

Global Text



Introductory
Business Statistics

Introductory Business Statistics

Thomas K. Tiemann

Copyright © 2010 by Thomas K. Tiemann

For any questions about this text, please email: drexel@uga.edu

Editor-In-Chief: Thomas K. Tiemann

Associate Editor: Marisa Drexel

Editorial Assistants: Jaclyn Sharman, LaKwanzaa Walton

The Global Text Project is funded by the Jacobs Foundation, Zurich, Switzerland.



[This book is licensed under a Creative Commons Attribution 3.0 License](https://creativecommons.org/licenses/by/3.0/)

Table of Contents

What is statistics?	5
1. Descriptive statistics and frequency distributions	10
Descriptive statistics	12
2. The normal and t-distributions	18
Normal things	18
The t-distribution	22
3. Making estimates	26
Estimating the population mean	26
Estimating the population proportion	27
Estimating population variance	29
4. Hypothesis testing	32
The strategy of hypothesis testing	33
5. The t-test	41
The t-distribution	41
6. F-test and one-way anova	52
Analysis of variance (ANOVA)	55
7. Some non-parametric tests	59
Do these populations have the same location? The Mann-Whitney U test	60
Testing with matched pairs: the Wilcoxon signed ranks test	63
Are these two variables related? Spearman's rank correlation	66
8. Regression basics	70
What is regression?	70
Correlation and covariance	79
Covariance, correlation, and regression	81

About the author

Author, Thomas K. Tiemann

Thomas K. Tiemann is Jefferson Pilot Professor of Economics at Elon University in North Carolina, USA. He earned an AB in Economics at Dartmouth College and a PhD at Vanderbilt University. He has been teaching basic business and economics statistics for over 30 years, and tries to take an intuitive approach, rather than a mathematical approach, when teaching statistics. He started working on this book 15 years ago, but got sidetracked by administrative duties. He hopes that this intuitive approach helps students around the world better understand the mysteries of statistics.

A note from the author: Why did I write this text?

I have been teaching introductory statistics to undergraduate economics and business students for almost 30 years. When I took the course as an undergraduate, before computers were widely available to students, we had lots of homework, and learned how to do the arithmetic needed to get the mathematical answer. When I got to graduate school, I found out that I did not have any idea of how statistics worked, or what test to use in what situation. The first few times I taught the course, I stressed learning what test to use in what situation and what the arithmetic answer meant.

As computers became more and more available, students would do statistical studies that would have taken months to perform before, and it became even more important that students understand some of the basic ideas behind statistics, especially the sampling distribution, so I shifted my courses toward an intuitive understanding of sampling distributions and their place in hypothesis testing. That is what is presented here—my attempt to help students understand how statistics works, not just how to “get the right number”.

What is statistics?

There are two common definitions of statistics. The first is "turning data into information", the second is "making inferences about populations from samples". These two definitions are quite different, but between them they capture most of what you will learn in most introductory statistics courses. The first, "turning data into information," is a good definition of descriptive statistics—the topic of the first part of this, and most, introductory texts. The second, "making inferences about populations from samples", is a good definition of inferential statistics—the topic of the latter part of this, and most, introductory texts.

To reach an understanding of the second definition an understanding of the first definition is needed; that is why we will study descriptive statistics before inferential statistics. To reach an understanding of how to turn data into information, an understanding of some terms and concepts is needed. This first chapter provides an explanation of the terms and concepts you will need before you can do anything statistical.

Before starting in on statistics, I want to introduce you to the two young managers who will be using statistics to solve problems throughout this book. Ann Howard and Kevin Schmidt just graduated from college last year, and were hired as "Assistants to the General Manager" at Foothill Mills, a small manufacturer of socks, stockings, and pantyhose. Since Foothill is a small firm, Ann and Kevin get a wide variety of assignments. Their boss, John McGrath, knows a lot about knitting hosiery, but is from the old school of management, and doesn't know much about using statistics to solve business problems. We will see Ann or Kevin, or both, in every chapter. By the end of the book, they may solve enough problems, and use enough statistics, to earn promotions.

Data and information; samples and populations

Though we tend to use data and information interchangeably in normal conversation, we need to think of them as different things when we are thinking about statistics. Data is the raw numbers before we do anything with them. Information is the product of arranging and summarizing those numbers. A listing of the score everyone earned on the first statistics test I gave last semester is data. If you summarize that data by computing the mean (the average score), or by producing a table that shows how many students earned A's, how many B's, etc. you have turned the data into information.

Imagine that one of Foothill Mill's high profile, but small sales, products is "Easy Bounce", a cushioned sock that helps keep basketball players from bruising their feet as they come down from jumping. John McGrath gave Ann and Kevin the task of finding new markets for Easy Bounce socks. Ann and Kevin have decided that a good extension of this market is college volleyball players. Before they start, they want to learn about what size socks college volleyball players wear. First they need to gather some data, maybe by calling some equipment managers from nearby colleges to ask how many of what size volleyball socks were used last season. Then they will want to turn that data into information by arranging and summarizing their data, possibly even comparing the sizes of volleyball socks used at nearby colleges to the sizes of socks sold to basketball players.

Some definitions and important concepts

It may seem obvious, but a population is all of the members of a certain group. A sample is some of the members of the population. The same group of individuals may be a population in one context and a sample in another. The women in your stat class are the population of "women enrolled in this statistics class", and they are also a sample of "all students enrolled in this statistics class". It is important to be aware of what sample you are using to make an inference about what population.

What is statistics?

How exact is statistics? Upon close inspection, you will find that statistics is not all that exact; sometimes I have told my classes that statistics is "knowing when its close enough to call it equal". When making estimations, you will find that you are almost never exactly right. If you make the estimations using the correct method however, you will seldom be far from wrong. The same idea goes for hypothesis testing. You can never be sure that you've made the correct judgement, but if you conduct the hypothesis test with the correct method, you can be sure that the chance you've made the wrong judgement is small.

A term that needs to be defined is **probability**. Probability is a measure of the chance that something will occur. In statistics, when an inference is made, it is made with some probability that it is wrong (or some confidence that it is right). Think about repeating some action, like using a certain procedure to infer the mean of a population, over and over and over. Inevitably, sometimes the procedure will give a faulty estimate, sometimes you will be wrong. The probability that the procedure gives the wrong answer is simply the proportion of the times that the estimate is wrong. The confidence is simply the proportion of times that the answer is right. The probability of something happening is expressed as the proportion of the time that it can be expected to happen. Proportions are written as decimal fractions, and so are probabilities. If the probability that Foothill Hosiery's best salesperson will make the sale is .75, three-quarters of the time the sale is made.

Why bother with stat?

Reflect on what you have just read. What you are going to learn to do by learning statistics is to learn the right way to make educated guesses. For most students, statistics is not a favorite course. Its viewed as hard, or cosmic, or just plain confusing. By now, you should be thinking: "I could just skip stat, and avoid making inferences about what populations are like by always collecting data on the whole population and knowing for sure what the population is like." Well, many things come back to money, and its money that makes you take stat. Collecting data on a whole population is usually very expensive, and often almost impossible. If you can make a good, educated inference about a population from data collected from a small portion of that population, you will be able to save yourself, and your employer, a lot of time and money. You will also be able to make inferences about populations for which collecting data on the whole population is virtually impossible. Learning statistics now will allow you to save resources later and if the resources saved later are greater than the cost of learning statistics now, it will be worthwhile to learn statistics. It is my hope that the approach followed in this text will reduce the initial cost of learning statistics. If you have already had finance, you'll understand it this way—this approach to learning statistics will increase the net present value of investing in learning statistics by decreasing the initial cost.

Imagine how long it would take and how expensive it would be if Ann and Kevin decided that they had to find out what size sock every college volleyball player wore in order to see if volleyball players wore the same size socks as basketball players. By knowing how samples are related to populations, Ann and Kevin can quickly and inexpensively get a good idea of what size socks volleyball players wear, saving Foothill a lot of money and keeping John McGrath happy.

There are two basic types of inferences that can be made. The first is to estimate something about the population, usually its mean. The second is to see if the population has certain characteristics, for example you might want to infer if a population has a mean greater than 5.6. This second type of inference, hypothesis testing, is what we will concentrate on. If you understand hypothesis testing, estimation is easy. There are many applications,

especially in more advanced statistics, in which the difference between estimation and hypothesis testing seems blurred.

Estimation

Estimation is one of the basic inferential statistics techniques. The idea is simple; collect data from a sample and process it in some way that yields a good inference of something about the population. There are two types of estimates: point estimates and interval estimates. To make a point estimate, you simply find the single number that you think is your best guess of the characteristic of the population. As you can imagine, you will seldom be exactly correct, but if you make your estimate correctly, you will seldom be very far wrong. How to correctly make these estimates is an important part of statistics.

To make an interval estimate, you define an interval within which you believe the population characteristic lies. Generally, the wider the interval, the more confident you are that it contains the population characteristic. At one extreme, you have complete confidence that the mean of a population lies between $-\infty$ and $+\infty$ but that information has little value. At the other extreme, though you can feel comfortable that the population mean has a value close to that guessed by a correctly conducted point estimate, you have almost no confidence ("zero plus" to statisticians) that the population mean is exactly equal to the estimate. There is a trade-off between width of the interval, and confidence that it contains the population mean. How to find a narrow range with an acceptable level of confidence is another skill learned when learning statistics.

Hypothesis testing

The other type of inference is hypothesis testing. Though hypothesis testing and interval estimation use similar mathematics, they make quite different inferences about the population. Estimation makes no prior statement about the population; it is designed to make an educated guess about a population that you know nothing about. Hypothesis testing tests to see if the population has a certain characteristic—say a certain mean. This works by using statisticians' knowledge of how samples taken from populations with certain characteristics are likely to look to see if the sample you have is likely to have come from such a population.

A simple example is probably the best way to get to this. Statisticians know that if the means of a large number of samples of the same size taken from the same population are averaged together, the mean of those sample means equals the mean of the original population, and that most of those sample means will be fairly close to the population mean. If you have a sample that you suspect comes from a certain population, you can test the hypothesis that the population mean equals some number, m , by seeing if your sample has a mean close to m or not. If your sample has a mean close to m , you can comfortably say that your sample is likely to be one of the samples from a population with a mean of m .

Sampling

It is important to recognize that there is another cost to using statistics, even after you have learned statistics. As we said before, you are never sure that your inferences are correct. The more precise you want your inference to be, either the larger the sample you will have to collect (and the more time and money you'll have to spend on collecting it), or the greater the chance you must take that you'll make a mistake. Basically, if your sample is a good representation of the whole population—if it contains members from across the range of the population in proportions similar to that in the population—the inferences made will be good. If you manage to pick a sample that is not a good representation of the population, your inferences are likely to be wrong. By choosing samples

What is statistics?

carefully, you can increase the chance of a sample which is representative of the population, and increase the chance of an accurate inference.

The intuition behind this is easy. Imagine that you want to infer the mean of a population. The way to do this is to choose a sample, find the mean of that sample, and use that sample mean as your inference of the population mean. If your sample happened to include all, or almost all, observations with values that are at the high end of those in the population, your sample mean will overestimate the population mean. If your sample includes roughly equal numbers of observations with "high" and "low" and "middle" values, the mean of the sample will be close to the population mean, and the sample mean will provide a good inference of the population mean. If your sample includes mostly observations from the middle of the population, you will also get a good inference. Note that the sample mean will seldom be exactly equal to the population mean, however, because most samples will have a rough balance between high and low and middle values, the sample mean will usually be close to the true population mean. The key to good sampling is to avoid choosing the members of your sample in a manner that tends to choose too many "high" or too many "low" observations.

There are three basic ways to accomplish this goal. You can choose your sample randomly, you can choose a stratified sample, or you can choose a cluster sample. While there is no way to insure that a single sample will be representative, following the discipline of random, stratified, or cluster sampling greatly reduces the probability of choosing an unrepresentative sample.

The sampling distribution

The thing that makes statistics work is that statisticians have discovered how samples are related to populations. This means that statisticians (and, by the end of the course, you) know that if all of the possible samples from a population are taken and something (generically called a "statistic") is computed for each sample, something is known about how the new population of statistics computed from each sample is related to the original population. For example, if all of the samples of a given size are taken from a population, the mean of each sample is computed, and then the mean of those sample means is found, statisticians know that the mean of the sample means is equal to the mean of the original population.

There are many possible sampling distributions. Many different statistics can be computed from the samples, and each different original population will generate a different set of samples. The amazing thing, and the thing that makes it possible to make inferences about populations from samples, is that there are a few statistics which all have about the same sampling distribution when computed from the samples from many different populations.

You are probably still a little confused about what a sampling distribution is. It will be discussed more in the chapter on the Normal and t-distributions. An example here will help. Imagine that you have a population—the sock sizes of all of the volleyball players in the South Atlantic Conference. You take a sample of a certain size, say six, and find the mean of that sample. Then take another sample of six sock sizes, and find the mean of that sample. Keep taking different samples until you've found the mean of all of the possible samples of six. You will have generated a new population, the population of sample means. This population is the sampling distribution. Because statisticians often can find what proportion of members of this new population will take on certain values if they know certain things about the original population, we will be able to make certain inferences about the original population from a single sample.

Univariate and multivariate statistics statistics and the idea of an observation.

A population may include just one thing about every member of a group, or it may include two or more things about every member. In either case there will be one observation for each group member. Univariate statistics are concerned with making inferences about one variable populations, like "what is the mean shoe size of business students?" Multivariate statistics is concerned with making inferences about the way that two or more variables are connected in the population like, "do students with high grade point averages usually have big feet?" What's important about multivariate statistics is that it allows you to make better predictions. If you had to predict the shoe size of a business student and you had found out that students with high grade point averages usually have big feet, knowing the student's grade point average might help. Multivariate statistics are powerful and find applications in economics, finance, and cost accounting.

Ann Howard and Kevin Schmidt might use multivariate statistics if Mr McGrath asked them to study the effects of radio advertising on sock sales. They could collect a multivariate sample by collecting two variables from each of a number of cities—recent changes in sales and the amount spent on radio ads. By using multivariate techniques you will learn in later chapters, Ann and Kevin can see if more radio advertising means more sock sales.

Conclusion

As you can see, there is a lot of ground to cover by the end of this course. There are a few ideas that tie most of what you learn together: populations and samples, the difference between data and information, and most important, sampling distributions. We'll start out with the easiest part, descriptive statistics, turning data into information. Your professor will probably skip some chapters, or do a chapter toward the end of the book before one that's earlier in the book. As long as you cover the chapters "Descriptive Statistics and frequency distributions", "The normal and the t-distributions", "Making estimates" and that is alright.

You should learn more than just statistics by the time the semester is over. Statistics is fairly difficult, largely because understanding what is going on requires that you learn to stand back and think about things; you cannot memorize it all, you have to figure out much of it. This will help you learn to use statistics, not just learn statistics for its own sake.

You will do much better if you attend class regularly and if you read each chapter at least three times. First, the day before you are going to discuss a topic in class, read the chapter carefully, but do not worry if you understand everything. Second, soon after a topic has been covered in class, read the chapter again, this time going slowly, making sure you can see what is going on. Finally, read it again before the exam. Though this is a great statistics book, the stuff is hard, and no one understands statistics the first time.

1. Descriptive statistics and frequency distributions

This chapter is about describing populations and samples, a subject known as descriptive statistics. This will all make more sense if you keep in mind that the information you want to produce is a description of the population or sample as a whole, not a description of one member of the population. The first topic in this chapter is a discussion of "distributions", essentially pictures of populations (or samples). Second will be the discussion of descriptive statistics. The topics are arranged in this order because the descriptive statistics can be thought of as ways to describe the picture of a population, the distribution.

Distributions

The first step in turning data into information is to create a distribution. The most primitive way to present a distribution is to simply list, in one column, each value that occurs in the population and, in the next column, the number of times it occurs. It is customary to list the values from lowest to highest. This simple listing is called a "frequency distribution". A more elegant way to turn data into information is to draw a graph of the distribution. Customarily, the values that occur are put along the horizontal axis and the frequency of the value is on the vertical axis.

Ann Howard called the equipment manager at two nearby colleges and found out the following data on sock sizes used by volleyball players. At Piedmont State last year, 14 pairs of size 7 socks, 18 pairs of size 8, 15 pairs of size 9, and 6 pairs of size 10 socks were used. At Graham College, the volleyball team used 3 pairs of size 6, 10 pairs of size 7, 15 pairs of size 8, 5 pairs of size 9, and 11 pairs of size 10. Ann arranged her data into a distribution and then drew a graph called a Histogram:

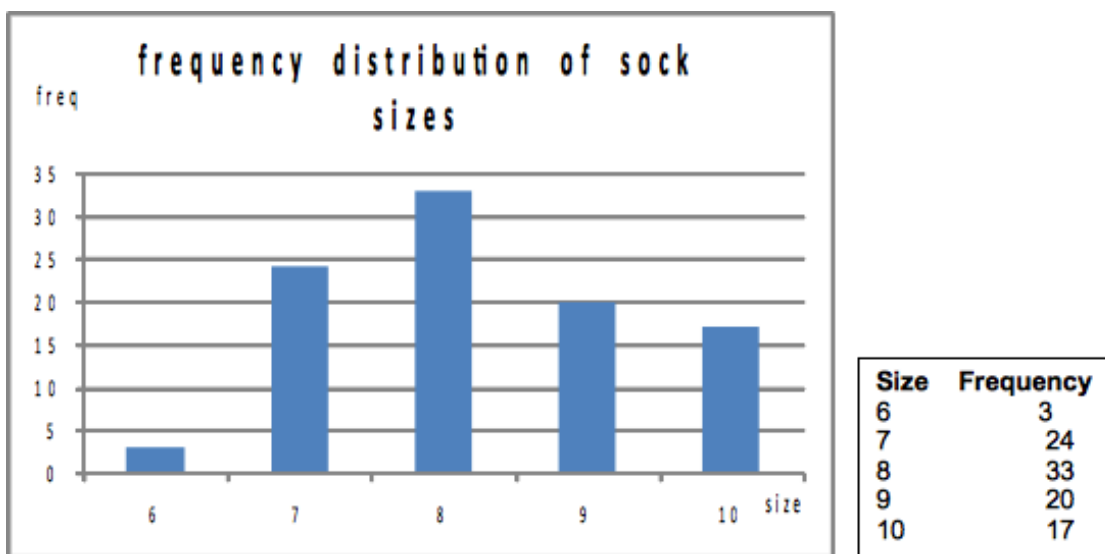


Exhibit 1: Frequency graph of sock sizes

1. Descriptive statistics and frequency distributions

Ann could have created a relative frequency distribution as well as a frequency distribution. The difference is that instead of listing how many times each value occurred, Ann would list what proportion of her sample was made up of socks of each size:

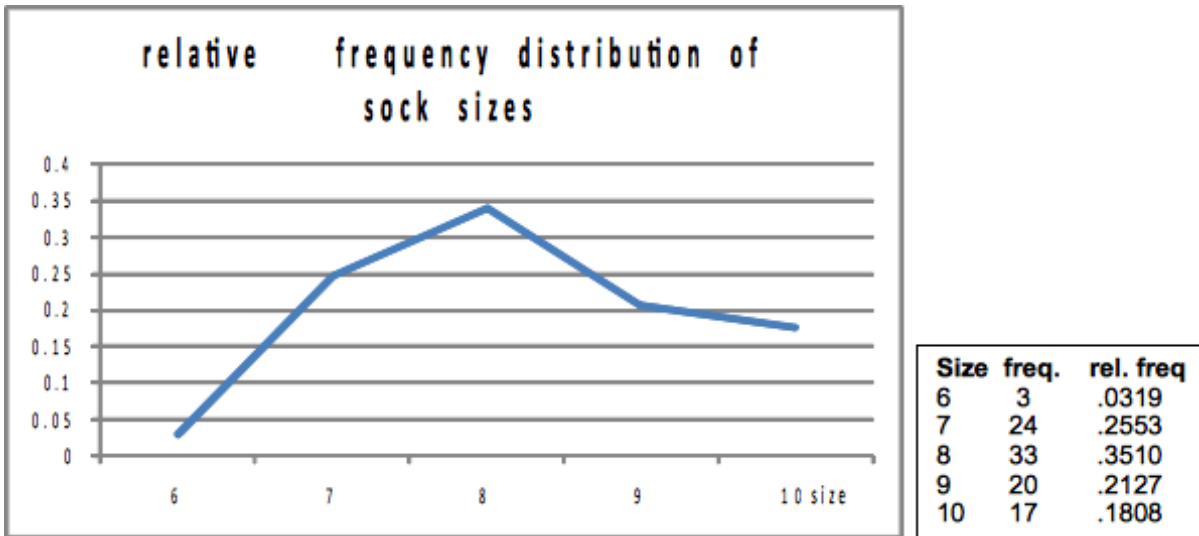


Exhibit 2: Relative frequency graph of sock sizes

Notice that Ann has drawn the graphs differently. In the first graph, she has used bars for each value, while on the second, she has drawn a point for the relative frequency of each size, and then "connected the dots". While both methods are correct, when you have values that are continuous, you will want to do something more like the "connect the dots" graph. Sock sizes are **discrete**, they only take on a limited number of values. Other things have **continuous** values, they can take on an infinite number of values, though we are often in the habit of rounding them off. An example is how much students weigh. While we usually give our weight in whole pounds in the US ("I weigh 156 pounds."), few have a weight that is exactly so many pounds. When you say "I weigh 156", you actually mean that you weigh between 155 1/2 and 156 1/2 pounds. We are heading toward a graph of a distribution of a continuous variable where the relative frequency of any **exact** value is very small, but the relative frequency of observations between two values is measurable. What we want to do is to get used to the idea that the total area under a "connect the dots" relative frequency graph, from the lowest to the highest possible value is one. Then the part of the area under the graph between two values is the relative frequency of observations with values within that range. The height of the line above any particular value has lost any direct meaning, because it is now the area under the line between two values that is the relative frequency of an observation between those two values occurring.

You can get some idea of how this works if you go back to the bar graph of the distribution of sock sizes, but draw it with relative frequency on the vertical axis. If you arbitrarily decide that each bar has a width of one, then the area "under the curve" between 7.5 and 8.5 is simply the height times the width of the bar for sock size 8: 0.3510×1 . If you wanted to find the relative frequency of sock sizes between 6.5 and 8.5, you could simply add together the area of the bar for size 7 (that's between 6.5 and 7.5) and the bar for size 8 (between 7.5 and 8.5).

Descriptive statistics

Now that you see how a distribution is created, you are ready to learn how to describe one. There are two main things that need to be described about a distribution: its location and its shape. Generally, it is best to give a single measure as the description of the location and a single measure as the description of the shape.

Mean

To describe the location of a distribution, statisticians use a "typical" value from the distribution. There are a number of different ways to find the typical value, but by far the most used is the "arithmetic mean", usually simply called the "mean". You already know how to find the arithmetic mean, you are just used to calling it the "average". Statisticians use average more generally—the arithmetic mean is one of a number of different averages. Look at the formula for the arithmetic mean:

$$\mu = \frac{\sum x}{N}$$

All you do is add up all of the members of the population, $\sum x$, and divide by how many members there are, N . The only trick is to remember that if there is more than one member of the population with a certain value, to add that value once for every member that has it. To reflect this, the equation for the mean sometimes is written :

$$\mu = \frac{\sum f_i x_i}{N}$$

where f_i is the frequency of members of the population with the value x_i .

