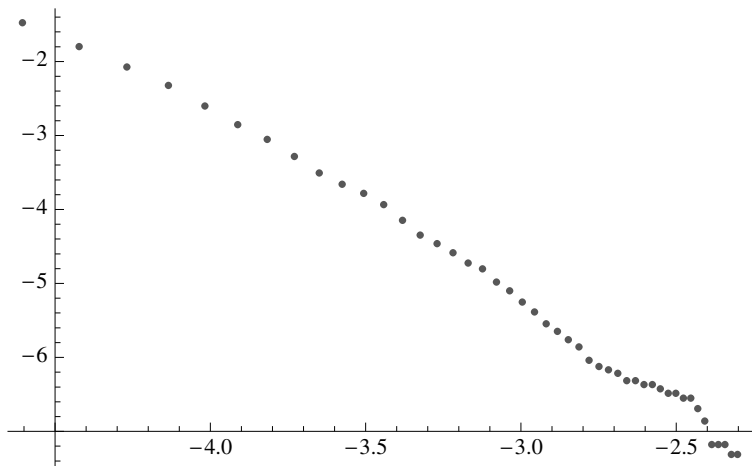


Figure 12. Log of tail probability versus log of input for positive returns of the Dow Jones

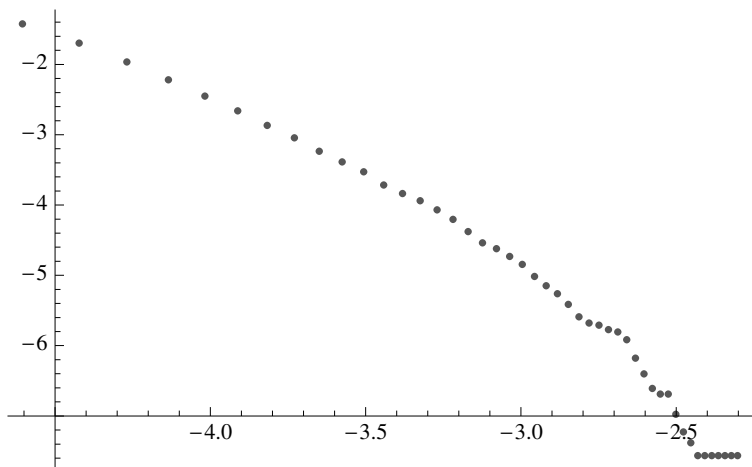


Figure 13. Log of tail probability versus log of input for negative returns of the Dow Jones

The best-fit line shown in Figure 14 corresponds to the tail probability function

$$g(x) = (1.33577 * 10^{-6}) * x^{-2.75552}.$$

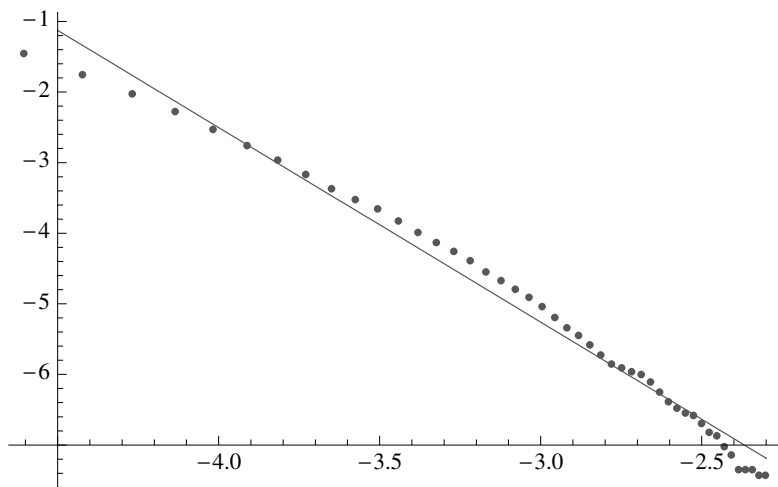


Figure 14. Log of tail probability versus log of input for absolute value of returns of the Dow Jones, together with a best-fit line

One can use this function to see how many large changes one would expect to see, under this power law model, in a set of 20202 days. For example, if we let $x = .1$, we find that

$$g(.1) = .000761 ,$$

which means in 20202 days, we would expect to see about 15 returns whose absolute values exceed .1. In fact, there were 12 such returns. Similarly, if we let $x = .05$, we would expect to see about 104 returns whose absolute values exceed .05; in fact, there were 131. Finally, what is the expected number of returns whose absolute value exceeds .256 (the largest absolute value among all returns)? This expected number is 1.15.

The point of these calculations is not that the model is accurate enough to predict the number of returns exceeding any bound with great accuracy, but rather that under this model, very large returns are not highly improbable. We shall see below that this statement is not true when we try to model the returns using a normal distribution.

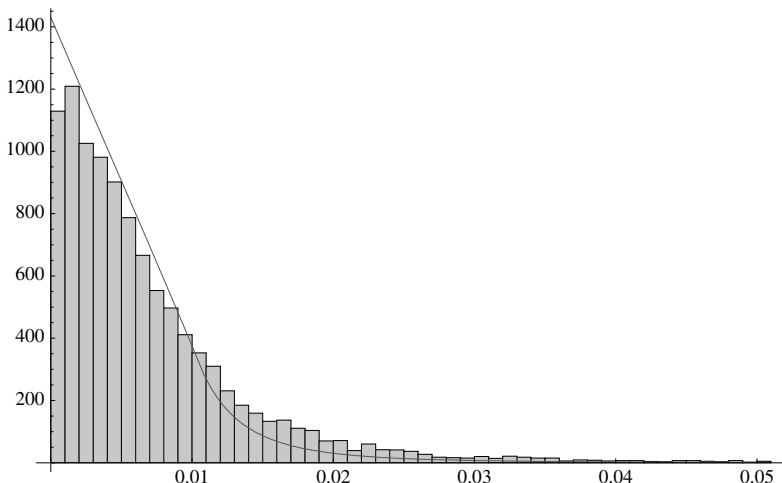


Figure 15. Fitting a power law with parameters $\alpha = 2.64$ and $a = .0111$ to the positive returns of the Dow Jones Industrial Average

It is worth noting that the estimated values for α for positive and negative returns are somewhat different. For positive returns, the best-fit line derived from Figure 12 yields a value of $\alpha = 2.64$, while for negative returns, we obtain a value of $\alpha = 2.91$. To see how well one of our power law distributions fits the data, we return to Figure 12, and note that the “linear” portion of the data begins for values of the return very close to the origin. At the origin in this figure, we find that $\log x = -4.5$, which is equivalent to the statement $x = .0111$. So, we will choose $\alpha = 2.64$ and $a = .0111$ as our parameters in the power law density $f(t)$ described in the appendix. Figure 15 shows the power law density using these parameters and the histogram of positive returns. We note that if we change the value of a to $a = .012$, the resulting power law, shown in Figure 16, seems to fit the data better.

Now let’s try to fit a normal distribution to the Dow Jones returns. If we use all of the data, we find that the mean is 0.000169685

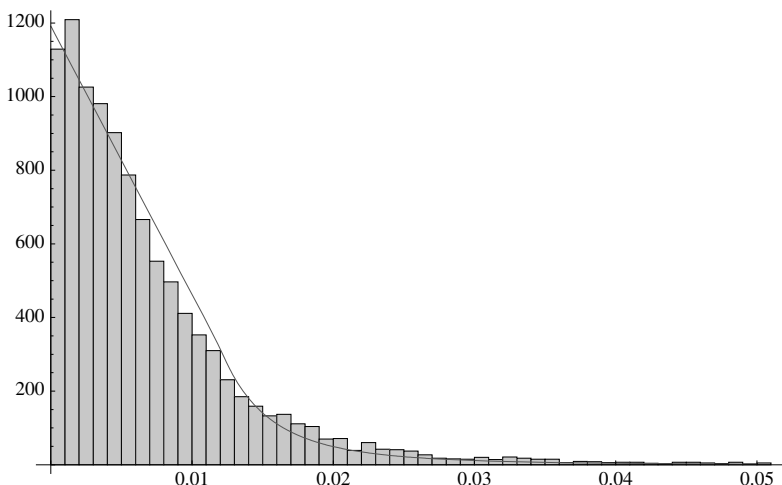


Figure 16. Fitting a power law with parameters $\alpha = 2.64$ and $a = .012$ to the positive returns of the Dow Jones Industrial Average

and the standard deviation is 0.0116385. This leads to a normal distribution with these parameters; a scaled graph of this normal density function is shown in Figure 17, along with a histogram of the returns, in the interval $[-.04, .04]$. The reason that the normal density function in this graph appears wider than the histogram is that there are large returns (outside the interval shown) that inflate standard deviation. For example, if we truncate the data to the interval $[-.04, .04]$, we find that the standard deviation drops to .00946. Of course, we shouldn't use this value, since we have thrown away 251 data points (which, in some sense, are the most interesting points in the set) in order to calculate this value. If we superimpose the resulting normal density on the graph in Figure 17, we obtain Figure 18. Note that neither of these graphs can be said to fit the data very closely.

Now, let's estimate the expected number of large returns we should see in 20202 trading days, using the normal distribution model. We use the value $\sigma = 0.0116385$, computed above, for the standard deviation. In terms of σ , the computed value of the mean is only

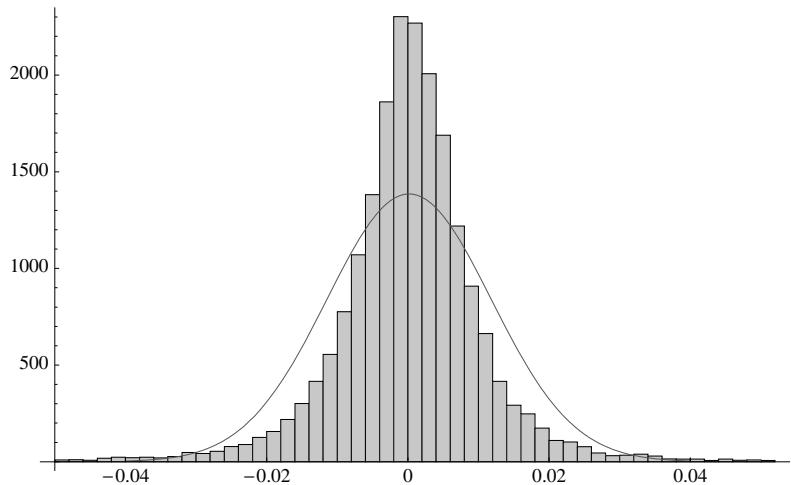


Figure 17. Fitting a normal distribution to the returns of the Dow Jones Industrial Average

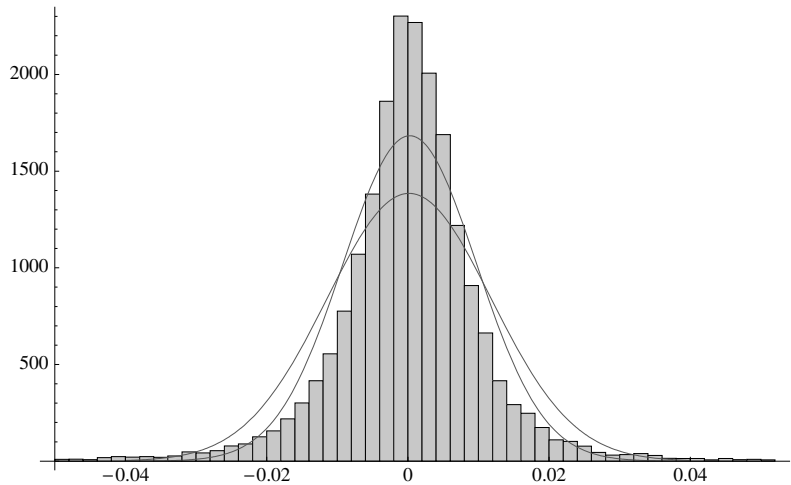


Figure 18. Fitting a normal distribution to truncated returns of the Dow Jones Industrial Average

$.028\sigma$, i.e. it is a very small fraction of the standard deviation. Hence, we will not be committing a big error if, for ease of calculation, we set the mean equal to 0. In fact, setting the mean to 0 corresponds to the situation involving power laws, where we assumed the density is symmetric about 0.

Using these parameter values, how many times should we see a return whose absolute value exceeds .1? We see that $.1 = 8.59\sigma$, so we are asking for the probability that a normally distributed experiment takes on a value that is farther than 8.59σ from the mean. The probability of this is 8.7×10^{-18} , meaning that we should expect a return of this magnitude once every 4.4×10^{14} years. (The reader will recall that there were 12 such returns in the last 80 years.) If we perform a similar calculation for returns whose absolute value exceed .05, we find that the normal distribution model predicts that in 20202 trading days, there should be about 0.35 returns of this size; this should be compared with the actual number of 131.

We have now seen that power laws do a much better job than normal distributions do, in some cases, of modeling the distribution of returns of certain stocks and the Dow Jones Industrial Average. But we have also seen that if we use power laws instead of normal distributions to model returns, there will be a much higher probability of large values. In finance, it is of great importance to understand the risks of various investments. Many measures have been proposed to quantify risk in this area. One measure, called value at risk, or VAR, can be described as follows. One begins by choosing a time interval, say 10 days, and a probability, say .99. Then one computes, using a model of returns on the investment, an interval $J = [a, \infty)$ having the property that with probability .99, the value of the investment at the end of the time period will be inside the interval J . The number VAR is the difference between the value of the investment at the beginning of the time interval and the left-hand endpoint of J . In other words, it is an estimate, under the given model, of the largest amount by which the investment's value could decrease in 10 days, *in 99 cases out of 100*.

The larger the value of VAR is, the riskier the investment is said to be. If one uses normal distributions to model stock prices, then for a given time interval and probability, the interval J will tend to be much smaller than the interval one would obtain by using power laws to model the prices. Thus, if one uses normal distributions, one will tend to drastically underestimate the risk associated with these stocks. Risk underestimates of this type have been cited as one of the causes of the stock market crash in 2008. See, for example, the news article [32].

The reader will recall that Mandelbrot first saw a connection between financial prices and power laws after looking at cotton prices. Figure 19 shows the logarithm of the tail probability versus the logarithm of the input for absolute values of the monthly changes in cotton prices from January 1784 to January 2009. These prices come from the Global Financial Data database. If we use the value of $a = .091$, which corresponds to the value of -2.4 on the horizontal axis in this figure, we obtain a value of $\alpha = 2.64$. The line corresponding to this power law is also shown in the figure. In Figure 20, we have plotted the logarithms of the tail probabilities for the changes in butter prices from January 1, 1890 to the present. It is not clear that a power law fits this graph very well.

5. Independence of Returns

In the last section, we studied the distribution of the geometric returns and price changes of various financial instruments. But a collection of such returns is more than just a set of numbers. We know the order in which these numbers occurred. Experiments in which the order of the outcomes is studied are called time series. It is frequently the case that some very interesting questions can be asked concerning the order of outcomes of an experiment. This is certainly the case when dealing with financial data.

As an example of this situation, we return to our coin-tossing game. We flip a fair coin n times and record each flip. If we simply

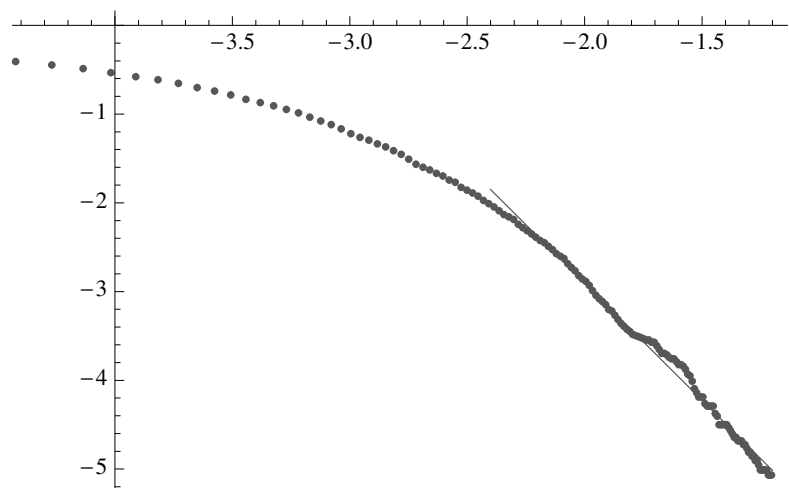


Figure 19. Log of tail probability versus log of input for absolute value of cotton price changes, together with a best-fit line

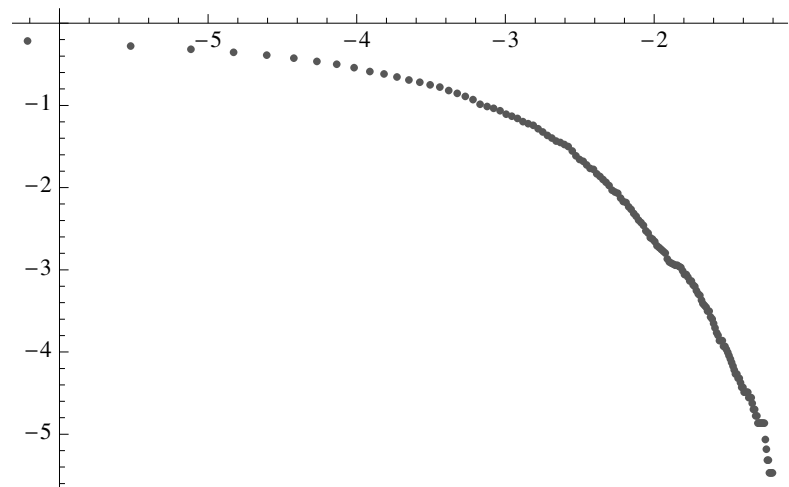


Figure 20. Log of tail probability versus log of input for absolute value of monthly butter price changes

ask for the distribution of the number of heads that occur, the answer is the well-known binomial distribution with parameters n and $1/2$.

Suppose that we define an equalization to mean that at a certain time, the numbers of heads and tails thrown up to that time are equal. We are now considering a concept that depends on the *order* in which the heads and tails occur. One can ask for the distribution of the number of equalizations in n tosses. There are some very interesting statements one can make concerning this distribution. For example, if we double the number of tosses in the experiment, we do not double the average number of equalizations; instead, this average number increases by a factor of about $\sqrt{2}$. As another example, it can be shown that in a sequence of n tosses, the *last* equalization is just as likely to occur in the first half of the sequence as in the last half. These examples show that by considering the order in which a set of outcomes occurs, a richer theory can be developed. For much more on this subject, the reader is urged to consult the first volume [12] of William Feller's masterpiece on probability theory.

The reader will recall that one of the early models proposed for sequences of stock prices was Brownian motion (or perhaps a discrete analogue thereof). In this model, the successive changes in a stock's price are mutually independent, and obey the normal distribution, with some mean and variance. In the previous section, we saw that in the real world, the geometric returns in many cases are not modeled well by a normal distribution. However, it is still possible that by suitably modifying this model, one can obtain a model that fits the data fairly well, at least in some cases. Models in which the changes are identically distributed and the set of changes is mutually independent are traditionally called random walks.

The intuitive meanings of these properties are as follows. We imagine a process that produces a value X_t for each non-negative integer t . One can think of these values as the logarithms of the price of a stock. The sequence of changes for this process is just the sequence

$$X_1 - X_0, X_2 - X_1, \dots$$

In the case of a stock, this is the sequence of geometric returns. For such a process to have changes that are identically distributed means that if we consider many such sequences, coming from this process, then the distributions of $X_{s+1} - X_s$ and $X_{t+1} - X_t$ are the same for any s and t . If, for example, the changes in a stock's price on Mondays tended to be larger than the changes on Fridays, then we would not be able to model the price sequence of this stock with a random walk.

A process has mutually independent changes if the knowledge of a finite set of changes does not give any information about any of the changes that are not in the set. Suppose, for example, that for a certain stock, two successive changes of size at least .01 are never followed by a third change of at least this size. In this case, we would not be able to model the stock's price sequence with a random walk.

There is some evidence that in the case of stocks, the successive changes may not be identically distributed. In other words, it seems that the variance of a stock's returns may change over time. One can guess that models in which the distribution of the increments is allowed to change over time are more complicated to analyze than are random walks. So for now, we will concentrate our attention on the identically distributed case.

We now turn to the question of whether it is roughly true for returns of stocks that the set of changes is mutually independent. The idea of the variance ratio was introduced in Chapter 1; we will remind the reader here of the idea. Suppose that the logarithm of the price of a stock at the end of day t is denoted by X_t . We see that

$$X_{t+2} - X_t = (X_{t+2} - X_{t+1}) + (X_{t+1} - X_t).$$

Under the assumption of mutual independence of the changes, the variance of the left-hand expression is the sum of the variances of the right-hand summands. Thus, the variance of the two-day changes should be twice the variance of the one-day changes. By the same reasoning, if the changes are mutually independent, then the variance of k -day changes should equal k times the variance of one day changes.

This idea has been used by several authors in the following way. One chooses a value of k and then computes the ratio of the variance of the k -day changes to k times the variance of the one-day changes. This statistic is called a variance ratio. We saw above that under the assumption of mutual independence of the changes, this ratio should be about one, for all values of k .

We are now at a typical point in a statistical test of hypothesis. In this case, the hypothesis being tested is mutual independence of the geometric returns of a stock. If the hypothesis is true, then the statistic being computed will have a particular distribution. This distribution may or may not be possible to calculate analytically, but one can usually use simulation to get good estimates of the distribution. It is important to have such values, because one tests the hypothesis for a given data set by computing the statistic for the data set. If the computed value of the statistic is very far from the mean of the distribution of the statistic, this is taken as evidence that the hypothesis is not satisfied by the data. If the distribution of the test statistic cannot be calculated theoretically, one can usually use simulation to obtain good estimates of this distribution.

We illustrate this with an example. Suppose that we have observed that the geometric returns of a certain stock are distributed according to a power law with $\alpha = 2.5$ and $a = .01$. We wish to test the hypothesis that the returns are mutually independent. We begin by choosing a value of k ; we'll choose $k = 2$. Next, we simulate the variance ratio for many simulated sequences of returns distributed according to this power law. These sequences are constructed so that they obey the hypothesis. In our simulation, the length of each sequence is taken to be 1000. We calculate the variance ratio for 10000 sequences and plot the values in Figure 21.

One uses this figure in the following way. We calculate the variance ratio for our stock data. Recall that the data for this stock was used to estimate α and a . Suppose that the variance ratio equals 1.08. One notes that this is quite far from the mean of the simulated distribution, which equals 0.999. In fact, only 108 of the 10000 data

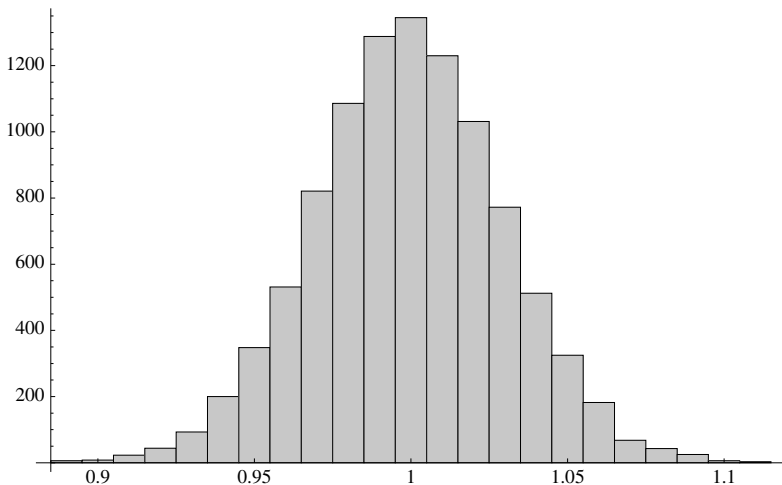


Figure 21. Simulated variance ratio distribution, with $k = 2$, for a power law with $\alpha = 2.5$ and $a = .01$

values, or about 1%, are as far or farther from the mean, in either direction, than the variance ratio of our stock. In other words, the return sequence for the stock is unlike “most” of the sequences that are governed by this power law and obey the assumption of mutual independence of individual returns. This result is taken as evidence that the assumption does not hold for this sequence of returns for this stock.