This is really the same formula as above. If there are seven members with a value of ten, the first formula would have you add seven ten times. The second formula simply has you multiply seven by ten—the same thing as adding together ten sevens.

Other measures of location are the median and the mode. The median is the value of the member of the population that is in the middle when the members are sorted from smallest to largest. Half of the members of the population have values higher than the median, and half have values lower. The median is a better measure of location if there are one or two members of the population that are a lot larger (or a lot smaller) than all the rest. Such extreme values can make the mean a poor measure of location, while they have little effect on the median. If there are an odd number of members of the population, there is no problem finding which member has the median value. If there are an even number of members of the population, then there is no single member in the middle. In that case, just average together the values of the two members that share the middle.

The third common measure of location is the mode. If you have arranged the population into a frequency or relative frequency distribution, the mode is easy to find because it is the value that occurs most often. While in some sense, the mode is really the most typical member of the population, it is often not very near the middle of the population. You can also have multiple modes. I am sure you have heard someone say that "it was a bimodal distribution". That simply means that there were two modes, two values that occurred equally most often.

If you think about it, you should not be surprised to learn that for bell-shaped distributions, the mean, median, and mode will be equal. Most of what statisticians do with the describing or inferring the location of a population is done with the mean. Another thing to think about is using a spreadsheet program, like Microsoft Excel when arranging data into a frequency distribution or when finding the median or mode. By using the sort and

1. Descriptive statistics and frequency distributions

distribution commands in 1-2-3, or similar commands in Excel, data can quickly be arranged in order or placed into value classes and the number in each class found. Excel also has a function, =AVERAGE(...), for finding the arithmetic mean. You can also have the spreadsheet program draw your frequency or relative frequency distribution.

One of the reasons that the arithmetic mean is the most used measure of location is because the mean of a sample is an "unbiased estimator" of the population mean. Because the sample mean is an unbiased estimator of the population mean, the sample mean is a good way to make an inference about the population mean. If you have a sample from a population, and you want to guess what the mean of that population is, you can legitimately guess that the population mean is equal to the mean of your sample. This is a legitimate way to make this inference because the mean of all the sample means equals the mean of the population, so, if you used this method many times to infer the population mean, on average you'd be correct.

All of these measures of location can be found for samples as well as populations, using the same formulas. Generally, μ is used for a population mean, and \bar{x} is used for sample means. Upper-case N, really a Greek "nu", is used for the size of a population, while lower case n is used for sample size. Though it is not universal, statisticians tend to use the Greek alphabet for population characteristics and the Roman alphabet for sample characteristics.

Measuring population shape

Measuring the shape of a distribution is more difficult. Location has only one dimension ("where?"), but shape has a lot of dimensions. We will talk about two, and you will find that most of the time, only one dimension of shape is measured. The two dimensions of shape discussed here are the width and symmetry of the distribution. The simplest way to measure the width is to do just that—the range in the distance between the lowest and highest members of the population. The range is obviously affected by one or two population members which are much higher or lower than all the rest.

The most common measures of distribution width are the standard deviation and the variance. The standard deviation is simply the square root of the variance, so if you know one (and have a calculator that does squares and square roots) you know the other. The standard deviation is just a strange measure of the mean distance between the members of a population and the mean of the population. This is easiest to see if you start out by looking at the formula for the variance:

$$\sigma^2 = \frac{\sum(x - \mu)^2}{N}$$

Look at the numerator. To find the variance, the first step (after you have the mean, μ) is to take each member of the population, and find the difference between its value and the mean; you should have N differences. Square each of those, and add them together, dividing the sum by N, the number of members of the population. Since you find the mean of a group of things by adding them together and then dividing by the number in the group, the variance is simply the "mean of the squared distances between members of the population and the population mean".

Notice that this is the formula for a population characteristic, so we use the Greek σ and that we write the variance as σ^2 , or "sigma square" because the standard deviation is simply the square root of the variance, its symbol is simply "sigma", σ .

One of the things statisticians have discovered is that 75 per cent of the members of any population are within two standard deviations of the mean of the population. This is known as Chebyshev's Theorem. If the mean of a

population of shoe sizes is 9.6 and the standard deviation is 1.1, then 75 per cent of the shoe sizes are between 7.4 (two standard deviations below the mean) and 11.8 (two standard deviations above the mean). This same theorem can be stated in probability terms: the probability that anything is within two standard deviations of the mean of its population is .75.

It is important to be careful when dealing with variances and standard deviations. In later chapters, there are formulas using the variance, and formulas using the standard deviation. Be sure you know which one you are supposed to be using. Here again, spreadsheet programs will figure out the standard deviation for you. In Excel, there is a function, =STDEVP(...), that does all of the arithmetic. Most calculators will also compute the standard deviation. Read the little instruction booklet, and find out how to have your calculator do the numbers before you do any homework or have a test.

The other measure of shape we will discuss here is the measure of "skewness". Skewness is simply a measure of whether or not the distribution is symmetric or if it has a long tail on one side, but not the other. There are a number of ways to measure skewness, with many of the measures based on a formula much like the variance. The formula looks a lot like that for the variance, except the distances between the members and the population mean are cubed, rather than squared, before they are added together:

$$sk = \frac{\sum(x - \mu)^3}{N}$$

At first it might not seem that cubing rather than squaring those distances would make much difference. Remember, however, that when you square either a positive or negative number you get a positive number, but that when you cube a positive, you get a positive and when you cube a negative you get a negative. Also remember that when you square a number, it gets larger, but that when you cube a number, it gets a whole lot larger. Think about a distribution with a long tail out to the left. There are a few members of that population much smaller than the mean, members for which $(x - \mu)$ is large and negative. When these are cubed, you end up with some really big negative numbers. Because there are not any members with such large, positive $(x - \mu)$, there are not any corresponding really big positive numbers to add in when you sum up the $(x - \mu)^3$, and the sum will be negative. A negative measure of skewness means that there is a tail out to the left, a positive measure means a tail to the right. Take a minute and convince yourself that if the distribution is symmetric, with equal tails on the left and right, the measure of skew is zero.

To be really complete, there is one more thing to measure, "kurtosis" or "peakedness". As you might expect by now, it is measured by taking the distances between the members and the mean and raising them to the fourth power before averaging them together.

Measuring sample shape

Measuring the location of a sample is done in exactly the way that the location of a population is done. Measuring the shape of a sample is done a little differently than measuring the shape of a population, however. The reason behind the difference is the desire to have the sample measurement serve as an unbiased estimator of the population measurement. If we took all of the possible samples of a certain size, n , from a population and found the variance of each one, and then found the mean of those sample variances, that mean would be a little smaller than the variance of the population.

1. Descriptive statistics and frequency distributions

You can see why this is so if you think it through. If you knew the population mean, you could find $\sum (x - \mu)^2 / n$ for each sample, and have an unbiased estimate for σ^2 . However, you do not know the population mean, so you will have to infer it. The best way to infer the population mean is to use the sample mean \bar{x} . The variance of a sample will then be found by averaging together all of the $\sum (x - \bar{x})^2 / n$.

The mean of a sample is obviously determined by where the members of that sample lie. If you have a sample that is mostly from the high (or right) side of a population's distribution, then the sample mean will almost for sure be greater than the population mean. For such a sample, $\sum (x - \bar{x})^2 / n$ would underestimate σ^2 . The same is true for samples that are mostly from the low (or left) side of the population. If you think about what kind of samples will have $\sum (x - \bar{x})^2 / n$ that is greater than the population σ^2 , you will come to the realization that it is only those samples with a few very high members and a few very low members—and there are not very many samples like that. By now you should have convinced yourself that $\sum (x - \bar{x})^2 / n$ will result in a biased estimate of σ^2 . You can see that, on average, it is too small.

How can an unbiased estimate of the population variance, σ^2 , be found? If $\sum (x - \bar{x})^2 / n$ on average too small, we need to do something to make it a little bigger. We want to keep the $\sum (x - \bar{x})^2$, but if we divide it by something a little smaller, the result will be a little larger. Statisticians have found out that the following way to compute the sample variance results in an unbiased estimator of the population variance:

$$s^2 = \frac{\sum (x - \bar{x})^2}{n - 1}$$

If we took all of the possible samples of some size, n , from a population, and found the sample variance for each of those samples, using this formula, the mean of those sample variances would equal the population variance, σ^2 .

Note that we use s^2 instead of σ^2 , and n instead of N (really "nu", not "en") since this is for a sample and we want to use the Roman letters rather than the Greek letters, which are used for populations.

There is another way to see why you divide by $n-1$. We also have to address something called "degrees of freedom" before too long, and it is the degrees of freedom that is the key of the other explanation. As we go through this explanation, you should be able to see that the two explanations are related.

Imagine that you have a sample with 10 members ($n=10$), and you want to use it to estimate the variance of the population from which it was drawn. You write each of the 10 values on a separate scrap of paper. If you know the population mean, you could start by computing all 10 $(x - \mu)^2$. In the usual case, you do not know μ , however, and you must start by finding \bar{x} from the values on the 10 scraps to use as an estimate of μ . Once you have found \bar{x} , you could lose any one of the 10 scraps and still be able to find the value that was on the lost scrap from the other 9 scraps. If you are going to use \bar{x} in the formula for sample variance, only 9 (or $n-1$), of the x 's are free to take on any value. Because only $n-1$ of the x 's can vary freely, you should divide $\sum (x - \bar{x})^2$ by $n-1$, the number of (x 's) that are really free. Once you use \bar{x} in the formula for sample variance, you use up one "degree of freedom", leaving only $n-1$. Generally, whenever you use something you have previously computed from a sample within a formula, you use up a degree of freedom.

A little thought will link the two explanations. The first explanation is based on the idea that \bar{x} , the estimator of μ , varies with the sample. It is because \bar{x} varies with the sample that a degree of freedom is used up in the second explanation.

The sample standard deviation is found simply by taking the square root of the sample variance:

$$s = \sqrt{[\sum(x - \bar{x})^2 / (n - 1)]}$$

While the sample variance is an unbiased estimator of population variance, the sample standard deviation is not an unbiased estimator of the population standard deviation—the square root of the average is not the same as the average of the square roots. This causes statisticians to use variance where it seems as though they are trying to get at standard deviation. In general, statisticians tend to use variance more than standard deviation. Be careful with formulas using sample variance and standard deviation in the following chapters. Make sure you are using the right one. Also note that many calculators will find standard deviation using both the population and sample formulas. Some use σ and s to show the difference between population and sample formulas, some use s_n and s_{n-1} to show the difference.

If Ann Howard wanted to infer what the population distribution of volleyball players' sock sizes looked like she could do so from her sample. If she is going to send volleyball coaches packages of socks for the players to try, she will want to have the packages contain an assortment of sizes that will allow each player to have a pair that fits. Ann wants to infer what the distribution of volleyball players sock sizes looks like. She wants to know the mean and variance of that distribution. Her data, again, is:

size	frequency
6	3
7	24
8	33
9	20
10	17

The mean sock size can be found:

$$= [(3 \times 6) + (24 \times 7) + (33 \times 8) + (20 \times 9) + (17 \times 10)] / 97 = 8.25.$$

To find the sample standard deviation, Ann decides to use Excel. She lists the sock sizes that were in the sample in column A, and the frequency of each of those sizes in column B. For column C, she has the computer find for each of $\sum(x - \bar{x})^2$ the sock sizes, using the formula $= (A1 - 8.25)^2$ in the first row, and then copying it down to the other four rows. In D1, she multiplies C1, by the frequency using the formula $= B1 * C1$, and copying it down into the other rows. Finally, she finds the sample standard deviation by adding up the five numbers in column D and dividing by $n - 1 = 96$ using the Excel formula $= \text{sum}(D1:D5) / 96$. The spreadsheet appears like this when she is done:

A	B	C	D	E
1	6	3	5.06	15.19

1. Descriptive statistics and frequency distributions

2	7	24	1.56	37.5
3	8	33	0.06	2.06
4	9	20	0.56	11.25
5	10	17	3.06	52.06
6	n=	97		Var = 1.217139
7				Std.dev = 1.103.24
8				

Ann now has an estimate of the variance of the sizes of socks worn by college volleyball players, 1.22. She has inferred that the population of college volleyball players' sock sizes has a mean of 8.25 and a variance of 1.22.

Summary

To describe a population you need to describe the picture or graph of its distribution. The two things that need to be described about the distribution are its location and its shape. Location is measured by an average, most often the arithmetic mean. The most important measure of shape is a measure of dispersion, roughly width, most often the variance or its square root the standard deviation.

Samples need to be described, too. If all we wanted to do with sample descriptions was describe the sample, we could use exactly the same measures for sample location and dispersion that are used for populations. We want to use the sample describers for dual purposes, however: (a) to describe the sample, and (b) to make inferences about the description of the population that sample came from. Because we want to use them to make inferences, we want our sample descriptions to be "unbiased estimators". Our desire to measure sample dispersion with an unbiased estimator of population dispersion means that the formula we use for computing sample variance is a little difference than the used for computing population variance.

2. The normal and t-distributions

The normal distribution is simply a distribution with a certain shape. It is "normal" because many things have this same shape. The normal distribution is the "bell-shaped distribution" that describes how so many natural, machine-made, or human performance outcomes are distributed. If you ever took a class when you were "graded on a bell curve", the instructor was fitting the class' grades into a normal distribution—not a bad practice if the class is large and the tests are objective, since human performance in such situations is normally distributed. This chapter will discuss the normal distribution and then move onto a common sampling distribution, the t-distribution. The t-distribution can be formed by taking many samples (strictly, all possible samples) of the same size from a normal population. For each sample, the same statistic, called the t-statistic, which we will learn more about later, is calculated. The relative frequency distribution of these t-statistics is the t-distribution. It turns out that t-statistics can be computed a number of different ways on samples drawn in a number of different situations and still have the same relative frequency distribution. This makes the t-distribution useful for making many different inferences, so it is one of the most important links between samples and populations used by statisticians. In between discussing the normal and t-distributions, we will discuss the central limit theorem. The t-distribution and the central limit theorem give us knowledge about the relationship between sample means and population means that allows us to make inferences about the population mean.

The way the t-distribution is used to make inferences about populations from samples is the model for many of the inferences that statisticians make. Since you will be learning to make inferences like a statistician, try to understand the general model of inference making as well as the specific cases presented. Briefly, the general model of inference-making is to use statisticians' knowledge of a sampling distribution like the t-distribution as a guide to the probable limits of where the sample lies relative to the population. Remember that the sample you are using to make an inference about the population is only one of many possible samples from the population. The samples will vary, some being highly representative of the population, most being fairly representative, and a few not being very representative at all. By assuming that the sample is at least fairly representative of the population, the sampling distribution can be used as a link between the sample and the population so you can make an inference about some characteristic of the population.

These ideas will be developed more later on. The immediate goal of this chapter is to introduce you to the normal distribution, the central limit theorem, and the t-distribution.

Normal things

Normal distributions are bell-shaped and symmetric. The mean, median, and mode are equal. Most of the members of a normally distributed population have values close to the mean—in a normal population 96 per cent of the members (much better than Chebyshev's 75 per cent), are within 2σ of the mean.

2. The normal and t-distributions

Statisticians have found that many things are normally distributed. In nature, the weights, lengths, and thicknesses of all sorts of plants and animals are normally distributed. In manufacturing, the diameter, weight, strength, and many other characteristics of man- or machine-made items are normally distributed. In human performance, scores on objective tests, the outcomes of many athletic exercises, and college student grade point averages are normally distributed. The normal distribution really is a normal occurrence.

If you are a skeptic, you are wondering how can GPAs and the exact diameter of holes drilled by some machine have the same distribution—they are not even measured with the same units. In order to see that so many things have the same normal shape, all must be measured in the same units (or have the units eliminated)—they must all be "standardized." Statisticians standardize many measures by using the STANDARD deviation. All normal distributions have the same shape because they all have the same relative frequency distribution *when the values for their members are measured in standard deviations above or below the mean.*

Using the United States customary system of measurement, if the weight of pet cats is normally distributed with a mean of 10.8 pounds and a standard deviation of 2.3 pounds and the daily sales at The First Brew Espresso Cafe are normally distributed with $\mu = \$341.46$ and $\sigma = \$53.21$, then the same proportion of pet cats weigh between 8.5 pounds ($\mu - 1\sigma$) and 10.8 pounds (μ) as the proportion of daily First Brew sales which lie between $\mu - 1\sigma$ (\$288.25) and μ (\$341.46). Any normally distributed population will have the same proportion of its members between the mean and one standard deviation below the mean. Converting the values of the members of a normal population so that each is now expressed in terms of standard deviations from the mean makes the populations all the same. This process is known as "standardization" and it makes all normal populations have the same location and shape.

This standardization process is accomplished by computing a "z-score" for every member of the normal population. The z-score is found by:

$$z = (x - \mu) / \sigma$$

This converts the original value, in its original units, into a standardized value in units of "standard deviations from the mean." Look at the formula. The numerator is simply the difference between the value of this member of the population, x , and the mean of the population μ . It can be measured in centimeters, or points, or whatever. The denominator is the standard deviation of the population, σ , and it is also measured in centimeters, or points, or whatever. If the numerator is 15cm and the standard deviation is 10cm, then the z will be 1.5. This particular member of the population, one with a diameter 15cm greater than the mean diameter of the population, has a z -value of 1.5 because its value is 1.5 standard deviations greater than the mean. Because the mean of the x 's is μ , the mean of the z -scores is zero.

We could convert the value of every member of *any* normal population into a z -score. If we did that for any normal population and arranged those z -scores into a relative frequency distribution, they would all be the same. Each and every one of those standardized normal distributions would have a mean of zero and the same shape. There are many tables which show what proportion of any normal population will have a z -score less than a certain value. Because the standard normal distribution is symmetric with a mean of zero, the same proportion of the population that is less than some positive z is also greater than the same negative z . Some values from a "standard normal" table appear below:

Proportion below	.75	.90	.95	.975	.99	.995
-------------------------	-----	-----	-----	------	-----	------

z-score	0.674	1.282	1.645	1.960	2.326	2.576
----------------	-------	-------	-------	-------	-------	-------

John McGrath has asked Kevin Schmidt "How much does a pair of size 11 mens dress socks usually weigh?" Kevin asks the people in quality control what they know about the weight of these socks and is told that the mean weight is 4.25 ounces with a standard deviation of .021 ounces. Kevin decides that Mr. McGrath probably wants more than the mean weight, and decides to give his boss the range of weights within which 95% of size 11 men's dress socks falls. Kevin sees that leaving 2.5% (.025) in the left tail and 2.5% (.025) in the right tail will leave 95% (.95) in the middle. He assumes that sock weights are normally distributed, a reasonable assumption for a machine-made product, and consulting a standard normal table, sees that .975 of the members of any normal population have a z-score less than 1.96 and that .975 have a z-score greater than -1.96, so .95 have a z-score between ± 1.96 .

Now that he knows that 95% of the socks will have a weight with a z-score between ± 1.96 , Kevin can translate those z's into ounces. By solving the equation for both +1.96 and -1.96, he will find the boundaries of the interval within which 95% of the weights of the socks fall:

$$1.96 = (x - 4.25)/.021$$

solving for x, Kevin finds that the upper limit is 4.29 ounces. He then solves for $z = -1.96$:

$$-1.96 = (x - 4.25)/.021$$

and finds that the lower limit is 4.21 ounces. He can now go to John McGrath and tell him: "95% of size 11 mens' dress socks weigh between 4.21 and 4.29 ounces."

The central limit theorem

If this was a statistics course for math majors, you would probably have to prove this theorem. Because this text is designed for business and other non-math students, you will only have to learn to understand what the theorem says and why it is important. To understand what it says, it helps to understand why it works. Here is an explanation of why it works.

The theorem is about sampling distributions and the relationship between the location and shape of a population and the location and shape of a sampling distribution generated from that population. Specifically, the central limit theorem explains the relationship between a population and the distribution of sample means found by taking all of the possible samples of a certain size from the original population, finding the mean of each sample, and arranging them into a distribution.

The sampling distribution of means is an easy concept. Assume that you have a population of x's. You take a sample of n of those x's and find the mean of that sample, giving you one \bar{x} . Then take another sample of the same size, n, and find its \bar{x} ...Do this over and over until you have chosen all possible samples of size n. You will have generated a new population, a population of \bar{x} 's. Arrange this population into a distribution, and you have the sampling distribution of means. You could find the sampling distribution of medians, or variances, or some other sample statistic by collecting all of the possible samples of some size, n, finding the median, variance, or other statistic about each sample, and arranging them into a distribution.

The central limit theorem is about the sampling distribution of means. It links the sampling distribution of \bar{x} 's with the original distribution of x's. It tells us that:

2. The normal and t-distributions

(1) The mean of the sample means equals the mean of the original population, $\mu_{\bar{x}} = \mu$. This is what makes \bar{x} an unbiased estimator of μ .

(2) The distribution of \bar{x} 's will be bell-shaped, no matter what the shape of the original distribution of x's.

This makes sense when you stop and think about it. It means that only a small portion of the samples have means that are far from the population mean. For a sample to have a mean that is far from μ_x , almost all of its members have to be from the right tail of the distribution of x's, or almost all have to be from the left tail. There are many more samples with most of their members from the middle of the distribution, or with some members from the right tail and some from the left tail, and all of those samples will have an \bar{x} close to μ_x .

(3a) The larger the samples, the closer the sampling distribution will be to normal, and

(3b) if the distribution of x's is normal, so is the distribution of \bar{x} 's.

These come from the same basic reasoning as 2), but would require a formal proof since "normal distribution" is a mathematical concept. It is not too hard to see that larger samples will generate a "more-bell-shaped" distribution of sample means than smaller samples, and that is what makes 3a) work.

(4) The variance of the \bar{x} 's is equal to the variance of the x's divided by the sample size, or:

$$\sigma_{\bar{x}}^2 = \sigma^2 / n$$

therefore the standard deviation of the sampling distribution is:

$$\sigma_{\bar{x}} = \sigma / \sqrt{n}$$

While it is a difficult to see why this exact formula holds without going through a formal proof, the basic idea that larger samples yield sampling distributions with smaller standard deviations can be understood intuitively. If $\sigma_{\bar{x}} = \sigma_x / \sqrt{n}$ then $\sigma_{\bar{x}} < \sigma_x$. Furthermore, when the sample size, n, rises, $\sigma_{\bar{x}}^2$ gets smaller. This is because it becomes more unusual to get a sample with an \bar{x} that is far from μ as n gets larger. The standard deviation of the sampling distribution includes an $(\bar{x} - \mu)$ for each, but remember that there are not many \bar{x} 's that are as far from μ as there are x's that are far from μ , and as n grows there are fewer and fewer samples with an \bar{x} far from μ . This means that there are not many $(\bar{x} - \mu)$ that are as large as quite a few $(x - \mu)$ are. By the time you square everything, the average $(\bar{x} - \mu)^2$ is going to be much smaller than the average $(x - \mu)^2$, so, $\sigma_{\bar{x}}$ is going to be smaller than σ_x . If the mean volume of soft drink in a population of 12 ounce cans is 12.05 ounces with a variance of .04 (and a standard deviation of .2), then the sampling distribution of means of samples of 9 cans will have a mean of 12.05 ounces and a variance of .04/9=.0044 (and a standard deviation of .2/3=.0667).