To understand this idea better, suppose that we generate a sequence of 1000 returns using the above power law, but instead of assuming mutual independence, we create the sequence by using the additional rule that with probability p , each return has the same sign as its predecessor. If $p = 1/2$, then since the power law distribution is symmetric about 0, this is the same as assuming that the returns are mutually independent. But if $p = 2/3$, say, then one might expect that the variance ratio is not close to 1. In fact, if $p > 1/2$, one can see that pairs of consecutive returns are more likely to be of the same sign than if $p = 1/2$. This means that there is less cancellation in the two-day returns, leading to a variance that is more than twice as

large as the one-day variance. Thus, the variance ratio is likely to be greater than 1.

In Figure 22, we show the results of the following simulation. For each value of p of the form $j/10$, for j an integer between 0 and 10, we created ten sequences of 1000 returns using this value of p and using a power law with $\alpha = 2.5$ and $a = .01$. For each sequence, we calculated the variance ratio and computed the average of these ten ratios. The average is what is shown in the figure on the vertical axis.

As expected, we see that if p is substantially different than $1/2$, the variance ratio is not close to 1. Recalling from Figure 21 that almost all of the variance ratios for sequences with mutually independent entries satisfying the given power law are in the interval $[.93, 1.07]$, we see that it would be easy to reject the hypothesis of independence for the sequences in our simulation for which $p \neq 1/2$. (There is an exception to this statement that can be seen in the figure. When $p = 1$, the variance ratio is just about 1. Can the reader supply an intuitive argument as to why this should be the case?)

The point of this simulation is not to show that one can discern non-randomness in the sign sequence of a set of returns using the variance ratio; this can be discerned in easier ways. Rather, we are simply showing that sometimes, if the returns are not independent, the variance ratio of the sequence will provide evidence that this is the case.

In Poterba and Summers [35], applications of the variance ratio are given for testing independence of returns in various settings. We will briefly summarize some of their results here. The authors created two portfolios of stocks, called value-weighted and equal-weighted. The time interval under consideration was the period from 1926 to 1985. At the beginning of this time interval, an equal dollar value of each stock on the New York Stock exchange was put into the first portfolio. Similarly, an equal number of shares of each stock was put into the second portfolio. They computed variance ratios for various pairs of time interval lengths, as we did above.

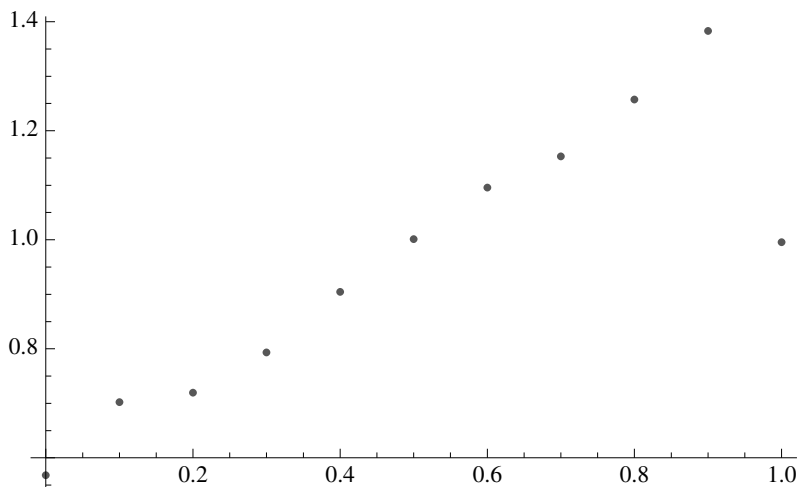


Figure 22. Average variance ratios for various values of p for the simulation described in the text

When they computed the variance ratio for eight-year returns versus one-year returns, they found it to be .509 for the first portfolio and .290 for the second portfolio. These values are quite far from the value of 1 that obtains under the assumption of mutual independence of returns. Comparing these values to the distributions of the variance ratio (obtained by simulation, as we did above) allowed them to reject the hypothesis of mutual independence of returns at the .08 level for the first portfolio and the .005 level for the second portfolio. Among their other results, they also computed the variance ratio for both portfolios for one-year returns versus one-month returns. For both portfolios, they found that the variance ratio was about 1.27.

These results suggest that there is much less variation in the price of a stock over an eight-year period than would be suggested by looking at the average variation of that stock over a one-year period. However, the opposite statement seems to be true for one-month versus one-year periods; the variance over a 12-month period tends to be greater than 12 times the variance over a one-month period. The authors state that the stocks “revert to the mean,” by which they mean

that a stock that grows, for a period of a year, more rapidly than does the market, will tend to have slower growth than the market for some time afterwards.

6. Is the Power Law Exponent Intrinsic?

In Section 4, we gave some examples showing that the geometric returns of the stock prices of certain companies appear to obey a power law. This information is useful in understanding the variability in a given company's stock price. We recall that the smaller the value of α in a power law, the more likely it is that there will be large values in a process that is governed by the power law. Thus one can see that the smaller the value of α is, the more volatile the stock prices are. So one might think that a good estimate of α for a given stock is useful information to have.

There is a possible problem with this idea, and it can be illustrated by considering an example. We use as our data set the geometric returns for IBM stock between January 2, 1962 and May 20, 2009. There are 11925 returns in this time period. If we split this time period in half, and for each half, plot the log of the tail probabilities versus the log of the input (as we have done numerous times above), we obtain the plots shown in Figure 23. Both plots have been put in this figure. The steeper-sloping set of points corresponds to the first time period.

The two values of α that correspond to these two time periods are different; the estimates are $\alpha = 3.07$ and $\alpha = 2.90$. The estimate for α when the entire data set is used is $\alpha = 2.97$. There are two possibilities that one might consider at this point. First, it is possible that there is an 'intrinsic' value of α for IBM stock, and the variations we see above are simply natural variations that can arise when one takes different samples from a distribution. (For example, if one takes samples of a certain size from a given normal distribution, and calculates their means, these means will vary.) Second, it might be the case that over time, the value of α changes for IBM stock. In this case, unless we can

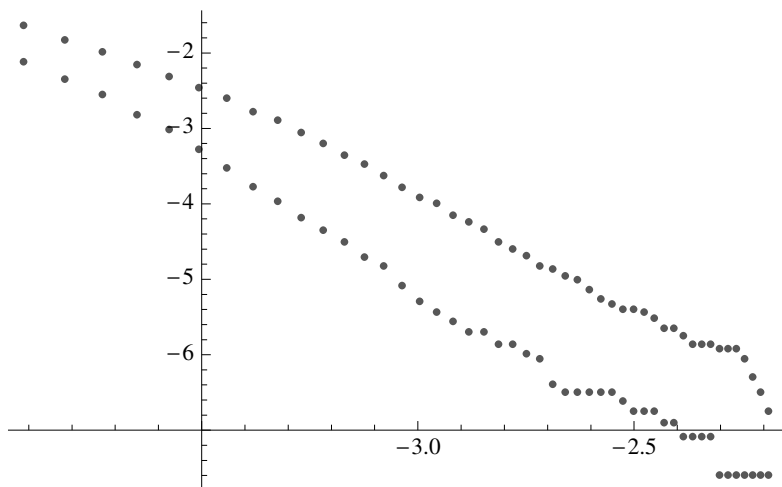


Figure 23. Log of tail probability versus log of input for returns of IBM for the first and second halves of the time period 1-2-62 through 5-20-09

understand how it changes, our will be of little use in understanding how IBM stock might behave in the future.

To investigate whether the first of the above possibilities is reasonable, we carried out the following simulation. We created 1000 samples from a power law distribution with $\alpha = 3.14$ and $a = .02$. Each sample contained 5900 values, as did each of the two IBM data sets corresponding to the two time periods described above. For each of these 1000 samples, we estimated the value of α . The results are shown in Figure 24.

One sees from this histogram that the observed values of α for the actual returns of IBM stock over the two time periods described above (these values are $\alpha = 3.07$ and $\alpha = 2.90$) are not unusual values at all. In other words, these values do not, by themselves, discredit the hypothesis that α does not change over time for a particular stock.

In Figure 25, we graph the estimated values of α for consecutive blocks of IBM stock prices of length 6000, where each block is offset

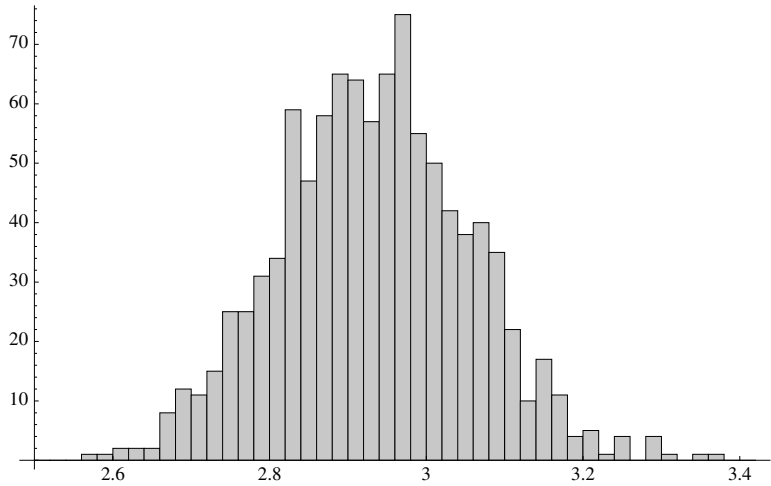


Figure 24. Estimated values of α for samples taken from a power law with $\alpha = 2.97$

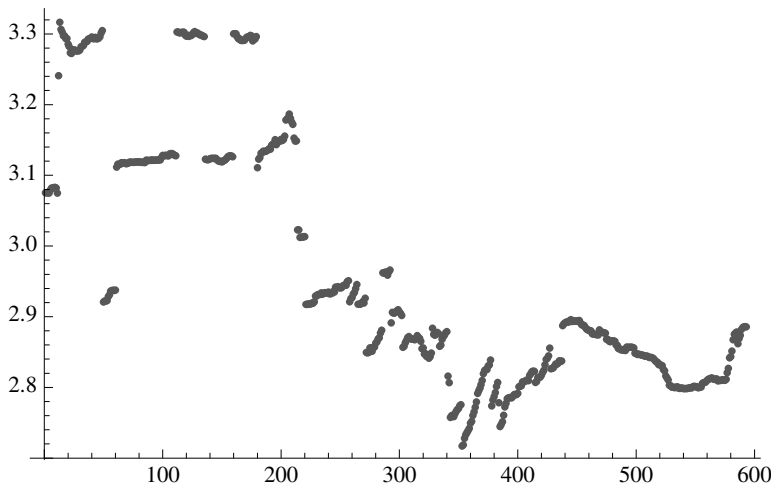


Figure 25. Estimated values of α for blocks of consecutive prices of IBM stock

from the previous block by ten days. Figure 26 shows the results of the same calculation applied to General Electric stock over the same

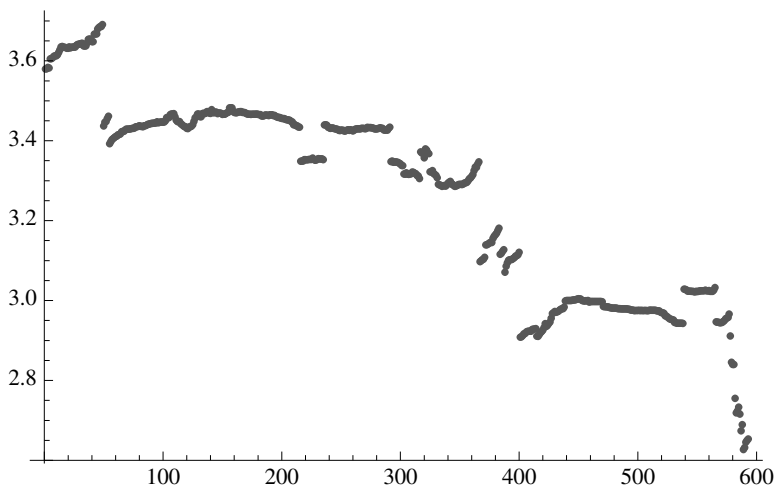


Figure 26. Estimated values of α for blocks of consecutive prices of General Electric stock

time period. The abrupt jumps in the graphs are probably artifacts of the algorithm we use to estimate the values of α for the blocks. Even if we ignore the jumps, these two figures show a clear downward trend in the sequences of estimates of α , suggesting that the value does in fact change over time.

7. Appendix

We begin with a brief description of the power law distributions that were used in the simulations in this chapter. We recall that a power law has a density function that is roughly of the form

$$f(t) = C_1 |t|^{-k},$$

where t is assumed to be not too close to the origin. Below, we will define f for values of t near the origin. The parameter of greatest interest in such a power law is the number α , defined to equal $k - 1$. If L represents a quantity that is governed by this power law, then

we saw above that for large positive x ,

$$(4) \quad P(L \geq x) = C_3 x^{-\alpha}$$

for some positive constant α . Since $f(t)$ is symmetric with respect to the origin, a similar equation holds for $P(L \leq x)$ where x is large and negative.

There are various other parameters in our distributions which we now define. The number a is the bound above which Equation 4 holds. The constant C_3 in Equation 4 must be chosen so that the area under $f(t)$ and above the interval $[a, \infty)$ is less than $1/2$ (since by symmetry, the area above the interval $(-\infty, -a]$ will be the same as this area). In fact, we wish to make $f(t)$ unimodal. This is not a requirement of a power law, but it is a reasonable assumption to make about a density that models stock price changes. If this is so, then $f(t)$ decreases on the interval $[0, a]$, so the area under $f(t)$ and above $[0, a]$ must be at least $af(a)$.

Finally, it would be nice to be able to simulate, using a computer, the quantity L . One method for accomplishing this involves inverting the cumulative distribution function corresponding to the density $f(t)$. If we write

$$F(x) = \int_{-\infty}^x f(t) dt,$$

then

$$F(x) = P(L \leq x).$$

This quantity is always between 0 and 1. Also, F is an increasing function on the real line. Hence it has an inverse function G i.e. G is a function from $[0, 1]$ to the real line such that for all $r \in [0, 1]$, we have

$$F(G(r)) = r.$$

Suppose that r is uniformly chosen (say, by a computer) in the interval $[0, 1]$. We let L be given by $G(r)$. To see that L is distributed according to the density function $f(t)$, it suffices to show that for every real value of x , the fraction of simulated values of L not exceeding x equals $F(x)$.

So let x be a fixed real number. Let r be the real number in $[0, 1]$ with the property that $G(r) = x$. If $r' < r$, then $G(r') < x$, and vice versa. So the fraction of the simulated values of L not exceeding x equals the fraction of randomly chosen real numbers r' not exceeding r . This equals r , since the real numbers are chosen uniformly in $[0, 1]$. But we know that $G(r) = x$ and $F(G(r)) = r$, so we see that this fraction equals $F(x)$, as desired.

In the domain $[a, \infty)$, we have

$$\begin{aligned} F(x) &= P(L \leq x) \\ &= 1 - P(L \geq x) \\ &= 1 - C_3 x^{-\alpha}, \end{aligned}$$

and this last expression is easily inverted. If we make $f(t)$ linear on $[0, a]$, say

$$f(t) = C_4 - C_5 t,$$

then for $x \in [0, a]$, we have

$$\begin{aligned} F(x) &= \int_{-\infty}^x f(t) dt \\ &= \frac{1}{2} + \int_0^x f(t) dt \\ &= \frac{1}{2} + C_4 x - \frac{1}{2} C_5 x^2, \end{aligned}$$

and this expression is also easy to invert.

One (arbitrary) choice of parameters, given that we have already chosen α and a , is as follows:

$$\begin{aligned} C_3 &= \frac{a^\alpha}{4(1 + \alpha)}, \\ C_4 &= \frac{2 + 3\alpha}{4a(1 + \alpha)}, \\ C_5 &= \frac{1}{2a^2}. \end{aligned}$$

For example, if we choose $\alpha = 2$ and $a = 1$, then the graph of the resulting density function $f(t)$ is as in Figure 27.

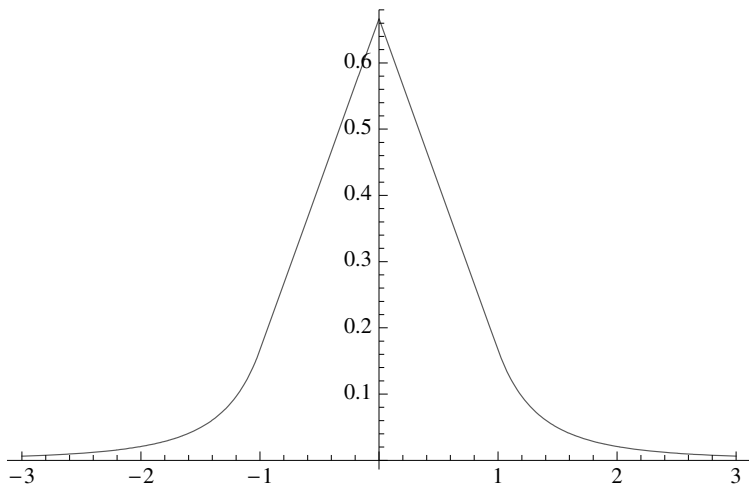


Figure 27. A power law density with $\alpha = 2$ and $a = 1$

We created these power law densities for use in illustrating some of the ideas in the modeling of stock prices. However, these densities are not the ones that are typically used. Ours are simpler than the ones in use and are adequate for our purposes. We will now give a brief explanation of the densities that are used in more advanced treatments.

Recall that the model of Bachelier used normal distributions to model arithmetic returns and assumed mutual independence of returns in non-overlapping time intervals. We saw above that instead of the normal distribution, power laws should be used to model the returns. We have also seen that the assumption of independence may not be justified. The second of these statements was demonstrated somewhat later (in the 1980's) than the first (first discussed by Mandelbrot in the early 1960's).

Because of this time difference, models were created that mimicked Brownian motion in one way, namely that non-overlapping increments were independent, but used power laws instead of normal distributions. A crucial property enjoyed by the family of normal distributions is the property of stability. A set of distributions is said to be stable if, when we have two independent experiments whose outcomes are governed by distributions from the set, then the distribution of the sum of the outcomes is also in the set. It can easily be shown by induction that in a stable family, the sum of any finite number of outcomes from independent experiments is distributed according to a member of the stable family.

In terms of geometric returns, if the distribution of daily returns, say, is a member of a stable family, and if we assume that the daily returns are mutually independent then, for example, the distribution of weekly returns will also be a member of the same stable family. This is a pleasing property for a set of distributions to possess, and it seems to be at least approximately true for geometric returns of stocks. In other words, the observed distributions of weekly returns look similar to those of daily returns.

We can easily perform a simulation of the sum of two independent experiments, each distributed according to a power law. In our simulation, we chose $\alpha = 2.5$ and $a = .5$. We constructed two lists of outcomes, each with 100,000 entries. Then we added the corresponding entries together and computed the tail probabilities. The results are shown in Figure 28. If we throw out the leftmost seven points in the figure and fit a line to the remaining points, we find the line has slope -2.48 , leading to a value of $\alpha = 2.48$ for the sum of the two experiments.

Given two independent experiments whose outcomes obey two (possibly different) distributions, it is typically quite difficult to find the distribution of their sum. This distribution is called the convolution of the distributions of the two summands. The most common way to proceed is to use Fourier transforms, a method that is beyond the scope of this book. Nevertheless, we will explain the basic idea

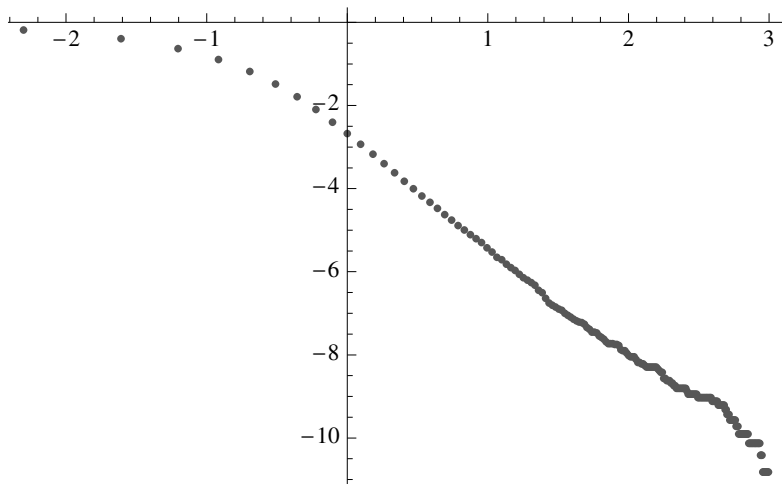


Figure 28. Log of tail probability versus log of input for the sum of two experiments each distributed according to a power law with $\alpha = 2.5$ and $a = .5$

here, since it helps in gaining a bit of understanding of the set of power law distributions that are commonly used.

If $f_1(t)$ and $f_2(t)$ denote the density functions of the two independent experiments, and if $f_3(t)$ denotes the density function of their sum, then it is fairly easy to show that

$$(5) \quad f_3(t) = \int_{-\infty}^{\infty} f_1(s)f_2(t-s) \, ds.$$

However, evaluating this integral is typically quite hard, or even impossible, by direct methods. (We will briefly discuss below the case when f_1 and f_2 are chosen to be two of our power law distributions.) The Fourier transform (or characteristic function) of a density function $f(t)$ is another function $f^*(w)$ with the following very helpful property:

$$f_3^*(w) = f_1^*(w)f_2^*(w).$$

In words, the Fourier transform of the density function of the sum of two independent experiments is the product of the Fourier transforms of the density functions of the summands. It is much easier to compute this product than it is to compute the above integral.

Of course, very little is free in life, and this situation is no exception. One has replaced the calculation of the above integral by the calculation of two Fourier transforms (one for each summand). Furthermore, after the product of the two transforms has been computed, one must determine the density function whose transform equals this product. In some cases, this last step is easy, since the product of the transforms is recognizable as the transform of some known density. This would typically be true when one is dealing with stable sets of distributions. In other cases, one computes the inverse Fourier transform of the product, thereby obtaining the required density.

The power law densities that are in use have relatively benign Fourier transforms (see McCulloch [30]), and it is relatively easy to show that the product of two such transforms, under suitable conditions concerning their parameters, is another function of the same form, thereby showing that the set of such densities forms a stable set. However, these densities all have values of α not exceeding two. The reason for this is if one sets α to a value larger than two in the defining formula, then the resulting function is not a density function; its integral over the positive reals exceeds one. This is somewhat troublesome, as we have seen above that the returns of some stocks seem to obey a power law with α near three.

It is a (hard) exercise to apply Equation 5 directly to the case when f_1 and f_2 are taken to be two of our power law density functions. Since we are interested in summing the returns of two successive time periods for the same stock, it makes sense to choose the same values of α and a for both density functions. One quickly finds that the integral in Equation 5 is intractable unless α is an integer. If one tries the value $\alpha = 3$, it is possible to show that for large $|t|$, the convolution density $f_3(t)$ is asymptotic to ct^{-4} , where $c = (3/32)a^3$. This shows that the sum of two independent experiments, each distributed according to

one of our power laws with $\alpha = 3$, is also distributed (for large $|t|$) according to a power law with the same value of α .

To those readers with knowledge of the Central Limit Theorem, the last statement might seem to contradict this theorem. Suppose that we say that an experiment is of class 3 if it is distributed according to a power law with $\alpha = 3$. Then the contradiction appears to stem from the fact that if the sum of two independent class 3 experiments is of class 3, then by induction, any finite sum of mutually independent class 3 experiments is also of class 3. But since the density of a class 3 experiment possesses a variance, the Central Limit Theorem implies that the normalized sum of mutually independent class 3 experiments approaches a standard normal.

This apparent contradiction can be resolved by noting that the power law property is asymptotic in terms of the input, i.e. if X is a random variable that is of class 3, then

$$P(X \geq x) \sim Cx^{-\alpha},$$

for some constant C . This condition says nothing about what happens for small values of $|x|$. In fact, the reader will recall that our power law densities do not behave according to a power law for small $|x|$. The sum of many class 3 experiments will behave like a normal density over a large x -interval. In fact, this interval will increase in size as the number of summands increases. Nevertheless, this result means that one should be very wary of using the normal distribution to model n -day stock price changes, even for moderate values of n .

The following theorem shows that the above statement concerning power laws of class 3 actually holds for power laws of any class (see [13], p. 278). In fact, one can say something even in the case where the summands are power laws of different classes (see [10]). Here is a statement of the result.

Let X be a real-valued random variable. We say that the right tail of X is asymptotic to $f(x)$ if

$$\lim_{x \rightarrow \infty} \frac{P(X > x)}{f(x)} = 1.$$

A similar definition holds for left tails. For example, the Cauchy distribution with density $1/(\pi(1+x^2))$ has right and left tails asymptotic to $1/(\pi x)$. Then the following holds.

Theorem 2. *Let X and Y be independent random variables with right tails asymptotic to c/x^α and d/x^β , where $0 < \alpha \leq \beta$ and $c, d > 0$. If $\alpha = \beta$ then $X + Y$ has right tail asymptotic to $(c + d)/x^\alpha$. If $\alpha < \beta$ then $X + Y$ has right tail asymptotic to c/x^α .*

Chapter 3

Lotteries

1. Rules of the Powerball Lottery

Lotteries are discussed frequently in the news. They are the most popular form of gambling and an increasingly important way for states to obtain revenue. In this chapter, we will use the Powerball lottery to illustrate some of the statistical ideas associated with lotteries. Our calculations are based on the rules in that were in effect for the February 18, 2006 drawing, whose \$365,000,000 jackpot is still the largest jackpot in the history of Powerball. This example will feature prominently throughout the chapter. We provide some further discussion on the evolution of the game below.

The Powerball Lottery is a multi-state lottery, a format which is gaining popularity because of the potential for large prizes. It is currently available in 42 states, Washington, D.C., and the U.S. Virgin Islands. It is run by the Multi-State Lottery Association, and we shall use information from the Powerball homepage. We found their “Frequently Asked Questions,” (hereafter abbreviated FAQ) to be particularly useful. These are compiled by Charles Strutt, the executive director of the association.

A Powerball lottery ticket costs \$1. For each ticket you are asked to mark your choice of numbers in two boxes displayed as shown in

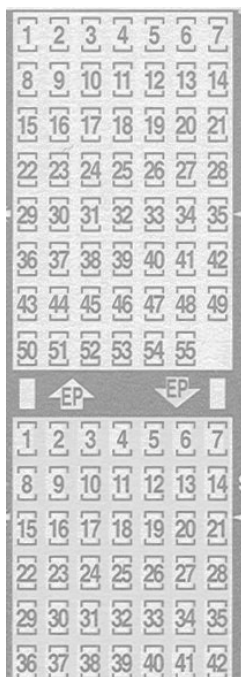


Figure 1. Picking your numbers

Table 1. You are asked to select five numbers from the top box and one from the bottom box. The latter number is called the “Powerball.” If you check EP (Easy Pick) at the top of either box, the computer will make the selections for you. You also must select “cash” or “annuity” to determine how the jackpot will be paid should you win. Finally, there is another option called the “Power Play.” In what follows, we will refer to a particular selection of five plus one numbers as a “pick.”

The Powerball Lottery was started in 1992. In the history of this lottery, there have been at least five versions. Before November 2, 1997, there were 45 numbers in the top box and 45 in the bottom box. On that date, these numbers were changed to 49 and 42, respectively. They were changed again on October 9, 2002 to 53 and 42, respectively. On August 28, 2005, they were changed to 55 and 42.

You Match	You Win	Win Probability
5 white balls and the red ball	JACKPOT	1/146,107,962
5 white balls but not the red ball	\$200,000	1/3,563,609
4 white balls and the red ball	\$10,000	1/584,432
4 white balls but not the red ball	\$100	1/14,254
3 white balls and the red ball	\$100	1/11,927
3 white balls but not the red ball	\$7	1/291
2 white balls and the red ball	\$7	1/745
1 white ball and the red ball	\$4	1/127
0 white balls and the red ball	\$3	1/69

Table 1. Possible prizes and their probabilities

On January 7, 2009, they were changed to 59 and 39. Unless stated otherwise, the calculations below will refer to the 2005 version. The interested reader can easily modify the calculations in this chapter to take into account any subsequent changes in the lottery.

Every Wednesday and Saturday night at 10:59 P.M. Eastern Time, lottery officials draw five white balls out of a drum with 55 balls and one red ball from a drum with 42 red balls. Players win prizes when the numbers on their ticket match some or all of the numbers drawn. The order in which the numbers are drawn does not matter. There are 9 different prizes. In Table 1 we give the possible prizes and their probabilities. We show how to calculate the probabilities in the next section.

The jackpot starts at \$15,000,000 and increases based on sales each time there is no winner. Each time a new record jackpot is reached, there is enormous media attention. Beginning in 2002, a new rule was added to control the growth of jackpots by spreading the prize money. Specifically, once a new record is reached, subsequent unsuccessful drawings will increase the jackpot by at most \$25,000,000. The remaining funds are added to a Bonus Prize Pool. When someone wins the record jackpot, the Bonus pool is divided

among the winners of the second prize (for matching 5 white balls but not the red ball), in addition to the usual \$200,000.

For other prizes, the player wins the amount described above, unless she selected the Power Play option (which costs an additional \$1 dollar per ticket). With the Power Play feature, all prize amounts except the jackpot are multiplied by either 2, 3, 4, or 5. This multiplier is the same for all players in a given lottery and is determined by spinning a Power Play wheel at the time the other balls are drawn. The wheel has sixteen equally likely slots, with four occurrences each of the multipliers 2, 3, 4, and 5.

If the player wins the jackpot, he or she must share it equally with all other players (if there are any) who have also won the jackpot. A few other comments concerning the jackpot are in order. First, the winning players have 60 days after they have won to declare whether they want a lump sum payment or a series of 30 graduated annuity payments. The first payment is made immediately and the others are made at the end of each subsequent year. The lump sum is the present value of the annuity, which is typically only about 50% of the announced value of the jackpot.

2. Calculating the Probabilities of Winning

The first question we ask is: how are the probabilities of the prizes determined? In what follows, we will calculate various probabilities assuming the player has bought one ticket. This is a counting problem that requires that you understand one simple counting rule: if you can do one task in n ways and, for each of these, another task in m ways, the number of ways the two tasks can be done is mn . A simple tree diagram makes this principle very clear.

When you watch the numbers being drawn on television, you see that, as the five winning white balls come out of the drum, they are lined up in a row. The first ball could be any one of 55. For each of these possibilities the next ball could be any of 54, etc. Hence the number of possibilities for the way the five white balls can come out

in the order drawn is

$$55 \cdot 54 \cdot 53 \cdot 52 \cdot 51 = 417,451,320.$$

But to win a prize, the order of these 5 white balls does not count. Thus, for a particular set of 5 balls all possible orders are considered the same. Again by our counting principle, there are $5 \cdot 4 \cdot 3 \cdot 2 \cdot 1 = 120$ possible orders. So, in the above count of the number of possibilities, each one has been counted 120 times. Thus, the number of possible sets of 5 white balls not counting order is

$$\frac{417,451,320}{120} = 3,478,761.$$

This is the familiar problem of choosing a set of 5 objects out of 55, and we denote the answer by $C(55, 5)$. Such numbers are called binomial coefficients. The general formula for these numbers is the following:

$$C(n, k) = \frac{(n) \cdot (n-1) \cdot \dots \cdot (n-k+1)}{(k) \cdot (k-1) \cdot \dots \cdot (1)}.$$

This number equals the number of possible sets of k objects chosen from a set of n objects, where the order of the chosen objects in the set is disregarded. The derivation of the above formula follows the same line of argument as the one given above.

Applying this formula to the example above gives

$$C(55, 5) = \binom{55}{5} = \frac{(55) \cdot (54) \cdot \dots \cdot (51)}{(5) \cdot (4) \cdot \dots \cdot (1)} = 3,478,761.$$

Now for each pick of five white numbers there are 42 possibilities for the red Powerball, so the total number of ways a set of six numbers can be chosen is

$$42 \cdot C(55, 5) = 146,107,962.$$

We will need this number often and denote it by b (for big).

The lottery officials go to great pains to make sure that all b possibilities are equally likely. So, a player has one chance in 146,107,962

of winning the jackpot. However, the player will have to share the prize with anyone else who has picked the same set of numbers.

We note that on the Powerball website, the column corresponding to the last column in Table 1 is labeled “odds.” The numbers in the column are in fact probabilities, not odds. The media prefers to use odds, and textbooks prefer to use probability or chance. Because the probabilities are small, there is not much difference between odds and probabilities. However, this is a good excuse to get the difference between the two concepts straightened out.

Suppose we are dealing with a chance situation in which there are f favorable outcomes and u unfavorable outcomes. Suppose in addition that all of the outcomes are equally likely. Then the probability of a favorable outcome is $f/(f + u)$, i.e. it is the fraction of all of the possible outcomes that are favorable. Odds are quoted in two different ways. The odds in favor of a favorable outcome are f to u , and the odds against a favorable outcome are u to f . Thus the chance of winning the jackpot is 1 in 146,107,962, whereas the odds are 1 to 146,107,961 in favor, or 146,107,961 to 1 against.

To win the second prize (\$200,000 plus a possible share of any bonus), the player must get the 5 white numbers correct but miss the red Powerball number. How many ways can this be accomplished? There is only one way to get the set of five white numbers, but the player’s Powerball pick can be any of the 41 numbers different from the red number that was drawn. Thus, the chance of winning second prize is 41 in 146,107,962; rounded to the nearest integer this is 1 in 3,563,609.

When there were 45 white balls and 45 red balls, the ticket listed the chances of getting only the red ball as 1 in 84. This often seemed wrong to players who have had elementary probability, as the following exchange from the Powerball FAQ¹ illustrates:

¹From the Multi-State Lottery Association website at <http://www.musl.com/>

COULD YOUR ODDS BE WRONG?

I have a simple question. You list the odds of matching only the powerball as one in 84 on the powerball “9 Ways to Win” page. From my understanding of statistics (I could be wrong, but I got an A), the odds of selecting one number out of a group is simply one over the number of choices. Since there are not 84 choices for the powerball, may I assume the balls are somehow “fixed” so that some are more common than others? Otherwise, the listed odds are somehow defying the laws of statistics. I am really very eager to hear your explanation, so please return my message. Thank you.

Susan G., via the Internet.

This is one of the most common questions we get about the statistics of the game. If you could play only the red Powerball, then your odds of matching it would indeed be 1 in 45. But to win the \$1 prize for matching the red Powerball alone, you must do just that: match the red Powerball ALONE. When you bet a dollar and play the game, you might match one white ball and the red Powerball. You might match three white balls and the red Powerball. To determine the probability of matching the red Powerball alone, you have to factor in the chances of matching one or more of the white balls too.

Charles Strutt

To win this last prize you must choose your six numbers so that only the Powerball number is correct. In the old version of the Powerball lottery this would be done as follows: there are $45 \cdot C(45, 5) = 54,979,155$ ways to choose your six numbers. But here your first 5

Number of Ways	Match
$n(1) = 1$	all six balls
$n(2) = 41$	5 white balls but not the red ball
$n(3) = C(5, 4) \cdot C(50, 1)$	4 white balls and the red ball
$n(4) = n(3) \cdot 41$	4 white balls but not the red ball
$n(5) = C(5, 3) \cdot C(50, 2)$	3 white balls and the red ball
$n(6) = n(5) \cdot 41$	3 white balls but not the red ball
$n(7) = C(5, 2) \cdot C(50, 3)$	2 white balls and the red ball
$n(8) = C(5, 1) \cdot C(50, 4)$	1 white ball and the red ball
$n(9) = C(50, 5)$	only the red ball

Table 2. Number of ways that a particular prize can be won

numbers must come from the 40 numbers not drawn by the lottery. This can happen in $C(40, 5) = 658,008$ ways. Now there is only one way to match the Powerball number, so overall you have 658,008 chances out of 54,979,155 to win this prize. This reduces to 1 chance in 83.55, or about 1 chance in 84, in agreement with the official lottery pronouncement.

The same kind of reasoning carries over to the present version of the game. To find the chance of winning any one of the prizes we need only count the number of ways to win the prize and divide this by the total number of possible picks b . Let $n(i)$ be the number of ways to win the i th prize. Then the values of $n(i)$ are shown in Table 2. Dividing these numbers by b , we obtain the chance of winning the corresponding prizes given in Table 1. Adding all the $n(i)$ values gives a total of 3,991,302 ways to win something. Thus we get an overall chance of winning of $3,991,302/b = 0.1732$, which is about 1 in 36.61.

In a textbook, we would be apt to give the results of Table 1 in the form shown in Table 3. As noted earlier, rounding the reciprocals of these probabilities to the nearest integer gives the numbers reported as “odds” on the lottery ticket.

You Match	You Win	Win Probability
5 white balls and the red ball	JACKPOT	0.0000000068
5 white balls but not the red ball	\$200,000	0.0000002806
4 white balls and the red ball	\$10,000	0.0000017111
4 white balls but not the red ball	\$100	0.0000701536
3 white balls and the red ball	\$100	0.0000838421
3 white balls but not the red ball	\$7	0.0034375266
2 white balls and the red ball	\$7	0.0013414738
1 white ball and the red ball	\$4	0.0078811585
0 white balls and the red ball	\$3	0.0145013316

Table 3. The probabilities of winning

Exercise.

1. Which of the two methods for presenting the chances of winning, Table 1 or Table 2, do you think is better understood by the general public? Which do you prefer?

3. What is Your Expected Winning for a \$1 Ticket?

The value of a gambling game is usually expressed in terms of the player's expected, or average, winning. If there are n prizes and $p(i)$ is the probability of winning the i th prize $w(i)$, then your expected winning is:

$$E = w(1) \cdot p(1) + w(2) \cdot p(2) + \dots + w(n) \cdot p(n).$$

As a simple example to illustrate the above formula, consider the following game. A fair coin is tossed twice. The number of dollars won in one play of this game is defined to be the square of the number of heads that appear. Thus, the possible payoffs are 0, 1, and 4 dollars. The expected winning is the average amount won per play. We can estimate this average amount by imagining that we have played the game 100 times. We would expect to obtain the 1 dollar payoff about

50 times, because if we toss a coin twice, one head will appear about half of the time. Similarly, we would expect to obtain the 4 dollar payoff about 25 times, because two heads occur about one quarter of the time in this game. Thus, we would expect our total payoff in 100 plays to be about \$150. Therefore, we estimate the average winning to be \$1.50. The above formula for the expected winning gives

$$E = 0 \cdot \frac{1}{4} + 1 \cdot \frac{1}{2} + 4 \cdot \frac{1}{4} = 1.5,$$

in agreement with our estimate.

We will first discuss the case when the Power Play option is not chosen and the bonus mechanism is not active (that is, we have not exceeded the record jackpot by more than \$25,000,000). Later, we will summarize the corresponding case when this option is chosen. For all prizes, except the jackpot, the value of the prize is known. However, since the size of the jackpot differs significantly from drawing to drawing, we will want to find the expected winning for different jackpot sizes. In the 508 drawings from the beginning of the lottery on April 22, 1992 through March 1, 1997 the jackpot was won 75 times. It was shared with one other winner 11 times. During this period the jackpot prize ranged from \$2,000,000 to \$314,900,000.

If x is the amount of the jackpot and $p(i)$ the probability of winning the i th prize, the expected winning is:

$$\begin{aligned} E &= x \cdot p(1) + 200,000 \cdot p(2) + 10,000 \cdot p(3) + 100 \cdot p(4) \\ &\quad + 100 \cdot p(5) + 7 \cdot p(6) + 7 \cdot p(7) + 4 \cdot p(8) + 3 \cdot p(9) \\ (6) \quad &= \frac{x}{b} + 0.197, \end{aligned}$$

where $b = 146,107,962$. The last expression above says that the expected value of a ticket has two components, the amount attributable to the jackpot, and the amount attributable to all of the other prizes. As we will soon see, the amount won by a jackpot winner is affected by many things. The second component is not changed by any of these things except for a possible bonus for the second prize.

We start, therefore, in the simplest case by assuming that hitting the jackpot gives the player the full amount, ignoring taxes and the

x = Jackpot (in millions of dollars) E = Expected winnings (in dollars)

20	0.334
40	0.471
60	0.608
80	0.745
100	0.882
120	1.018
140	1.155
160	1.292
180	1.429
200	1.566
220	1.703
240	1.840
260	1.977
280	2.114
300	2.250
320	2.387
340	2.524
360	2.661
380	2.798
400	2.935

Table 4. Expected winnings for different size jackpots

possibility of sharing the prize. Table 4 shows the expected winning E for various values of the jackpot.

A game is said to be *favorable* if the expected winning is greater than the cost of playing. Here we compare the expected winning with the \$1 cost of buying a ticket. Looking at Table 3, we see that the lottery appears to be a favorable game as soon as x gets up to \$100 million.

The jackpot for the Powerball lottery for February 8, 2006 built up to \$365 million, as hordes of players lined up at ticket outlets for a

shot at what had become the largest prize for any lottery in history. At first glance, it certainly looks as if this was a favorable bet!

However, the reader will recall that the winner must choose, within 60 days of winning the jackpot, whether to take a lump sum payment of cash or to take an annuity. In fact, the Powerball website regularly updates the value of the jackpot for each choice. For example, a few months after the record, the website was showing the estimated jackpot for the May 24, 2006 drawing as \$20 million, with a cash value of \$8.9 million.

The \$20 million here corresponds to the \$365 million from the record lottery, and is the number that the media likes to hype. But note that this corresponds to the annuity amount to be paid out over time, not the immediate cash value. You are not going to get this money tomorrow; in fact the lottery doesn't even have it on hand! This is explained further in an earlier excerpt from the FAQ:

When we advertise a prize of \$100 million paid over 29 years (30 payments), we actually have less than \$50 million in cash. When someone wins the jackpot and wants cash, we give them all of the cash in the jackpot prize pool. If the winner wants the annuity, we invest the \$50 million in cash to fund the annuity payments. The winner gets the cash plus the interest earned. When you see an estimated jackpot annuity prize, we are guessing both what the sales will be and what the market's bond prices will be. The annuity jackpot amount and the cash jackpot amount that we announce are always estimates.

The cash option on the record \$365 million jackpot was \$177.3 million. (The ratio of the cash option to the jackpot size depends upon current interest rates. It is usually about $1/2$.) From Table 3, we see that a \$1 ticket still looks like a favorable bet, with expected value of about \$1.40. We have been assuming that the player has elected the lump

sum cash payment, and treating the annuity as equivalent in present value terms. You may want to think harder about this. An article² in the *Star Tribune* discusses the question of lump sum or annuity. The article is based on an interview with Linda Crouse, a financial planner and certified public accountant in Portland, Oregon. (Note that this article was based on an earlier version of the lottery.) From this article we read:

Crouse ran numbers to help determine whether it's better to take a windfall in payments over time—an annuity—or in a lump sum.

Crouse used the Powerball jackpot as an example to determine which pays off in the long run: the ticket that pays the winner \$104.3 million now or pays \$7.7 million annually for 25 years. (Both are before taxes.)

The annuity represents a 5.4 percent return. That sounds easy to beat if you take the lump sum and invest it—until you consider the huge negative effect of paying all the taxes up front instead of over 25 years. Figure 45 percent of the payout—\$46.9 million—goes to state and federal taxes right off the bat. If you invest the remaining \$57.4 million and receive an average return of 8 percent, you still can't beat the annuity. After all taxes are paid, you receive \$4,235,000 annually for the annuity vs. \$3,991,000 for the lump sum you invested at 8 percent.

Beyond about a 9 percent return, you start to beat the annuity.

Of course, one should consider the fact that the annuity is a guaranteed payment while your investments are subject to the volatility of the way you invest your money.

²Julie Trip, *Star Tribune*, June 7, 1998, Metro Section, p. 1D

Well, at least with the lump sum above, we convinced ourselves that we had a favorable game. Alas, there is another rub. We have been implicitly assuming that if we hold the lucky numbers, we will get the whole prize! But if other ticket holders have selected the same numbers, the jackpot will be split. This will be a particularly important factor when large number of tickets are sold. As the jackpot grows, an increasing number of tickets are sold. For example, for the February 8, 2006 Powerball lottery, about 175.3 million tickets were sold.

The chance of having to share the jackpot depends upon several factors. First, it depends upon the number of tickets sold. Second, it depends upon whether the numbers are all roughly equally likely to be chosen. In the Easy Pick method, the numbers are chosen by the computer, and we can therefore assume that they are all equally likely. We are told that about 70% of tickets sold in a typical lottery are chosen by the Easy Pick method. Probably this percentage is even larger when the jackpot is large since people tend to buy a number of tickets and would be more likely to use the Easy Pick method when they do this.

The remaining tickets have numbers that are chosen by the ticket buyers. Figure 7 (displayed later in this chapter) shows that the numbers are not chosen with equal probabilities by the buyers. In a given winning set of numbers, some of them will be more likely than others to have been chosen by these buyers. For the reasons stated above, we expect that this will have little effect on the number of jackpot winners.