You can follow this same line of reasoning once again, and see that as the sample size gets larger, the variance and standard deviation of the sampling distribution will get smaller. Just remember that as sample size grows, samples with an \bar{x} that is far from μ get rarer and rarer, so that the average $(\bar{x} - \mu)^2$ will get smaller. The average $(\bar{x} - \mu)^2$ is the variance. If larger samples of soft drink bottles are taken, say samples of 16, even fewer of the samples will have means that are very far from the mean of 12.05 ounces. The variance of the sampling distribution when n=16 will therefore be smaller. According to what you have just learned, the variance will be only .04/16=.0025 (and the standard deviation will be .2/4=.05). The formula matches what logically is happening; as the samples get bigger, the probability of getting a sample with a mean that is far away from the population mean

gets smaller, so the sampling distribution of means gets narrower and the variance (and standard deviation) get smaller. In the formula, you divide the population variance by the sample size to get the sampling distribution variance. Since bigger samples means dividing by a bigger number, the variance falls as sample size rises. If you are using the sample mean as to infer the population mean, using a bigger sample will increase the probability that your inference is very close to correct because more of the sample means are very close to the population mean.. There is obviously a trade-off here. The reason you wanted to use statistics in the first place was to avoid having to go to the bother and expense of collecting lots of data, but if you collect more data, your statistics will probably be more accurate.

The t-distribution

The central limit theorem tells us about the relationship between the sampling distribution of means and the original population. Notice that if we want to know the variance of the sampling distribution we need to know the variance of the original population. You do not need to know the variance of the sampling distribution to make a point estimate of the mean, but other, more elaborate, estimation techniques require that you either know or estimate the variance of the population. If you reflect for a moment, you will realize that it would be strange to know the variance of the population when you do not know the mean. Since you need to know the population mean to calculate the population variance and standard deviation, the only time when you would know the population variance without the population mean are examples and problems in textbooks. The usual case occurs when you have to estimate both the population variance and mean. Statisticians have figured out how to handle these cases by using the sample variance as an estimate of the population variance (and being able to use that to estimate the variance of the sampling distribution). Remember that s^2 is an unbiased estimator of σ^2 . Remember, too, that the variance of the sampling distribution of means is related to the variance of the original population according to the equation:

$$\sigma_{\bar{x}}^2 = \sigma^2 / n$$

so, the estimated standard deviation of a sampling distribution of means is:

$$\text{estimated } \sigma_{\bar{x}} = s / \sqrt{n}$$

Following this thought, statisticians found that if they took samples of a constant size from a normal population, computed a statistic called a "t-score" for each sample, and put those into a relative frequency distribution, the distribution would be the same for samples of the same size drawn from any normal population. The shape of this sampling distribution of t's varies somewhat as sample size varies, but for any n it's always the same. For example, for samples of 5, 90% of the samples have t-scores between -1.943 and +1.943, while for samples of 15, 90% have t-scores between ± 1.761 . The bigger the samples, the narrower the range of scores that covers any particular proportion of the samples. That t-score is computed by the formula:

$$t = (\bar{x} - \mu) / (s/\sqrt{n})$$


By comparing the formula for the t-score with the formula for the z-score, you will be able to see that the t is just an estimated z. Since there is one t-score for each sample, the t is just another sampling distribution. It turns out that there are other things that can be computed from a sample that have the same distribution as this t. Notice that we've used the sample standard deviation, s, in computing each t-score. Since we've used s, we've used up one degree of freedom. Because there are other useful sampling distributions that have this same shape, but use up various numbers of degrees of freedom, it is the usual practice to refer to the t-distribution not as the distribution

2. The normal and t-distributions

for a particular sample size, but as the distribution for a particular number of degrees of freedom. There are published tables showing the shapes of the t-distributions, and they are arranged by degrees of freedom so that they can be used in all situations.

Looking at the formula, you can see that the mean t-score will be zero since the mean \bar{x} equals μ . Each t-distribution is symmetric, with half of the t-scores being positive and half negative because we know from the central limit theorem that the sampling distribution of means is normal, and therefore symmetric, when the original population is normal.

An excerpt from a typical t-table is printed below. Note that there is one line each for various degrees of freedom. Across the top are the proportions of the distributions that will be left out in the tail--the amount shaded in the picture. The body of the table shows which t-score divides the bulk of the distribution of t's for that df from the area shaded in the tail, which t-score leaves that proportion of t's to its right. For example, if you chose all of the possible samples with 9 df, and found the t-score for each, .025 (2 1/2 %) of those samples would have t-scores greater than 2.262, and .975 would have t-scores less than 2.262.



df	prob = .10	prob. = .05	prob. = .025	prob. = .01	prob. = .005
1	3.078	6.314	12.70	13.81	63.65
5	1.476	2.015	2.571	3.365	4.032
6	1.440	1.943	2.447	3.143	3.707
7	1.415	1.895	2.365	2.998	3.499
8	1.397	1.860	2.306	2.896	3.355
9	1.383	1.833	2.262	2.821	3.250
10	1.372	1.812	2.228	2.764	3.169
20	1.325	1.725	2.086	2.528	2.845
30	1.310	1.697	2.046	2.457	2.750
40	1.303	1.684	2.021	2.423	2.704
Infinity	1.282	1.645	1.960	2.326	2.58

Exhibit 3: A sampling of a student's t-table. The table shows the probability of exceeding the value in the body. With 5 df, there is a .05 probability that a sample will have a t-score > 2.015.

Since the t-distributions are symmetric, if 2 1/2% (.025) of the t's with 9df are greater than 2.262, then 2 1/2% are less than -2.262. The middle 95% (.95) of the t's, when there are 9df, are between -2.262 and +2.262. The

middle .90 of t=scores when there are 14df are between ± 1.761 , because -1.761 leaves .05 in the left tail and $+1.761$ leaves .05 in the right tail. The t-distribution gets closer and closer to the normal distribution as the number of degrees of freedom rises. As a result, the last line in the t-table, for infinity df, can also be used to find the z-scores that leave different proportions of the sample in the tail.

What could Kevin have done if he had been asked "about how much does a pair of size 11 socks weigh?" and he could not easily find good data on the population? Since he knows statistics, he could take a sample and make an inference about the population mean. Because the distribution of weights of socks is the result of a manufacturing process, it is almost certainly normal. The characteristics of almost every manufactured product are normally distributed. In a manufacturing process, even one that is precise and well-controlled, each individual piece varies slightly as the temperature varies some, the strength of the power varies as other machines are turned on and off, the consistency of the raw material varies slightly, and dozens of other forces that affect the final outcome vary slightly. Most of the socks, or bolts, or whatever is being manufactured, will be very close to the mean weight, or size, with just as many a little heavier or larger as there are that are a little lighter or smaller. Even though the process is supposed to be producing a population of "identical" items, there will be some variation among them. This is what causes so many populations to be normally distributed. Because the distribution of weights is normal, he can use the t-table to find the shape of the distribution of sample t-scores. Because he can use the t-table to tell him about the shape of the distribution of sample t-scores, he can make a good inference about the mean weight of a pair of socks. This is how he could make that inference:

STEP 1. Take a sample of n, say 15, pairs size 11 socks and carefully weigh each pair.

STEP 2. Find \bar{x} and s for his sample.

STEP 3 (where the tricky part starts). Look at the t-table, and find the t-scores that leave some proportion, say .95, of sample t's with n-1df in the middle.

STEP 4 (the heart of the tricky part). Assume that his sample has a t-score that is in the middle part of the distribution of t-scores.

STEP 5 (the arithmetic). Take his \bar{x} , s, n, and t's from the t-table, and set up two equations, one for each of his two table t-values. When he solves each of these equations for μ , he will find a interval that he is 95% sure (a statistician would say "with .95 confidence") contains the population mean.

Kevin decides this is the way he will go to answer the question. His sample contains pairs of socks with weights of:

4.36, 4.32, 4.29, 4.41, 4.45, 4.50, 4.36, 4.35, 4.33, 4.30, 4.39, 4.41, 4.43, 4.28, 4.46 oz.

He finds his sample mean, $\bar{x} = 4.376$ ounces, and his sample standard deviation (remembering to use the sample formula), $s = .067$ ounces. The t-table tells him that .95 of sample t's with 14df are between ± 2.145 . He solves these two equations for μ :

$$+2.145 = (4.376 - \mu)/(.067/\sqrt{14}) \quad \text{and} \quad -2.145 = (4.376 - \mu)/(.067/\sqrt{14})$$

finding $\mu = 4.366$ ounces and $\mu = 4.386$. With these results, Kevin can report that he is "95 per cent sure that the mean weight of a pair of size 11 socks is between 4.366 and 4.386 ounces". Notice that this is different from when he knew more about the population in the previous example.

2. The normal and t-distributions

Summary

A lot of material has been covered in this chapter, and not much of it has been easy. We are getting into real statistics now, and it will require care on your part if you are going to keep making sense of statistics.

The chapter outline is simple:

- Many things are distributed the same way, at least once we've standardized the members' values into z-scores.
- The central limit theorem gives users of statistics a lot of useful information about how the sampling distribution of \bar{x} is related to the original population of x's.
- The t-distribution lets us do many of the things the central limit theorem permits, even when the variance of the population, s_x^2 , is not known.

We will soon see that statisticians have learned about other sampling distributions and how they can be used to make inferences about populations from samples. It is through these known sampling distributions that most statistics is done. It is these known sampling distributions that give us the link between the sample we have and the population that we want to make an inference about.

3. Making estimates

The most basic kind of inference about a population is an estimate of the location (or shape) of a distribution. The central limit theorem says that the sample mean is an unbiased estimator of the population mean and can be used to make a single point inference of the population mean. While making this kind of inference will give you the correct estimate on average, it seldom gives you exactly the correct estimate. As an alternative, statisticians have found out how to estimate an interval that almost certainly contains the population mean. In the next few pages, you will learn how to make three different inferences about a population from a sample. You will learn how to make interval estimates of the mean, the proportion of members with a certain characteristic, and the variance. Each of these procedures follows the same outline, yet each uses a different sampling distribution to link the sample you have chosen with the population you are trying to learn about.

Estimating the population mean

Though the sample mean is an unbiased estimator of the population mean, very few samples have a mean exactly equal to the population mean. Though few samples have a mean, exactly equal to the population mean, m , the central limit theorem tells us that most samples have a mean that is close to the population mean. As a result, if you use the central limit theorem to estimate μ , you will seldom be exactly right, but you will seldom be far wrong. Statisticians have learned how often a point estimate will be how wrong. Using this knowledge you can find an interval, a range of values, which probably contains the population mean. You even get to choose how great a probability you want to have, though to raise the probability, the interval must be wider.

Most of the time, estimates are interval estimates. When you make an interval estimate, you can say "I am z per cent sure that the mean of this population is between x and y ". Quite often, you will hear someone say that they have estimated that the mean is some number " \pm so much". What they have done is quoted the midpoint of the interval for the "some number", so that the interval between x and y can then be split in half with $+$ "so much" above the midpoint and $-$ "so much" below. They usually do not tell you that they are only " z per cent sure". Making such an estimate is not hard— it is what Kevin Schmidt did at the end of the last chapter. It is worth your while to go through the steps carefully now, because the same basic steps are followed for making any interval estimate.

In making any interval estimate, you need to use a sampling distribution. In making an interval estimate of the population mean, the sampling distribution you use is the t -distribution.

The basic method is to pick a sample and then find the range of population means that would put your sample's t -score in the central part of the t -distribution. To make this a little clearer, look at the formula for t :

$$t = \frac{\bar{x} - \mu}{\frac{s}{\sqrt{n}}}$$

n is your sample's size and \bar{x} and s are computed from your sample. μ is what you are trying to estimate. From the t -table, you can find the range of t -scores that include the middle 80 per cent, or 90 per cent, or whatever per

3. Making estimates

cent, for $n-1$ degrees of freedom. Choose the percentage you want and use the table. You now have the lowest and highest t-scores, \bar{x} , s and n . You can then substitute the lowest t-score into the equation and solve for μ to find one of the limits for μ if your sample's t-score is in the middle of the distribution. Then substitute the highest t-score into the equation, and find the other limit. Remember that you want two μ 's because you want to be able to say that the population mean is between two numbers.

The two t-scores are almost always \pm the same number. The only heroic thing you have done is to assume that your sample has a t-score that is "in the middle" of the distribution. As long as your sample meets that assumption, the population mean will be within the limits of your interval. The probability part of your interval estimate, "I am z per cent sure that the mean is between...", or "with z confidence, the mean is between...", comes from how much of the t-distribution you want to include as "in the middle". If you have a sample of 25 (so there are 24df), looking at the table you will see that .95 of all samples of 25 will have a t-score between ± 2.064 ; that also means that for any sample of 25, the probability that its t is between ± 2.064 is .95.

As the probability goes up, the range of t-scores necessary to cover the larger proportion of the sample gets larger. This makes sense. If you want to improve the chance that your interval contains the population mean, you could simply choose a wider interval. For example, if your sample mean was 15, sample standard deviation was 10, and sample size was 25, to be .95 sure you were correct, you would need to base your mean on t-scores of ± 2.064 . Working through the arithmetic gives you an interval from 10.872 to 19.128. To have .99 confidence, you would need to base your interval on t-scores of ± 2.797 . Using these larger t-scores gives you a wider interval, one from 9.416 to 20.584. This trade-off between precision (a narrower interval is more precise) and confidence (probability of being correct), occurs in any interval estimation situation. There is also a trade-off with sample size. Looking at the t-table, note that the t-scores for any level of confidence are smaller when there are more degrees of freedom. Because sample size determines degrees of freedom, you can make an interval estimate for any level of confidence more precise if you have a larger sample. Larger samples are more expensive to collect, however, and one of the main reasons we want to learn statistics is to save money. There is a three-way trade-off in interval estimation between precision, confidence, and cost.

At Foothill Hosiery, John McGrath has become concerned that the hiring practices discriminate against older workers. He asks Kevin to look into the age at which new workers are hired, and Kevin decides to find the average age at hiring. He goes to the personnel office, and finds out that over 2,500 different people have worked at Foothill in the past fifteen years. In order to save time and money, Kevin decides to make an interval estimate of the mean age at date of hire. He decides that he wants to make this estimate with .95 confidence. Going into the personnel files, Kevin chooses 30 folders, and records the birth date and date of hiring from each. He finds the age at hiring for each person, and computes the sample mean and standard deviation, finding $\bar{x} = 24.71$ years and $s = 2.13$ years. Going to the t-table, he finds that .95 of t-scores with 29df are between ± 2.045 . He solves two equations:

$$\pm 2.045 = (24.71 - \mu) / (2.13 / \sqrt{30})$$

and finds that the limits to his interval are 23.91 and 25.51. Kevin tells Mr McGrath: "With .95 confidence, the mean age at date of hire is between 23.91 years and 25.51 years."

Estimating the population proportion

There are many times when you, or your boss, will want to estimate the proportion of a population that has a certain characteristic. The best known examples are political polls when the proportion of voters who would vote

for a certain candidate is estimated. This is a little trickier than estimating a population mean. It should only be done with large samples and there are adjustments that should be made under various conditions. We will cover the simplest case here, assuming that the population is very large, the sample is large, and that once a member of the population is chosen to be in the sample, it is replaced so that it might be chosen again. Statisticians have found that, when all of the assumptions are met, there is a sample statistic that follows the standard normal distribution. If all of the possible samples of a certain size are chosen, and for each sample, p , the proportion of the sample with a certain characteristic, is found, and for each sample a z -statistic is computed with the formula:

$$z = \frac{p - \pi}{\sqrt{\frac{(p)(1-p)}{n}}}$$

where π = proportion of population with the characteristic these will be distributed normally. Looking at the bottom line of the t -table, .90 of these z 's will be between ± 1.645 , .99 will be between ± 2.326 , etc.

Because statisticians know that the z -scores found from sample will be distributed normally, you can make an interval estimate of the proportion of the population with the characteristic. This is simple to do, and the method is parallel to that used to make an interval estimate of the population mean: (1) choose the sample, (2) find the sample p , (3) assume that your sample has a z -score that is not in the tails of the sampling distribution, (4) using the sample p as an estimate of the population π in the denominator and the table z -values for the desired level of confidence, solve twice to find the limits of the interval that you believe contains the population proportion, p .

At Foothill Hosiery, Ann Howard is also asked by John McGrath to look into the age at hiring at the plant. Ann takes a different approach than Kevin, and decides to investigate what proportion of new hires were at least 35. She looks at the personnel records and, like Kevin, decides to make an inference from a sample after finding that over 2,500 different people have worked at Foothill at some time in the last fifteen years. She chooses 100 personnel files, replacing each file after she has recorded the age of the person at hiring. She finds 17 who were 35 or older when they first worked at Foothill. She decides to make her inference with .95 confidence, and from the last line of the t -table finds that .95 of z -scores lie between ± 1.96 . She finds her upper and lower bounds:

$$+1.96 = \frac{.17 - \pi}{\sqrt{\frac{(.17)(1-.17)}{100}}}$$

$$\pi = .17 - (.038)(1.96) = .095$$

and, she finds the other boundary:

$$-1.96 = \frac{.17 - p}{\sqrt{\frac{(.17)(1-.17)}{100}}}$$

$$\pi = .17 - (.038)(1.96) = .245$$

and concludes, that with .95 confidence, the proportion of people who have worked at Foothills Hosiery who were over 35 when hired is between .095 and .245. This is a fairly wide interval. Looking at the equation for constructing the interval, you should be able to see that a larger sample size will result in a narrower interval, just as it did when estimating the population mean.

3. Making estimates

Estimating population variance

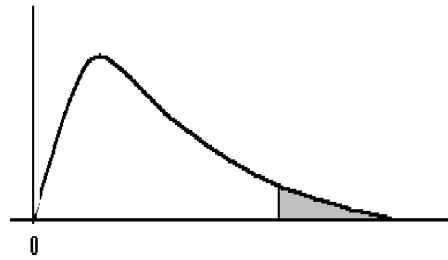
Another common interval estimation task is to estimate the variance of a population. High quality products not only need to have the proper mean dimension, the variance should be small. The estimation of population variance follows the same strategy as the other estimations. By choosing a sample and assuming that it is from the middle of the population, you can use a known sampling distribution to find a range of values that you are confident contains the population variance. Once again, we will use a sampling distribution that statisticians have discovered forms a link between samples and populations.

Take a sample of size n from a normal population with known variance, and compute a statistic called " χ^2 " (pronounced "chi square") for that sample using the following formula:

$$\chi^2 = \frac{\sum(x - \bar{x})^2}{\sigma^2}$$

You can see that χ^2 will always be positive, because both the numerator and denominator will always be positive. Thinking it through a little, you can also see that as n gets larger, χ^2 will generally be larger since the numerator will tend to be larger as more and more $(x - \bar{x})^2$ are summed together. It should not be too surprising by now to find out that if all of the possible samples of a size n are taken from any normal population, that when χ^2 is computed for each sample and those χ^2 are arranged into a relative frequency distribution, the distribution is always the same.

Because the size of the sample obviously affects χ^2 , there is a different distribution for each different sample size. There are other sample statistics that are distributed like χ^2 , so, like the t-distribution, tables of the χ^2 distribution are arranged by degrees of freedom so that they can be used in any procedure where appropriate. As you might expect, in this procedure, $df = n-1$. A portion of a χ^2 table is reproduced below.



The χ^2 distribution

	p	.95	.90	.10	.05
n	df				
2	1	0.004	0.02	2.706	3.841
10	9	3.33	4.17	14.68	19.92
15	14	6.57	7.79	21.1	23.7
20	19	10.12	11.65	27.2	30.1
30	19	17.71	19.77	39.1	42.6

Exhibit 4: The χ^2 distribution

Variance is important in quality control because you want your product to be consistently the same. John McGrath has just returned from a seminar called "Quality Socks, Quality Profits". He learned something about variance, and has asked Kevin to measure the variance of the weight of Foothill's socks. Kevin decides that he can fulfill this request by using the data he collected when Mr McGrath asked about the average weight of size 11 men's dress socks. Kevin knows that the sample variance is an unbiased estimator of the population variance, but he decides to produce an interval estimate of the variance of the weight of pairs of size 11 men's socks. He also decides that .90 confidence will be good until he finds out more about what Mr McGrath wants.

Kevin goes and finds the data for the size 11 socks, and gets ready to use the χ^2 distribution to make a .90 confidence interval estimate of the variance of the weights of socks. His sample has 15 pairs in it, so he will have 14 df. From the χ^2 table he sees that .95 of χ^2 are greater than 6.57 and only .05 are greater than 23.7 when there are 14df. This means that .90 are between 6.57 and 23.7. Assuming that his sample has a χ^2 that is in the middle .90, Kevin gets ready to compute the limits of his interval. He notices that he will have to find $\sum (x - \bar{x})^2$ and decides to use his spreadsheet program rather than find $(x - \bar{x})^2$ fifteen times. He puts the original sample values in the first column, and has the program compute the mean. Then he has the program find $(x - \bar{x})^2$ fifteen times. Finally, he has the spreadsheet sum up the squared differences and finds 0.062.

3. Making estimates

Kevin then takes the χ^2 formula, and solves it twice, once by setting χ^2 equal to 6.57:

$$\chi^2 = 6.57 = .062/\sigma^2$$

Solving for σ^2 , he finds one limit for his interval is .0094. He solves the second time by setting $\chi^2 = 23.6$:

$$23.6 = .062/\sigma^2$$

and find that the other limit is .0026. Armed with his data, Kevin reports to Mr McGrath that "with .90 confidence, the variance of weights of size 11 men's socks is between .0026 and .0094."

What is this confidence stuff mean anyway?

In the example we just did, Ann found "that with .95 confidence..." What exactly does "with .95 confidence" mean? The easiest way to understand this is to think about the assumption that Ann had made that she had a sample with a z-score that was not in the tails of the sampling distribution. More specifically, she assumed that her sample had a z-score between ± 1.96 ; that it was in the middle 95 per cent of z-scores. Her assumption is true 95% of the time because 95% of z-scores *are* between ± 1.96 . If Ann did this same estimate, including drawing a new sample, over and over, in .95 of those repetitions, the population proportion *would* be within the interval because in .95 of the samples the z-score would be between ± 1.96 . In .95 of the repetitions, her estimate would be right.

4. Hypothesis testing

Hypothesis testing is the other widely used form of inferential statistics. It is different from estimation because you start a hypothesis test with some idea of what the population is like and then test to see if the sample supports your idea. Though the mathematics of hypothesis testing is very much like the mathematics used in interval estimation, the inference being made is quite different. In estimation, you are answering the question "what is the population like?" While in hypothesis testing you are answering the question "is the population like this or not?"