Because the effect of the non-Easy Pick tickets is small, we will use $n = 175,300,000$ as the number of tickets sold. We will assume that they were all chosen by the Easy Pick method. The probability that a particular ticket is the winning ticket is $1/146,107,962$. The probability of k winners can be obtained from a binomial distribution with $p = 1/146,107,962$ and $n = 175,300,000$. The expected number of

k	$P(k \text{ winners})$
0	.3013
1	.3614
2	.2168
3	.0867
4	.1600
5	.0062
6	.0012
7	.0002

Table 5. Poisson probabilities for numbers of jackpot winners

winning tickets is

$$np = \frac{175,300,000}{146,107,962} = 1.20.$$

Since p is small and n is large we can use the Poisson approximation to the binomial distribution:

$$P(k \text{ winners}) = e^{-m} \frac{m^k}{k!},$$

where $m = 1.20$ is the mean of the binomial distribution that we are approximating. The results of these calculations are shown in Table 5. From this table we find that the conditional probability that there is only one winner, given that there is at least one winner, is .517. Thus the probability that the winner has to share the prize is .483.

Recall that the cash value of the February 6 jackpot was \$177,300,000. Letting $s(k)$ denote the probabilities in Table 5, we see that $s(k)/(1 - s(0))$ is the probability that a winning ticket gets fraction $1/k$ of the jackpot. We can now find the expected amount that a winning player will end up with, by summing the values $(177,300,000/k) * s(k)/(1 - s(0))$. Carrying out this calculation we find that the expected cash value of the record jackpot is \$128,600,000. Using this number for x in Equation 6 for the expected value of a ticket gives \$1.077.

This still sounds like a favorable bet. Unfortunately, the government isn't about to let a lucky winner just walk off with a lump sum without paying taxes. Since 2003 the top marginal federal tax bracket rate has been 35%, so this percentage must be withheld in federal tax. In fact, the situation is even worse than this, since some states take out additional money to cover state tax! Here in New Hampshire (at this writing at least!), there is no state income tax. Thus, if we win the jackpot, we are left with 65% of our expected winning of \$128,585,771, obtaining \$83,580,751. Thus the expected value of a ticket, after taxes, is $(.65)(\$1.078) = \$.700$.

What happens if the player chooses the Power Play option (described in Section 1)? The multiplier is equally likely to be 2, 3, 4, or 5, so the expected multiplier is 3.5. In this case, the second component of the expected value formula gets multiplied by 3.5, so the expression for the expected value of the ticket changes to

$$E = \frac{x}{b} + 0.690,$$

where x is still \$128,585,771, since the Power Play option does not multiply the jackpot. The other award values have been changed by multiplying by the expected value of the multiplier.

The expected value of the ticket is now \$1.57, which leaves \$1.02 after taxes. Alas, we recall that we paid an extra \$1 for this option, so this is a \$2 bet. Thus the expected return per dollar is only \$0.51.

Perhaps we have now explained the famous quote:

“The lottery: A tax on people who flunked math.”

– Monique Lloyd

Exercise.

1. Suppose that someone buys two lottery tickets and fills them out with the same six numbers. How does this affect his expected value?

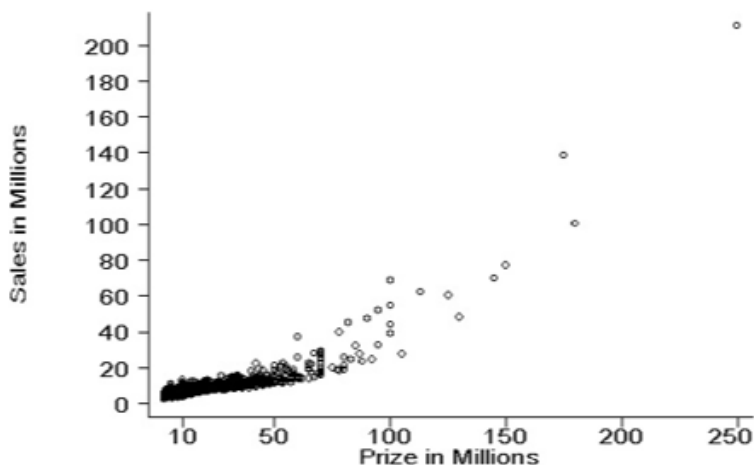


Figure 2. Powerball sales vs. jackpot size

4. Does a Ticket's Expected Value Ever Exceed \$1?

We saw above that even in the record jackpot drawing of February 18, 2006, the expected value of a ticket was significantly less than \$1. We now consider whether, for certain jackpot sizes, the expected value of a ticket ever exceeds \$1. To do this, we need to be able to estimate the number of tickets sold as a function of the jackpot size. This has been done for the Powerball lottery by Emily Oster in her senior thesis at Harvard.³ Figure 2 shows the data for the Powerball lottery from 1992 to 2000.

Oster used a log-linear fit to arrive at the following equation. In this equation, s denotes the number of tickets sold, in millions, and p denotes the announced size of the jackpot, in millions, for a given drawing:

$$\log(s) = 15.558 + .016p,$$

³Emily Oster, "Dreaming Big: Why Do People Play the Powerball?", Harvard University, March 2002.

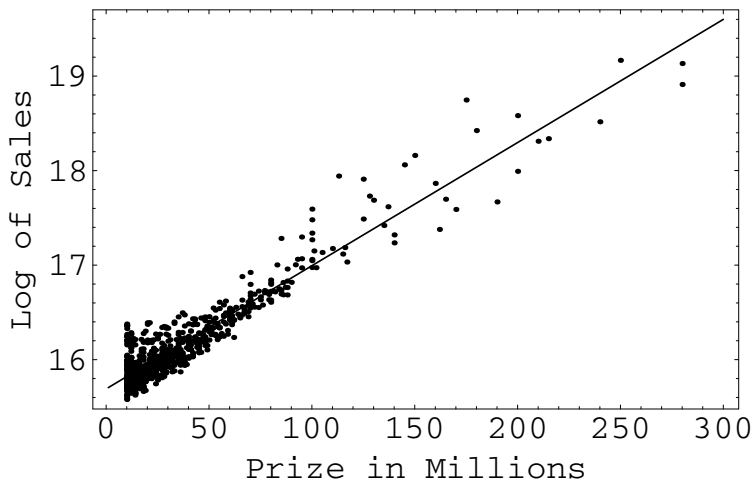


Figure 3. Log-linear fit of Powerball sales vs. jackpot size

or equivalently,

$$s = (5,712,000)(1.016)^p.$$

Thus, for example, if we let p equal 100, then the predicted value of s is approximately $s = 28,829,000$. Figure 3 shows this log-linear fit.

It is interesting to compare Figure 2 with the corresponding data from the United Kingdom lottery, shown in Figure 4. In the latter figure, the association between size of the jackpot and the number of tickets sold seems much weaker. In part, however, this is attributable to the high ticket sales that accompany unusually large Powerball jackpots. Indeed, the idea of the Powerball format was to create such exciting opportunities.

There is no reason to think that the above equation will be accurate for announced jackpot sizes larger than about \$300,000,000, since there are no data for jackpots larger than this size. In fact, if the announced jackpot size were \$500,000,000, say, then the above formula predicts that the number of tickets sold would be more than 17 billion, or more than 54 tickets for each person in the United States. This is

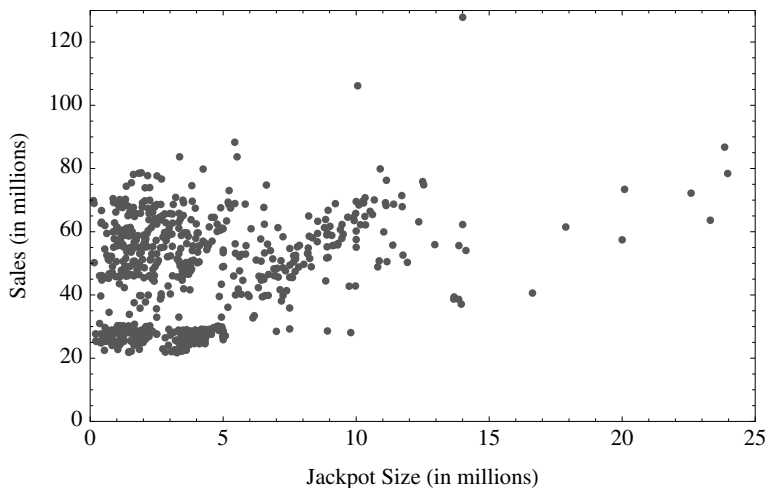


Figure 4. United Kingdom sales vs. jackpot size

clearly unreasonable. Thus, in what follows, we will initially assume that the announced jackpot size is no larger than \$300,000,000.

We can now proceed as we did in the previous section. Given an announced jackpot size j , we can estimate, using Oster's equation above, the number of tickets that will be sold. Then, using the Poisson distribution, we can calculate the expected pre-tax value of a jackpot-winning ticket. Here we assume that the cash value of the jackpot is one-half of the announced jackpot size. This is actually slightly generous. Next, we add 19.7 cents to this value; this represents the contribution to the ticket's expected value by the prizes other than the jackpot. Finally, we multiply the result by 0.65, obtaining the after-tax expected value of the ticket. Figure 5 shows the after-tax expected value of a ticket, as a function of the announced jackpot size. Note that the expected after-tax value never exceeds 48 cents.

As we said above, we cannot extrapolate our estimate of the number of tickets sold past \$300,000,000 with any assurance of accuracy. Suppose that we assume the number of tickets sold has a certain maximum value t , no matter how large the announced jackpot size is.

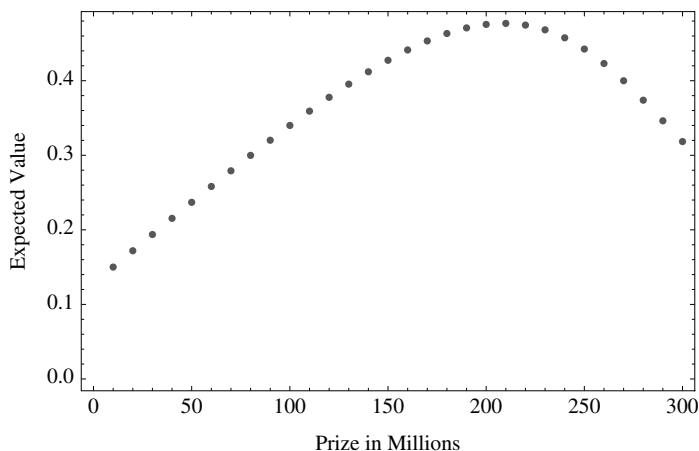


Figure 5. Expected after-tax value of a ticket vs. announced jackpot size

Given t , how large must the announced jackpot size be so that the expected after-tax value of a ticket will exceed \$1?

This question is easy to answer, using the same methods as we used above. Note that for fixed t , the expected number of jackpot winners is constant, so the expected after-tax value increases linearly in the announced jackpot size. If $t = 200,000,000$, then the expected after-tax value of a ticket equals \$1 if the announced jackpot size is \$567,000,000; if $t = 500,000,000$, the corresponding announced jackpot size is \$1,019,000,000. These examples have jackpot sizes that far exceed any Powerball jackpot that has ever occurred. Thus it is safe to say that under almost no conceivable circumstance would the after-tax value of a ticket exceed \$1.

5. What Kind of Numbers Do Lottery Buyers Choose?

We have suggested that we might at least be able to avoid sharing the jackpot with people who choose their own numbers if we choose

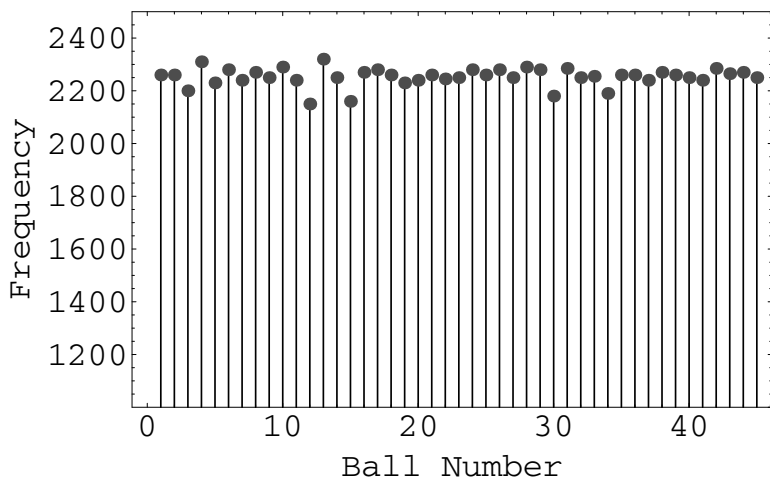


Figure 6. Frequencies of numbers chosen by computer

our own cleverly. It is well known that people who choose their own numbers do not choose randomly. They use their child's birthday, their "lucky" numbers, arithmetic progressions such as 2-4-6-8-10-12, numbers from sports, etc.

To see what players actually do, we obtained the numbers chosen by players in the Powerball Lottery in one state on May 3, 1996. Recall that at this time the game was played by selecting five of 45 white balls and one of 45 red balls. On this day, 17,001 of the picks were chosen by the buyers, and 56,496 were chosen by the computer (Easy Pick). Thus only about 23% of the picks were chosen by the buyers.

We first compare the distribution of the individual numbers from the picks chosen by the computer and those chosen by the buyers. To make the two sets the same size, we use only the first 17,001 picks produced by the Easy Pick method. Each pick contributed 6 numbers, so in both cases we have 102,006 numbers between 1 and 45. Figure 6 is a plot of the number of times each of the numbers from 1 to 45 occurred for the picks chosen by the computer. There does not

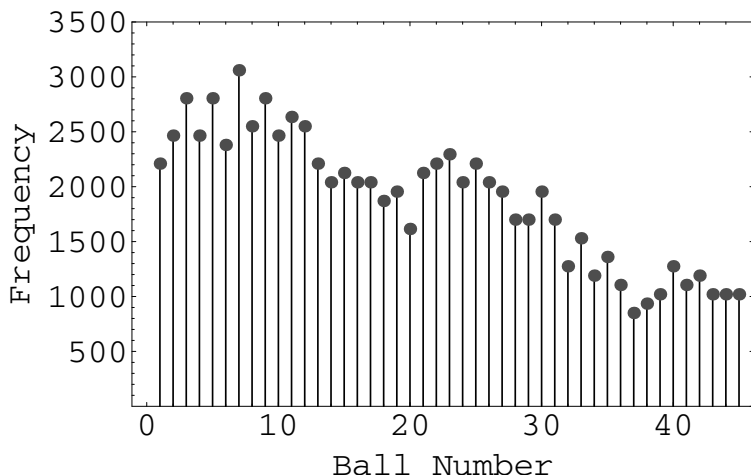


Figure 7. Frequencies of numbers chosen by players

seem to be very much variation, but it is worth checking how much variation we would expect if the numbers were, in fact, randomly chosen. If they were randomly chosen, the numbers of occurrences of a particular number, say 17, would have a binomial distribution with $n = 102,006$ and $p = 1/45$. Such a distribution has mean np and standard deviation \sqrt{npq} , where $q = 1 - p$. This gives values for the mean and standard deviation of 2267 and 47.

It is hard to tell the actual differences from the graph, so we looked at the actual data. We found that, for all but two numbers, the results were within two standard deviations of the mean. For the other 2 numbers the results were within 3 standard deviations of the mean. Thus the picks chosen by the computer do not appear to be inconsistent with the random model. A chi-square test would show how to proceed with a formal test of this hypothesis.

We look next at the picks chosen by the players. Recall that we have the same number 17,001 of picks, so we again have 85,005 individual numbers. The observed frequencies are plotted in Figure 7. You don't have to do any fancy tests to see that these are not ran-

Number Probability Number Probability Number Probability

37	0.010	29	0.10	1	0.16
38	0.011	28	0.10	22	0.16
43	0.012	31	0.10	13	0.16
45	0.012	18	0.12	23	0.17
39	0.012	30	0.13	6	0.18
44	0.012	19	0.13	2	0.19
41	0.013	27	0.13	10	0.19
36	0.013	24	0.14	4	0.19
42	0.014	14	0.14	8	0.030
34	0.014	26	0.14	12	0.030
40	0.015	16	0.14	11	0.030
32	0.015	17	0.14	3	0.033
35	0.016	21	0.15	5	0.033
33	0.018	15	0.15	9	0.033
20	0.019	25	0.16	7	0.036

Table 6. Observed probabilities for numbers chosen by players

domly chosen numbers. The most popular number 7 was chosen 3,176 times, which would be 19 standard deviations above the mean if the numbers were randomly chosen!

It is often observed that people use birthdays to choose their numbers. If they did, we would expect numbers from 1 to 12 to be most likely to be picked since such numbers can occur both in the month and the day. The next most likely numbers to be picked would be those from 13 to 31 where the remaining days of the months could occur. The least likely numbers would then be those from 32 to 45 where the year of the birthday could occur but only for those at least 61 years old (at this writing). Note that this is indeed what appears to be happening.

Finally, we look at the winning numbers to see if they could be considered to be randomly chosen. Recall that the lottery officials try to ensure that they are. Here we have many fewer numbers so we

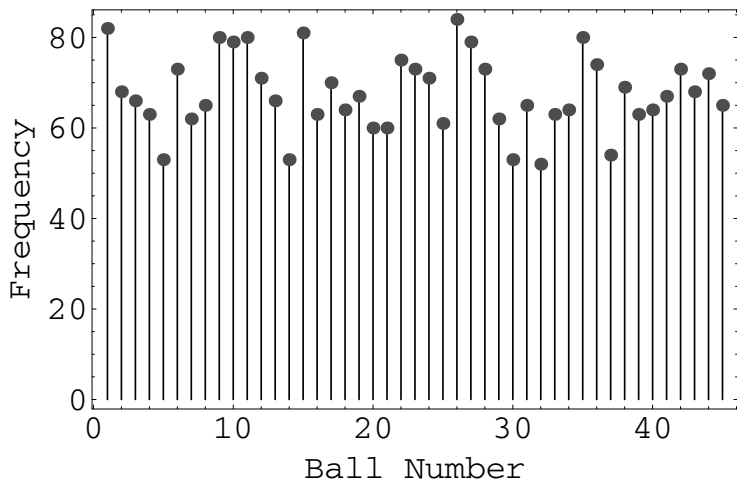


Figure 8. Winning number frequencies

expect more variation even if they are randomly chosen. Since there were 508 drawings in the period we are considering, we have $6 \cdot 508 = 3048$ numbers. Now, if the numbers are randomly chosen, the number of times a particular numbers occurs has a binomial distribution with $n = 3048$ and $p = 1/45$. Such a distribution has a mean of 67.7 and standard deviation 8.14. The biggest deviations from the mean are about 2 standard deviations so this looks consistent with the hypothesis that the numbers were randomly chosen. Again, we could check this with a chi-square test.

Exercises.

1. We have seen that the numbers picked by the players fall into three sets, with the property that numbers in the same set are approximately equally likely to be chosen, but numbers in different sets are not equally likely to be chosen. Let us denote the three sets by

$$S_1 = \{1, 2, \dots, 12\},$$

$$S_2 = \{13, 14, \dots, 31\},$$

$$S_3 = \{32, 33, \dots, 45\}.$$

The collection of all sets of five unequal numbers, between 1 and 45, written in increasing order, serves as the sample space of all possible choices by the players. The numbers in S_1 , S_2 , and S_3 occur with average frequencies .0306, .134, and .0134, respectively. Using these frequencies, show that if a chosen set of numbers has all of its numbers in S_1 , then it occurs with probability 2.7×10^{-8} , while if all of its numbers are in S_3 , then it occurs with probability 4.3×10^{-10} . How does this affect the expected value of these two tickets, i.e. is there a difference between the expected values of tickets of the above two types?

2. Using the above data from the December 25, 2002 lottery, we saw that the expected value of one lottery ticket, after taxes, is \$.607. Suppose someone buys two tickets and puts the same numbers on both tickets. What is the expected value of each ticket?

6. Finding Patterns

Recall that players choose their first five numbers to match the white balls separately from their choice of the Powerball number to match the red ball. Thus, if we are looking for patterns in the way people choose their numbers, it is best to consider the first five numbers by themselves. We recall that our data set consists of two sets of picks of size 17,001; the first set contains picks chosen by the Easy Pick method, and the second set contains picks chosen by the players. For the Easy Picks, we found that 136 of these were represented twice and 2 were represented 3 times.

To see if we should have expected to find 3 picks the same, we use the solution of the birthday problem, which is a well-known problem in probability. The most basic form of this problem asks for the probability that at least two people among a set of k people will

share a birthday. Part of the reason that this problem is interesting stems from the fact that if one asks how many people are needed in a room in order to have a favorable bet that at least two of these people have the same birthday, then the surprising answer is that only 23 people are needed.

In our case, we are asking a more difficult question: Given 17,001 choices from a set of size $C(45, 5) = 1,221,759$ (this is the number of possible choices of five numbers from 45 numbers), what is the probability that at least three of the choices are equal? So, instead of 366 possible birthdays, there are now 1,221,759 possible birthdays. It can be calculated that in this case, the probability of finding three or more choices the same is about .42. A formula for this calculation can be found at Wolfram MathWorld⁴. One can also calculate that there is only a probability of .002 of finding 4 or more the same birthday. Thus we should not be surprised at finding 3 picks the same and should not expect to find 4 the same. Again, the computer picks seem to conform to random choices.

We look next at the 17,001 picks of 5 numbers chosen by the lottery players. We found 966 sets of numbers that were represented more than once (compared to 138 for the Easy Pick numbers). The largest number of times a particular set of numbers was chosen was 24. This occurred for the pick 02-14-18-21-39. Looking at the order in which the picks were given to us, we noticed that these occurred consecutively in blocks of 5, with the blocks themselves close together. The ticket on which you mark your numbers allows room for 5 sets of numbers. We concluded that one player had made 24 picks all with the same five numbers for the white balls. He at least chose different Powerball numbers. The same explanation applied to the next most popular pick 08-12-24-25-27, which occurred 16 times.

The third most popular set 03-13-23-33-43 was picked by 13 people and was more typical of the patterns that people chose. In this

⁴MathWorld, "Birthday Problem", <http://mathworld.wolfram.com/BirthdayProblem.html>

version of Powerball, the numbers were arranged on the ticket as shown below:

Pick 5									EP---
01	02	03	04	05	06	07	08	09	
10	11	12	13	14	15	16	17	18	
19	20	21	22	23	24	25	26	27	
28	29	30	31	32	33	34	35	36	
37	38	39	40	41	42	43	44	45	

Note that the pick 03-13-23-33-43 is an arithmetic progression obtained by going down a diagonal starting with 03. Similarly, the set of numbers 01-11-21-31-41, which was chosen 10 times, corresponds to going down a diagonal starting with 01 and the set 06-15-24-33-42, chosen 9 times, corresponds to going down a column starting with 06. The most interesting pattern noticed was 01-09-23-37-45, occurring 8 times, which results from choosing the corner points and the middle point. Since we do not expect repetitions of 4 or more to occur by chance, we looked at all those that occurred 4 or more times. We could explain all but three such sets of 5 numbers. These were:

01-03-09-30-34 (occurred 5 times, always with Powerball number 40)
 05-06-16-18-23 (occurred 4 times, always with Powerball number 31)
 02-05-20-26-43 (occurred 4 times, with different Powerball numbers) .

Here are two letters that appeared in *The Times* (London) related to the problem of people choosing popular numbers. The letters followed an article in *The Times* stating that the inaugural drawing of the new British Lottery had five times the number of winners expected, including seven people who had to share the jackpot. They blamed this on the fact that the six winning numbers 03-05-14-22-30-44 had five numbers under 31 and most people chose low numbers. In this lottery, you choose 6 numbers between 1 and 49 and have to get them all correct to win the jackpot. If you get three numbers correct you win £10. The amount you win for any other prize depends on the number of other people who win this prize.

Here is the first letter.

The Times, 24 November 1994, letters to the editor.

Slim pickings in National Lottery

From Mr. George Coggan

Sir, With random choices, the odds against there being seven or more jackpot winners in the National Lottery when only 44 million tickets have been sold are 23-1. This suggests that those who predicted that low numbers would be popular were right as the slightly disproportionate number of single digits (3 and 5 came up) would combine to produce more winners than would be produced by entirely random selections.

Mildly interesting, one might think, but then one suddenly realizes that there is a lurking danger that the rules create the possibility that when (as will happen sooner or later) three single digit numbers come up the prize fund may not be enough to cover the Pounds 10 guaranteed minimum prize, never mind a jackpot. I estimate that if the number 7 had come up instead of say 44 the prize fund in this first lottery would have been about Pounds 5 million short of the guarantee. What then panic?

Yours sincerely,

GEORGE COGGAN,

14 Cavendish Crescent North,

The Park, Nottingham.

Here is the response.

The Times, 29 November 1994, letters to the editor.

No need to fear a lottery shortfall

From the Director General of the National Lottery

Sir, Mr. George Coggan (letter, November 24) raises concerns about the National Lottery Prize Fund's ability to pay winners when "popular" numbers are selected in the weekly draw.

We are aware that players do not choose numbers randomly but use birthdays, sequences or other lucky numbers. This causes the number of winners to deviate each week from the number predicted by statistical theory. Experience from other lotteries shows that the number of winners of the lower prizes can vary by up to 30 per cent from the theoretical expectation.

In the first National Lottery game there were many more Pounds 10 prize-winners than theory predicted. It is just as likely that future draws will produce fewer than expected winners and, because each higher prize pool is shared between the winners, prize values will rise accordingly.

Mr. Coggan suggests a pessimistic scenario in which the cost of paying the fixed Pounds 10 prizes to those who choose three correct numbers exceeds the prize fund. Best advice, and observations from other lotteries around the world, is that, even after allowing for the concentration on "popular" numbers, the chances of this happening are extremely remote.

Your readers will be reassured to know, however, that I have not relied totally upon statistics or evidence from other lotteries. Camelot's license to operate the National Lottery also requires them to provide substantial additional funds by way of deposit in trust and by guarantee to protect the

interests of the prize-winners in unexpected circumstances.

Yours faithfully,

PETER A. DAVIS

Director General, National Lottery

PO Box 4465

London SW1Y 5XL

Of course, it is interesting to look at this problem for the Powerball lottery. We noted that, in our sample of 17,001 sets of five numbers where players picked their own numbers, there were particular sets of five numbers for the white balls that were chosen as many as 10 times. For example the set of numbers 01-11-21-31-41 obtained by going down a diagonal starting at 1 in the box where you mark your numbers was chosen 10 times in our sample of 17,001.

For the July 29, 1998 drawing there were 210,800,000 tickets sold. If we assume that about 30% of the players pick their own numbers, then using the above example, it is possible that there exists a set of white numbers that was picked by

$$\frac{.3(210,800,000)}{1,700} = 37,200$$

players. For the prize schedule in force at that time, a player who picked all five white numbers won \$100,000. If the lottery officials had the bad luck to also choose this set of five numbers this would have cost them 3.72 billion dollars! The new boxes are not as symmetric as the old ones to which our data applied. This may help them with this potential problem. It is also the case that the lottery computer is very unlikely to choose any particular set of numbers. However, the lottery is protected in still another way. In the fine print that describes the lottery, it is stated that if there are too many awards won, the officials reserve the right to divide a certain amount of money among the prize winners, just like what is done with the jackpot. This last event has never happened in the history of the Powerball lottery.

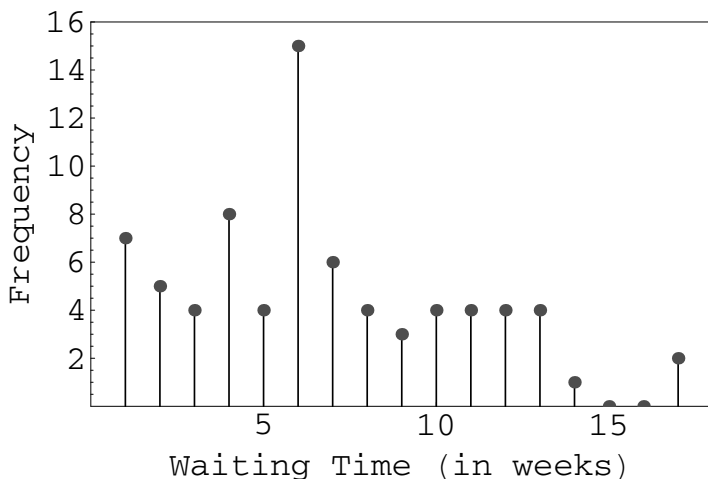


Figure 9. Waiting times between jackpots

7. How Often is the Jackpot Won?

The size of the jackpot changes from one drawing to the next. If, on a given drawing, no one chooses the winning numbers, the jackpot is increased several million dollars for the next drawing. When there is a winner, the next jackpot goes back to the minimum amount, which currently is 15 million dollars. The size of the increase when there is no winner depends upon the number of tickets sold for the previous drawing. We investigate the size of the jackpots through the years of the original rules.

According to the data from Emily Oster, the jackpots from the beginning of the Powerball lottery on April 22, 1992 until March 1, 1997, ranged from \$2 million to \$111,240,463. The jackpot was won in 75 of these 805 drawings. Eleven of these times there were two winners and never more than two winners. The total of all these jackpots was \$2,206,159,204 with an average of \$29,415,456. The average number of drawings between jackpots being won was 6.72 or,

Millions of tickets sold	Probability no one wins
10	.834
20	.695
30	.579
40	.483
50	.403
60	.336
70	.280

Table 7. Probability no one wins the jackpot

since there are two drawing a week, about 3 weeks. The distribution of times between jackpots is shown in Figure 9.

It is interesting to ask what kind of a probability model would be appropriate to describe the time between winning the jackpot. The probability of the jackpot being won depends upon the number of tickets sold. (Actually, it depends upon the number of different picks chosen.) If the same number of tickets were sold for each drawing then the appropriate model would be: toss a coin with probability p for heads until the first time heads turns up where $1/p$ is the average time between heads. Unfortunately, it is not at all reasonable to assume the same number of tickets will be sold. Here is what the Powerball FAQ says about this.

For a \$10 million jackpot draw we sell about \$11 million. For a \$20 million jackpot we sell about \$13 million. With a \$100 million jackpot we sell \$50 to \$70 million for the draw (depending on time of year and other factors).

Let's assume that, for a particular jackpot, n tickets are sold. Then the probability that a particular person does not win the jackpot is $(a - 1)/a$, where, for the old version of the game,

$$a = 45 \cdot C(45, 5) = 54,979,155.$$

The probability that none of the tickets sold wins the jackpot is

$$\left(\frac{a-1}{a}\right)^n.$$

Table 7 contains these probabilities for some different values of the number of tickets sold.

Exercise.

1. In Figure 9, the mode 6 seems rather over-represented. Can you think of any explanation for this, or is it just chance?

8. Other Lotteries Pose New Questions

There are many other interesting questions that can be explored about lotteries. The questions that one asks depend, to some extent, on the nature of the lottery. For example, in September 1996 the Multi-State lottery introduced a new lottery called Daily Millions, where the amount of the jackpot is always \$1 million and, if you win, you don't have to share it with another person who also has the same winning pick. Actually, if there are more than 10 such winners they share a \$10 million prize. An article appeared in the *Star Tribune* shortly after this lottery was introduced⁵. The article began as follows:

The lottery wizards said it wasn't supposed to work this way.

Nearly five months after the Daily Millions lottery began, none of the 34 million tickets sold has won the \$1 million jackpot. The probability of such a long losing streak is 1 in 38.

The drought has lasted so long, "We even have solid believers in statistics questioning the wisdom of numbers," quipped Charles Strutt, executive director of the Multi-State Lottery Association, which runs Daily Millions and Powerball.

⁵Doyle, Pat, "Daily Millions Beats Odds: No One Wins—5-Month Losing Streak Puzzles Even Statisticians", *Star Tribune*, Minneapolis, 7 Feb. 1997, p. 1B.

They're not the only ones.

At the East Grand Forks, Minn., Holiday store, "People are starting to get a little disgusted with it," said cashier Steve Nelson. The store has been among the top sellers of Powerball tickets, but sales of Daily Millions tickets have declined.

The article states that the ticket sales in the first week of the lottery were \$2.75 million, but five months later, they had declined to \$1.23 million.

The day after the above article appeared, the Daily Millions lottery had its first winner.

9. Using Lottery Stories to Discuss Coincidences

James Hanley⁶ has discussed how stories about lottery winners provide good examples to discuss the meaning of apparent coincidences. Here is his first example⁷.

Lottery officials say that there is 1 chance in 100 million that the same four-digit lottery numbers would be drawn in Massachusetts and New Hampshire on the same night. That's just what happened Tuesday.

The number 8092 came up, paying \$5,482 in Massachusetts and \$4,500 in New Hampshire. "There is a 1-in-10,000 chance of any four digit number being drawn at any given time," Massachusetts Lottery Commission official David Ellis said. "But the odds of it happening with two states at any one time are just fantastic," he said.

What is the probability that the same four-digit lottery number would be drawn in Massachusetts and New Hampshire on the same

⁶Hanley, James A. "Jumping to Coincidences", *American Statistician*, Vol 46, No. 3, pp. 197-201.

⁷"Same Number 2-State Winner", *Montreal Gazette*, September 10, 1981.

night? What is the probability that some two such lotteries have the same two numbers during a given period of time? Is this different from a reporter noticing that the number that turned up in the lottery in New Hampshire on Wednesday happened also to occur in the Massachusetts lottery on Saturday?

Here is another of Hanley's examples⁸.

Defying odds in the realm of the preposterous—1 in 17 million—a woman who won \$3.9 million in the New Jersey state lottery last October has hit the jackpot again and yesterday laid claim to an additional \$1.5 million prize...

She was the first two time million-dollar winner in the history of New Jersey's lottery, state officials said. They added that they had never before heard of a person winning two million-dollar prizes in any of the nation's 22 state lotteries.

For aficionados of miraculous odds, the numbers were mind boggling: In winning her first prize last Oct. 24, Mrs. Adams was up against odds of 1 in 3.2 million. The odds of winning last Monday, when numbers were drawn in a somewhat modified game, were 1 in 5.2 million.

And after due consultation with a professor of statistics at Rutgers University, lottery officials concluded that the odds of winning the top lottery prize twice in a lifetime were 1 in about 17.3 trillion—that is, 17,300,000,000,000.

Exercise.

1. Does it matter that she played the lottery many times, often buying more than one ticket? Again, are we talking about

⁸ "Odds-Defying Jersey Woman Hits Lottery Jackpot 2nd Time", *New York Times*, February 14, 1986.

this happening somewhere, sometime? Should we ever believe that something with these odds has happened?

10. Lottery Systems

Richard Paulson⁹ observes that claims made about systems for improving your chances at lotteries illustrate important statistical concepts. For example, people claim that it is possible to predict future winners by analyzing the historical data of winning numbers. Indeed, lottery sites encourage this by making this historical data available. Sometimes the argument is simply that, when a particular number has not turned up as often as would be expected, then this number is more likely to come up in the future. This is often called the “gambler’s fallacy” and all too many people believe it. The fact that it is not true is the basis for many beautiful applications of chance processes called martingales.

Paulson remarks that he particularly enjoys discussing the following system. Consider, the six winning numbers in the Powerball Lottery. If they occur randomly their sum should be approximately normally distributed with mean $6(1+45)/2 = 138$ and standard deviation approximately 32. Thus, sets of six numbers whose sum is more than 64 away from the mean 138 are not likely to occur as winning sets and should be avoided. It is better to pick six numbers whose sum is near 138. We leave the reader to ponder this last system.

One of our teachers, the well-known probabilist Joseph Doob, was often asked for a system to use when playing roulette. His advice ran as follows. Play red once. Then wait until there have been two blacks and play red again. Then wait until there have been three blacks and play red again. Continue in this manner. You will enjoy playing and not lose too much.

⁹Paulson, Richard A. “Using Lottery Games to Illustrate Statistical Concepts and Abuses”, *American Statistician*, Vol. 45, No. 3, pp. 202-204.

11. Lottery Stories from Chance News

As remarked earlier, lotteries often make the news, so we now include some commentaries from Chance News. We traditionally start Chance News with a quotation. Here is a sample of lottery quotes.

“To be born an Englishman is to win first prize in the lottery of life?” Cecil Rhodes

A Lottery is a Taxation, Upon all the Fools in Creation;
And Heav’n be prais’d, It is easily rais’d,
Credulity’s always in Fashion:
For, Folly’s a Fund, Will never lose Ground,
While Fools are so rife in the Nation.

Henry Fielding’s play *The Lottery*, a farce (1724)

Here is a better known and simpler version of this idea.

Lotteries are a tax on stupidity.

This is often attributed to Voltaire, but we have not been able to find a reference.

Lotteries have also been the subject of Forsooths in our Chance News. Here is one of these.

In theory, if you were to buy 50 tickets and your neighbor bought one, neither of you would have a better or worse chance of winning. We like to say it only takes one ticket to win.

Brian Rockey, Nebraska Lottery Spokesman.
Omaha World-Herald, February 18, 2005

Here are some stories from Chance News:

Math professor shares \$15 million lotto jackpot.
Denver Post, 9 April 1996, B3
Peter G. Chronis

Math professor Celestino Mendez was discussing expected value in his class at Metropolitan State College and remarked that, in a lottery, the expected winning increases when the jackpot gets higher. He told his students that they ought to buy a ticket in the current Colorado Lottery because the expected winning was positive (14 cents when you buy a \$1 ticket). Professor Mendez thought he should put his money where his mouth is, and so, on the way home, he stopped and bought ten tickets. One of these had the lucky numbers, and he shared the \$15 million prize with one other winner.

Is your ticket a winner? Odds are, it's not.

Valley News 29 July 1998, A2

Sarah M. Earle

At the time of the record 50-million-dollar Powerball jackpot, most newspapers felt that they had to help the public think about what the 1-in-80 million chance of winning the jackpot really means. We got a call for help with this from a local newspaper, the *Valley News*. We gave two suggestions for thinking about these odds. The first was suggested by Fred Hoppe: if you toss a coin 26 times your chance of getting all heads is greater than your chance of winning the Powerball jackpot.

The second we learned from Arnold Barnett in his “Chance Lecture: Risk in Everyday Life.”¹⁰ Arnold discussed how to help people understand the chance of being killed on an airplane flight. He estimated that if you go on a randomly chosen airplane flight you have a 1-in-7 million chance of being killed. He said his first attempt to explain these odds was to you had a three times greater chance of winning the Massachusetts lottery. He said most people think they will win the lottery, so this “juxtaposition of people’s greatest hope with their worst fear” did not work.

After some experimentation, he found more success saying that if they took a randomly selected airplane flight every single day, on

¹⁰<http://www.dartmouth.edu/~chance/ChanceLecture/AudioVideo.html#Videos98>.

average it would take 19,000 years to experience a fatal accident. They were happy with this since they said they weren't going to live this long. By the same logic, you would expect to have to buy a lottery ticket twice a week for 1.4 million years in order to experience a jackpot.

An article¹¹ in the *New York Times* reports unexpected winnings in the Powerball lottery of March 30, 2005 which lottery officials thought might be fraudulent but which had a much simpler explanation. At the time of this article, for the Powerball lottery, players chose five numbers from 1 to 52 and an additional powerpoint number chosen from 1 to 42. On the March 30, 2005 drawing of the Powerball lottery, 110 players made a \$1 bet, choosing as their five numbers 22,28,32,33,39 and as their Powerball number 40. The lottery chose the same five basic numbers but chose 42 for the Powerball number. It turns out that 89 of these winners did not choose the Power Play and so each won \$100,000; the other 21 players chose the Power Play and, the multiplier was 5, so these winners won \$500,000. Thus the lottery paid out \$19.4 million to these winners.

Powerball officials stated that, considering the number of tickets sold in the 29 states, they expected 4 or 5 winners. The article quotes Chuck Strutt, executive director of the Multi-State Lottery Association as saying: "Panic began at 11:30 P.M. March 30 when I got a call from a worried staff member. We didn't sleep a lot that night. Is there someone trying to cheat the system?"

The lottery authorities tried a number of theories about how people choose their numbers. For example, many players pick their numbers following a geometric design on the ticket. Nothing worked. But then the first three winners said that they had obtained the numbers from a fortune cookie.

With this lead, they just had to find the fortune cookie maker who had the winning numbers. They found that many different brands of fortune cookies come from the same Long Island City factory owned

¹¹Jennifer Lee, "Who needs Giacomo? Bet on the fortune cookie", *New York Times*, May 11, 2005, p. 1.

by Wonton Food. This company turns out four million fortune cookies a day, which are delivered to dealers over the entire country. When shown the numbers, Derrick Wong, of Wonton Food, verified that they had used them. The numbers for the fortune cookies were chosen from a bowl but the company plans to switch to having them chosen by a computer and Derrick plans to start playing the lottery.

We turn to a letter¹² to the editor of the *New York Times*, from John P. Rash. Rash starts his letter with:

Your July 13 Week in Review article on lottery advertising repeats stereotypes about lottery players being poor and uneducated and swept up into gambling addictions. No doubt many are. But Gov. George E. Pataki's statement that "It has always bothered me to hold up the prospect of instant riches" could also be recast as, "I want to take away the only prospect poor people have of getting out of their rut."

Rash goes on to say that before lotteries, there were other ways to improve your life, but now graduate education is expensive and required for many of the better jobs. He explains his experience playing the lottery and concludes with the remark:

Yes, I have lost more than I've won. But in the tedious world I inhabit along with so many other New Yorkers, I've bought a fantasy. If I ever win the Jackpot, I'll wave to you from Sutton Place.

Our next article supports Rash's argument. An article¹³ in the *Los Angeles Times* concerns the winners of the record jackpot of \$365 million on October 22, 2005. The author writes:

On Wednesday morning in Lincoln, Neb., after four days of speculation about who had won the

¹²John P. Rash, "Seen from a rut, the lottery is essential", *New York Times*, July 16, 1996.

¹³Meghan Daum, "Who's the idiot now?", *Los Angeles Times*, Feb. 25, 2005.

biggest jackpot in Powerball history, eight employees of a ConAgra ham processing plant came forward and identified themselves as the winners of the \$365 million purse. As lottery stories go, this is about as heartwarming as it gets. Two of the winners are immigrants from Vietnam and one is a political refugee from the Republic of Congo—and all worked the second and third shifts, some clocking as many as 70 hours a week. There is probably no jobsite as gruesome as a meatpacking house. If anyone deserves an express ticket to a new life, it's these folks.

12. Lottery Questions from John Haigh

The questions below were asked by John Haigh in a commentary that appeared in January 2003 in the RSS News. You can learn the rules for the UK Lottery at the National Lottery homepage.¹⁴ Haigh has also written an expository piece entitled “The UK National Lottery - a guide for beginners.”¹⁵ A collection of data on the UK Lottery, including volume of tickets sales (this is not available for Powerball) can be found at Richard Lloyd's website.¹⁶

In order to try to answer these questions, the reader should understand how the UK Lottery works. Before March 1996, there was one drawing per week, and since that time, there have been two drawings per week. Each drawing consists of picking seven numbers out of the set of integers between 1 and 49. The first six of these numbers form the set of “main” numbers for that drawing, and the seventh number is called the “bonus” number. Winning the jackpot prize requires getting all six main numbers. (The bonus number is used to determine some of the other prize amounts.) To answer questions 1, 4, 6, and 7, one needs to know that we are considering a set of

¹⁴<http://www.national-lottery.co.uk/player/p/home/home.do>

¹⁵<http://plus.maths.org/issue29/features/haigh/index.html>

¹⁶<http://lottery.merseyworld.com/>

721 drawings, most of which occurred after March 1996. To answer questions 5, 6, and 7, one needs to understand more about the distribution of tickets sold. Since March 1996 (i.e. for most of the time interval covering the data), the Saturday sales averaged roughly 40 million and Wednesday sales averaged roughly half of that.

Here are Haigh's questions (answers are given beginning on page 225):

1. How often would you expect the most frequent and the least frequent main numbers to arise? (The mean and variance of the frequency of any one integer should be about 88.3 and 77.5.)
2. How many draws should be needed until each number has arisen at least once as a bonus number?
3. How many draws should be needed until we achieve a complete collection of main numbers? How many draws until we achieve a complete collection of pairs of main numbers?
4. What is your guess for the length of the maximum run to date? A run is a consecutive set of drawings with the property that a given number appears in each set of main numbers.
5. Over the set of 721 weekly draws, guess the frequency of no jackpot winners.
6. Guess the modal number of jackpot winners (i.e. the most common number of jackpot winners in a drawing).
7. Guess the smallest number n for which there has not yet been exactly n jackpot winners.

Chapter 4

Fingerprints

1. Introduction

On January 7, 2002, in the case *U.S. v. Llera Plaza*, Louis H. Pollack, a federal judge in the United States District Court in Philadelphia, barred any expert testimony on fingerprinting that asserted that a particular print gathered at the scene of a crime is or is not the print of a particular person. As might be imagined, this decision was met with much interest, since it seemed to call into question whether fingerprinting can be used to help prove the guilt or innocence of an accused person.

In this chapter, we will consider the ways in which fingerprints have been used by society and show how the current quandary was reached. We will also consider what probability and statistics have to say about certain questions concerning fingerprints.

2. History of Fingerprinting

It seems that the first use of fingerprints in human society was to give evidence of authenticity to certain documents in seventh-century China, although it is possible that they were used even earlier than

this. Fingerprints were used in a similar way in Japan, Tibet, and India. In Simon Cole's excellent book on the history of fingerprinting, the Persian historian Rashid-eddin is quoted as having declared in 1303 that "Experience shows that no two individuals have fingers exactly alike."¹ This statement is one with which the reader is no doubt familiar. A little thought will show that unless all the fingerprints in the world are observed, it is impossible to verify this statement. Thus, one might turn to a probability model to help understand how likely it is that this statement is true. We will consider such models below.

In the Western world, fingerprints were not discussed in any written work until 1685, when an illustration of the papillary ridges of a thumb was placed in an anatomy book written by the Dutch scientist Govard Bidloo. A century later, the statement that fingerprints are unique appeared in a book by the German anatomist J. C. A. Mayer.

In 1857, a group of Indian conscripts rebelled against the British. After this rebellion had been put down, the British government decided that it needed stricter law enforcement in its colonies. William Herschel, the grandson of the discoverer of the planet Uranus, was the chief administrator of a district in Bengal. Herschel noted that the unrest in his district had given rise to a great amount of perjury and fraud. For example, it was believed that many people were impersonating deceased officers to collect their pensions. Such impersonation was hard to prove, since there was no method that could be used to decide whether a person was who he or she claimed to be.