A hypothesis is essentially an idea about the population that you think might be true, but which you cannot prove to be true. While you usually have good reasons to think it is true, and you often hope that it is true, you need to show that the sample data supports your idea. Hypothesis testing allows you to find out, in a formal manner, if the sample supports your idea about the population. Because the samples drawn from any population vary, you can never be positive of your finding, but by following generally accepted hypothesis testing procedures, you can limit the uncertainty of your results.

As you will learn in this chapter, you need to choose between two statements about the population. These two statements are the hypotheses. The first, known as the "null hypothesis", is basically "the population is like this". It states, in formal terms, that the population is no different than usual. The second, known as the "alternative hypothesis", is "the population is like something else". It states that the population is different than the usual, that something has happened to this population, and as a result it has a different mean, or different shape than the usual case. Between the two hypotheses, all possibilities must be covered. Remember that you are making an inference about a population from a sample. Keeping this inference in mind, you can informally translate the two hypotheses into "I am almost positive that the sample came from a population like this" and "I really doubt that the sample came from a population like this, so it probably came from a population that is like something else". Notice that you are never entirely sure, even after you have chosen the hypothesis which is best. Though the formal hypotheses are written as though you will choose with certainty between the one that is true and the one that is false, the informal translations of the hypotheses, with "almost positive" or "probably came", is a better reflection of what you actually find.

Hypothesis testing has many applications in business, though few managers are aware that that is what they are doing. As you will see, hypothesis testing, though disguised, is used in quality control, marketing, and other business applications. Many decisions are made by thinking as though a hypothesis is being tested, even though the manager is not aware of it. Learning the formal details of hypothesis testing will help you make better decisions and better understand the decisions made by others.

The next section will give an overview of the hypothesis testing method by following along with a young decision-maker as he uses hypothesis testing. The rest of the chapter will present some specific applications of hypothesis tests as examples of the general method.

4. Hypothesis testing

The strategy of hypothesis testing

Usually, when you use hypothesis testing, you have an idea that the world is a little bit surprising, that it is not exactly as conventional wisdom says it is. Occasionally, when you use hypothesis testing, you are hoping to confirm that the world is not surprising, that it is like conventional wisdom predicts. Keep in mind that in either case you are asking "is the world different from the usual, is it surprising?" Because the world is usually not surprising and because in statistics you are never 100 per cent sure about what a sample tells you about a population, you cannot say that your sample implies that the world is surprising unless you are almost positive that it does. The dull, unsurprising, usual case not only wins if there is a tie, it gets a big lead at the start. You cannot say that the world is surprising, that the population is unusual, unless the evidence is very strong. This means that when you arrange your tests, you have to do it in a manner that makes it difficult for the unusual, surprising world to win support.

The first step in the basic method of hypothesis testing is to decide what value some measure of the population would take if the world was unsurprising. Second, decide what the sampling distribution of some sample statistic would look like if the population measure had that unsurprising value. Third, compute that statistic from your sample and see if it could easily have come from the sampling distribution of that statistic if the population was unsurprising. Fourth, decide if the population your sample came from is surprising because your sample statistic could not easily have come from the sampling distribution generated from the unsurprising population.

That all sounds complicated, but it is really pretty simple. You have a sample and the mean, or some other statistic, from that sample. With conventional wisdom, the null hypothesis that the world is dull and not surprising, tells you that your sample comes from a certain population. Combining the null hypothesis with what statisticians know tells you what sampling distribution your sample statistic comes from if the null hypothesis is true. If you are "almost positive" that the sample statistic came from that sampling distribution, the sample supports the null. If the sample statistic "probably came" from a sampling distribution generated by some other population, the sample supports the alternative hypothesis that the population is "like something else".

Imagine that Thad Stoykov works in the marketing department of Pedal Pushers, a company that makes clothes for bicycle riders. Pedal Pushers has just completed a big advertising campaign in various bicycle and outdoor magazines, and Thad wants to know if the campaign has raised the recognition of the Pedal Pushers brand so that more than 30 per cent of the potential customers recognize it. One way to do this would be to take a sample of prospective customers and see if at least 30 per cent of those in the sample recognize the Pedal Pushers brand. However, what if the sample is small and just barely 30 per cent of the sample recognizes Pedal Pushers? Because there is variance among samples, such a sample could easily have come from a population in which less than 30 percent recognize the brand—if the population actually had slightly less than 30 per cent recognition, the sampling distribution would include quite a few samples with sample proportions a little above 30 per cent, especially if the samples are small. In order to be comfortable that more than 30 per cent of the **population** recognizes Pedal Pushers, Thad will want to find that a bit more than 30 per cent of the **sample** does. How much more depends on the size of the sample, the variance within the sample, and how much chance he wants to take that he'll conclude that the campaign did not work when it actually did.

Let us follow the formal hypothesis testing strategy along with Thad. First, he must explicitly describe the population his sample could come from in two different cases. The first case is the unsurprising case, the case where there is no difference between the population his sample came from and most other populations. This is the case where the ad campaign did not really make a difference, and it generates the null hypothesis. The second case is the

surprising case when his sample comes from a population that is different from most others. This is where the ad campaign worked, and it generates the alternative hypothesis. The descriptions of these cases are written in a formal manner. The null hypothesis is usually called " H_o ". The alternative hypothesis is called either " H_1 :" or " H_a :". For Thad and the Pedal Pushers marketing department, the null will be :

H_o : proportion of the population recognizing Pedal Pushers brand $\leq .30$ and the alternative will be:

H_a : proportion of the population recognizing Pedal Pushers brand $>.30$.

Notice that Thad has stacked the deck against the campaign having worked by putting the value of the population proportion that means that the campaign was successful in the alternative hypothesis. Also notice that between H_o : and H_a : all possible values of the population proportion— $>$, $=$, and $< .30$ — have been covered.

Second, Thad must create a rule for deciding between the two hypotheses. He must decide what statistic to compute from his sample and what sampling distribution that statistic would come from if the null hypothesis,

H_o :, is true. He also needs to divide the possible values of that statistic into "usual" and "unusual" ranges if the null is true. Thad's decision rule will be that if his sample statistic has a "usual" value, one that could easily occur if

H_o : is true, then his sample could easily have come from a population like that described in H_o :. If his sample's statistic has a value that would be "unusual" if H_o : is true, then the sample probably comes from a population like that described in H_a :. Notice that the hypotheses and the inference are about the original population while the decision rule is about a sample statistic. The link between the population and the sample is the sampling distribution. Knowing the relative frequency of a sample statistic when the original population has a proportion with a known value is what allows Thad to decide what are "usual" and "unusual" values for the sample statistic.

The basic idea behind the decision rule is to decide, with the help of what statisticians know about sampling distributions, how far from the null hypothesis' value for the population the sample value can be before you are uncomfortable deciding that the sample comes from a population like that hypothesized in the null. Though the hypotheses are written in terms of descriptive statistics about the population—means, proportions, or even a distribution of values—the decision rule is usually written in terms of one of the standardized sampling distributions—the t, the normal z, or another of the statistics whose distributions are in the tables at the back of statistics books. It is the sampling distributions in these tables that are the link between the sample statistic and the population in the null hypothesis. If you learn to look at how the sample statistic is computed you will see that all of the different hypothesis tests are simply variations on a theme. If you insist on simply trying to memorize how each of the many different statistics is computed, you will not see that all of the hypothesis tests are conducted in a similar manner, and you will have to learn many different things rather than learn the variations of one thing.

Thad has taken enough statistics to know that the sampling distribution of sample proportions is normally distributed with a mean equal to the population proportion and a standard deviation that depends on the population proportion and the sample size. Because the distribution of sample proportions is normally distributed, he can look at the bottom line of a t-table and find out that only .05 of all samples will have a proportion more than 1.645 standard deviations above .30 if the null hypothesis is true. Thad decides that he is willing to take a 5 per cent chance that he will conclude that the campaign did not work when it actually did, and therefore decides that he will

4. Hypothesis testing

conclude that the sample comes from a population with a proportion that has heard of Pedal Pushers that is greater than .30 if the sample's proportion is more than 1.645 standard deviations above .30. After doing a little arithmetic (which you'll learn how to do later in the chapter), Thad finds that his decision rule is to decide that the campaign was effective if the sample has a proportion which has heard of Pedal Pushers that is greater than .375. Otherwise the sample could too easily have come from a population with a proportion equal to or less than .30.

alpha	0.1	0.05	0.03	0.01
df infinity	1.28	1.65	1.96	2.33

Exhibit 5: The bottom line of a t-table, showing the normal distribution

The final step is to compute the sample statistic and apply the decision rule. If the sample statistic falls in the usual range, the data supports H_o ;, and the world is probably unsurprising and the campaign did not make any difference. If the sample statistic is outside the usual range, the data supports H_a ;, and the world is a little surprising, the campaign affected how many people have heard of Pedal Pushers. When Thad finally looks at the sample data, he finds that .39 of the sample had heard of Pedal Pushers. The ad campaign was successful!

A straight-forward example: testing for "goodness-of-fit"

There are many different types of hypothesis tests, including many that are used more often than the "goodness-of-fit" test. This test will be used to help introduce hypothesis testing because it gives a clear illustration of how the strategy of hypothesis testing is put to use, not because it is used frequently. Follow this example carefully, concentrating on matching the steps described in previous sections with the steps described in this section; the arithmetic is not that important right now.

We will go back to Ann Howard's problem with marketing "Easy Bounce" socks to volleyball teams. Remember that Ann works for Foothills Hosiery, and she is trying to market these sports socks to volleyball teams. She wants to send out some samples to convince volleyball players that wearing "Easy Bounce" socks will be more comfortable than wearing other socks. Her idea is to send out a package of socks to volleyball coaches in the area, so the players can try them out. She needs to include an assortment of sizes in those packages and is trying to find out what sizes to include. The Production Department knows what mix of sizes they currently produce, and Ann has collected a sample of 97 volleyball players' sock sizes from nearby teams. She needs to test to see if her sample supports the hypothesis that volleyball players have the same distribution of sock sizes as Foothills is currently producing—is the distribution of volleyball players' sock sizes a "good fit" to the distribution of sizes now being produced?

Ann's sample, a sample of the sock sizes worn by volleyball players, as a frequency distribution of sizes:

size	frequency
6	3
7	24
8	33
9	20
10	17

From the Production Department, Ann finds that the current relative frequency distribution of production of "Easy Bounce" socks is like this:

size	re. frequency
6	0.06
7	0.13
8	0.22
9	0.3
10	0.26
11	0.03

If the world in "unsurprising", volleyball players will wear the socks sized in the same proportions as other athletes, so Ann writes her hypotheses:

H_o : Volleyball players' sock sizes are distributed just like current production.

H_a : Volleyball players' sock sizes are distributed differently.

Ann's sample has $n=97$. By applying the relative frequencies in the current production mix, she can find out how many players would be "expected" to wear each size if her sample was perfectly representative of the distribution of sizes in current production. This would give her a description of what a sample from the population in the null hypothesis would be like. It would show what a sample that had a "very good fit" with the distribution of sizes in the population currently being produced would look like.

Statisticians know the sampling distribution of a statistic which compares the "expected" frequency of a sample with the actual, or "observed" frequency. For a sample with c different classes (the sizes here), this statistic is distributed like χ^2 with $c-1$ df. The χ^2 is computed by the formula:

$$\text{sample } \chi^2 = \sum \frac{(O - E)^2}{E}$$

where:

O = observed frequency in the sample in this class

E = expected frequency in the sample in this class.

The expected frequency, E, is found by multiplying the relative frequency of this class in the H_o :hypothesized population by the sample size. This gives you the number in that class in the sample if the relative frequency distribution across the classes in the sample exactly matches the distribution in the population.

Notice that χ^2 is always ≥ 0 and equals 0 only if the observed is equal to the expected in each class. Look at the equation and make sure that you see that a larger value of goes with samples with large differences between the observed and expected frequencies.

Ann now needs to come up with a rule to decide if the data supports H_o : or H_a :. She looks at the table and sees that for 5 df (there are 6 classes—there is an expected frequency for size 11 socks), only .05 of samples drawn from a given population will have a $\chi^2 > 11.07$ and only .10 will have a $\chi^2 > 9.24$. She decides that it

4. Hypothesis testing

would not be all that surprising if volleyball players had a different distribution of sock sizes than the athletes who are currently buying "Easy Bounce", since all of the volleyball players are women and many of the current customers are men. As a result, she uses the smaller .10 value of 9.24 for her decision rule. Now she must compute her sample χ^2 . Ann starts by finding the expected frequency of size 6 socks by multiplying the relative frequency of size 6 in the population being produced by 97, the sample size. She gets $E = .06 \cdot 97 = 5.82$. She then finds $O - E = 3 - 5.82 = -2.82$, squares that and divides by 5.82, eventually getting 1.37. She then realizes that she will have to do the same computation for the other five sizes, and quickly decides that a spreadsheet will make this much easier. Her spreadsheet looks like this:

sock size	frequency in sample	population relative frequency	expected frequency = $97 \cdot C$	$(O - E)^2 / E$
6	3	0.06	5.82	1.3663918
7	24	0.13	12.61	10.288033
8	33	0.22	21.34	6.3709278
9	20	0.3	29.1	2.8457045
10	17	0.26	25.22	2.6791594
11	0	0.03	2.91	2.91
	97			$\chi^2 =$ 26.460217

Exhibit 6: Ann's Excel sheet

Ann performs her third step, computing her sample statistic, using the spreadsheet. As you can see, her sample $\chi^2 = 26.46$, which is well into the "unusual" range which starts at 9.24 according to her decision rule. Ann has found that her sample data supports the hypothesis that the distribution of sock sizes of volleyball players is different from the distribution of sock sizes that are currently being manufactured. If Ann's employer, Foothill Hosiery, is going to market "Easy Bounce" socks to volleyball players, they are going to have to send out packages of samples that contain a different mix of sizes than they are currently making. If "Easy Bounce" are successfully marketed to volleyball players, the mix of sizes manufactured will have to be altered.

Now, review what Ann has done to test to see if the data in her sample supports the hypothesis that the world is "unsurprising" and that volleyball players have the same distribution of sock sizes as Foothill Hosiery is currently producing for other athletes. The essence of Ann's test was to see if her sample χ^2 could easily have come from the sampling distribution of χ^2 's generated by taking samples from the population of socks currently being produced. Since her sample χ^2 would be way out in the tail of that sampling distribution, she judged that her sample data supported the other hypothesis, that there is a difference between volleyball players and the athletes who are currently buying "Easy Bounce" socks.

Formally, Ann first wrote null and alternative hypotheses, describing the population her sample comes from in two different cases. The first case is the null hypothesis; this occurs if volleyball players wear socks of the same sizes in the same proportions as Foothill is currently producing. The second case is the alternative hypothesis; this occurs if volleyball players wear different sizes. After she wrote her hypotheses, she found that there was a sampling

distribution that statisticians knew about that would help her choose between them. This is the χ^2 distribution. Looking at the formula for computing χ^2 and consulting the tables, Ann decided that a sample χ^2 value greater than 9.24 would be unusual if her null hypothesis was true. Finally, she computed her sample statistic, and found that her χ^2 , at 26.46, was well above her cut-off value. Ann had found that the data in her sample supported the alternative, H_a ; that the distribution of volleyball players' sock sizes is different from the distribution that Foothill is currently manufacturing. Acting on this finding, Ann will send a different mix of sizes in the sample packages she sends volleyball coaches.

Testing population proportions

As you learned in the chapter "Making estimates", sample proportions can be used to compute a statistic that has a known sampling distribution. Reviewing, the z-statistic is:

$$z = \frac{p - \pi}{\sqrt{\frac{(\pi)(1-\pi)}{n}}}$$

where: p = the proportion of the sample with a certain characteristic

π = the proportion of the population with that characteristic.

These sample z-statistics are distributed normally, so that by using the bottom line of the t table, you can find what portion of all samples from a population with a given population proportion, π , have z-statistics within different ranges. If you look at the table, you can see that .95 of all samples from any population have a z-statistics between ± 1.96 , for instance.

If you have a sample that you think is from a population containing a certain proportion, π , of members with some characteristic, you can test to see if the data in your sample supports what you think. The basic strategy is the same as that explained earlier in this chapter and followed in the "goodness-of-fit" example: (a) write two hypotheses, (b) find a sample statistic and sampling distribution that will let you develop a decision rule for choosing between the two hypotheses, and (c) compute your sample statistic and choose the hypothesis supported by the data.

Foothill Hosiery recently received an order for children's socks decorated with embroidered patches of cartoon characters. Foothill did not have the right machinery to sew on the embroidered patches and contracted out the sewing. While the order was filled and Foothill made a profit on it, the sewing contractor's price seemed high, and Foothill had to keep pressure on the contractor to deliver the socks by the date agreed upon. Foothill's CEO, John McGrath has explored buying the machinery necessary to allow Foothill to sew patches on socks themselves. He has discovered that if more than a quarter of the children's socks they make are ordered with patches, the machinery will be a sound investment. Mr McGrath asks Kevin Schmidt to find out if more than 25 per cent of children's socks are being sold with patches.

Kevin calls the major trade organizations for the hosiery, embroidery, and children's clothes industries, and no one can answer his question. Kevin decides it must be time to take a sample and to test to see if more than 25 per cent of children's socks are decorated with patches. He calls the sales manager at Foothill and she agrees to ask her salespeople to look at store displays of children's socks, counting how many pairs are displayed and how many of

4. Hypothesis testing

those are decorated with patches. Two weeks later, Kevin gets a memo from the sales manager telling him that of the 2,483 pairs of children's socks on display at stores where the salespeople counted, 716 pairs had embroidered patches.

Kevin writes his hypotheses, remembering that Foothill will be making a decision about spending a fair amount of money based on what he finds. To be more certain that he is right if he recommends that the money be spent, Kevin writes his hypotheses so that the "unusual" world would be the one where more than 25 per cent of children's socks are decorated:

$$H_0: \pi_{\text{decorated socks}} \leq .25$$

$$H_a: \pi_{\text{decorated socks}} > .25$$

When writing his hypotheses, Kevin knows that if his sample has a proportion of decorated socks well below .25, he will want to recommend against buying the machinery. He only wants to say the data supports the alternative if the sample proportion is well above .25. To include the low values in the null hypothesis and only the high values in the alternative, he uses a "one-tail" test, judging that the data supports the alternative only if his z-score is in the upper tail. He will conclude that the machinery should be bought only if his z-statistic is too large to have easily have come from the sampling distribution drawn from a population with a proportion of .25. Kevin will accept H_a : only if his z is large and positive.

Checking the bottom line of the t-table, Kevin sees that .95 of all z-scores are less than 1.645. His rule is therefore to conclude that his sample data supports the null hypothesis that 25 per cent or less of children's socks are decorated if his sample z is less than 1.645. If his sample z is greater than 1.645, he will conclude that more than 25 per cent of children's socks are decorated and that Foothill Hosiery should invest in the machinery needed to sew embroidered patches on socks.

Using the data the salespeople collected, Kevin finds the proportion of the sample that is decorated:

$$p = \frac{716}{2483} = .288$$

Using this value, he computes his sample z-statistic:

$$\begin{aligned} z &= \frac{p - \pi}{\sqrt{\frac{(\pi)(1 - \pi)}{n}}} \\ &= \frac{.288 - .25}{\sqrt{\frac{(.25)(1 - .25)}{2483}}} \\ &= \frac{.0380}{.0087} = 4.368. \end{aligned}$$

Because his sample z-score is larger than 1.645, it is unlikely that his sample z came from the sampling distribution of z's drawn from a population where $\pi \leq .25$, so it is unlikely that his sample comes from a population with $\pi \leq .25$. Kevin can tell John McGrath that the sample the sales people collected supports the conclusion that

more than 25 per cent of children's socks are decorated with embroidered patches. John can feel comfortable making the decision to buy the embroidery and sewing machinery.

Summary

This chapter has been an introduction to hypothesis testing. You should be able to see the relationship between the mathematics and strategies of hypothesis testing and the mathematics and strategies of interval estimation. When making an interval estimate, you construct an interval around your sample statistic based on a known sampling distribution. When testing a hypothesis, you construct an interval around a hypothesized population parameter, using a known sampling distribution to determine the width of that interval. You then see if your sample statistic falls within that interval to decide if your sample probably came from a population with that hypothesized population parameter.

Hypothesis testing is a very widely used statistical technique. It forces you to think ahead about what you might find. By forcing you to think ahead, it often helps with decision-making by forcing you to think about what goes into your decision. All of statistics requires clear thinking, and clear thinking generally makes better decisions. Hypothesis testing requires very clear thinking and often leads to better decision-making.

5. The t-test

In Chapter 3 a sampling distribution, the t-distribution, was introduced. In Chapter 4 you learned how to use the t-distribution to make an important inference, an interval estimate of the population mean. Here you will learn how to use that same t-distribution to make more inferences, this time in the form of hypothesis tests. Before we start to learn about those tests, a quick review of the t-distribution is in order.

The t-distribution

The t-distribution is a sampling distribution. You could generate your own t-distribution with $n-1$ degrees of freedom by starting with a normal population, choosing all possible samples of one size, n , computing a t-score for each sample:

$$t = \frac{\bar{x} - \mu}{\frac{s}{\sqrt{n}}}$$

where: \bar{x} = the sample mean

μ = the population mean

s = the sample standard deviation

n = the size of the sample.

When you have all of the samples' t-scores, form a relative frequency distribution and you will have your t-distribution. Luckily, you do not have to generate your own t-distributions because any statistics book has a table that shows the shape of the t-distribution for many different degrees of freedom. Exhibit 7 reproduces a portion of a typical t-table. See below.

5. The t-test



df	$p = .10$	$p = .05$	$p = .01$
	upper tail = .05	upper tail = 0.025	upper tail = .005
5	2.02	2.57	4.03
10	1.81	2.23	3.17
15	1.75	2.13	2.95
20	1.72	2.09	2.85
25	1.71	2.06	2.79
30	1.70	2.04	2.75
35	1.69	2.03	2.72
40	1.68	2.02	2.70
45	1.68	2.01	2.69
50	1.68	2.01	2.68
55	1.67	2.00	2.67
60	1.67	2.00	2.66
65	1.67	2.00	2.65
70	1.67	1.99	2.65
75	1.67	1.99	2.64
80	1.66	1.99	2.64
85	1.66	1.99	2.63
90	1.66	1.99	2.63
95	1.66	1.99	2.63
100	1.66	1.98	2.63
Infinite	1.65	1.96	2.58

Exhibit 7: A portion of a typical t-table

When you look at the formula for the t-score, you should be able to see that the mean t-score is zero because the mean of the \bar{x} 's is equal to μ . Because most samples have \bar{x} 's that are close to μ , most will have t-scores that are close to zero. The t-distribution is symmetric, because half of the samples will have \bar{x} 's greater than μ , and half less. As you can see from the table, if there are 10 df, only .005 of the samples taken from a normal population will have a t-score greater than +3.17. Because the distribution is symmetric, .005 also have a t-score less than -3.17. Ninety-nine per cent of samples will have a t-score between ± 3.17 . Like the example in Exhibit 7, most t-tables have a picture showing what is in the body of the table. In Exhibit 7, the shaded area is in the right tail, the body of the table shows the t-score that leaves the α in the right tail. This t-table also lists the two-tail α above the one-tail where it has $p = .xx$. For 5 df, there is a .05 probability that a sample will have a t-score greater than 2.02, and a .10 probability that a sample will have a t score either $> +2.02$ or < -2.02 .