In 1858, Herschel asked a road contractor for a handprint, to deter the contractor from trying to contest, at a later date, the authenticity of a certain contract. A few years subsequent to this, Herschel began using fingerprints. It is interesting to note that in India, as in China, the first use of fingerprints was in civil, not criminal, identification.

At about the same time, the British were increasingly concerned about crime in India. One of the main problems was to determine

¹Cole, Simon A., *Suspect Identities: A History of Fingerprinting and Criminal Identification*, Harvard University Press, Cambridge, MA, 2001, pp. 60-61.

whether a person arrested and tried for a crime was a habitual offender. Of course, to determine this required that some method be used to identify people who had been convicted of crimes. Presumably, a list would be created by the authorities, and if a person was arrested, this list would be consulted to determine whether the person in question had prior convictions. In order for such a method to be useful, it would have to possess two properties. First, there would have to be a way to store, in written form, enough information about a person so as to uniquely identify that person. Second, the list containing this information would have to be in a form that would allow quick and accurate searches.

Although, in hindsight, it might seem obvious that one should use fingerprints to help with the formation of such a list, this method was not the first to be used. Instead, a system based on anthropometry was developed. Anthropometry is the study and measurement of the size and proportions of the human body. It was naturally thought that once adulthood is reached, the lengths of bones do not change. In the 1880s Alphonse Bertillon, a French police official, developed a system in which eleven different measurements were taken and recorded. In addition to these measurements, a detailed physical description, including information on such things as eyes, ears, hair color, general demeanor, and many other attributes, was recorded. Finally, descriptions of any “peculiar marks” were recorded. This system was called Bertillonage and was widely used in Europe, India, and the United States, as well as other locations, for several decades.

One of the main problems encountered in the use of Bertillonage was inconsistency in measurement. The “operators,” as the measurers were called, were well trained, and many measurements of each person were taken. Nonetheless, if a criminal suspect was measured in custody, and the suspect’s measurements were already in the list, the two sets of measurements might vary enough so that no match would be made.

Another problem was the amount of time required to search the list of known offenders, in order to determine whether a person in

custody had been arrested before. In some places in India, the lists grew to contain many thousands of records. Although these records were certainly stored in a logical way, the variations in measurements made it necessary to look at many records that were “near” the place that the searched-for record should be.

The chief problem at that time with the use of fingerprints for identification was that no good classification system had been developed. In this regard, fingerprints were not thought to be as useful as Bertillonage, since the latter method did involve numerical records that could be sorted. In the 1880s, Henry Faulds, a British physician who was serving in a Tokyo hospital at the time, devised a method for classifying fingerprints. This method consisted of identifying each major type of print (like those shown in Figure 1) with a certain written syllable, followed by other syllables representing different features in the print. Once a set of syllables for a given print was determined, the set was added to a alphabetical list of stored sets of syllables representing other prints.

Faulds wrote to Charles Darwin about his ideas, and Darwin forwarded them to his cousin, Francis Galton. Galton was one of the giants among British scientists in the late 19th century. His interests included meteorology, statistics, psychology, genetics, and geography. Early in his adulthood, he spent two years exploring southwest Africa. He is known as the first modern-day proponent of eugenics; in fact, this word is due to Galton.

Galton became interested in fingerprints for several reasons. He was interested in the heritability of certain traits, and one such trait that could easily be tested were fingerprint patterns. He was concerned with ethnology, and sought to compare the various races. One question that he considered in this vein was whether the proportions of the various types of fingerprints differed among the races. He also tried to determine whether any other traits were related to fingerprints. Finally, he understood the value that such a system would have in helping the police and the courts identify recidivists.

To carry out such research, it was necessary for Galton to have access to many fingerprints. By the early 1890s, he had amassed a collection of thousands of prints. This collection contained prints from people belonging to many different ethnic groups. He also collected fingerprints from certain types of people, such as criminals. He was able to show that fingerprints are partially controlled by heredity. For example, it was found that a peculiarity in a pattern in a fingerprint of a parent might pass to the same finger of a child, or, with less probability, to another finger of that child. Nonetheless, it must be stated that his work in this area did not lead to any discoveries of great import.

One of Galton's most fundamental contributions to the study of fingerprints consisted of his publishing of material, much of which was due to William Herschel, that fully established the fact that fingerprint patterns persist over the lifetime of an individual. Of at least equal importance was his development of a method to classify fingerprints. An important attribute of his method was that it allowed fingerprint records to be quickly searched to determine if a given fingerprint were present.

Very shortly thereafter, a committee consisting of various high officials in British law enforcement was formed to compare Bertillonage and the Galton fingerprint method. The goal was to decide which method to adopt. Although Bertillonage was in use in continental Europe, India, and elsewhere, it had not yet been used in Britain. The committee also considered whether it might be still better to use both methods at once.

In their deliberations, the committee noted that the Galton fingerprint method is a much easier process than the one that is used by Bertillonage operators. In addition, a fingerprint, if it is properly taken (i.e. if the resulting impression is legible), is a true and accurate rendition of the patterns on the finger. Both of these statements lead to the conclusion that Galton's method is more accurate than Bertillonage.

Given these remarks, it might seem strange that the committee did not recommend that fingerprints be the method of choice. However, there was still some concern about the accuracy of the classification method used by Galton. It was recommended that identification be made by fingerprints, with indexing by Bertillonage. The committee did foresee that the problems with fingerprint indexing could be overcome, and that in this case, the fingerprint method might be the sole system in use.

Galton continued to work on his method of classification, and in 1895, he published a system that greatly improved his previous attempts. Edward Henry, a magistrate of a district in India, worked on and modified Galton's classification method between 1898 and 1900. This modification was adopted by Scotland Yard. Regarding credit for the method, a letter from Sir George Darwin to the *London Times* had this to say: "Sir Edward Henry undoubtedly deserves great credit in recognising the merits of the system and in organising its use in a practical manner in India, the Cape and England, but it would seem that the yet greater credit is due to Mr. Francis Galton."²

In 1902, Galton published a letter entitled "Finger-Print Evidence" in the journal *Nature*, in which he discusses a pair of enlarged photographs, sent to him by Scotland Yard, of fingerprints. The first photograph came from the scene of a burglary, and the second came from the fingerprint files at Scotland Yard. Galton discusses how the use of his system allows the prosecution to explain the similarities in the two prints. The question of accuracy in matching prints obtained from a crime scene with those in a database is one that is still being considered today. Before turning to this question, we will describe Galton's method.

Galton begins by noting that in the center of most fingerprints there is a "core," which consists of patterns that he calls loops and whorls. (See Figure 1.) If no such core exists, the pattern is said to be an arch. Next, he defines a delta as the region where the parallel

²George Darwin, quoted in Karl Pearson, "Life, Letters, and Labors of Francis Galton," Cambridge University Press.

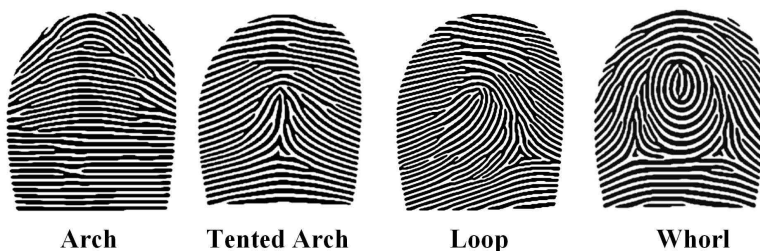


Figure 1. Four types of fingerprints. (A modified version of this figure appeared in E. Keogh, “An Overview of the Science of Fingerprints”, *Anil Aggrawal’s Internet Journal of Forensic Medicine and Toxicology*, 2001, vol. 2, no. 1 (January-June 2001).)

ridges begin to diverge to form the core. Loops have one delta, and whorls have two. These deltas serve as axes of reference for the rest of the classification. By tracing the ridges as they leave the delta(s) and cross the core, and keeping track of certain aspects such as the direction in which the loops open up, one can partition fingerprints into ten classes. We will not describe these ten classes in detail here, as the specifics are not important in what follows. Since each finger would be in one of the ten classes, there are 10^{10} possible sets of ten classes. Even though the ten classes do not occur with equal frequency among all recorded fingerprints, this first level of classification already serves to distinguish between most pairs of people.

Of the ten classes, only two correspond to loops, as opposed to arches and whorls. However, about half of all observed fingerprints are loops, which suggests that the scheme is not yet precise enough. Galton was aware of this and added two other types of information to the process. The first involved using the axes of reference arising from the deltas to count ridges in certain directions. The second involved the counting and classification of what he termed “minutiae.” This term refers to places in the print where a ridge bifurcates or ends. The idea of minutiae is still in use today, although the minutiae are now sometimes referred to as “Galton points” or “points.”

There are many different types of points, and the places that they occur in a given fingerprint seems to be somewhat random. In addition, a typical fingerprint has many such points. These observations imply that if one can accurately write down where the points occur and which types of points occur, then one has a very powerful way to distinguish two fingerprints. The method is even more powerful when it is used to compare sets of ten fingerprints from two people.

3. Models of Fingerprints

We shall investigate some probabilistic models for fingerprints that incorporate the idea of points. The two most basic questions that one might use such models to help answer are as follows. First, in a given model, what is the probability that no two fingerprints, among all people who are now alive, are exactly alike? Second, suppose that we have a partial fingerprint, such as one that might have been recovered from a crime scene. Such partial prints are called latent prints. What is the probability that this latent print exactly matches more than one fingerprint, among all fingerprints in the world? The reason that we are interested in whether the latent print matches more than one fingerprint is that it clearly matches one print, namely the one belonging to the person who left the latent print. It is typically the case that the latent print, if it is to be of any use, will identify a suspect, i.e. someone who has a fingerprint that matches the latent print. It is obviously of great interest in a court of law as to how likely it is that someone other than the suspect has a fingerprint that matches the latent print. We will see that this second question is of central importance in the discussions going on today about the accuracy of fingerprinting as a crimefighting tool.

Galton seems to have been the first person to consider a probabilistic model that might shed some light on the answer to the first question. He began by imagining a fingerprint as a random set of ridges, with roughly 24 ridge intervals across the finger and 36 ridge intervals along the finger. Next, he imagined covering up an n by n ridge interval square on a fingerprint and attempting to recreate the

ridge pattern in the area that was covered. Galton maintained that if n were small, say at most 4, then most of the time, the pattern could be recreated by using the information in the rest of the fingerprint. However, if n were 6, he found that he was wrong more often than right when he carried out this experiment.

He then let $n = 5$ and claimed that he would be right about one-half of the time in reconstructing the fingerprint. This led him to consider the fingerprint as consisting of a set of non-overlapping n by n squares, which he considered to be independent random variables. In Pearson's account, Galton used $n = 6$, although his argument is more understandable had he used $n = 5$. Galton claimed that any of the reconstructions, both the correct and incorrect ones, might have occurred in nature, so each random variable has two possible values, given the way that the ridges leave and enter the square, and given how many ridges leave and enter. Pearson says that Galton "proceeds to give a rough approximation to two other chances, which he considers to be involved: the first concerns guessing correctly the general course of the ridges adjacent to each square and the second of guessing rightly the number of ridges that enter and issue from the square."³ Finally, Galton multiplies all of these probabilities together, under the assumption of independence, and arrives at the number 1 out of 64 billion. At the time, there were about 16 billion fingerprints in the world. (Galton claims that the odds are roughly 39 to 1 against any particular fingerprint occurring anywhere in the world. It seems to us that the odds should be 3 to 1 against.)

We will soon see other models of fingerprints that arrive at much different answers. However, it should be remembered that we are trying to estimate the probability that no two fingerprints, among all people who are now alive, are exactly alike. Suppose, as Galton did, that there are 16 billion fingerprints among the people of the world, and there are 64 billion possible fingerprints. Does the reader think that these assumptions make it very likely or very unlikely that there are two fingerprints that are the same? To answer this

³Pearson, *ibid.*, p. 182.

question, we can proceed as follows. Consider an urn with 64 billion labeled balls in it. We choose, one at a time, 16 billion balls from the urn, replacing the balls after each choice. We are asking for the probability that we never choose the same ball more than once. This is the celebrated birthday problem, in a world where there are 64 billion days in a year and 16 billion people. The birthday problem asks what is the probability that at least two people share a birthday. The complementary probability, i.e. the probability that no two people share a birthday, is

$$\left(1 - \frac{0}{n}\right) \left(1 - \frac{1}{n}\right) \left(1 - \frac{2}{n}\right) \cdots \left(1 - \frac{k-1}{n}\right),$$

where $n = 64$ billion and $k = 16$ billion. This can be seen by considering the people one at a time. If 6 people, say, have already been considered and if they all have different birthdays, then the probability that the seventh person has a birthday that is different than all of the first 6 people equals

$$\left(1 - \frac{6}{n}\right).$$

(One way to obtain a rough upper bound on this product is to use the inequality

$$1 - x < e^{-x},$$

which is valid for $x > 0$.) For the values given by Galton, the product is less than

$$\frac{1}{10^{10^9}}.$$

This means that in Galton's model, with his estimates, it is extremely likely that there are two fingerprints that are the same.

In fact, to our knowledge, no two fingerprints from different people have ever been found that are identical. Of course, it is not the case that all fingerprints on Earth have been recorded or compared, but the FBI has a database with more than 10 million fingerprints in it, and we presume that no two fingerprints in it are exactly the same. (It must be said that it is not clear to us that all pairs of fingerprints

in this database have actually been compared. In addition, one wonders whether the FBI, if it found a pair of identical fingerprints, would announce this to the world.) In any case, if we use Galton's estimate for the number of possible fingerprints and let $k = 10$ million, the probability that no two are alike is still very small; it is less than

$$\frac{1}{10^{339}}.$$

We can turn the above question around and ask the following question. Suppose that there are 60 billion fingerprints in the world, and suppose that we imagine they are chosen from a set of n possible fingerprints. How large would n have to be in order that the probability that all of the chosen fingerprints are different exceeds .999? An approximate answer to this question is that it would suffice for n to be at least 10^{25} . Although this is quite a bit larger than Galton's estimate, there have been other, more sophisticated models of fingerprints, some of which we will now describe, have come up with estimates for n that are much larger than 10^{25} . Thus, if these models are at all accurate, it is extremely unlikely that there exist two fingerprints in the world that are exactly alike.

In 1933, T. Roxburgh described a model for fingerprint classification that is much more intricate than Galton's model. This model, and many others, are described and compared in an article in the *Journal of Forensic Sciences*, written by D. A. Stoney and J. I. Thornton.⁴ In Roxburgh's model, a vertical ray is drawn upwards from the center of the fingerprint. This idea must be accurately defined, but for our purposes, we can take it to mean the center of the loop or whorl or the top of the arch. This ray is defined to be 0 degrees. Another ray, with endpoint at the center, is revolved clockwise from the first ray. As this ray passes over minutiae, the types of the minutiae are recorded, along with the ridge numbers on which the minutiae lie. If a fingerprint has R concentric ridges, n minutiae, and there are T

⁴Stoney, D. A. and J. I. Thornton, "A Critical Analysis of Quantitative Fingerprint Individuality Models", *Journal of Forensic Sciences*, vol. 31, no. 4 (1986), pp. 1187-1216.

minutia types, then the number of possible patterns equals

$$(RT)^n ,$$

since as the second ray revolves clockwise, the next minutia encountered could be on any of the R ridges and be of any of the T minutia types. Roxburgh also introduces a factor of P that corresponds to the number of different overall patterns and core types that might be encountered. Thus, he estimates the number of possible fingerprints to be

$$P(RT)^n .$$

He takes $P = 1000$, $R = 10$, $T = 4$, and $n = 35$; this last value is Galton's estimate for the typical number of minutia in a fingerprint. If we calculate the above expression with these values, we obtain the number

$$1.18 \times 10^{59} .$$

Roxburgh modified the above expression for the number of possible fingerprints to attempt to account for ambiguities between various types of minutiae. For example, it is possible that a fork in a ridge might be seen as a ridge ending, depending upon whether the ridges in question meet each other or not. Roxburgh suggested using a number Q which would vary depending upon the quality of the fingerprint under examination. The value of Q ranges from 1.5 to 3, with the smaller value corresponding to a higher quality fingerprint. For each minutia, Roxburgh replaced the factor RT by the factor RT/Q . This leads to the expression

$$P((RT)/Q)^n$$

as an estimate for the number of discernable types of fingerprints, assuming their quality corresponds to a particular value of Q . Note that even if $Q = 3$, so that $RT/Q = 1.33R$, the number of discernable types of fingerprints in this model is

$$2.16 \times 10^{42} .$$

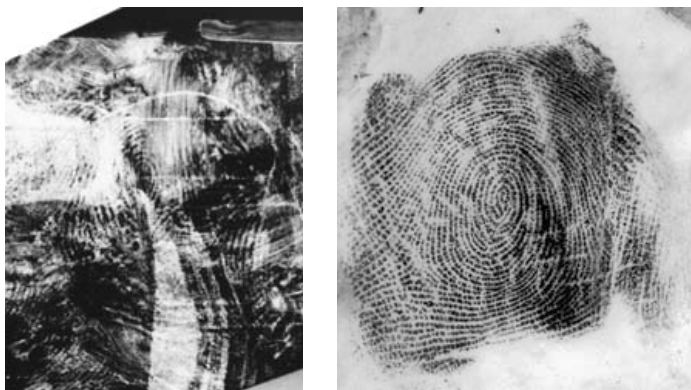


Figure 2. Examples of latent and rolled prints

Stoney and Thornton note that although this is a very interesting, sophisticated model, it has been “totally ignored by the forensic science community.”⁵

4. Latent Fingerprints

According to a government expert who testified at a recent trial, the average size of a latent fingerprint fragment is about one-fifth the size of a full fingerprint. Since a typical fingerprint contains between 75 and 175 minutiae⁶, this means that a typical latent print has between 15 and 35 minutiae, assuming that minutiae are roughly evenly distributed across the print. In addition, the latent print recovered from a crime scene is frequently of poor quality, which tends to increase the likelihood of mistaking the types of minutiae being observed.

In a criminal case, the latent print is compared with a high quality print taken from the hand of the accused or from a database of fingerprints. Figure 2 shows a latent print and the corresponding rolled print to which the latent print was matched. Figure 3 shows another

⁵ibid., p. 1192.

⁶“An Analysis of Standards in Fingerprint Identification 1,” *Federal Bureau of Investigation Department of Justice Law Enforcement Bulletin*, vol. 1 (June 1972).

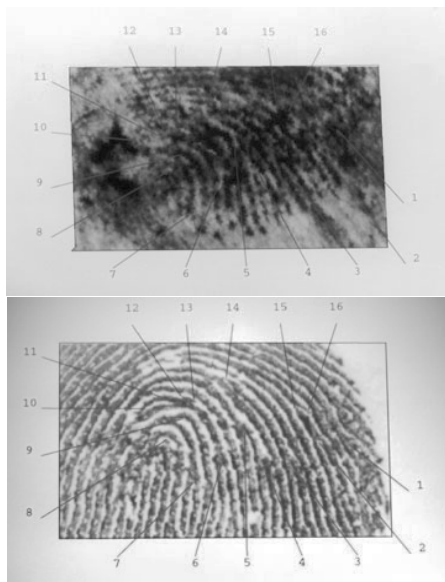


Figure 3. Minutiae matches

pair of prints, one latent and one rolled, from the same case. The figure also shows the claimed matching minutiae in the two prints.

The person making the comparison states that there is a match if he or she believes that there are a sufficient number of common minutiae, both in type and location, in the two prints. There have been many criminal cases in which an identification was made with fewer than fifteen matching minutiae⁷. There is no general agreement among various law enforcement agencies or among various countries on the number of matching minutiae that must exist in order for a match to be declared. In fact, according to Robert Epstein⁸, “many examiners ... including those at the FBI, currently believe that there

⁷see footnote 25 in Epstein, Robert, “Fingerprints Meet Daubert: The Myth of Fingerprint ‘Science’ is Revealed,” *Southern California Law Review*, vol. 75 (2002), pp. 605-658.

⁸*ibid.*, p. 610.

should be no minimum standard whatsoever and that the determination of whether there is a sufficient basis for an identification should be left to the subjective judgment of the individual examiner.” It is quite understandable that a law enforcement agency might object to constraints on its ability to claim matches between fingerprints, as this could only serve to decrease the number of matches obtained.

In some countries, fingerprint matches can be declared with as few as eight minutiae matches. However, there are examples of fingerprints from different people that have seven matching minutiae. In a California bank robbery trial, *U.S. v. Parks*, in 1991, the prosecution introduced evidence that showed that the suspect’s fingerprint and the latent print had ten points. The trial judge, Spencer Letts, asked the prosecution expert what the minimum standard was for points in order to declare a match. The expert announced that the minimum was eight. Judge Letts had seen fingerprint evidence entered in other trials. He said “If you only have ten points, you’re comfortable with eight; if you have twelve, you’re comfortable with ten; if you have fifty, you’re comfortable with twenty.”⁹ Later in the same trial, the following exchange occurred between Judge Letts and another prosecution fingerprint expert:

The Witness: The thing you have there is that each department has their own goals or their own rules as far as the number of points being a make [an identification]. ...that number really just varies from department to department.

The Court: I don’t think I’m ever going to use fingerprint testimony again; that simply won’t do...

The Witness: That just may be one of the problems of the field, but I think if there was [a] survey taken, you would probably get a different number from every department that has a fingerprint

⁹Cole, op. cit., p. 272.

section as to their lowest number of points for a comparison and make.

The Court: That's the most incredible thing I've ever heard of.¹⁰

According to Simon Cole, no scientific study has been carried out to estimate the probability of two different prints sharing a given number of minutiae. David Stoney and John Thornton claim that none of the fingerprint models proposed during the past century "even approaches theoretical accuracy ..., and none has been subjected to empirical validations."¹¹ In fact, latent print examiners are prohibited by their primary professional association, the International Association for Identification (IAI), from offering opinions of identification using probabilistic terminology. A resolution, passed by the IAI at one of its meetings, states that "any member, officer, or certified latent print examiner who provides oral or written reports, or gives testimony of possible, probable, or likely friction ridge identification shall be deemed to be engaged in [unbecoming] conduct... and charges may be brought."¹²

In 1993, the Supreme Court rendered a decision in the case *Daubert v. Merrell Dow Pharmaceuticals, Inc.*¹³ The Court described certain factors that courts needed to consider when deciding whether to admit expert testimony. In this decision, the Court concentrated on scientific expert testimony; it considered the issue of expert testimony of a non-scientific nature in the case *Kumho Tire Co. v. Carmichael*¹⁴, a few years later. In the first decision, the Court interpreted the Federal Rule of Evidence 702, which defines the term "expert witness" and states when such witnesses are allowed, as requiring trial judges to determine whether the opinion of an expert witness lacks sufficient reliability, and if so, to exclude this testimony. The Daubert decision listed five factors that could be considered when determining

¹⁰ibid., pp. 272-273.

¹¹Stoney and Thornton, op. cit., p. 1187.

¹²Epstein, op. cit., p. 611, footnote 32.

¹³509 U.S. 579 (1993).