There are other sample statistics which follow this same shape and which can be used as the basis for different hypothesis tests. You will see the t-distribution used to test three different types of hypotheses in this chapter and that the t-distribution can be used to test other hypotheses in later chapters.

Though t-tables show how the sampling distribution of t-scores is shaped if the original population is normal, it turns out that the sampling distribution of t-scores is very close to the one in the table even if the original population is not quite normal, and most researchers do not worry too much about the normality of the original population. An even more important fact is that the sampling distribution of t-scores is very close to the one in the table even if the original population is not very close to being normal as long as the samples are large. This means that you can safely use the t-distribution to make inferences when you are not sure that the population is normal as long as you are sure that it is bell-shaped. You can also make inferences based on samples of about 30 or more using the t-distribution when you are not sure if the population is normal. Not only does the t-distribution describe

the shape of the distributions of a number of sample statistics, it does a good job of describing those shapes when the samples are drawn from a wide range of populations, normal or not.

A simple test: does this sample come from a population with that mean?

Imagine that you have taken all of the samples with $n=10$ from a population that you knew the mean of, found the t-distribution for 9 df by computing a t-score for each sample and generated a relative frequency distribution of the t's. When you were finished, someone brought you another sample ($n=10$) wondering if that new sample came from the original population. You could use your sampling distribution of t's to test if the new sample comes from the original population or not. To conduct the test, first hypothesize that the new sample comes from the original population. With this hypothesis, you have hypothesized a value for μ , the mean of the original population, to use to compute a t-score for the new sample. If the t for the new sample is close to zero—if the t-score for the new sample could easily have come from the middle of the t-distribution you generated—your hypothesis that the new sample comes from a population with the hypothesized mean seems reasonable and you can conclude that the data supports the new sample coming from the original population. If the t-score from the new sample was far above or far below zero, your hypothesis that this new sample comes from the original population seems unlikely to be true, for few samples from the original population would have t-scores far from zero. In that case, conclude that the data gives support to the idea that the new sample comes from some other population.

This is the basic method of using this t-test. Hypothesize the mean of the population you think a sample might come from. Using that mean, compute the t-score for the sample. If the t-score is close to zero, conclude that your hypothesis was probably correct and that you know the mean of the population from which the sample came. If the t-score is far from zero, conclude that your hypothesis is incorrect, and the sample comes from a population with a different mean.

Once you understand the basics, the details can be filled in. The details of conducting a "hypothesis test of the population mean", testing to see if a sample comes from a population with a certain mean—are of two types. The first type concerns how to do all of this in the formal language of statisticians. The second type of detail is how to decide what range of t-scores implies that the new sample comes from the original population.

You should remember from the last chapter that the formal language of hypothesis testing always requires two hypotheses. The first hypothesis is called the "null hypothesis", usually denoted H_o :. It states that there is no difference between the mean of the population from which the sample is drawn and the hypothesized mean. The second is the "alternative hypothesis", denoted H_1 : or H_a :. It states that the mean of the population from which the sample comes is different from the hypothesized value. If your question is simply "does this sample come from a population with **this** mean?", your H_a : is simply $\mu \neq \text{the hypothesized value}$. If your question is "does this sample come from a population with a mean **greater than** some value", then your H_a : becomes $\mu > \text{the hypothesized value}$.

The other detail is deciding how "close to zero" the sample t-score has to be before you conclude that the null hypothesis is probably correct. How close to zero the sample t-score must be before you conclude that the data supports H_o : depends on the df and how big a chance you want to take that you will make a mistake. If you decide to conclude that the sample comes from a population with the hypothesized mean only if the sample t is very, very close to zero, there are many samples actually from the population that will have t-scores that would lead

5. The t-test

you to believe they come from a population with some other mean—it would be easy to make a mistake and conclude that these samples come from another population. On the other hand, if you decide to accept the null hypothesis even if the sample t-score is quite far from zero, you will seldom make the mistake of concluding that a sample from the original population is from some other population, but you will often make another mistake—concluding that samples from other populations are from the original population. There are no hard rules for deciding how much of which sort of chance to take. Since there is a trade-off between the chance of making the two different mistakes, the proper amount of risk to take will depend on the relative costs of the two mistakes. Though there is no firm basis for doing so, many researchers use a 5 per cent chance of the first sort of mistake as a default. The level of chance of making the first error is usually called "alpha" (α) and the value of alpha chosen is usually written as a decimal fraction—taking a 5 per cent chance of making the first mistake would be stated as " $\alpha = .05$ ". When in doubt, use $\alpha = .05$.

If your alternative hypothesis is "not equal to", you will conclude that the data supports H_a : if your sample t-score is either well below or well above zero and you need to divide α between the two tails of the t-distribution. If you want to use $\alpha = .05$, you will support H_a : if the t is in either the lowest .025 or the highest .025 of the distribution. If your alternative is "greater than", you will conclude that the data supports H_a : only if the sample t-score is well above zero. So, put all of your α in the right tail. Similarly, if your alternative is "less than", put the whole α in the left tail.

The table itself can be confusing even after you know how many degrees of freedom you have and if you want to split your α between the two tails or not. Adding to the confusion, not all t-tables look exactly the same. Look at the typical t-table above and notice that it has three parts: column headings of decimal fractions, row headings of whole numbers, and a body of numbers generally with values between 1 and 3. The column headings are labeled p or "area in the right tail," and sometimes are labeled " α ." The row headings are labeled "df," but are sometimes labeled "v" or "degrees of freedom". The body is usually left unlabeled and it shows the t-score which goes with the " α " and "degrees of freedom" of that column and row. These tables are set up to be used for a number of different statistical tests, so they are presented in a way that is a compromise between ease of use in a particular situation and the ability to use the same table for a wide variety of tests. My favorite t tables are available online at <http://www.itl.nist.gov/div898/handbook/eda/section3/eda3672.htm>

In order to use the table to test to see if "this sample comes from a population with a certain mean" choose α and find the number of degrees of freedom. The number of degrees of freedom in a test involving one sample mean is simply the size of the sample minus one ($df = n-1$). The α you choose may not be the α in the column heading. The column headings show the "right tail areas"—the chance you'll get a t-score *larger* than the one in the body of the table. Assume that you had a sample with ten members and chose $\alpha = .05$. There are nine degrees of freedom, so go across the 9 df row to the .025 column since this is a two-tail test, and find the t-score of 2.262. This means that in any sampling distribution of t-scores, with samples of ten drawn from a normal population, only 2.5 per cent (.025) of the samples would have t-scores greater than 2.262—any t-score greater than 2.262 probably occurs because the sample is from some other population with a larger mean. Because the t-distributions are symmetrical, it is also true that only 2.5 per cent of the samples of ten drawn from a normal population will have t-scores less than -2.262. Putting the two together, 5 per cent of the t-scores will have an absolute value greater the 2.262. So if

you choose $\alpha = .05$, you will probably be using a t-score in the .025 column. The picture that is at the top of most t-tables shows what is going on. Look at it when in doubt.

LaTonya Williams is the plant manager for Eileen's Dental Care Company (EDC) which makes dental floss. EDC has a good, stable work force of semi-skilled workers who work packaging floss, paid by piece-work, and the company wants to make sure that these workers are paid more than the local average wage. A recent report by the local Chamber of Commerce shows an average wage for "machine operators" of USD 8.71. LaTonya needs to decide if a raise is needed to keep her workers above the average. She takes a sample of workers, pulls their work reports, finds what each one earned last week and divides their earnings by the hours they worked to find average hourly earnings.

That data appears below:

Smith 9.01

Wilson 8.67

Peterson 8.90

Jones 8.45

Gordon 8.88

McCoy 9.13

Bland 8.77

LaTonya wants to test to see if the mean of the average hourly earnings of her workers is greater than USD 8.71. She wants to use a one-tail test because her question is "greater than" not "unequal to". Her hypotheses are:

$$H_0: \mu \leq 8.71 \quad \text{and} \quad H_a: \mu > 8.71$$

As is usual in this kind of situation, LaTonya is hoping that the data supports H_a ; but she wants to be confident that it does before she decides her workers are earning above average wages. Remember that she will compute a t-score for her sample using USD 8.71 for μ . If her t-score is negative or close to zero, she will conclude that the data supports H_0 . Only if her t-score is large and positive will she go with H_a . She decides to use $\alpha = .025$ because she is unwilling to take much risk of saying the workers earn above average wages when they really do not. Because her sample has $n=7$, she has 6 df. Looking at the table, she sees that the data will support H_a ; the workers earn more than average, only if the sample t-score is greater than 2.447.

Finding the sample mean and standard deviation, $\bar{x} = \$8.83$ and $s = .225$, LaTonya computes her sample t-score:

$$t = \frac{\bar{x} - \mu}{\frac{s}{\sqrt{n}}} = \frac{8.83 - 8.71}{\frac{.225}{\sqrt{7}}}$$

$$t = \frac{.12}{.085} = 1.41$$

Because her sample t is not greater than +2.447, LaTonya concludes that she will have to raise the piece rates EDC pays in order to be really sure that mean hourly earnings are above the local average wage.

5. The t-test

If LaTonya had simply wanted to know if EDC's workers earned the same as other workers in the area, she would have used a two-tail test. In that case her hypotheses would have been:

$$H_o: \mu = 8.71 \quad \text{and} \quad H_a: \mu \neq 8.71$$

Using $\alpha = .10$, LaTonya would split the .10 between the two tails since the data supports H_a : if the sample t-score is either large and negative or large and positive. Her arithmetic is the same, her sample t-score is still 1.41, but she now will decide that the data supports H_a : only if it outside ± 1.943 .

An alternative to choosing an α

Many researchers now report how unusual the sample t-score would be if the null hypothesis was true rather than choosing an α and stating whether the sample t-score implies the data supports one or the other of the hypotheses based on that α . When a researcher does this, he is essentially letting the reader of his report decide how much risk to take of making which kind of mistake. There are even two ways to do this. If you look at the portion of the t-table reproduced above, you will see that it is not set up very well for this purpose; if you wanted to be able to find out what part of a t-distribution was above any t-score, you would need a table that listed many more t-scores. Since the t-distribution varies as the df changes, you would really need a whole series of t-tables, one for each df.

The old-fashioned way of making the reader decide how much of which risk to take is to not state an α in the body of your report, but only give the sample t-score in the main text. To give the reader some guidance, you look at the usual t-table and find the smallest α , say it is .01, that has a t-value less than the one you computed for the sample. Then write a footnote saying "the data supports the alternative hypothesis for any $\alpha > .01$ ".

The more modern way uses the capability of a computer to store lots of data. Many statistical software packages store a set α detailed t-tables, and when a t-score is computed, the package has the computer look up exactly what proportion of samples would have t-scores larger than the one for your sample. Exhibit 2 shows the computer output for LaTonya's problem from a typical statistical package. Notice that the program gets the same t-score that LaTonya did, it just goes to more decimal places. Also notice that it shows something called the "P value". The P value is the proportion of t-scores that are larger than the one just computed. Looking at the example, the computed t statistic is 1.41188 and the P value is 0.1038. This means that if there are 6 df, a little over 10 per cent of samples will have a t-score greater than 1.41188. Remember that LaTonya used an $\alpha = .025$ and decided that the data supported H_o ;, the P value of .1038 means that H_o : would be supported for any α less than .1038. Since LaTonya had used $\alpha = .025$, this p value means she does not find support for H_o .

Hypothesis test: Mean
Null Hypothesis: Mean = 8.71
Alternative: greater than
Computed t statistic = 1.41188
P value = 0.1038

Exhibit 8: Output from typical statistical software for LaTonya's problem

The P-value approach is becoming the preferred way to The P-value presents research results to audiences of professional researchers. Most of the statistical research conducted for a business firm will be used directly for decision making or presented to an audience of executives to aid them in making a decision. These audiences will generally not be interested in deciding for themselves which hypothesis the data supports. When you are making a presentation of results to your boss, you will want to simply state which hypothesis the evidence supports. You may decide by using either the traditional α approach or the more modern P-value approach, but deciding what the evidence says is probably your job.

Another t-test: do these two (independent) samples come from populations with the same mean?

One of the other statistics that has a sampling distribution that follows the t-distribution is the difference between two sample means. If samples of one size (n_1) are taken from one normal population and samples of another size (n_2) are taken from another normal population (and the populations have the same standard deviation), then a statistic based on the difference between the sample means and the difference between the population means is distributed like t with $n_1 + n_2 - 2$ degrees of freedom. These samples are independent because the members in one sample do not affect which members are in the other sample. You can choose the samples independently of each other, and the two samples do not need to be the same size. The t- statistic is:

$$t = \frac{(\bar{x}_1 - \bar{x}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{s^2}{n_1} + \frac{s^2}{n_2}}}$$

where: \bar{x}_i = the mean of sample i

μ_i = the mean of population i

s^2 = the pooled variance

n_i = the size of sample i.

The usual case is to test to see if the samples come from populations with the same mean, the case where $(\mu_1 - \mu_2) = 0$. The pooled variance is simply a weighted average of the two sample variances, with the weights based on the sample sizes. This means that you will have to calculate the pooled variance before you calculate the t-score. The formula for pooled variance is:

$$s^2 = \frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{n_1 + n_2 - 2}$$

To use the pooled variance t-score, it is necessary to assume that the two populations have equal variances. If you are wondering about why statisticians make a strong assumption in order to use such a complicated formula, it is because the formula that does not need the assumption of equal variances is even more complicated, and reduces the degrees of freedom in the final statistic. In any case, unless you have small samples, the amount of arithmetic needed means that you will probably want to use a statistical software package for this test. You should also note

5. The t-test

that you can test to see if two samples come from populations that are any hypothesized distance apart by setting $(\mu_1 - \mu_2)$ equal to that distance.

An article in *U. S. News and World Report* (Nov. 1993) lamenting grade inflation in colleges states that economics grades have not been inflated as much as most other grades. Nora Alston chairs the Economics Department at Oaks College, and the dean has sent her a copy of the article with a note attached saying "Is this true here at Oaks? Let me know." Dr Alston is not sure if the Dean would be happier if economics grades were higher or lower than other grades, but the article claims that economics grades are lower. Her first stop is the Registrar's office.

She has the clerk in that office pick a sample of 10 class grade reports from across the college spread over the past three semesters. She also has the clerk pick out a sample of 10 reports for economics classes. She ends up with a total of 38 grades for economics classes and 51 grades for other classes. Her hypotheses are:

$$H_o : \mu_{econ} - \mu_{other} \geq 0$$

$$H_a : \mu_{econ} - \mu_{other} < 0$$

She decides to use $\alpha = .05$.

This is a lot of data, and Dr Alston knows she will want to use the computer to help. She initially thought she would use a spreadsheet to find the sample means and variances, but after thinking a minute, she decided to use a statistical software package. The one she is most familiar with is one called SAS. She loads SAS onto her computer, enters the data, and gives the proper SAS commands. The computer gives her the output in Exhibit 8.

The SAS System						
TTFST Procedure						
Variable: GRADE						
Dept	N	Mean	Dev	Std Error	Minimum	Maximum
Econ	38	2.28947	1.01096	0.16400	0	4.00000
Variance	T	DF	Prob>[T]			
Unequal	-2.3858	85.1	0.0193			
Equal	-2.3345	87.0	0.0219			
For HO: Variances are equal, F=1.35 DF=[58.37] Prob>F=0.3485						

Exhibit 9: SAS system software output for Dr Alston's grade study

Dr Alston has 87 df, and has decided to use a one-tailed, left tail test with $\alpha = .05$. She goes to her t-table and finds that 87 df does not appear, the table skipping from 60 to 120 df. There are two things she could do. She could try to interpolate the t-score that leaves .05 in the tail with 87 df, or she could choose between the t-value for 60 and 120 in a conservative manner. Using the conservative choice is the best initial approach, and looking at her table she sees that for 60 df .05 of t-scores are less than -1.671, and for 120 df, .05 are less than -1.658. She does not want to conclude that the data supports economics grades being lower unless her sample t-score is far from zero, so she

decides that she will accept H_a : if her sample t is to the left of -1.671 . If her sample t happens to be between -1.658 and -1.671 , she will have to interpolate.

Looking at the SAS output, Dr Alston sees that her t -score for the equal variances formula is -2.3858 , which is well below -1.671 . She concludes that she will tell the dean that economics grades are lower than grades elsewhere at Oaks College.

Notice that SAS also provides the t -score and df for the case where equal variances are not assumed in the "unequal" line. SAS also provides a P value, but it is for a two-tail test because it gives the probability that a t with a larger absolute value, $>|T|$, occurs. Be careful when using the p values from software: notice if they are one-tail or two-tail p -values before you make your report!

A third t-test: do these (paired) samples come from the sample population?

Managers are often interested in "before and after" questions. As a manager or researcher you will often want to look at "longitudinal" studies, studies that ask about what has happened to an individual as a result of some treatment or across time. Are they different after than they were before? For example, if your firm has conducted a training program you will want to know if the workers who participated became more productive. If the work area has been re-arranged, do workers produce more than before? Though you can use the difference of means test developed earlier, this is a different situation. Earlier, you had two samples that were chosen independently of each other; you might have a sample of workers who received the training and a sample of workers who had not. The situation for this test is different; now you have a sample of workers and for each worker you have measured their productivity before the training or re-arrangement of the work space and you have measured their productivity after. For each worker you have a pair of measures, before and after. Another way to look at this is that for each member of the sample you have a difference between before and after.

You can test to see if these differences equal zero, or any other value, because a statistic based on these differences follows the t -distribution for $n-1$ df when you have n matched pairs. That statistic is:

$$t = \frac{\bar{D} - \delta}{\frac{s_D}{\sqrt{n}}}$$

where: \bar{D} = the mean of the differences in the pairs in the sample

δ = the mean of the differences in the pairs in the population

s_D = the standard deviation of the differences in the sample

n = the number of pairs in the sample.

It is a good idea to take a minute and figure out this formula. There are paired samples and the differences in those pairs, the D 's, are actually a population. The mean of those D 's is δ . Any sample of pairs will also yield a sample of D 's. If those D 's are normally distributed, then the t -statistic in the formula above will follow the t -distribution. If you think of the D 's as the same as x 's in the t -formula at the beginning of the chapter, and think of δ as the population mean, you should realize that this formula is really just that basic t formula.

Lew Podolsky is division manager for Dairyland Lighting, a manufacturer of outdoor lights for parking lots, barnyards, and playing fields. Dairyland Lighting organizes its production work by teams. The size of the team

5. The t-test

varies somewhat with the product being assembled, but there are usually three to six in a team, and a team usually stays together for a few weeks assembling the same product. Dairyland Lighting has a branch plant in the US state of Arizona that serves their west coast customers and Lew has noticed that productivity seems to be lower in Arizona during the summer, a problem that does not occur at the main plant in the US city of Green Bay, Wisconsin. After visiting the Arizona plant in July, August, and November, and talking with the workers during each visit, Lew suspects that the un-air conditioned plant just gets too hot for good productivity. Unfortunately, it is difficult to directly compare plant-wide productivity at different times of the year because there is quite a bit of variation in the number of employees and product mix across the year. Lew decides to see if the same workers working on the same products are more productive on cool days than hot days by asking the local manager, Dave Mueller, to find a cool day and a hot day from last fall and choose ten work teams who were assembling the same products on the two days. Dave sends Lew the following data:

Team leader	Output—cool day	Output—hot day	Difference (cool-hot)
	October 14	October 20	
Martinez	153	149	4
McAlan	167	170	-3
Wilson	164	155	9
Burningtree	183	179	4
Sanchez	177	167	10
Lilly	162	150	12
Cantu	165	158	7

Exhibit 10: Lew Podolsky's data for the air-conditioning decision

Lew decides that if the data support productivity being higher of cool days, he will call in a heating/air-conditioning contractor to get some cost estimates so that he can decide if installing air conditioning in the Arizona plant is cost-effective. Notice that he has matched pairs data--for each team he has production on October 14, a cool day, and on October 20, a hot day. His hypotheses are:

$$H_0: \delta \leq 0 \text{ and } H_a: \delta > 0$$

Using $\alpha = .05$ in this one-tail test, Lew will decide to call the engineer if his sample t-score is greater than 1.943, since there are 6 df. This sample is small, so it is just as easy to do the computations on a calculator. Lew finds:

$$\bar{D} = 6.1428$$

$$s_D = 5.0142$$

and his sample t-score is:

$$t = \frac{\bar{D} - \delta}{\frac{s_D}{\sqrt{n}}} = \frac{6.14 - 0}{\frac{5.01}{\sqrt{7}}}$$

$$t = \frac{6.14}{1.89} = 3.24$$

Because his sample t-score is greater than 1.943, Lew gets out the telephone book and looks under air conditioning contractors to call for some estimates.

Summary

The t-tests are commonly used hypothesis tests. Researchers often find themselves in situations where they need to test to see if a sample comes from a certain population, and therefore test to see if the sample probably came from a population with that certain mean. Even more often, researchers will find themselves with two samples and want to know if the samples come from the same population, and will test to see if the samples probably come from populations with the same mean. Researchers also frequently find themselves asking if two sets of paired samples have equal means. In any case, the basic strategy is the same as for any hypothesis test. First, translate the question into null and alternative hypotheses, making sure that the null hypothesis includes an equal sign. Second, choose α . Third, compute the relevant statistics, here the t-score, from the sample or samples. Fourth, using the tables, decide if the sample statistic leads you to conclude that the sample came from a population where the null hypothesis is true or a population where the alternative is true.

The t-distribution is also used in testing hypotheses in other situations since there are other sampling distributions with the same t-distribution shape. So, remember how to use the t-tables for later chapters.

Statisticians have also found how to test to see if three or more samples come from populations with the same mean. That technique is known as "one-way analysis of variance". The approach used in analysis of variance is quite different from that used in the t-test. It will be covered in chapter, "The F-test and One-Way ANOVA".

6. F-test and one-way anova

Years ago, statisticians discovered that when pairs of samples are taken from a normal population, the ratios of the variances of the samples in each pair will always follow the same distribution. Not surprisingly, over the intervening years, statisticians have found that the ratio of sample variances collected in a number of different ways follow this same distribution, the F-distribution. Because we know that sampling distributions of the ratio of variances follow a known distribution, we can conduct hypothesis tests using the ratio of variances.

The F-statistic is simply:

$$F = \frac{S_1^2}{S_2^2}$$

where s_i^2 is the variance of sample i . Remember that the sample variance is:

$$s^2 = \frac{\sum(x - \bar{x})^2}{(n - 1)}$$

Think about the shape that the F-distribution will have. If s_1^2 and s_2^2 come from samples from the same population, then if many pairs of samples were taken and F-scores computed, most of those F-scores would be close to one. All of the F-scores will be positive since variances are always positive—the numerator in the formula is the sum of squares, so it will be positive, the denominator is the sample size minus one, which will also be positive. Thinking about ratios requires some care. If s_1^2 is a lot larger than s_2^2 , F can be quite large. It is equally possible for s_2^2 to be a lot larger than s_1^2 , and then F would be very close to zero. Since F goes from zero to very large, with most of the values around one, it is obviously not symmetric; there is a long tail to the right, and a steep descent to zero on the left.

There are two uses of the F-distribution which will be discussed in this chapter. The first is a very simple test to see if two samples come from populations with the same variance. The second is one-way Analysis Of Variance (ANOVA) which uses the F-distribution to test to see if three or more samples come from populations with the same **mean**.

A simple test: Do these two samples come from populations with the same variance?

Because the F-distribution is generated by drawing two samples from the same normal population, it can be used to test the hypothesis that two samples come from populations with the same variance. You would have two samples (one of size n_1 and one of size n_2), and the sample variance from each. Obviously, if the two variances are very close to being equal the two samples could easily have come from populations with equal variances. Because the F-statistic is the ratio of two sample variances, when the two sample variances are close to

6. F-test and one-way anova

equal, the F-score is close to one. If you compute the F-score, and it is close to one, you accept your hypothesis that the samples come from populations with the same variance.

This is the basic method of the F-test. Hypothesize that the samples come from populations with the same variance. Compute the F-score by finding the ratio of the sample variances. If the F-score is close to one, conclude that your hypothesis is correct and that the samples do come from populations with equal variances. If the F-score is far from one, then conclude that the populations probably have different variances.

The basic method must be fleshed out with some details if you are going to use this test at work. There are two sets of details: first, formally writing hypotheses, and second, using the F-distribution tables so that you can tell if your F-score is close to one or not. Formally, two hypotheses are needed for completeness. The first is the null hypothesis that there is no difference (hence "null"). It is usually denoted as H_o . The second is that there is a difference, and it is called the alternative, and is denoted H_1 : or H_a .

Using the F-tables to decide how close to one is close enough to accept the null hypothesis (truly formal statisticians would say "fail to reject the null"), is fairly tricky, for the F-distribution tables are fairly tricky. Before using the tables, the researcher must decide how much chance he or she is willing to take that the null will be rejected when it is really true. The usual choice is 5 per cent, or as statisticians say, " $\alpha = .05$ ". If more, or less, chance is wanted, α can be varied. Choose your α and go to the F-tables. First notice that there are a number of F-tables, one for each of several different levels of α (or at least a table for each two α 's with the F-values for one α in **bold** type and the values for the other in regular type). There are rows and columns on each F-table, and both are for degrees of freedom. Because two separate samples are taken to compute an F-score and the samples do not have to be the same size, there are two separate degrees of freedom—one for each sample. For each sample, the number of degrees of freedom is $n-1$, one less than the sample size. Going to the table, how do you decide which sample's degrees of freedom (df) is for the row and which is for the column? While you could put either one in either place, you can save yourself a step if you put the sample with the **larger variance** (not necessarily the larger sample) **in the numerator**, and then that sample's df determines the column and the other sample's df determines the row. The reason that this saves you a step is that the tables only show the values of F that leave α in the right tail where $F > 1$, the picture at the top of most F-tables shows that. Finding the critical F value for left tails requires another step, which is outlined in the box below.

Table 1: F Distribution $\alpha = .05$ (rows are df in the numerator, columns are df in denominator)

df	10	20	30	120	infinity
10	2.98	2.77	2.70	2.58	2.54
20	2.35	2.12	2.04	1.90	1.84
30	2.16	1.93	1.84	1.68	1.62
120	1.91	1.66	1.55	1.35	1.25
infinity	1.83	1.57	1.46	1.22	1.00

F-tables are virtually always printed as "one-tail" tables, showing the critical F-value that separates the right tail from the rest of the distribution. In most statistical applications of the F-distribution, only the right tail is of interest, because most applications are testing to see if the variance from a certain source is greater than the variance from another source, so the researcher is interested in finding if the F-score is greater than one. In the test of equal variances, the researcher is interested in finding out if the F-score is *close* to one, so that either a large F-score or a small F-score would lead the researcher to conclude that the variances are not equal. Because the critical F-value that separates the left tail from the rest of the distribution is not printed, and not simply the negative of the printed value, researchers often simply divide the larger sample variance by the smaller sample variance, and use the printed tables to see if the quotient is "larger than one", effectively rigging the test into a one-tail format. For purists, and occasional instances, the left tail critical value can be computed fairly easily.

The left tail critical value for x, y degrees of freedom (df) is simply the inverse of the right tail (table) critical value for y, x df. Looking at an F-table, you would see that the F-value that leaves $\alpha = .05$ in the right tail when there are **10, 20** df is $F=2.35$. To find the F-value that leaves $\alpha = .05$ in the left tail with 10, 20 df, look up $F=2.77$ for $\alpha = .05$, **20, 10** df. Divide one by 2.77, finding .36. That means that 5 per cent of the F-distribution for **10, 20** df is below the critical value of .36, and 5 per cent is above the critical value of 2.35.

Putting all of this together, here is how to conduct the test to see if two samples come from populations with the same variance. First, collect two samples and compute the sample variance of each, s_1^2 and s_2^2 . Second, write your hypotheses and choose α . Third find the F-score from your samples, dividing the larger s^2 by the smaller so that $F > 1$. Fourth, go to the tables, find the table for $\alpha/2$, and find the critical (table) F-score for the proper degrees of freedom ($n-1$ and $n-1$). Compare it to the samples' F-score. If the samples' F is larger than the critical F, the samples' F is not "close to one", and H_a : the population variances are not equal, is the best hypothesis. If the samples' F is less than the critical F, H_o : that the population variances are equal should be accepted.

An example

A young banker has recently been promoted and made the manager of her own branch. After a few weeks, she has discovered that maintaining the correct number of tellers seems to be more difficult than it was when she was assistant manager of a larger branch. Some days, the lines are very long, and other days, the tellers seem to have little to do. She wonders if the number of customers at her new branch is simply more variable than the number of customers at the branch where she used to work. Because tellers work for a whole day or half a day (morning or afternoon), she collects the following data on the number of transactions in a half day from her branch and the branch where she used too work:

Her branch: 156, 278, 134, 202, 236, 198, 187, 199, 143, 165, 223

Old branch: 345, 332, 309, 367, 388, 312, 355, 363, 381

She hypothesizes:

$$H_o: \sigma_h^2 = \sigma_o^2$$

$$H_a: \sigma_h^2 \neq \sigma_o^2$$

She decides to use $\alpha = .05$. She computes the sample variances and finds:

$$s_h^2 = 2027.1$$

6. F-test and one-way anova

$$s_o^2 = 795.2$$

Following the rule to put the larger variance in the numerator, so that she saves a step, she finds:

$$F = \frac{s_k^2}{s_o^2} = \frac{2027.1}{795.2} = 2.55$$

From the table, (remembering to use the $\alpha = .025$ table because the table is one-tail and the test is two-tail) she finds that the critical F for 10,8 df is 4.30. Because her F-score is less than the critical score, she concludes that her F-score is "close to one", and that the variance of customers in her office is the same as it was in the old office. She will need to look further to solve her staffing problem.

Analysis of variance (ANOVA)

The importance of ANOVA

A more important use of the F-distribution is in analyzing variance to see if three or more samples come from populations with equal means. This is an important statistical test, not so much because it is frequently used, but because it is a bridge between univariate statistics and multivariate statistics and because the strategy it uses is one which is used in many multivariate tests and procedures.

One-way ANOVA: Do these three (or more) samples all come from populations with the same mean?

This seems wrong—we will test a hypothesis about means by "analyzing variance". It is not wrong, but rather a really clever insight that some statistician had years ago. This idea—looking at variance to find out about differences in means—is the basis for much of the multivariate statistics used by researchers today. The ideas behind ANOVA are used when we look for relationships between two or more variables, the big reason we use multivariate statistics.

Testing to see if three or more samples come from populations with the same mean can often be a sort of multivariate exercise. If the three samples came from three different factories or were subject to different treatments, we are effectively seeing if there is a difference in the results because of different factories or treatments—is there a relationship between factory (or treatment) and the outcome?

Think about three samples. A group of x's have been collected, and for some good reason (other than their x value) they can be divided into three groups. You have some x's from group (sample) 1, some from group (sample) 2, and some from group (sample) 3. If the samples were combined, you could compute a "grand mean" and a "total variance" around that grand mean. You could also find the mean and (sample) "variance within" each of the groups. Finally, you could take the three sample means, and find the "variance between" them. ANOVA is based on analyzing where the "total" variance comes from. If you picked one x, the source of its variance, its distance from the grand mean would have two parts: (1) how far it is from the mean of its sample, and (2) how far its sample's mean is from the grand mean. If the three samples really do come from populations with different means, then for most of the x's, the distance between the sample mean and the grand mean will probably be greater than the distance between the x and its group mean. When these distances are gathered together and turned into variances, you can see that if the population means are different, the variance between the sample means is likely to be greater than the variance within the samples.

By this point in the book, it should not surprise you to learn that statisticians have found that if three or more samples are taken from a normal population, and the variance between the samples is divided by the variance within the samples, a sampling distribution formed by doing that over and over will have a known shape. In this case it will be distributed like F with $m-1$, $n-m$ df, where m is the number of samples and n is the size of the m samples altogether. "Variance between" is found by:

$$s_b^2 = \frac{\sum_{j=1}^m n_j (\bar{x}_j - \bar{X})^2}{m-1}$$

where \bar{x}_j is the mean of sample j , and \bar{X} is the "grand mean".

The numerator of the variance between is the sum of the squares of the distance between each \bar{x} 's sample mean and the grand mean. It is simply a summing of one of those sources of variance across all of the observations.

The "variance within" is found by:

$$s_w^2 = \frac{\sum_{j=1}^m \sum_{i=1}^{n_j} (x_{ij} - \bar{x}_j)^2}{n - m}$$

Double sums need to be handled with care. First (operating on the "inside" or second sum sign) find the mean of each sample and the sum of the squares of the distances of each x in the sample from its mean. Second (operating on the "outside" sum sign), add together the results from each of the samples.

The strategy for conducting a one-way analysis of variance is simple. Gather m samples. Compute the variance between the samples, the variance within the samples, and the ratio of between to within, yielding the F-score. If the F-score is less than one, or not much greater than one, the variance between the samples is no greater than the variance within the samples and the samples probably come from populations with the same mean. If the F-score is much greater than one, the variance between is probably the source of most of the variance in the total sample, and the samples probably come from populations with different means.

The details of conducting a one-way ANOVA fall into three categories: (1) writing hypotheses, (2) keeping the calculations organized, and (3) using the F-tables. The null hypothesis is that all of the population means are equal, and the alternative is that not all of the means are equal. Quite often, though two hypotheses are really needed for completeness, only H_o is written:

$$H_o : \mu_1 = \mu_2 = \dots = \mu_m$$

Keeping the calculations organized is important when you are finding the variance within. Remember that the variance within is found by squaring, and then summing, the distance between each observation and the **mean of its sample**. Though different people do the calculations differently, I find the best way to keep it all straight is to find the sample means, find the squared distances in each of the samples, and then add those together. It is also important to keep the calculations organized in the final computing of the F-score. If you remember that the goal is to see if the variance **between** is large, then it's easy to remember to divide variance between by variance within.

6. F-test and one-way anova

Using the F-tables is the second detail. Remember that F-tables are one-tail tables and that ANOVA is a one-tail test. Though the null hypothesis is that all of the means are equal, you are testing that hypothesis by seeing if the variance between is less than or equal to the variance within. The number of degrees of freedom is $m-1$, $n-m$, where m is the number of samples and n is the total size of all the samples together.

An example

The young bank manager in the last example is still struggling with finding the best way to staff her branch. She knows that she needs to have more tellers on Fridays than on other days, but is trying to find if the need for tellers is constant across the rest of the week. She collects data for the number of transactions each day for two months. Here is her data:

Mondays: 276, 323, 298, 256, 277, 309, 312, 265, 311

Tuesdays: 243, 279, 301, 285, 274, 243, 228, 298, 255

Wednesdays: 288, 292, 310, 267, 243, 293, 255, 273

Thursdays: 254, 279, 241, 227, 278, 276, 256, 262

She tests the null hypothesis:

$$H_o: m_m = m_{tu} = m_w = m_{th}$$

and decides to use $\alpha = .05$. She finds:

$$m = 291.8$$

$$tu = 267.3$$

$$w = 277.6$$

$$th = 259.1$$

and the grand mean, = 274.3

She computes variance within:

$$[(276-291.8)^2 + (323-291.8)^2 + \dots + (243-267.6)^2 + \dots + (288-277.6)^2 + \dots + (254-259.1)^2] / [34-4] = 15887.6 / 30 = 529.6$$

Then she computes variance between:

$$[9(291.8-274.3)^2 + 9(267.3-274.3)^2 + 8(277.6-274.3)^2 + 8(259.1-274.3)^2] / [4-1]$$

$$= 5151.8 / 3 = 1717.3$$

She computes her F-score:

$$F = \frac{1717.3}{529.6} = 3.24$$

Consulting the F-tables for $\alpha = .05$ and 3, 30 df, she finds that the critical F-value is 2.92. Because her F-score is larger than the critical F-value, she concludes that the mean number of transactions is not the equal on different days of the week. She will want to adjust her staffing so that she has more tellers on some days than on others.

Summary

The F-distribution is the sampling distribution of the ratio of the variances of two samples drawn from a normal population. It is used directly to test to see if two samples come from populations with the same variance. Though you will occasionally see it used to test equality of variances, the more important use is in analysis of variance,

(ANOVA). ANOVA, at least in its simplest form as presented in this chapter, is used to test to see if three or more samples come from populations with the same mean. By testing to see if the variance of the observations comes more from the variation of each observation from the mean of its sample or from the variation of the means of the samples from the grand mean, ANOVA tests to see if the samples come from populations with equal means or not.

Connections

ANOVA has more elegant forms that appear in later chapters. It also forms the basis for regression analysis, a statistical technique that has many business applications; it is covered in later chapters. The F-tables are also used in testing hypotheses about regression results.

This is also the beginning of multivariate statistics. Notice that in the one-way ANOVA, each observation is for two variables: the "x" variable and the group of which the observation is a part. In later chapters, observations will have two, three, or more variables.

The F-test for equality of variances is sometimes used before using the t-test for equality of means because the t-test, at least in the form presented in this text, requires that the samples come from populations with equal variances. You will see it used along with t-tests when the stakes are high or the researcher is a little compulsive.

7. Some non-parametric tests

Remember that you use statistics to make inferences about populations from samples. Most of the techniques statisticians use require that two assumptions are met. First, the population that the sample comes from is normal. Second, whenever means and variances were computed, the numbers in the data are "cardinal" or "interval", meaning that the value given an observation not only tells you which observation is larger or smaller, but how much larger or smaller. There are many situations when these assumptions are not met, and using the techniques developed so far will not be appropriate. Fortunately, statisticians have developed another set of statistical techniques, non-parametric statistics, for these situations. Three of these tests will be explained in this chapter. These three are the Mann-Whitney U-Test, which tests to see if two independently chosen samples come from populations with the same location; the Wilcoxon Rank Sum Test, which tests to see if two paired samples come from populations with the same location; and Spearman's Rank Correlation, which tests to see if two variables are related.

What does "non-parametric" mean?

To a statistician, a parameter is a measurable characteristic of a population. The population characteristics that usually interest statisticians are the location and the shape. Non-parametric statistics are used when the parameters of the population are not measurable or do not meet certain standards. In cases when the data only orders the observations, so that the interval between the observations is unknown, neither a mean nor a variance can be meaningfully computed. In such cases you need to use non-parametric tests. Because your sample does not have cardinal, or interval, data you cannot use it to estimate the mean or variance of the population, though you can make other inferences. Even if your data is cardinal, the population must be normal before the shape of the many sampling distributions are known. Fortunately, even if the population is not normal, such sampling distributions are usually close to the known shape if large samples are used. In that case, using the usual techniques is acceptable. However, if the samples are small and the population is not normal, you have to use non-parametric statistics. As you know, "there is no such thing as a free lunch". If you want to make an inference about a population without having cardinal data, or without knowing that the population is normal, or with very small samples, you will have to give up something. In general, non-parametric statistics are less precise than parametric statistics. Because you know less about the population you are trying to learn about, the inferences you make are less exact.

When either (1) the population is not normal and the samples are small, or (2) when the data is not cardinal, the same non-parametric statistics are used. Most of these tests involve ranking the members of the sample, and most involve comparing the ranking of two or more samples. Because we cannot compute meaningful sample statistics to compare to a hypothesized standard, we end up comparing two samples.

7. Some non-parametric tests

Do these populations have the same location? The Mann-Whitney U test.

In the chapter “T-test”, you learned how to test to see if two samples came from populations with the same mean by using the t-test. If your samples are small and you are not sure if the original populations are normal, or if your data does not measure intervals, you cannot use that t-test because the sample t-scores will not follow the sampling distribution in the t-table. Though there are two different data problems that keep you from using the t-test, the solution to both problems is the same, the non-parametric Mann-Whitney U test. The basic idea behind the test is to put the samples together, rank the members of the combined sample, and then see if the two samples are mixed together in the common ranking.

Once you have a single ranked list containing the members of both samples, you are ready to conduct a Mann-Whitney U test. This test is based on a simple idea. If the first part of the combined ranking is largely made up of members from one sample, and the last part is largely made up of members from the other sample, then the two samples are probably from populations with different “averages” and therefore different locations. You can test to see if the members of one sample are lumped together or spread through the ranks by adding up the ranks of each of the two groups and comparing the sums. If these “rank sums” are about equal, the two groups are mixed together. If these rank sums are far from equal, each of the samples is lumped together at the beginning or the end of the overall ranking.

Willy Senn works for Old North Gadgets, a maker and marketer of computer peripherals aimed at scientists, consultants, and college faculty. Old North's home office and production facilities are in a small town in the US state of Maine. While this is a nice place to work, the firm wants to expand its sales and needs a sales office in a location closer to potential customers and closer to a major airport. Willy has been given the task of deciding where that office should be. Before he starts to look at office buildings and airline schedules, he needs to decide if Old North's potential customers are in the east or the west. Willy finds an article in Fortune magazine that lists the best cities for finding “knowledge workers”, Old North's customers. That article lists the ten best cities in the United States.

Rank	Metro Area	Region
1	Raleigh-Durham	East
2	New York	East
3	Boston	East
4	Seattle	West
5	Austin	West
6	Chicago	East
7	Houston	West
8	San Jose	West
9	Philadelphia	East
10	Minnesota-St Paul	East

Exhibit 11: Data for Willy's problem. From Kenneth Labich, "The Best Cities for Knowledge Workers," *Fortune*, 128:12, Nov. 15, 1993, pp. 50 ff.

Six of the top ten are in the east and four are in the west, but these ten represent only a sample of the market. It looks like the eastern places tend to be higher in the top ten, but is that really the case? If you add up the ranks, the six eastern cities have a "rank sum" of 31 while the western cities have a rank sum of 24, but there are more eastern cities and even if there were the same number would that difference be due to a different "average" in the rankings, or is it just due to sampling? The Mann-Whitney U test can tell you if the rank sum of 31 for the eastern cities is significantly less than would be expected if the two groups really were about the same and six of the ten in the sample happened to be from one group. The general formula for computing the Mann-Whitney U for the first of two groups is:

$$U_1 = n_1 n_2 + [n_1(n_1+1)]/2 - T_1$$

where:

T_1 = the sum of the ranks of group 1.

n_1 = the number of members of the sample from group 1

n_2 = the number of members of the sample from group 2.

This formula seems strange at first, but a little careful thought will show you what is going on. The last third of the formula, $-T_1$, subtracts the rank sum of the group from the rest of the formula. What is the first two-thirds of the formula? The bigger the total of your two samples, and the more of that total that is in the first group, the bigger you would expect T_1 to be, everything else equal. Looking at the first two-thirds of the formula, you can see that the only variables in it are n_1 and n_2 , the sizes of the two samples. The first two-thirds of the formula depends on the how big the total group is and how it is divided between the two samples. If either n_1 or n_2 gets larger, so does this part of the formula. The first two-thirds of the formula is the maximum value for T_1 , the rank sum of group 1. T_1 will be at its maximum if the members of the first group were all at the bottom of the rankings for the combined samples. The U_1 score then is the difference between the actual rank sum and the maximum possible. A bigger U_1

7. Some non-parametric tests

means that the members of group 1 are bunched more at the top of the rankings and a smaller U_1 means that the members of group 1 are bunched near the bottom of the rankings so that the rank sum is close to its maximum. Obviously, a U-score can be computed for either group, so there is always a U_1 and a U_2 . If U_1 is larger, U_2 is smaller for a given n_1 and n_2 because if T_1 is smaller, T_2 is larger.

What should Willy expect if the best cities are in one region rather than being evenly distributed across the country? If the best cities are evenly distributed, then the eastern group and the western group should have U's that are close together since neither group will have a T that is close to either its minimum or its maximum. If the one group is mostly at the top of the list, then that group will have a large U since its T will be small, and the other group will have a smaller U since its T will be large. $U_1 + U_2$ is always equal to $n_1 n_2$ so either one can be used to test the hypothesis that the two groups come from the same population. Though there is always a pair of U-scores for any Mann-Whitney U-test, the published tables only show the smaller of the pair. Like all of the other tables you have used, this one shows what the sampling distribution of U's is like.

The sampling distribution, and this test, were first described by HB Mann and DR Whitney in 1947.¹ While you have to compute both U-scores, you only use the smaller one to test a two-tailed hypothesis. Because the tables only show the smaller U, you need to be careful when conducting a one-tail test. Because you will accept the alternative hypothesis if U is very small, you use the U computed for that sample which H_a : says is farther down the list. You are testing to see if one of the samples is located to the right of the other, so you test to see if the rank sum of that sample is large enough to make its U small enough to accept H_a :. If you learn to think through this formula, you will not have to memorize all of this detail because you will be able to figure out what to do.

Let us return to Willy 's problem. He needs to test to see if the best cities in which to locate the sales office, the best cities for finding "knowledge workers", are concentrated in one part of the country or not. He can attack his problem with a hypothesis test using the Mann-Whitney U-test. His hypotheses are:

H_0 : The distributions of eastern and western city rankings among the "best places to find knowledge workers" are the same.

H_a : The distributions are different.

Looking at the table of Mann-Whitney values, he finds the following if one of the n's is 6:

	n_1			
U_0	1	2	3	4
0	0.1429	0.0357	0.0119	0.0005
1	0.2857	0.0714	0.0238	0.0095
2	0.4286	0.1429	0.0476	0.0190
3	0.5714	0.2143	0.0833	0.0333
4		0.4286	0.1310	0.0571
5		0.5714	0.1905	0.0857
6			0.2738	0.1286
7			0.3571	0.1762
8			0.4524	0.2381

1 ["On a test of whether one or two random variables is stochastically larger than the other." *Annals of Mathematical Statistics*, 18, 50-60.].