¹⁴526 U.S. 137 (1999).

whether scientific expert testimony should be retained or excluded. These factors are as follows:

1. “A preliminary assessment of whether the reasoning or methodology underlying the testimony is scientifically valid and of whether that reasoning or methodology properly can be applied to the facts in issue.”¹⁵
2. “The court ordinarily should consider the known or potential rate of error... .”¹⁶
3. The court should consider “the existence and maintenance of standards controlling the technique’s operation... .”¹⁷
4. “‘General acceptance’ can ... have a bearing on the inquiry. A reliability assessment does not require, although it does permit, explicit identification of a relevant scientific community and an express determination of a particular degree of acceptance within that community.”¹⁸
5. “A pertinent consideration is whether the theory or technique has been subjected to peer review and publication... .”¹⁹

In the *Kumho* case, the Court held that a trial court’s obligation to decide whether to admit expert testimony applies to all experts, not just scientific experts. The Court also held that the factors listed above may be used by a court in assessing nonscientific expert testimony.

In the case (*U.S. v. Llera Plaza*) mentioned at the beginning of the chapter, the presiding judge, Louis Pollack, applied the Daubert criteria to the fingerprint identification process, as he was instructed

¹⁵509 U.S. 579 (1993), note 593.

¹⁶*ibid.*, note 594.

¹⁷*ibid.*

¹⁸*ibid.*, quoted from *U.S. v. Downing*, 753 F.2d 1224, 1238 (3d Cir. 1985).

¹⁹*ibid.*, note 593.

to do by the *Kumho* case. In particular, he discussed the problem with the current process employed by the FBI (and other law enforcement agencies), which is called the ACE-V system. This name is an acronym that stands for analysis, comparison, evaluation, and verification. Judge Pollack ruled that the third part of this process, in which a fingerprint expert states his or her opinion that the latent print and the comparison print (either a rolled print from a suspect or a print from a database) either match or do not match, did not measure up to several of the Daubert criteria.

With regard to the first criterion, the government (the plaintiff in the case) argued that the method of fingerprint matching had been tested empirically over a period of 100 years. It also argued that in any particular case, the method can be tested through the testimony of a fingerprint expert other than the one whose testimony is being heard. The judge rejected this argument, saying that neither of these actions could be considered as scientific tests of the method. He further noted that in the second case, the strength of the second examiner's "test" of a claimed match is diluted by the fact that in many cases, the second examiner has been advised of the first examiner's claims in advance.

On the point of testing, it is interesting to note that in 2000, the National Institute of Justice (NIJ), which is an arm of the Department of Justice, had solicited proposals for research projects to study the reliability of fingerprinting. This solicitation was mentioned by the judge in his ruling and was also taken as evidence by the defense that the government did not know whether fingerprinting was reliable.

The second Daubert criterion concerns the "known or potential rate of error" of the method. In their arguments before the court, the government contended that there were two types of error—methodology error and practitioner error. One of the government's witnesses, when asked to explain methodology error, stated that "an error rate is a wispy thing like smoke, it changes over time..."²⁰ The

²⁰ *U.S. v. Llera Plaza*, January 7, 2002, at 47.

judge said that “the full import of [this] testimony is not easy to grasp.” He summarizes this testimony as saying that if a method, together with its limitations, has been defined, then there is no methodology error. All of the error is practitioner error. The other government witness, Stephen Meagher, a supervisory fingerprint specialist with the FBI, also testified that if the scientific method is followed, then the methodology error rate will be zero, i.e. all of the error is practitioner error. We will have more to say about practitioner error below.

Judge Pollack also found problems concerning the third Daubert criterion, which deals with standards controlling a technique’s operation. There are three types of standards discussed in the judge’s ruling. The first is whether there is a minimum number of Galton points that must be matched before an overall match is declared. In the ACE-V process, no minimum number is prescribed, and in fact, in some jurisdictions, there is no minimum. The second type of standard concerns the evaluation of whether a match exists. The government and defense witnesses agreed that this decision is subjective. The judge concluded that “it is difficult to see how fingerprint identification—the matching of a latent print to a known fingerprint—is controlled by any clearly describable set of standards to which most examiners subscribe.”²¹ Finally, there is the issue of the qualifications of examiners. There are no mandatory qualification standards that must be attained in order for someone to become a fingerprint examiner; nor are there any uniform certification processes.

Regarding the fourth Daubert criterion, the judge had this to say:

General acceptance by the fingerprint examiner community does not ... meet the standard... . First, there is the difficulty that fingerprint examiners, while respected professionals, do not constitute a ‘scientific community’ in the Daubert sense... . Second, the Court cautioned in *Kumho*

²¹ibid. at 58.

Tire that general acceptance does not help show that an expert's testimony is reliable where the discipline itself lacks reliability. The failure of fingerprint identifications fully to satisfy the first three Daubert factors militates against heavy reliance on the general acceptance factor. Thus, while fingerprint examinations conducted under the general ACE-V rubric are generally accepted as reliable by fingerprint examiners, this by itself cannot sustain the government's burden in making the case for the admissibility of fingerprint testimony under Federal Rule of Evidence 702.²²

The conclusion of the judge's ruling was as follows:

For the foregoing reasons:

A. This court will take judicial notice of the uniqueness and permanence of fingerprints.

B. The parties will be able to present expert fingerprint testimony (1) describing how any latent and rolled prints at issue in this case were obtained, (2) identifying, and placing before the jury, such fingerprints and any necessary magnifications, and (3) pointing out any observed similarities and differences between a particular latent print and a particular rolled print alleged by the government to be attributable to the same persons. But the parties will not be permitted to present testimony expressing an opinion of an expert witness that a particular latent print matches, or does not match, the rolled print of a particular person and hence is, or is not, the fingerprint of that person.²³

²²ibid. at 61.

²³ibid. at 69.

The government asked for a reconsideration of this ruling. Not surprisingly, it felt that its effectiveness in both the trial at hand and in future trials would be seriously compromised if witnesses were not allowed to express an opinion on whether or not a latent print matches a rolled print. The government asked to be allowed to submit evidence that would show the accuracy of FBI fingerprint examiners. The defendants argued that the judge should decline to reconsider his ruling, and Judge Pollack stated that their argument was solid: “Neither of the circumstances conventionally justifying reconsideration—new, or hitherto unavailable facts or new controlling law—was present here.”²⁴ Nonetheless, the judge decided to grant a reconsideration hearing, arguing that the record on which he made his previous ruling was testimony presented two years earlier in another courtroom. “It seemed prudent to hear such live witnesses as the government wished to present, together with any rebuttal witnesses the defense would elect to present.”²⁵

At this point in our narrative, it makes sense to consider the various attempts to measure error rates in the field of fingerprint analysis. Lyn and Ralph Haber, who are consultants at a private company in California and also adjuncts at the University of California at Santa Cruz, have obtained and analyzed relevant data from many sources.²⁶ These data include both results on crime laboratories and individual practitioners. We will summarize some of their findings here.

The American Society of Crime Laboratory Directors (ASCLD) is an organization that provides leadership in the management of forensic science. It is in their interest to evaluate and improve the quality of operations of crime laboratories. In 1977, the ASCLD began developing an accreditation program for crime laboratories. By 1999, 182 labs had been accredited. One requirement for a lab to be accredited

²⁴ *U.S. v. Llera Plaza*, March 13, 2002, at 11.

²⁵ *ibid.*

²⁶ Haber, Lynn, and Ralph Norman Haber, “Error Rates for Human Latent Fingerprint Examiners,” in *Advances in Automatic Fingerprint Recognition*, Nalini K. Ratha, ed., New York, Springer-Verlag, 2003.

is that the examiners working in the lab must pass an externally administered proficiency test. We note that since it is the lab, and not the individual examiners, that is being tested, these proficiency tests are taken by all of the examiners as a group in a given lab.

Beginning in 1983, the ASCLD began administering such a test in the area of fingerprint identification. The test, which is given each year to all labs requesting accreditation, consists of pictures of 12 or more latent prints and a set of ten-print (rolled print) cards. The set of latent prints contains a range of quality and is supposed to be representative of what is actually seen in practice. For each latent print, the lab was asked to decide whether it is “scorable,” i.e. whether it is of sufficient quality to attempt to match it with a rolled print. If it is judged to be scorable, then the lab is asked to decide whether or not it matches one of the prints on the ten-print cards. There are “correct” answers for each latent print on the test, i.e. the ASCLD has decided, in each case, whether or not a latent print is scorable, and if so, whether or not it matches any of the rolled prints.

The Habers report on results from 1983 to 1991. During this time, the number of labs that took the exam increased from 24 to 88; many labs took the tests more than once. A new test was constructed each year. Assuming that in many cases the labs have more than one fingerprint expert, this means that hundreds of these experts participated in the test at least once during this period.

Each lab returned one answer for each question. There are four types of errors that can be made on each question of each test. A scorable print can be ruled unscorable or vice versa. If a print is correctly judged to be scorable, it can be erroneously matched to a rolled print, or it can fail to be matched at all, even though a match exists. Of these four types of errors, the second and third are more serious than the others, assuming that we take the point of view that erroneous evidence against an innocent person should be strenuously guarded against.

The percentage of answers with errors of each of the four types were 8%, 2%, 2%, and 8%, respectively. What should we make of these error rates? We see that the more serious types of errors had lower rates, but we must remember that these answers are consensus answers of the experts in a given lab. For purposes of illustration, suppose that there are two experts in a given lab and they agree on an answer that turns out to be incorrect. Presumably they consulted each other on their answers, so we cannot multiply their individual error rates to obtain their group error rate, since their answers were not independent events. However, we can certainly suppose that if a lab error rate is 2%, say, then the individual error rate of at least one of the experts at the lab who took the test is at least 2%.

In 1994, the ASCLD asked the IAI for assistance in creating and reviewing future tests. The IAI asked a company called Collaborative Testing Services (CTS) to design and administer these tests. The format of these tests is similar to the earlier ones, but all of the latent prints are scorable, so there are only two possible types of errors for each question. In addition, individual fingerprint examiners who wish to do so may take the exam by themselves. The Habers report on the error rates for the examinations given from 1995 through 2001. Of the 1685 tests that were graded by CTS, 95 of them, or more than 5%, had at least one erroneous identification, and 502 of the tests, or more than 29%, had at least one missed identification.

Since 1995, the FBI has administered its own examinations to all of its fingerprint examiners. These examinations are similar in nature to the ones described above, but there are a few differences worthy of note. These differences were described in Judge Pollack's reconsideration ruling, in the testimony of Allan Bayle, a fingerprint examiner for 25 years at Scotland Yard.²⁷

Mr. Bayle had reviewed copies of the internal FBI proficiency tests before taking the stand. He found the latent prints utilized in those tests to be, on

²⁷ *U.S. v. Llera Plaza*, March 13, 2002, at 24.

the whole, markedly unrepresentative of the latent prints that would be lifted at a crime scene. In general, Mr. Bayle found the test latent prints to be far clearer than the prints an examiner would routinely deal with. The prints were too clear—they were, according to Mr. Bayle, lacking in the “background noise” and “distortion” one would expect in latent prints that were not identifiable; according to Mr. Bayle, at a typical crime scene only about ten per cent of the lifted latent prints will turn out to be matched. In Mr. Bayle’s view the paucity of non-identifiable prints: “makes the test too easy. It’s not testing their ability. It doesn’t test their expertise. I mean I’ve set these tests to trainees and advanced technicians. And if I gave my experts these tests, they’d fall about laughing.”

Approximately 60 FBI fingerprint examiners took the FBI test each year in the period from 1995 to 2001. On these tests, virtually all of the latent prints had matches among the rolled prints. Since many of the examiners took the tests most or all of these years, it is reasonable to suppose that they knew this fact, and hence would hardly ever claim that a latent print had no match. The results of these tests are as follows: there were no erroneous matches, and only three cases where an examiner claimed there was no match when there was one. Thus, the error rates for the two types of error were 0% and 1%.

It seems clear that the error rates of the crime labs for the various types of error are small, but not negligible, and the FBI’s rates are suspect for the reasons given above. Given that in many criminal cases fingerprint evidence forms a crucial part of the prosecution’s case, it is reasonable to ask whether the above data, were it to be submitted to a jury, would make it difficult for the jury to find the defendant guilty “beyond a reasonable doubt,” which is the standard that must be met in such cases.

The question of what this last phrase means is a fascinating one, and the answers show how hard it is to use probabilistic language in the legal world. The U.S. Supreme Court recently weighed in on this issue, and the majority opinion is thorough in its attempt to explicate the history of the usage of this phrase. The Court agreed to review two cases involving instructions given to juries by judges. Standard instructions to juries state that “guilt beyond a reasonable doubt” means that the jurors need to be convinced “to a moral certainty” of the defendant’s guilt. In one case, “California defended the use of the moral-certainty language as a “commonsense and natural” phrase that conveys an “extraordinarily high degree of certainty.”²⁸ In the second case, a judge in Nebraska “included not only the moral-certainty language but also a definition of reasonable doubt as ‘an actual and substantial doubt.’ The jurors were instructed that ‘you may find an accused guilty upon the strong probabilities of the case, provided such probabilities are strong enough to exclude any doubt of his guilt that is reasonable.’”²⁹ The Supreme Court upheld both sets of instructions. The decision regarding the first set was unanimous, while in the second case, two justices dissented, noting that “the jury was likely to have interpreted the phrase ‘substantial doubt’ to mean that ‘a large as opposed to a merely reasonable doubt is required to acquit a defendant.’”³⁰

The Court went on to note that the meaning of the phrase “moral certainty” has changed over time. In the mid-19th century, the phrase generally meant a high degree of certainty, whereas today, some dictionaries define the phrase to mean “based on strong likelihood or firm conviction, rather than on the actual evidence.”³¹ Although the Court upheld both sets of instructions, the majority opinion stated that the Court did not condone the use of the phrase “moral certainty.”

²⁸Linda Greenhouse, “High Court Warns About Test for Reasonable Doubt,” *New York Times*, March 22, 1994.

²⁹*ibid.*

³⁰*ibid.*

³¹*American Heritage Dictionary of the English Language*, 1992.

In a concurring opinion, Justice Ruth Bader Ginsburg noted that some Federal appellate circuit courts have instructed trial judges not to provide any definition of the phrase “beyond a reasonable doubt.” Justice Ginsburg said that it would be advisable to construct a better definition than the one used in the instructions in the cases under review. She cited one suggested in 1987 by the Federal Judicial Center, a research arm of the Federal judiciary. Making no reference to moral certainty, that definition says in part, “Proof beyond a reasonable doubt is proof that leaves you firmly convinced of the defendant’s guilt.”³²

It may very well be the case that after wading through the above verbiage, the reader has no clearer an idea (and perhaps even has a less clear idea) than before of what the phrase “beyond a reasonable doubt” means. However, juries are given this phrase as part of their instructions, and in the case of fingerprint evidence, they deserve to be educated about error rates involved. We leave it to the reader to ponder whether evidence produced by a technique whose error rate seems to be at least 2% is strong enough to be beyond a reasonable doubt.

On March 13, 2002, Judge Pollack filed his second decision in the Llera Plaza case. The judge’s ruling was a partial reversal of the original one. His ruling allowed FBI fingerprint examiners to state in court whether there is a match between a latent and a rolled print, but nothing was said in the ruling about examiners not in the employ of the FBI. The judge’s mind was changed primarily because of the testimony of Mr. Bayle who, ironically, was a witness for the defense. Although, as noted above and in the judge’s decision, there are shortcomings in the FBI’s proficiency testing of its examiners, the judge was convinced by the facts that the ACE-V system used by the FBI is essentially the same as the system used in Great Britain and that Mr. Bayle believes in this system without reservation.

³²Greenhouse, loc. cit.

As an interesting footnote to this case, after Judge Pollack announced his second ruling, the NIJ cancelled its original solicitation, described above, and replaced it by a “General Forensic Research and Development” solicitation. In the guidelines for this proposal under “what will not be funded,” we find the phrase “proposals to evaluate, validate, or implement existing forensic technologies.” This is a somewhat strange way to respond to the judge’s worries about whether the method has been adequately tested in a scientific manner.

5. The 50K Study

At the beginning of Section 3, we stated that in order to decide whether fingerprints are useful in forensics, it is of central importance to be able to estimate how likely it is that a latent print will be incorrectly matched to a rolled print. In 1999, the FBI asked the Lockheed Martin Company to carry out a study of fingerprints. In a pre-trial hearing in the case *U.S. v. Mitchell*³³, Stephen Meagher, whom we have introduced earlier, explained why he commissioned the study. The primary reason for carrying out this study, he said, was to use the FBI database of over 34 million sets of 10 rolled prints to see how well the automatic fingerprint recognition computer programs distinguished between prints of different fingers. The results of the study could also be used, he reasoned, to strengthen the claim that no two fingerprints are alike. Thus, this study was not originally conceived as a test of the accuracy of matching latent and rolled prints. Nonetheless, as we shall see, this study touched on this second issue.

Together with Bruce Budlowe, a statistician who works for the FBI, Meagher came up with the following design for the experiment. The overall idea was to compare every pair of rolled prints in the database, to see if the computer algorithms could distinguish among different prints with high accuracy. It was decided that carrying this

³³ *U.S. v. Mitchell*, July 7, 1999.

out for the whole database was not reasonable (about 5.8×10^{16} comparisons would be required), so they instead chose 50,000 rolled fingerprints from the FBI's master file. These prints were not chosen at random; rather, they were the first 50,000 that were of the pattern "left loop" from white males. It was decided to restrict the fingerprints in this way because according to Meagher, race and gender have some effect on the size and types of fingerprints. By restricting in this way, the resulting set of fingerprints are probably more homogeneous than a set of randomly chosen fingerprints would be, thereby making it harder to distinguish between pairs from the set. If the study could show that each pair could be distinguished, then the result is more impressive than a similar result accomplished using a set of randomly chosen prints.

At this point, Meagher turned the problem of design over to the Lockheed group, where the design and implementation of the study were carried out by Donald Zeisig, an applied mathematician and software designer, and James O'Sullivan, a statistician. Much of what follows comes from testimony that Zeisig gave at the pre-trial hearing in *U.S. v. Mitchell*.

Two experiments with this data were performed. The first began by using two different software programs that each generated a measure of similarity between two fingerprints based on their minutiae patterns. A third program was used to merge these two measures. A paper by David Kaye³⁴ delved into various difficulties presented by this study. Information about this study was also provided by the fascinating transcripts of the pre-trial hearing mentioned above³⁵.

We follow Kaye in denoting the measure of similarity between fingerprints f_i and f_j by $x(f_i, f_j)$. Each of the fingerprints was compared with itself, and the function x was normalized. Although this normalization is not explicitly defined in either the court testimony or the Lockheed summary of the test, we will proceed as best we

³⁴Kaye, David, "Questioning a Courtroom Proof of the Uniqueness of Fingerprints", *International Statistical Review*, vol. 71, no. 3 (2003), pp. 521-533.

³⁵Daubert Hearing Transcripts, at www.clpex.com/Mitchell.htm

can. It seems that the values of $x(f_i, f_j)$ were all multiplied by a constant, so that $x(f_i, f_i) \leq 1$ for all i , and there is an i such that $x(f_i, f_i) = 1$. One would expect that a measure of similarity would be symmetric, i.e. that $x(f_i, f_j) = x(f_j, f_i)$, but this is never mentioned in the report, and in fact there is evidence that this is not true for this measure.

The value of $x(f_i, f_j)$ is then computed for all 2.5×10^9 ordered pairs of fingerprints. If this measure of similarity is of any value, it should be very small for all pairs of non-identical fingerprints and large (i.e. close to 1) for all pairs of identical fingerprints.

Next, for each rolled print f_i , the 500 largest values of $x(f_i, f_j)$ are recorded. One of these values, namely when $j = i$, will presumably be very close to 1, but the other 499 values will probably be very close to 0. At this point, the Lockheed group calculated the mean and standard deviation of this set of 500 values (for each fixed value of i). Presumably, the mean and the standard deviation are both positive and very close to 0 (since all but one of the values is very small and positive).

Next, Zeisig and O'Sullivan assume that the distribution, for each i , is normal, with the calculated mean and standard deviation. No reason is given for making this assumption, and we shall see that it gives rise to some amazing probabilities. Under this assumption, one can change the values of $x(f_i, f_j)$ into values of a standard normal distribution by subtracting the mean and dividing by the standard deviation. The Lockheed group calls these normalized values Z scores. The reader can see that if this is done for a typical set of 500 values of $x(f_i, f_j)$, with i fixed, one should obtain 499 Z scores that are fairly close to 0 and one Z score, corresponding to $x(f_i, f_i)$, that is quite large.

It is then pointed out that if one takes 500 independent values from the standard normal distribution, the expected value of the largest value obtained should be about 3. This value is easy to estimate by simulation; we simulated 50,000 repetitions of the maximum

of 500 independent standard normal values, and found the mean of the maximum values to be 3.04. Thus, Zeisig and O'Sullivan would be worried if any of the non-mate Z scores (i.e. Z scores corresponding to pairs (f_i, f_j) with $i \neq j$) were much greater than 3. In fact, except for three cases, which will be discussed below, all of the non-mate Z scores were less than 1.83. This fact casts much doubt on whether the distribution in question is normal.

The three non-mate Z scores that were larger than 1.83 corresponded to the (i, j) -pairs (48541, 48543), (48543, 48541), and (18372, 18373). The scores in these cases were 6.98, 6.95, and 3.41. When Zeisig and O'Sullivan found these high Z values, they discovered that in all three cases, the pairs were different rolled prints of the same finger. In other words, the sample of 50,000 fingerprints were from at most 49,998 different people. It is interesting to note that the ordered pair (18373, 18372) must have had a Z score of less than 1.83, even though the pair corresponds to two prints of the same finger. We'll have more to say about this below. This shows that it is possible for two different prints of the same finger to generate a Z score which is in the same range as a score generated by two prints of different fingers.

Now things get murky. The smallest Z score of any fingerprint paired with itself was stated to be 21.7. This high value is to be expected; the reader will recall that for any fingerprint f_i , the 500 values correspond to 499 small Z scores and one very large Z score. However, the conclusion drawn from this statement is far from clear. If one calculates the probability that a standard normal random variable will take on a value greater than 21.0, one obtains a value of less than 10^{-97} . The Lockheed group states its conclusion as follows³⁶: "The probability of a non-mate rolled fingerprint being identical to any particular fingerprint is less than 10^{-97} ."

David Kaye points out that the real question is not whether a computer program can detect copies of photographs of rolled prints,

³⁶Kaye, op. cit. , p. 530

as is done in this study when a rolled print is compared with itself. Rather, it is whether such a program can, for each finger in the world, put all rolled prints of that finger in one category and make sure that no rolled prints from any other finger fall into that same category. Kaye notes that although there were so few repeated fingers in the study that one cannot determine the answer to this question with any great degree of certainty, one of the three pairs noted above, of different rolled prints of the same finger, produced a Z score that would occur about once in every 3000 comparisons, assuming the comparisons generate scores that are normally distributed. This means that if one were to make millions of comparisons between pairs of rolled prints of different fingers, one would find thousands of Z scores as high as the one corresponding to the pair (18372, 18373). This would put the computer programmer in a difficult situation. To satisfy Kaye, the program would have to be assigned a number Z^* with the property that if a Z score were generated that was above this value, the program would state that the prints were of the same finger, while if the generated Z score were below this value, the program would state that the two prints were of different fingers. But we can see that there can be no such Z^* value that will always be right. If $Z^* > 3.41$ (the value corresponding to the pair (18372, 18373)) then the program would declare that the thousands of the pairs of prints of different fingers mentioned above are in fact prints of the same finger. If $Z^* < 3.41$, then the program would declare that the pair (18372, 18373) are prints of different fingers.