9	0.5476	0.3048
10		0.3810

Exhibit 12: Some lower-tail values for the Mann Whitney U statistic

The values in the table show what portion of the sampling distribution of U-statistics is in the lower tail, below the U value in the first column, if the null hypothesis is true. Willy decides to use an $\alpha = .10$. Since he will decide that the data supports H_a if either the east or the west has a small U, Willy has a two-tail test and needs to divide his α between the two tails. He will choose H_a if either U is in the lowest .05 of the distribution. Going down the column for the other n equal to 4, Willy finds that if the null hypothesis is true, the probability that the smaller of the two U-scores will be 4 or less is only .0571, and probability that the lower U-score will be 3 or less is .0333. His half α of .05 is between these two, so he decides to be conservative and use as a decision rule to conclude that the data supports H_a : The distributions are different, if his sample U is less than 3 and that the data supports H_0 : the distributions are the same, if his U is greater than or equal to 3. Now he computes his U, finding both U_e and U_w .

Remembering the formula from above, he finds his two U values::

For the eastern cities:

$$U_e = 6 \times 4 + \frac{6 \times 7}{2} - 31 = 14$$

For the western cities:

$$U_w = 6 \times 4 + \frac{4 \times 5}{2} - 24 = 10$$

The smaller of his two U-scores is $U_w = 10$. Because 10 is larger than 3, his decision rule tells him that the data supports the null hypothesis that eastern and western cities rank about the same. Willy decides that the sales office can be in either an eastern or western city, at least based on locating the office close to near large numbers of knowledge workers. The decision will depend on office cost and availability and airline schedules.

Testing with matched pairs: the Wilcoxon signed ranks test

During your career, you will often be interested in finding out if the same population is different in different situations. Do the same workers perform better after a training session? Do customers who used one of your products prefer the "new improved" version? Are the same characteristics important to different groups? When you are comparing the same group in two different situations, you have "matched pairs". For each member of the population or sample you have what happened under two different sets of conditions.

There is a non-parametric test using matched pairs that allows you to see if the location of the population is different in the different situations. This test is the Wilcoxon Signed Ranks Test. To understand the basis of this test, think about a group of subjects who are tested under two sets of conditions, A and B. Subtract the test score under B from the test score under A for each subject. Rank the subjects by the absolute size of that difference, and look to see if those who scored better under A are mostly lumped together at one end of your ranking. If most of the biggest absolute differences belong to subjects who scored higher under one of the sets of conditions, then the subjects probably perform differently under A than under B.

The details of how to perform this test were published by Frank Wilcoxon in 1945². Wilcoxon found a method to find out if the subjects who scored better under one of the sets of conditions were lumped together or not. He also

² "Individual comparisons by ranking methods", *Biometrics*, 1, 80-83

7. Some non-parametric tests

found the sampling distribution needed to test hypotheses based on the rankings. To use Wilcoxon's test, collect a sample of matched pairs. For each subject, find the difference in the outcome between the two sets of conditions and then rank the subjects according to the absolute value of the differences. Next, add together the ranks of those with negative differences and add together the ranks of those with positive differences. If these rank sums are about the same, then the subjects who did better under one set of conditions are mixed together with those who did better under the other condition, and there is no difference. If the rank sums are far apart, then there is a difference between the two sets of conditions.

Because the sum of the rank sums is always equal to $[N(N-1)]/2$, if you know the rank sum for either the positives or the negatives, you know it for the other. This means that you do not really have to compare the rank sums, you can simply look at the smallest and see if it is very small to see if the positive and negative differences are separated or mixed together. The sampling distribution of the smaller rank sums when the populations the samples come from are the same was published by Wilcoxon. A portion of a table showing this sampling distribution is in Exhibit 3. See below.

one-tail significance	0.05	0.025	0.01
two-tail significance	0.1	0.05	0.02
number of pairs, N			
5	0		
6	2	0	
7	3	2	0
8	5	3	1
9	8	5	3
10	10	8	5

Exhibit 13: Sampling distribution

Wendy Woodruff is the President of the Student Accounting Society at the University of North Carolina at Burlington (UNC-B). Wendy recently came across a study by Baker and McGregor ["Empirically Assessing the Utility of Accounting Student Characteristics", unpublished, 1993] in which both accounting firm partners and students were asked to score the importance of student characteristics in the hiring process. A summary of their findings is in Exhibit 11.

ATTRIBUTE	Mean: student rating	Mean: big firm rating
High Accounting GPA	2.06	2.56
High Overall GPA	0.08	-0.08
Communication Skills	4.15	4.25
Personal Integrity	4.27	7.5
Energy, drive, enthusiasm	4.82	3.15
Appearance	2.68	2.31

Exhibit 14: Data on importance of student attributes. From Baker and McGregor.

Wendy is wondering if the two groups think the same things are important. If the two groups think that different things are important, Wendy will need to have some society meetings devoted to discussing the differences. Wendy

has read over the article, and while she is not exactly sure how Baker and McGregor's scheme for rating the importance of student attributes works, she feels that the scores are probably not distributed normally. Her test to see if the groups rate the attributes differently will have to be non-parametric since the scores are not normally distributed and the samples are small. Wendy uses the Wilcoxon Signed Ranks Test.

Her hypotheses are:

H_0 : There is no true difference between what students and Big 6 partners think is important.

H_a : There is a difference.

She decides to use a level of significance of .05. Wendy's test is a two-tail test because she wants to see if the scores are different, not if the Big 6 partners value these things more highly. Looking at the table, she finds that, for a two-tail test, the smaller of the two sum of ranks must be less than or equal to 2 to accept H_a .

Wendy finds the differences between student and Big 6 scores, and ranks the absolute differences, keeping track of which are negative and which are positive. She then sums the positive ranks and sum the negative ranks. Her work is shown below:

ATTRIBUTE	Mean student rating	Mean big firm rating	Difference	Rank
High Accounting GPA	2.06	2.56	-0.5	-4
High Overall GPA	0.08	-0.08	0.16	2
Communication Skills	4.15	4.25	-0.1	-1
Personal Integrity	4.27	7.5	-2.75	-6
Energy, drive, enthusiasm	4.82	3.15	1.67	5
Appearance	2.68	2.31	0.37	3

sum of positive ranks = 4+5+3=10

sum of negative ranks = 4+1=6=11

number of pairs=6

Exhibit 15: The worksheet for the Wilcoxon Signed Ranks Test

Her sample statistic, T, is the smaller of the two sums of ranks, so $T=10$. According to her decision rule to accept H_a : if $T < 2$, she decides that the data supports H_0 : that there is no difference in what students and Big 6 firms think is important to look for when hiring students. This makes sense, because the attributes that students score as more important, those with positive differences, and those that the Big 6 score as more important, those with negative differences, are mixed together when the absolute values of the differences are ranked. Notice that using the rankings of the differences rather than the size of the differences reduces the importance of the large difference between the importance students and Big 6 partners place on Personal integrity. This is one of the costs of using non-parametric statistics. The Student Accounting Society at UNC-B does not need to have a major program on what accounting firms look for in hiring. However, Wendy thinks that the discrepancy in the importance in hiring placed on Personal Integrity by Big 6 firms and the students means that she needs to schedule a speaker on that subject. Wendy wisely tempers her statistical finding with some common sense.

7. Some non-parametric tests

Are these two variables related? Spearman's rank correlation

Are sales higher in those geographic areas where more is spent on advertising? Does spending more on preventive maintenance reduce down-time? Are production workers with more seniority assigned the most popular jobs? All of these questions ask how the two variables move up and down together; when one goes up, does the other also rise? when one goes up does the other go down? Does the level of one have no effect on the level of the other? Statisticians measure the way two variables move together by measuring the **correlation coefficient** between the two.

Correlation will be discussed again in the next chapter, but it will not hurt to hear about the idea behind it twice. The basic idea is to measure how well two variables are tied together. Simply looking at the word, you can see that it means co-related. If whenever variable X goes up by 1, variable Y changes by a set amount, then X and Y are perfectly tied together, and a statistician would say that they are perfectly correlated. Measuring correlation usually requires interval data from normal populations, but a procedure to measure correlation from ranked data has been developed. Regular correlation coefficients range from -1 to +1. The sign tells you if the two variables move in the same direction (positive correlation) or in opposite directions (negative correlation) as they change together. The absolute value of the correlation coefficient tells you how closely tied together the variables are; a correlation coefficient close to +1 or to -1 means they are closely tied together, a correlation coefficient close to 0 means that they are not very closely tied together. The non-parametric Spearman's Rank Correlation Coefficient is scaled so that it follows these same conventions.

The true formula for computing the Spearman's Rank Correlation Coefficient is complex. Most people using rank correlation compute the coefficient with a computer program, but looking at the equation will help you see how Spearman's Rank Correlation works. It is:

$$r_s = 1 - \left(\frac{6}{n(n^2 - 1)} \right) (\sum d^2)$$

where:

n = the number of observations

d = the difference between the ranks for an observation

Keep in mind that we want this non-parametric correlation coefficient to range from -1 to +1 so that it acts like the parametric correlation coefficient. Now look at the equation. For a given sample size, n, the only thing that will vary is $\sum d^2$. If the samples are perfectly positively correlated, then the same observation will be ranked first for both variables, another observation ranked second for both variables, etc. That means that each difference in ranks, d, will be zero, the numerator of the fraction at the end of the equation will be zero, and that fraction will be zero. Subtracting zero from one leaves one, so if the observations are ranked in the same order by both variables, the Spearman's Rank Correlation Coefficient is +1. Similarly, if the observations are ranked in exactly the opposite order by the two variables, there will many large d's, and $\sum d^2$ will be at its maximum. The rank correlation coefficient should equal -1, so you want to subtract 2 from 1 in the equation. The middle part of the equation, $6/n(n^2-1)$, simply scales $\sum d^2$ so that the whole term equals 2. As n grows larger, $\sum d^2$ will grow larger if the two variables produce exactly opposite rankings. At the same time, $n(n^2-1)$ will grow larger so that $6/n(n^2-1)$ will grow smaller.

Colonial Milling Company produces flour, corn meal, grits, and muffin, cake, and quickbread mixes. They are considering introducing a new product, Instant Cheese Grits mix. Cheese grits is a dish made by cooking grits, combining the cooked grits with cheese and eggs, and then baking the mixture. It is a southern favorite in the United States, but because it takes a long time to cook, is not served much anymore. The Colonial mix will allow someone to prepare cheese grits in 20 minutes in only one pan, so if it tastes right, it should be a good-selling product in the South. Sandy Owens is the product manager for Instant Cheese Grits, and is deciding what kind of cheese flavoring to use. Nine different cheese flavorings have been successfully tested in production, and samples made with each of those nine flavorings have been rated by two groups: first, a group of food experts, and second, a group of potential customers. The group of experts was given a taste of three dishes of "homemade" cheese grits and ranked the samples according to how well they matched the real thing. The customers were given the samples and asked to rank them according to how much they tasted like "real cheese grits should taste". Over time, Colonial has found that using experts is a better way of identifying the flavorings that will make a successful product, but they always check the experts' opinion against a panel of customers. Sandy must decide if the experts and customers basically agree. If they do, then she will use the flavoring rated first by the experts. The data from the taste tests is in Exhibit 13.

	Expert ranking	Consumer ranking
Flavoring		
NYS21	7	8
K73	4	3
K88	1	4
Ba4	8	6
Bc11	2	5
McA A	3	1
McA A	9	9
WIS 4	5	2
WIS 43	6	7

Exhibit 16: Data from two taste tests of cheese flavorings

Sandy decides to use the SAS statistical software that Colonial has purchased. Her hypotheses are:

H_0 : The correlation between the expert and consumer rankings is zero or negative.

H_a : The correlation is positive.

Sandy will decide that the expert panel does know best if the data supports H_a : that there is a positive correlation between the experts and the consumers. She goes to a table that shows what value of the Spearman's Rank Correlation Coefficient will separate one tail from the rest of the sampling distribution if there is no association in the population. A portion of such a table is in Exhibit 12.

7. Some non-parametric tests

n	α=.05	α=.025	α=.10
5	0.9		
6	0.829	0.886	0.943
7	0.714	0.786	0.893
8	0.643	0.738	0.833
9	0.6	0.683	0.783
10	0.564	0.648	0.745
11	0.523	0.623	0.736
12	0.497	0.591	0.703

Exhibit 17: Some one-tail critical values for Spearman's Rank Correlation Coefficient

Using $\alpha = .05$, going across the $n = 9$ row in Exhibit 12, Sandy sees that if H_0 is true, only .05 of all samples will have an r_s greater than .600. Sandy decides that if her sample rank correlation is greater than .600, the data supports the alternative, and flavoring K88, the one ranked highest by the experts, will be used. She first goes back to the two sets of rankings and finds the difference in the rank given each flavor by the two groups, squares those differences and adds them together:

	Expert ranking	Consumer ranking	difference	d²
Flavoring				
NYS21	7	8	-1	1
K73	4	3	1	1
K88	1	4	-3	9
Ba4	8	6	2	4
Bc11	2	5	-3	9
McA A	3	1	2	4
McA A	9	9	0	0
WIS 4	5	2	3	9
WIS 43	6	7	-1	1
			sum =	38

Exhibit 18: Sandy's worksheet

Then she uses the formula from above to find her Spearman rank correlation coefficient:

$$1 - [6/(9)(9^2-1)][38] = 1 - .3166 = .6834$$

Her sample correlation coefficient is .6834, greater than .600, so she decides that the experts are reliable, and decides to use flavoring K88. Even though Sandy has ordinal data that only ranks the flavorings, she can still perform a valid statistical test to see if the experts are reliable. Statistics has helped another manager make a decision.

Summary

Though they are less precise than other statistics, non-parametric statistics are useful. You will find yourself faced with small samples, populations that are obviously not normal, and data that is not cardinal. At those times, you can still make inferences about populations from samples by using non-parametric statistics.

Non-parametric statistical methods are also useful because they can often be used without a computer, or even a calculator. The Mann-Whitney U, and the t-test for the difference of sample means, test the same thing. You can usually perform the U-test without any computational help, while performing a t-test without at least a good calculator can take a lot of time. Similarly, the Wilcoxon Signed Ranks test and Spearman's Rank Correlation are easy to compute once the data has been carefully ranked. Though you should proceed on to the parametric statistics when you have access to a computer or calculator, in a pinch you can use non-parametric methods for a rough estimate.

Notice that each different non-parametric test has its own table. When your data is not cardinal, or your populations are not normal, the sampling distributions of each statistic is different. The common distributions, the t, the χ^2 , and the F, cannot be used.

Non-parametric statistics have their place. They do not require that we know as much about the population, or that the data measure as much about the observations. Even though they are less precise, they are often very useful.

8. Regression basics

Regression analysis, like most multivariate statistics, allows you to infer that there is a relationship between two or more variables. These relationships are seldom exact because there is variation caused by many variables, not just the variables being studied.

If you say that students who study more make better grades, you are really hypothesizing that there is a positive relationship between one variable, studying, and another variable, grades. You could then complete your inference and test your hypothesis by gathering a sample of (amount studied, grades) data from some students and use regression to see if the relationship in the sample is strong enough to safely infer that there is a relationship in the population. Notice that even if students who study more make better grades, the relationship in the population would not be perfect; the same amount of studying will not result in the same grades for every student (or for one student every time). Some students are taking harder courses, like chemistry or statistics, some are smarter, some will study effectively, some will get lucky and find that the professor has asked them exactly what they understood best. For each level of amount studied, there will be a distribution of grades. If there is a relationship between studying and grades, the location of that distribution of grades will change in an orderly manner as you move from lower to higher levels of studying.

Regression analysis is one of the most used and most powerful multivariate statistical techniques for it infers the existence and form of a functional relationship in a population. Once you learn how to use regression you will be able to estimate the parameters—the slope and intercept—of the function which links two or more variables. With that estimated function, you will be able to infer or forecast things like unit costs, interest rates, or sales over a wide range of conditions. Though the simplest regression techniques seem limited in their applications, statisticians have developed a number of variations on regression which greatly expand the usefulness of the technique. In this chapter, the basics will be discussed. In later chapters a few of the variations on, and problems with, regression will be covered. Once again, the t-distribution and F-distribution will be used to test hypotheses.

What is regression?

Before starting to learn about regression, go back to algebra and review what a function is. The definition of a function can be formal, like the one in my freshman calculus text: "A function is a set of ordered pairs of numbers (x,y) such that to each value of the first variable (x) there corresponds a unique value of the second variable (y)".³ More intuitively, if there is a regular relationship between two variables, there is usually a function that describes the relationship. Functions are written in a number of forms. The most general is "y = f(x)", which simply says that the value of y depends on the value of x in some regular fashion, though the form of the relationship is not specified. The simplest functional form is the linear function where

$$y = \alpha + \beta x$$

3 George B. Thomas, *Calculus and Analytical Geometry*, 3rd ed., Addison-Wesley, 1960.

8. Regression basics

α and β are parameters, remaining constant as x and y change. α is the intercept and β is the slope. If the values of α and β are known, you can find the y that goes with any x by putting the x into the equation and solving. There can be functions where one variable depends on the values of two or more other variables:

$$y = \alpha + \beta_1 x_1 + \beta_2 x_2$$

where x_1 and x_2 together determine the value of y . There can also be non-linear functions, where the value of the dependent variable ("y" in all of the examples we have used so far) depends on the values of one or more other variables, but the values of the other variables are squared, or taken to some other power or root or multiplied together, before the value of the dependent variable is determined. Regression allows you to estimate directly the parameters in linear functions only, though there are tricks which allow many non-linear functional forms to be estimated indirectly. Regression also allows you to test to see if there is a functional relationship between the variables, by testing the hypothesis that each of the slopes has a value of zero.

First, let us consider the simple case of a two variable function. You believe that y , the dependent variable, is a linear function of x , the independent variable— y depends on x . Collect a sample of (x, y) pairs, and plot them on a set of x, y axes. The basic idea behind regression is to find the equation of the straight line that "comes as close as possible to as many of the points as possible". The parameters of the line drawn through the sample are unbiased estimators of the parameters of the line that would "come as close as possible to as many of the point as possible" in the population, if the population had been gathered and plotted. In keeping with the convention of using Greek letters for population values and Roman letters for sample values, the line drawn through a population is

$$y = \alpha + \beta x$$

while the line drawn through a sample is

$$y = a + bx.$$

In most cases, even if the whole population had been gathered, the regression line would not go through every point. Most of the phenomena that business researchers deal with are not perfectly deterministic, so no function will perfectly predict or explain every observation.

Imagine that you wanted to study household use of laundry soap. You decide to estimate soap use as a function of family size. If you collected a large sample of (family size, soap use) pairs you would find that different families of the same size use different amounts of laundry soap—there is a distribution of soap use at each family size. When you use regression to estimate the parameters of soap use = $f(\text{family size})$, you are estimating the parameters of the line that connects the mean soap use at each family size. Because the best that can be expected is to predict the mean soap use for a certain size family, researchers often write their regression models with an extra term, the "error term", which notes that many of the members of the population of (family size, soap use) pairs will not have exactly the predicted soap use because many of the points do not lie directly on the regression line. The error term is usually denoted as " ε ", or "epsilon", and you often see regression equations written

$$y = \alpha + \beta x + \varepsilon$$

Strictly, the distribution of ε at each family size must be normal, and the distributions of ε for all of the family sizes must have the same variance (this is known as homoskedasticity to statisticians).

It is common to use regression to estimate the form of a function which has more than one independent, or explanatory, variable. If household soap use depends on household income as well as family size, then soap use = f(family size, income), or

$$y = \alpha + \beta_1 x_1 + \beta_2 x_2$$

where y is soap use, x_1 is family size and x_2 is income. This is the equation for a plane, the three-dimensional equivalent of a straight line. It is still a linear function because neither of the x's nor y is raised to a power nor taken to some root nor are the x's multiplied together. You can have even more independent variables, and as long as the function is linear, you can estimate the slope, β , for each independent variable.

Testing your regression: does y really depend upon x?

Understanding that there is a distribution of y (soap use) values at each x (family size) is the key for understanding how regression results from a sample can be used to test the hypothesis that there is (or is not) a relationship between x and y. When you hypothesize that $y = f(x)$, you hypothesize that the slope of the line (β in $y = \alpha + \beta x + \epsilon$) is not equal to zero. If β was equal to zero, changes in x would not cause any change in y. Choosing a sample of families, and finding each family's size and soap use, gives you a sample of (x, y). Finding the equation of the line that best fits the sample will give you a sample intercept, α , and a sample slope, β . These sample statistics are unbiased estimators of the population intercept, α , and slope, β . If another sample of the same size is taken another sample equation could be generated. If many samples are taken, a sampling distribution of sample β 's, the slopes of the sample lines, will be generated. Statisticians know that this sampling distribution of b's will be normal with a mean equal to β , the population slope. Because the standard deviation of this sampling distribution is seldom known, statisticians developed a method to estimate it from a single sample. With this estimated s_b , a t-statistic for each sample can be computed:

$$t = \frac{b - \beta}{\text{estimated } s_b} = \frac{b - \beta}{s_b}$$

where n = sample size

m = number of explanatory (x) variables

b = sample slope

β = population slope

s_b = estimated standard deviation of b's, often called the "standard error".

These t's follow the t-distribution in the tables with n-m-1 df.

Computing s_b is tedious, and is almost always left to a computer, especially when there is more than one explanatory variable. The estimate is based on how much the sample points vary from the regression line. If the points in the sample are not very close to the sample regression line, it seems reasonable that the population points are also widely scattered around the population regression line and different samples could easily produce lines with quite varied slopes. Though there are other factors involved, in general when the points in the sample are farther from the regression line s_b is greater. Rather than learn how to compute s_b , it is more useful for you

8. Regression basics

to learn how to find it on the regression results that you get from statistical software. It is often called the "standard error" and there is one for each independent variable. The printout in Exhibit 19 is typical.

Variable	DF	Parameter	Std Error	t-score
Intercept	1	27.01	4.07	6.64
TtB	1	-3.75	1.54	-2.43

Exhibit 19: Typical statistical package output for regression

You will need these standard errors in order to test to see if y depends upon x or not. You want to test to see if the slope of the line in the population, β , is equal to zero or not. If the slope equals zero, then changes in x do not result in any change in y . Formally, for each independent variable, you will have a test of the hypotheses:

$$H_o: \beta = 0$$

$$H_a: \beta \neq 0$$

if the t-score is large (either negative or positive), then the sample b is far from zero (the hypothesized β), and

H_a should be accepted. Substitute zero for b into the t-score equation, and if the t-score is small, b is close enough to zero to accept H_o . To find out what t-value separates "close to zero" from "far from zero", choose an α , find the degrees of freedom, and use a t-table to find the critical value of t . Remember to halve α when conducting a two-tail test like this. The degrees of freedom equal $n - m - 1$, where n is the size of the sample and m is the number of independent x variables. There is a separate hypothesis test for each independent variable. This means you test to see if y is a function of each x separately. You can also test to see if $\beta > 0$ (or $\beta < 0$) rather than simply if $\beta \neq 0$ by using a one-tail test, or test to see if his some particular value by substituting that value for β when computing the sample t-score.

Casper Gains has noticed that various stock market newsletters and services often recommend stocks by rating if this is a good time to buy that stock. Cap is cynical and thinks that by the time a newsletter is published with such a recommendation the smart investors will already have bought the stocks that are timely buys, driving the price up. To test to see if he is right or not, Cap collects a sample of the price-earnings ratio (P/E) and the "time to buy" rating (TtB) for 27 stocks. P/E measures the value of a stock relative to the profitability of the firm. Many investors search for stocks with P/E's that are lower than would be expected, so a high P/E probably means that the smart investors have discovered the stock. He decides to estimate the functional relationship between P/E and TtB using regression. Since a TtB of 1 means "excellent time to buy", and a TtB of 4 means "terrible time to buy", Cap expects that the slope, β , of the line $P/E = \alpha + \beta * TtB + \epsilon$ will be negative. Plotting out the data gives the graph in

Error: Reference source not found.

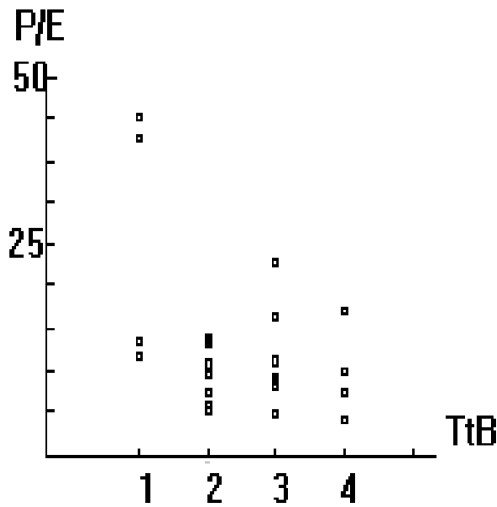


Exhibit 20: A plot of Cap's stock data

Entering the data into the computer, and using the SAS statistical software Cap has at work to estimate the function, yields the output given above.

Because Cap Gains wants to test to see if P/E is already high by the time a low TtB rating is published, he wants to test to see if the slope of the line, which is estimated by the parameter for TtB, is negative or not. His hypotheses are:

$$H_0: \beta \geq 0$$

$$H_a: \beta < 0$$

He should use a one-tail t-test, because the alternative is "less than zero", not simply "not equal to zero". Using an $\alpha = .05$, and noting that there are $n - m - 1$, $26 - 1 - 1 = 24$ degrees of freedom, Cap goes to the t-table and finds that he will accept H_a : if the t-score for the slope of the line with respect to TtB is smaller (more negative) than -1.711. Since the t-score from the computer output is -2.43, Cap should accept H_a : and conclude that by the time the TtB rating is published, the stock price has already been bid up, raising P/E. Buying stocks only on the basis of TtB is not an easy way to make money quickly in the stock market. Cap's cynicism seems to be well founded.

Both the laundry soap and Cap Gains's examples have an independent variable that is always a whole number. Usually, all of the variables are continuous, and to use the hypothesis test developed in this chapter all of the variables really should be continuous. The limit on the values of x in these examples is to make it easier for you to understand how regression works; these are not limits on using regression.

8. Regression basics

Testing your regression. Does this equation really help predict?

Returning to the laundry soap illustration, the easiest way to predict how much laundry soap a particular family (or any family, for that matter) uses would be to take a sample of families, find the mean soap use of that sample, and use that sample mean for your prediction, no matter what the family size. To test to see if the regression equation really helps, see how much of the error that would be made using the mean of all of the y's to predict is eliminated by using the regression equation to predict. By testing to see if the regression helps predict, you are testing to see if there is a functional relationship in the population.

Imagine that you have found the mean soap use for the families in a sample, and for each family you have made the simple prediction that soap use will be equal to the sample mean, \bar{y} . This is not a very sophisticated prediction technique, but remember that the sample mean is an unbiased estimator of population mean, so "on average" you will be right. For each family, you could compute your "error" by finding the difference between your prediction (the sample mean, \bar{y}) and the actual amount of soap used.

As an alternative way to predict soap use, you can have a computer find the intercept, α , and slope, β , of the sample regression line. Now, you can make another prediction of how much soap each family in the sample uses by computing:

$$\hat{y} = \alpha + \beta(\text{familysize})$$

Once again, you can find the error made for each family by finding the difference between soap use predicted using the regression equation, \hat{y} , and actual soap use, y . Finally, find how much using the regression improves your prediction by finding the difference between soap use predicted using the mean, \bar{y} , and soap use predicted using regression, \hat{y} . Notice that the measures of these differences could be positive or negative numbers, but that "error" or "improvement" implies a positive distance. There are probably a few families where the error from using the regression is greater than the error from using the mean, but generally the error using regression will be smaller.

If you use the sample mean to predict the amount of soap each family uses, your error is $(y - \bar{y})$ for each family. Squaring each error so that worries about signs are overcome, and then adding the squared errors together, gives you a measure of the total mistake you make if you use to predict y. Your total mistake is $\sum (y - \bar{y})^2$. The total mistake you make using the regression model would be $\sum (y - \hat{y})^2$. The difference between the mistakes, a raw measure of how much your prediction has improved, is $\sum (\hat{y} - \bar{y})^2$. To make this raw measure of the improvement meaningful, you need to compare it to one of the two measures of the total mistake. This means that there are two measures of "how good" your regression equation is. One compares the improvement to the mistakes still made with regression. The other compares the improvement to the mistakes that would be made if the mean was used to predict. The first is called an F-score because the sampling distribution of these measures follows the F-distribution seen in the "F-test and one-way anova" chapter. The second is called R^2 , or the "coefficient of determination".

All of these mistakes and improvements have names, and talking about them will be easier once you know those names. The total mistake made using the sample mean to predict, $\sum (y - \bar{y})^2$, is called the "sum of squares, total". The total mistake made using the regression, $\sum (y - \hat{y})^2$, is called the "sum of squares, residual" or the

"sum of squares, error". The total improvement made by using regression, $\sum (\hat{y} - \bar{y})^2$ is called the "sum of squares, regression" or "sum of squares, model". You should be able to see that:

Sum of Squares Total = Sum of Squares Regression + Sum of Squares Residual

$$\sum (y - \bar{y})^2 = \sum (\hat{y} - \bar{y})^2 + \sum (y - \hat{y})^2$$

The F-score is the measure usually used in a hypothesis test to see if the regression made a significant improvement over using the mean. It is used because the sampling distribution of F-scores that it follows is printed in the tables at the back of most statistics books, so that it can be used for hypothesis testing. There is also a good set of F-tables at <http://www.itl.nist.gov/div898/handbook/eda/section3/eda3673.htm>. It works no matter how many explanatory variables are used. More formally if there was a population of multivariate observations,

$(y, x_1, x_2, \dots, x_m)$, and there was no linear relationship between y and the x's, so that $y \neq f(x_1, x_2, \dots, x_m)$, if samples of n observations are taken, a regression equation estimated for each sample, and a statistic, F, found for each sample regression, then those F's will be distributed like those in the F-table with (m, n-m-1) df. That F is:

$$F = \frac{\frac{\text{Sum of Squares Regression}}{m}}{\frac{\text{Sum of Squares Residual}}{(n - m - 1)}}$$

$$= \frac{\frac{\text{improvement made}}{m}}{\frac{\text{mistakes still made}}{n - m - 1}}$$

$$F = \frac{\frac{\sum (\hat{y} - \bar{y})^2}{m}}{\frac{\sum (y - \hat{y})^2}{(n - m - 1)}}$$

where: n is the size of the sample

m is the number of explanatory variables (how many x's there are in the regression equation).

If, $\sum (\hat{y} - \bar{y})^2$ the sum of squares regression (the improvement), is large relative to $\sum (y - \hat{y})^2$, the sum of squares residual (the mistakes still made), then the F-score will be large. In a population where there is no functional relationship between y and the x's, the regression line will have a slope of zero (it will be flat), and the \hat{y} will be close to y. As a result very few samples from such populations will have a large sum of squares regression and large F-scores. Because this F-score is distributed like the one in the F-tables, the tables can tell you whether the F-score a sample regression equation produces is large enough to be judged unlikely to occur if $y \neq f(x_1, x_2, \dots, x_m)$. The sum of squares regression is divided by the number of explanatory variables to account for the fact that it always decreases when more variables are added. You can also look at this as finding the improvement per explanatory variable. The sum of squares residual is divided by a number very close to the

8. Regression basics

number of observations because it always increases if more observations are added. You can also look at this as the approximate mistake per observation.

$$H_o: y \neq f(x_1, x_2, \dots, x_m)$$

To test to see if a regression equation was worth estimating, test to see if there seems to be a functional relationship:

$$H_a: y = f(x_1, x_2, \dots, x_m)$$

This might look like a two-tailed test since H_o : has an equal sign. But, by looking at the equation for the F-score you should be able to see that the data supports H_a : only if the F-score is large. This is because the data supports the existence of a functional relationship if sum of squares regression is large relative to the sum of squares residual. Since F-tables are usually one-tailed tables, choose an α , go to the F-tables for that α and (m, n-m-1) df, and find the table F. If the computed F is greater than the table F, then the computed F is unlikely to have occurred if H_o : is true, and you can safely decide that the data supports H_a :. There is a functional relationship in the population.

The other measure of how good your model is, the ratio of the improvement made using the regression to the mistakes made using the mean is called "R-square", usually written R^2 . While R^2 is not used to test hypotheses, it has a more intuitive meaning than the F-score. R^2 is found by:

$$R^2 = \frac{\text{\textit{\sum of Squares Regression}}}{\text{\textit{\sum of Squares Total}}}$$

The numerator is the improvement regression makes over using the mean to predict, the denominator is the mistakes made using the mean, so R^2 simply shows what proportion of the mistakes made using the mean are eliminated by using regression.

Cap Gains, who in the example earlier in this chapter, was trying to see if there is a relationship between price-earnings ratio (P/E) and a "time to buy" rating (TtB), has decided to see if he can do a good job of predicting P/E by using a regression of TtB and profits as a percent of net worth (per cent profit) on P/E. He collects a sample of (P/E, TtB, per cent profit) for 25 firms, and using a computer, estimates the function

$$P/E = a + \beta_1 TtB + \beta_2 \textit{profit}$$

He again uses the SAS program, and his computer printout gives him the results in Figure 3. This time he notices that there are two pages in the printout.

The SAS System
Analysis of Variance

Source	DF	Sum of Squares	Mean Square	F Value	R Sq
Model	2	374.779	187.389	2.724	0.192
Error	23	1582.235	58.72		
Total	25	1957.015			

The SAS System
 Parameter Estimates

Variable	DF	Parameter Estimate	Standard Error	t
Intercept	1	27.281	6.199	4.401
TtB	1	-3.772	1.627	-2.318
Profit	1	-0.012	0.279	-0.042

Exhibit 21: Cap's SAS computer printout

The equation the regression estimates is:

$$P/E = 27.281 - 3.772TtB - 0.012 \text{ Profit}$$

Cap can now test three hypotheses. First, he can use the F-score to test to see if the regression model improves his ability to predict P/E. Second and third, he can use the t-scores to test to see if the slopes of TtB and Profit are different from zero.

To conduct the first test, Cap decides to choose an $\alpha = .10$. The F-score is the regression or model mean square over the residual or error mean square, so the df for the F-statistic are first the df for the model and second the df for the error. There are 2,23 df for the F-test. According to his F-table, with 2.23 degrees of freedom, the critical F-score for $\alpha = .10$ is 2.55. His hypotheses are:

$$H_0: P/E \neq f(TtB, Profit)$$

$$H_a: P/E = f(TtB, Profit)$$

Because the F-score from the regression, 2.724, is greater than the critical F-score, 2.55, Cap decides that the data supports H_a : and concludes that the model helps him predict P/E. There is a functional relationship in the population.

Cap can also test to see if P/E depends on TtB and Profit individually by using the t-scores for the parameter estimates. There are $(n-m-1)=23$ degrees of freedom. There are two sets of hypotheses, one set for β_1 , the slope for TtB, and one set for β_2 , the slope for Profit. He expects that β_1 , the slope for TtB, will be negative, but he does not

8. Regression basics

have any reason to expect that β_2 will be either negative or positive. Therefore, Cap will use a one-tail test on β_1 , and a two-tail test on β_2 :

$$H_0: \beta_1 \geq 0 \quad H_0: \beta_2 = 0$$

$$H_a: \beta_1 < 0 \quad H_a: \beta_2 \neq 0$$

Since he has one one-tail test and one two-tail test, the t-values he chooses from the t-table will be different for the two tests. Using $\alpha = .10$, Cap finds that his t-score for β_1 the one-tail test, will have to be more negative than -1.32 before the data supports P/E being negatively dependent on TtB. He also finds that his t-score for β_2 , the two-tail test, will have to be outside ± 1.71 to decide that P/E depends upon Profit. Looking back at his printout and checking the t-scores, Cap decides that Profit does not affect P/E, but that higher TtB ratings mean a lower P/E. Notice that the printout also gives a t-score for the intercept, so Cap could test to see if the intercept equals zero or not.

Though it is possible to do all of the computations with just a calculator, it is much easier, and more dependably accurate, to use a computer to find regression results. Many software packages are available, and most spreadsheet programs will find regression slopes. I left out the steps needed to calculate regression results without a computer on purpose, for you will never compute a regression without a computer (or a high end calculator) in all of your working years, and there is little most people can learn about how regression works from looking at the calculation method.

Correlation and covariance

The correlation between two variables is important in statistics, and it is commonly reported. What is correlation? The meaning of correlation can be discovered by looking closely at the word—it is almost co-relation, and that is what it means: how two variables are co-related. Correlation is also closely related to regression. The covariance between two variables is also important in statistics, but it is seldom reported. Its meaning can also be discovered by looking closely at the word—it is co-variance, how two variables vary together. Covariance plays a behind-the-scenes role in multivariate statistics. Though you will not see covariance reported very often, understanding it will help you understand multivariate statistics like understanding variance helps you understand univariate statistics.

There are two ways to look at correlation. The first flows directly from regression and the second from covariance. Since you just learned about regression, it makes sense to start with that approach.

Correlation is measured with a number between -1 and +1 called the correlation coefficient. The population correlation coefficient is usually written as the Greek "rho", ρ , and the sample correlation coefficient as r . If you have a linear regression equation with only one explanatory variable, the sign of the correlation coefficient shows whether the slope of the regression line is positive or negative, while the absolute value of the coefficient shows how close to the regression line the points lie. If ρ is +.95, then the regression line has a positive slope and the points in the population are very close to the regression line. If r is -.13 then the regression line has a negative slope and the points in the sample are scattered far from the regression line. If you square r , you will get R^2 , which is higher if the points in the sample lie very close to the regression line so that the sum of squares regression is close to the sum of squares total.

The other approach to explaining correlation requires understanding covariance, how two variables vary together. Because covariance is a multivariate statistic it measures something about a sample or population of observations where each observation has two or more variables. Think of a population of (x,y) pairs. First find the mean of the x's and the mean of the y's, μ_x and μ_y . Then for each observation, find $(x - \mu_x)(y - \mu_y)$. If the x and the y in this observation are both far above their means, then this number will be large and positive. If both are far below their means, it will also be large and positive. If you found $\sum (x - \mu_x)(y - \mu_y)$, it would be large and positive if x and y move up and down together, so that large x's go with large y's, small x's go with small y's, and medium x's go with medium y's. However, if some of the large x's go with medium y's, etc. then the sum will be smaller, though probably still positive. A $\sum (x - \mu_x)(y - \mu_y)$ implies that x's above μ_x are generally paired with y's above μ_y , and those x's below their mean are generally paired with y's below their mean. As you can see, the sum is a measure of how x and y vary together. The more often similar x's are paired with similar y's, the more x and y vary together and the larger the sum and the covariance.

The term for a single observation, $(x - \mu_x)(y - \mu_y)$, will be negative when the x and y are on opposite sides of their means. If large x's are usually paired with small y's, and vice-versa, most of the terms will be negative and the sum will be negative. If the largest x's are paired with the smallest y's and the smallest x's with the largest y's, then many of the $(x - \mu_x)(y - \mu_y)$ will be large and negative and so will the sum. A population with more members will have a larger sum simply because there are more terms to be added together, so you divide the sum by the number of observations to get the final measure, the covariance, or cov:

$$\text{population cov} = \frac{\sum (x - \mu_x)(y - \mu_y)}{(N)}$$

The maximum for the covariance is the product of the standard deviations of the x values and of the y values, $\sigma_x \sigma_y$. While proving that the maximum is exactly equal to the product of the standard deviations is complicated, you should be able to see that the more spread out the points are, the greater the covariance can be. By now you should understand that a larger standard deviation means that the points are more spread out, so you should understand that a larger σ_x or a larger σ_y will allow for a greater covariance.

Sample covariance is measured similarly, except the sum is divided by n-1 so that sample covariance is an unbiased estimator of population covariance:

$$\text{sample cov} = \frac{\sum (x - \bar{x})(y - \bar{y})}{(n-1)}$$

Correlation simply compares the covariance to the standard deviations of the two variables. Using the formula for population correlation:

$$\rho = \frac{\text{cov}}{\rho_x \rho_y} \quad \text{or} \quad \rho = \frac{\sum (x - \mu_x)(y - \mu_y) / N}{\sqrt{\sum (x - \mu_x)^2 / N} \sqrt{\sum (y - \mu_y)^2 / N}}$$

At its maximum, the absolute value of the covariance equals the product of the standard deviations, so at its maximum, the absolute value of r will be 1. Since the covariance can be negative or positive while standard deviations are always positive, r can be either negative or positive. Putting these two facts together, you can see that

8. Regression basics

r will be between -1 and $+1$. The sign depends on the sign of the covariance and the absolute value depends on how close the covariance is to its maximum. The covariance rises as the relationship between x and y grows stronger, so a strong relationship between x and y will result in r having a value close to -1 or $+1$.

Covariance, correlation, and regression

Now it is time to think about how all of this fits together and to see how the two approaches to correlation are related. Start by assuming that you have a population of (x, y) which covers a wide range of y -values, but only a narrow range of x -values. This means that σ_y is large while σ_x is small. Assume that you graph the (x, y) points and find that they all lie in a narrow band stretched linearly from bottom left to top right, so that the largest y 's are paired with the largest x 's and the smallest y 's with the smallest x 's. This means both that the covariance is large and a good regression line that comes very close to almost all the points is easily drawn. The correlation coefficient will also be very high (close to $+1$). An example will show why all these happen together.

Imagine that the equation for the regression line is $y=3+4x$, $\mu_y = 31$, and $\mu_x = 7$, and the two points farthest to the top right, $(10, 43)$ and $(12, 51)$, lie exactly on the regression line. These two points together contribute $\sum(x-\mu_x)(y-\mu_y) = (10-7)(43-31) + (12-7)(51-31) = 136$ to the numerator of the covariance. If we switched the x 's and y 's of these two points, moving them off the regression line, so that they became $(10, 51)$ and $(12, 43)$, μ_x, μ_y, σ_x , and σ_y would remain the same, but these points would only contribute $(10-7)(51-31) + (12-7)(43-31) = 120$ to the numerator. As you can see, covariance is at its greatest, given the distributions of the x 's and y 's, when the (x, y) points lie on a straight line. Given that correlation, r , equals 1 when the covariance is maximized, you can see that $r=+1$ when the points lie exactly on a straight line (with a positive slope). The closer the points lie to a straight line, the closer the covariance is to its maximum, and the greater the correlation.

As this example shows, the closer the points lie to a straight line, the higher the correlation. Regression finds the straight line that comes as close to the points as possible, so it should not be surprising that correlation and regression are related. One of the ways the "goodness of fit" of a regression line can be measured is by R^2 . For the simple two-variable case, R^2 is simply the correlation coefficient, r , squared.

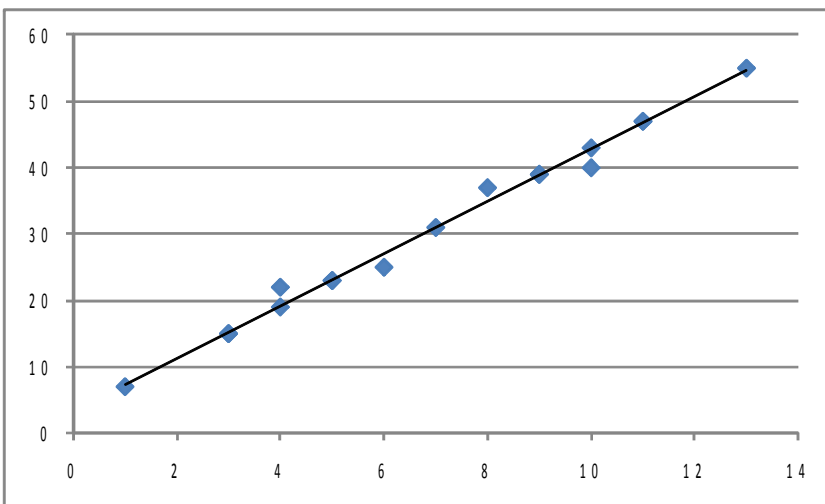


Exhibit 22: Plot of initial population

Correlation does not tell us anything about how steep or flat the regression line is, though it does tell us if the slope is positive or negative. If we took the initial population shown in Exhibit 20, and stretched it both left and right horizontally so that each point's x -value changed, but its y -value stayed the same, σ_x would grow while σ_y

stayed the same. If you pulled equally to the right and to the left, both μ_x and μ_y would stay the same. The covariance would certainly grow since the $(x - \mu_x)$ that goes with each point would be larger absolutely while the $(y - \mu_y)$'s would stay the same. The equation of the regression line would change, with the slope, b , becoming smaller, but the correlation coefficient would be the same because the points would be just as close to the regression line as before. Once again, notice that correlation tells you how well the line fits the points, but it does not tell you anything about the slope other than if it is positive or negative. If the points are stretched out horizontally, the slope changes but correlation does not. Also notice that though the covariance increases, correlation does not because σ_x increases causing the denominator in the equation for finding r to increase as much as covariance, the numerator.

The regression line and covariance approaches to understanding correlation are obviously related. If the points in the population lie very close to the regression line, the covariance will be large in absolute value since the x 's that are far from their mean will be paired with y 's which are far from theirs. A positive regression slope means that x and y rise and fall together, which also means that the covariance will be positive. A negative regression slope means that x and y move in opposite directions, which means a negative covariance.

Summary

Simple linear regression allows researchers to estimate the parameters—the intercept and slopes—of linear equations connecting two or more variables. Knowing that a dependent variable is functionally related to one or more independent or explanatory variables, and having an estimate of the parameters of that function, greatly improves the ability of a researcher to predict the values the dependent variable will take under many conditions. Being able to estimate the effect that one independent variable has on the value of the dependent variable in isolation from changes in other independent variables can be a powerful aid in decision making and policy design. Being able to test the existence of individual effects of a number of independent variables helps decision makers, researchers, and policy makers identify what variables are most important. Regression is a very powerful statistical tool in many ways.

The idea behind regression is simple, it is simply the equation of the line that "comes as close as possible to as many of the points as possible". The mathematics of regression are not so simple, however. Instead of trying to learn the math, most researchers use computers to find regression equations, so this chapter stressed reading computer printouts rather than the mathematics of regression.

Two other topics, which are related to each other and to regression, correlation and covariance, were also covered.

Something as powerful as linear regression must have limitations and problems. In following chapters those limitations, and ways to overcome some of them, will be discussed. There is a whole subject, econometrics, which deals with identifying and overcoming the limitations and problems of regression.