As we noted above, the pair (18373, 18372) was not flagged as having a large Z score. The reason for this is that when the three non-mate pairings mentioned above were flagged, it was not yet known that they corresponded to the same fingers. However, one does wonder whether the Lockheed group looked at the Z score of this pair, once the reversed pair was discovered to have a high Z score. In any event, the Z score of this pair is not given in the summary of the experiments. Robert Epstein, an attorney for the defense in *U.S.*

v. Mitchell, noticed this fact as well and asked Donald Zeisig, during cross-examination, what the Z score of this pair was. It turns out that the Z score was 1.79. This makes things still worse for the matching algorithm.

First, there were other non-mate pairs with larger Z scores. Second, one might expect that the Z score of a pair would be roughly the same in either order (although it isn't clear that this should be so). In any event, a Z score of 1.79 does not correspond to an extremely unlikely event; thus, the algorithm might fail, with some not-so-small probability, to detect an identification between two fingerprints (or else might, with some not-so-small probability, make false identifications). In fact Epstein, in his cross-examination, noted that the pair (12640, 21111) had the Z values 1.83 and 1.47 (depending upon the order), even though it was later discovered that both of this prints were of the same finger. When asked by Epstein, Zeisig agreed that there could possibly have been other pairs of different prints of the same finger (which must have had low Z values, since they were not flagged).

The second experiment that the Lockheed group performed was an attempt to find out how well their computer algorithms dealt with latent fingerprints. To that end, a set of "pseudo" latent fingerprints was made up, by taking the central 21.7% of each of the 50,000 rolled prints in the original data set. This percentage was arrived at by taking the average size of 300 latent prints from crime scenes versus the size of the corresponding rolled prints.

At this point, the experiment was carried out in essentially the same way as the first experiment. Each pseudo latent l_i was compared with all 50,000 rolled prints, and a score $y(l_i, f_j)$ was determined. For each latent l_i , the largest 500 values of $y(l_i, f_j)$ were used to construct Z scores. As before, the Z score corresponding to the pair (l_i, f_i) was expected to be the largest of these by far. Any non-mate Z scores that were high were a cause for concern.

The two pairs (48541, 48543) and (18372, 18373) did give high Z scores, but it was already known at this point that these pairs corresponded to different rolled images of the same finger. There were three other pairs whose Z scores were above 3.6. One pair, (21852, 21853) gave a Z score of 3.64. The latent and the rolled prints were of fingers 7 and 8 of the same person. Further examination of this pair determined that part of finger 8 had intruded into the box containing the rolled print of finger 7. The computer algorithm had found this intrusion, when the pseudo latent for finger 8 was compared with the rolled print of finger 7. This is a somewhat impressive achievement.

One other pair, (12640, 21111), generated large Z scores in both orders. At the time the summary was written, it had not yet been determined whether these two prints were of the same finger. The Lockheed group compared all 20 fingerprints (taken from the two sets of 10 rolled prints corresponding to this pair) with each other. Not surprisingly, the largest scores were generated by prints being compared with themselves. The second highest score for each print was generated when that print was compared with the corresponding print in the other set of 10 rolled prints, and these second-highest scores were quite a bit higher than any of the remaining scores. This is certainly strong evidence that the two sets of 10 rolled prints corresponded to the same person.

The second experiment does not get at one of the central issues concerning latent prints, namely how the quality of the latent print affects the ability of the fingerprint examiner (or a computer algorithm) to match this latent print with a rolled one. Figures 2 and 3 show that latent prints do not look much like the central 21.7% of a rolled print. Yet it is just these types of comparisons that are used as evidence in court. It would be interesting to conduct a third experiment with the Lockheed data set, in which care was taken to create a more realistic set of latent prints.

Exercise.

1. By the middle of the 20th century, the FBI had compiled a set of more than 10 million fingerprints. Suppose that there are n fingerprint patterns among all of the people on Earth. Thus, n is some number that does not exceed 10 times the number of people on Earth, and it equals this value if and only if no two fingerprints are exactly alike.
 - (a) Suppose that all n fingerprint patterns are equally likely. Estimate the number $f(n)$ of random fingerprints that must be observed in order that the probability that two of the same pattern are observed exceeds .5. Hint: To do this using a computer, try different values of n and guess an approximate relationship between n and $f(n)$.
 - (b) Under the supposition in part a), given that $f(n) = 10$ million, estimate n . Note that it is possible to show that if not all n fingerprint patterns are assumed to be equally likely, then the value of $f(n)$ decreases.
 - (c) Suppose that $n < 60$ billion (so that at least two fingerprints are alike). Estimate $f(n)$.
 - (d) Suppose that $n = 30$ billion, so that, on the average, every pattern appears twice among the people who are presently alive. Using a computer, choose 10 million patterns at random, thereby simulating the set compiled by the FBI. Was any pattern chosen more than once? Repeat this process many times, keeping track of whether or not at least one pattern is chosen twice. What percentage of the time was at least one pattern chosen at least twice?
 - (e) Do the above calculations convince you that no fingerprint pattern appears more than once among the people who are alive today?

Answers to John Haigh's Lottery Questions

1. In 721 drawings, there are 4326 ($= 721 \cdot 6$) main numbers. Thus, the average number of times a given number occurs as a main number is 88.3 ($= 4326/49$). To estimate the spread of the frequencies of the various numbers, we can use a normal approximation. The probability that a given number appears as a winning number in a given drawing is $6/49$. If we call the occurrence of this number a success and its non-occurrence a failure, then we have a Bernoulli trials process, with $n = 721$, $p = 6/49$, and $q = 43/49$. The variance of this process is $npq \approx 77.5$ (so the standard deviation is about 8.8). The frequencies of the various numbers can therefore be approximated by a normal distribution with mean 88.3 and variance 77.5. For a normal distribution, the probability that a given observation lies at least two standard deviations from the mean is .0456, and the corresponding probability for three standard deviations is .0026. Since there are 49 numbers, we see that the expected number of the frequencies lying more than two standard deviations from the mean is $2.23 (= 49 \cdot .0456)$

while the corresponding expected number for three standard deviations is $0.13 (= 49 \cdot .0026)$. Thus we should not be surprised to see a frequency lying slightly more than two standard deviations from the mean. Thus, we might estimate that the highest and lowest frequencies are around 106 and 70. The actual answers for this set of draws are 113 and 70.

2. This is an example of the coupon collector's problem. A description of this problem can be found in Grinstead and Snell [19] in Exercise 3.2.34. One finds that to make the probability at least $1/2$ of getting each number at least once as the bonus number, one needs to have about $n \log n + n \log 2$, where $n = 49$, drawings. This yields the estimate 225. In the actual set of 721 drawings, a complete set of bonus numbers was achieved in the first 262 drawings.

3. If we write the sets of main numbers in a sequence, we can imagine that this sequence is a sequence of independent draws. This isn't quite right, because in our sequence, it is not possible, for example, that the first and fifth numbers are equal, since they both occur in the first set of main numbers. In the answer to problem 2, we saw that we would expect to need around 225 draws of numbers between 1 and 49 to get, with probability $1/2$, each of them at least once. Thus, we would expect that it would take about $37.5 (= 225/6)$ drawings to get each number at least once as a main number. In the actual sequence of drawings, all numbers appeared at least once in the first 26 drawings.

There are $1176 (= 49 \cdot 48/2)$ pairs among the 49 numbers. In each drawing, $15 (= 6 \cdot 5/2)$ pairs of main numbers are drawn. Using the approximation in the answer to problem 2, and ignoring dependence of pairs within a single draw, we see that we would expect to have to draw about 9129 pairs to get each of the 1176 pairs at least once. Since 15 pairs are drawn in each drawing, we would expect to get all of the pairs in about $609 (= 9129/15)$ drawings. In the actual sequence, it took 591 drawings.

4. For any fixed number between 1 and 49, the probability that it is a main number in any given draw is $6/49$. If we use the formula stated in Chapter 1, Section 4 for the expected length of the longest success run, we find, by setting $n = 721$ and $p = 6/49$, that this expected length is about 2.85. But this isn't really what we want to estimate. We want the expected length of the longest success run for all 49 numbers. We would expect this to be somewhat longer than the expected length of the longest success run for a given number.

Although it might be possible to calculate the distribution of the longest success run for 49 numbers, a better way to proceed is to simulate the distribution. We simulated a sequence of 721 draws 1000 times. We found that the longest success run equaled four for 403 of the 1000 simulations and equaled five for 480 of the 1000 simulations. It was longer than five 117 times. In the actual data, there was one number that occurred in five consecutive draws.

John Haigh gave us the following analysis of this question. In any drawing after the first one, consider any one of the six numbers drawn. It will begin a streak of length k if (a) it was not drawn in the previous drawing, and (b) it is also drawn in the next $k - 1$ drawings. As there are six main numbers, the chance that one of them begins a streak of length at least n is (very slightly less than)

$$(6) \left(\frac{43}{49} \right) \left(\frac{6}{49} \right)^{k-1}.$$

So the mean number of streaks of length at least k is n times the above number, where n is the number of drawings. When $n = 721$, this expression equals 0.85 when $k = 5$ and 0.10 when $k = 6$. This means that we should not be surprised to see a streak of length five but we should be surprised to see a streak of length six.

5. The number of possible sets of main numbers is

$$\binom{49}{6} = 13983816.$$

Thus, the probability of winning the jackpot with one ticket is the reciprocal of this number. To estimate the probability of no winners in a given drawing, we can use the Poisson approximation to the binomial distribution. However, we have seen that the number of tickets sold for a given drawing varies quite a bit. This affects the probability of having no winners. For example, if the number of tickets sold for a given drawing is 27.7 million (the average number for the Wednesday drawings), then we use the Poisson distribution with

$$\lambda = \frac{27700000}{13983816} = 1.98.$$

In this case, the probability of having no winners is about

$$e^{-1.98} = 0.138.$$

If the number of tickets sold is 57.5 million (the average number for the Saturday drawings), then we use

$$\lambda = \frac{57500000}{13983816} = 4.11.$$

In this case, the probability of having no winners is about

$$e^{-4.11} = 0.016.$$

Of the 721 drawings in our data set, 303 occurred on Wednesday and 418 occurred on Saturday. So, we would expect about $(303)(0.138) = 42$ of these drawings would have no winners, and about $(418)(0.016) = 7$ of the Saturday drawings would have no winners. In fact, there were no winners in 107 of the 721 drawings.

6. Once again, we consider the Wednesday and Saturday drawings separately. In the Wednesday drawings, we would expect there to be about

$$\frac{27700000}{13983816} \approx 1.98$$

winners. In the Saturday drawings, we would expect there to be about

$$\frac{57500000}{13983816} \approx 4.11$$

winners. Haigh says the most common number of winners in the 721 draws was two, with 150 occurrences. He says there was one winner 149 times.

7. Again we consider the Wednesday and Saturday drawings separately. For the Wednesday drawings, using the Poisson approximation to the binomial distribution, the expected number of drawings (in 303) with exactly 7 winners is 0.99, and the corresponding number with exactly 8 winners is 0.25. For the Saturday drawings, the expected number of drawings with exactly 11 winners is 0.97 and the corresponding number with exactly 12 winners is 0.33. Thus, using only the average number of tickets sold, we might guess that the answer to this question is close to $n = 12$. In fact, the answer is $n = 17$.

Bibliography

- [1] Albert, Jim, and Jay Bennett. *Curve Ball*. New York: Copernicus Books, 2001.
- [2] Albright, S. Christian. "A Statistical Analysis of Hitting Streaks in Baseball." *Journal of the American Statistical Association*, 88.424 (December 1993): 1175–1183.
- [3] Arbesman, Samuel, and Steven Strogatz. "A Journey to Baseball's Alternate Universe." *New York Times* 30 Mar. 2008.
- [4] Berry, Scott. "The Summer of '41: A Probabilistic Analysis of DiMaggio's 'Streak' and Williams's Average of .406." *Chance*, 4.4 (1991): 8–11.
- [5] Camerer, Colin. "Does the Basketball Market Believe in the 'Hot Hand'?" *American Economic Review*, 79.5 (1989):1257–1261.
- [6] Champod, Christophe, and Ian W. Evett. "A Probabilistic Approach to Fingerprint Evidence." *Journal of Forensic Identification*, 51.2 (2001):101–122.
- [7] Chung, K. L. *A Course in Probability Theory, second edition*. London, Academic Press, 2001.
- [8] Cleary, Rick. "Surprising Streaks and Playoff Parity: Probability Problems in a Sports Context."
- [9] Coit, Michael. "Santa Rosa Woman Identified as Vegas Slaying Victim Turns Up Alive." *Press Democrat* 13 Sept. 2002.
- [10] Embrechts, Paul, Charles Goldie, and Noel Veraverbeke. "Subexponentiality and Infinite Divisibility." *Z. für Wahrscheinlichkeitstheorie*, 49 (1979): 335–347.

- [11] Epstein, Robert. "Fingerprints Meet Daubert: The Myth of Fingerprint 'Science' is Revealed." *California Law Review*, 75 (2002): 605–658.
- [12] Feller, William. *An Introduction to Probability Theory and Its Applications, Volume 1, third edition*. New York: John Wiley and Sons, 1968.
- [13] Feller, William. *An Introduction to Probability Theory and Its Applications, Volume 2*. New York: John Wiley and Sons, 1991.
- [14] Fitzgerald, Thomas J. "Fingerprints on File, Right From the Patrol Car." *New York Times* 23 Sept. 2004: G7.
- [15] Gerth Jeff. "Fingerprinting Glitches are Said to Hurt Antiterror Effort." *New York Times* 27 Oct. 2004: A10.
- [16] Gilovich, Thomas, Robert Vallone, and Amos Tversky. "The Hot Hand in Basketball: On the Misperception of Random Sequences." *Cognitive Psychology*, 17 (1995): 295–314.
- [17] Gordon, L., M. F. Schilling, and M. S. Waterman. "An extreme value theory for long head runs." *Probability Theory and Related Fields*, 72: 279–287, 1986.
- [18] Gould, Stephen Jay. "The Streak of Streaks." *Chance*, 2.2 (1989): 10–16.
- [19] Grinstead, Charles M., and J. Laurie Snell. *Introduction to Probability*. Providence: American Mathematical Society, 1997.
- [20] Haber, Lyn, and Ralph Norman Haber. "Error Rates for Human Latent Fingerprint Examiners." *Advances in Automatic Fingerprint Recognition*. Nalini K. Ratha, editor. New York: Springer-Verlag, 2003.
- [21] "Method for Fingerprint Identification." 29th European Regional Conference, May 17, 2000. Interpol European Expert Group on Fingerprint Identification.
- [22] Jackson, D. A., and K. Mosurski. "Heavy Defeats in Tennis: Psychological Momentum or Random Effect?" *Chance*, 10.2 (1997): 27–34.
- [23] Kay, David H. "Questioning a Courtroom Proof of the Uniqueness of Fingerprints." *International Statistical Review*, 71.3 (2003): 521–533.
- [24] Kemeny, John G., and J. Laurie Snell. *Finite Markov Chains*. New York: D. Van Nostrand and Company, 1960.
- [25] Kershaw, Sarah, and Eric Lichtblau. "Judge Rejects Bomb Case Against Oregon Lawyer." *New York Times*, 25 May 2004: A16.
- [26] Kingston, Charles R. "Probabilistic Analysis of Partial Fingerprint Patterns." Doctoral Thesis, University of California, Berkeley, 1964.
- [27] Lackritz, James R. "Two of Baseball's Great Marks: Can They Ever Be Broken?" *Chance*, 9.4 (1996): 12–18.

- [28] Langenburg, Glenn. "Defending Against the Critic's Curse." *the Print*, 19.3 (2003): 1–10.
- [29] Lo, Andrew, and Craig MacKinlay. *A Non-Random Walk Down Wall Street*. Princeton: Princeton University Press, 1999.
- [30] McCulloch, J. Huston. *Financial Applications of Stable Distributions*, in *Handbook of Statistics*, Vol. 14: 393–425. G. S. Maddala and C. R. Rao, eds. Amsterdam: Elsevier Science, 1996.
- [31] Mood, A. M. "The Distribution Theory of Runs." *Annals of Mathematical Statistics*, 11.4 (1940): 367–392.
- [32] Nocera, Joe. "Risk Mismanagement." *New York Times* 4 January 2009: MM24.
- [33] Osterberg, James W., T. Parthasarathy, T. E. S. Raghavan, and Stanley L. Sclove. "Development of a Mathematical Formula for the Calculation of Fingerprint Probabilities Based on Individual Characteristics." *Journal of the American Statistical Association*, 72.360 (1977): 772–778.
- [34] Pankanti, Sharath, Salil Prabhakar, and Anil K. Jain. "On the Individuality of Fingerprints." *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 24.8 (2002): 1010–1025.
- [35] Poterba, James M., and Lawrence H. Summers. "Mean Reversion in Stock Prices." *Journal of Financial Economics*, 22.1 (1988): 27–59.
- [36] Rockoff, David M., and Philip A. Yates. "Chasing DiMaggio: Streaks in Simulated Seasons Using Non-Constant At-Bats." *Journal of Quantitative Analysis in Sports*, 5.2 (2009).
- [37] Schilling, Mark F. "The Longest Run of Heads." *The College Mathematics Journal*, 21.3 (1990): 196–207.
- [38] Short, Tom, and Larry Wasserman. "Should We Be Surprised by the Streak of Streaks?" *Chance*, 2.2 (1989): 13.
- [39] Smith, Gary. "Horseshoe Pitchers' Hot Hands." *Psychonomic Bulletin & Review*, 10 (2003): 753–758.
- [40] Specter, Michael. "Do Fingerprints Lie?" *New Yorker* 27 May 2002: 1–17.
- [41] Stoney, David. "Measurement of Fingerprint Individuality." in *Advances in Fingerprint Technology*, 2nd ed. Lee and Gaensslen, editors. Boca Raton: CRC Press, 2001.
- [42] Stoney, David A., and John I. Thornton. "A Critical Analysis of Quantitative Fingerprint Individuality Models." *Journal of Forensic Sciences*, 31.4 (1986): 1187–1216.
- [43] Stoney, David A., and John I. Thornton. "A Method for the Description of Minutia Pairs in Epidermal Ridge Patterns." *Journal of Forensic Sciences*, 31.4 (1986): 1217–1234.

- [44] Stout, David. “Report Faults F.B.I.’s Fingerprint Scrutiny in Arrest of Lawyer.” *New York Times* 17 Nov. 2004: A18
- [45] Zaharov, V. K., and O. V. Sarmanov. “Distribution Law for the Number of Series in a Homogeneous Markov Chain.” *Soviet Math. Dokl.*, 9.2 (1968): 399–402.

Index

- χ^2 -test, 11
 m -independence of random
 variables, 72
- absorbing Markov chain, 54
absorbing state, 54
ACE-V system, 208
Albert, J., 40
Albright, S., 11
American Society of Crime
 Laboratory Directors, 211
annuity, 152
Arbesman, S., 37
arch, 196
- Bachelier, L., 103, 142
Barnett, A., 186
Bennett, J., 40
Bernoulli trials model, 4
Bertillonage, 193
beyond a reasonable doubt, 214
binomial coefficient, 153
binomial distribution, 104
block-Bernoulli process, 3, 14
Boston Celtics, 47
Brown, R., 106
Brownian motion, 105
Budlowe, B., 217
butter prices, 128
- Camerer, C., 48
- Central Limit Theorem, 106
characteristic function, 144
Cleary, R., 45
convolution, 143
core, fingerprint, 196
cotton prices, 108
- Daubert v. Merrell Dow*, 206
delta, 196
DiMaggio, Joe, 30
discrete return, 98
dividend, 101
Doob, J., 184
doubleton, 11
doubleton distribution
 asymptotics, in Markov chain
 model, 88
 exact, in Bernoulli trials model,
 90
Dow Jones Industrial Average, 120
Dropo, Walt, 44
- Easy Pick, 162
Einstein, A., 105
Epstein, R., 221
Euler's constant, 10
- fat tails, 113
Feller, W., 130
Fielding, Henry, 185

- fingerprint
 - latent, 198
 - rolled, 203
- Fourier transform, 144
- Galton points, 197
- Galton, Francis, 194
- General Electric, 118
- geometric return, 98
- Gilovich, T., 46
- Haber, L., 211
- Haber, R., 211
- Haigh, John, 189, 225
- Henry, Edward, 196
- Herschel, William, 192
- hitting streak, 30
- horseshoes, 51
- hypothesis test, 13
- hypothesis testing, 7
- IBM, 119, 136
- International Association for Identification, 206
- Jackson, D., 62
- Kaye, D., 218
- Kumho Tire v. Carmichael*, 206
- latent fingerprint, 198
- Llera Plaza, U.S. v.*, 191
- Lo, A., 75
- Lockheed Martin study, 217
- longest run
 - distribution of, 22
- loop, 196
- MacKinlay, C., 75
- Mandelbrot, B., 108
- Markov chain, 2
 - absorbing, 54
- Meagher, S., 209, 217
- minutiae, 197
- Mitchell, U.S. v.*, 217
- Mood, A., 79
- Mosurski, K., 62
- mutual fund, 103
- normal distribution, 106
- O'Sullivan, J., 218
- odds, 154
- odds model, 5, 63
- Oster, E., 165
- p-value, 7
- Philadelphia 76ers, 46
- points, Galton, 197
- Poisson approximation, 163
- Poterba, J., 134
- power law, 108
- power of a test, 14
- random walk, 103
- Rhodes, Cecil, 185
- Rockoff, D., 38
- rolled fingerprint, 203
- Roxburgh, T., 201
- run, 2, 6
 - distribution of longest, 22
- runs distribution
 - in Bernoulli trials model, 79
 - in Markov chain model, 84
- S&P 500 index, 102
- S&P 500 list, 72
- Sarmanov, O., 79
- Schilling, M., 10
- Smith, G., 52
- stable set of distributions, 143
- stock market, 71
- Stoney, D., 201
- Strogatz, S., 37
- Summers, L., 134
- tennis, 5, 54
- Thornton, J., 201
- transient state, 54
- Tversky, A., 46
- type I error, 14
- type II error, 14
- Vallone, R., 46
- value at risk, 127
- VAR, 127
- variance ratio, 75, 131
- vig, 47
- National Institute of Justice, 208

-
- whorl, 196
Wiener, N., 105

Yates, P., 38

Zaharov, V., 79
Zeisig, D., 218



Photo by James E. Shoridan

This book explores four real-world topics through the lens of probability theory. It can be used to supplement a standard text in probability or statistics. Most elementary textbooks present the basic theory and then illustrate the ideas with some neatly packaged examples. Here the authors assume that the reader has seen, or is learning, the basic theory from another book and concentrate in some depth on the following topics: streaks, the stock market, lotteries, and fingerprints. This extended format allows the authors to present multiple approaches to problems and to pursue promising side discussions in ways that would not be possible in a book constrained to cover a fixed set of topics.

To keep the main narrative accessible, the authors have placed the more technical mathematical details in appendices. The appendices can be understood by someone who has taken one or two semesters of calculus.

ISBN 978-0-8218-5261-3



9 780821 852613

STML/57



For additional information
and updates on this book, visit

www.ams.org/bookpages/stml-57

AMS on the Web
www.ams.org