Making Everything Easier!"

Statistics for Big Data

Learn to:

- Collect, clean, and interpret data
- Effectively communicate data analysis
- Make good predictions

Alan Anderson, PhD

Finance, economics, statistics, and math instructor

with David Semmelroth

Author of Data Driven Marketing For Dummies

Introduction

Welcome to *Statistics For Big Data For Dummies!* Every day, what has come to be known as *big data* is making its influence felt in our lives. Some of the most useful innovations of the past 20 years have been made possible by the advent of massive data-gathering capabilities combined with rapidly improving computer technology.

For example, of course, we have become accustomed to finding almost any information we need through the Internet. You can locate nearly anything under the sun immediately by using a search engine such as Google or DuckDuckGo. Finding information this way has become so commonplace that Google has slowly become a verb, as in "I don't know where to find that restaurant — I'll just Google it." Just think how much more efficient our lives have become as a result of search engines. But how does Google work? Google couldn't exist without the ability to process massive quantities of information at an extremely rapid speed, and its software has to be extremely efficient.

Another area that has changed our lives forever is ecommerce, of which the classic example is Amazon.com. People can buy virtually every product they use in their daily lives online (and have it delivered promptly, too). Often online prices are lower than in traditional "brickand-mortar" stores, and the range of choices is wider. Online shopping also lets people find the best available items at the lowest possible prices.

Another huge advantage to online shopping is the ability of the sellers to provide reviews of products and recommendations for future purchases. Reviews from other shoppers can give extremely important information that isn't available from a simple product description provided by manufacturers. And recommendations for future purchases are a great way for consumers to find new products that they might not otherwise have known about. Recommendations are enabled by one application of big data — the use of highly sophisticated programs that analyze shopping data and identify items that tend to be purchased by the same consumers.

Although online shopping is now second nature for many consumers, the reality is that e-commerce has only come into its own in the last 15–20 years, largely thanks to the rise of big data. A website such as Amazon.com must process quantities of information that would have been unthinkably gigantic just a few years ago, and that processing must be done quickly and efficiently. Thanks to rapidly improving technology, many traditional retailers now also offer the option of making purchases online; failure to do so would put a retailer at a huge competitive disadvantage.

In addition to search engines and e-commerce, big data is making a major impact in a surprising number of other areas that affect our daily lives:

- 🛩 Social media
- Online auction sites
- 🖊 Insurance
- 🖊 Healthcare
- 🛩 Energy
- Political polling
- Weather forecasting
- 🖊 Education
- 🖊 Travel

🛩 Finance

About This Book

This book is intended as an overview of the field of big data, with a focus on the statistical methods used. It also provides a look at several key applications of big data. Big data is a broad topic; it includes quantitative subjects such as math, statistics, computer science, and data science. Big data also covers many applications, such as weather forecasting, financial modeling, political polling methods, and so forth.

Our intentions for this book specifically include the following:

- Provide an overview of the field of big data.
- Introduce many useful applications of big data.
- Show how data may be organized and checked for bad or missing information.
- Show how to handle outliers in a dataset.
- Explain how to identify assumptions that are made when analyzing data.
- Provide a detailed explanation of how data may be analyzed with graphical techniques.
- Cover several key univariate (involving only one variable) statistical techniques for analyzing data.
- Explain widely used *multivariate* (involving more than one variable) statistical techniques.
- Provide an overview of modeling techniques such as regression analysis.
- Explain the techniques that are commonly used to analyze time series data.

- Cover techniques used to forecast the future values of a dataset.
- Provide a brief overview of software packages and how they can be used to analyze statistical data.

Because this is a *For Dummies* book, the chapters are written so you can pick and choose whichever topics that interest you the most and dive right in. There's no need to read the chapters in sequential order, although you certainly could. We do suggest, though, that you make sure you're comfortable with the ideas developed in <u>Chapters 4</u> and <u>5</u> before proceeding to the later chapters in the book. Each chapter also contains several tips, reminders, and other tidbits, and in several cases there are links to websites you can use to further pursue the subject. There's also an online Cheat Sheet that includes a summary of key equations for ease of reference.

As mentioned, this is a big topic and a fairly new field. Space constraints make possible only an introduction to the statistical concepts that underlie big data. But we hope it is enough to get you started in the right direction.

Foolish Assumptions

We make some assumptions about you, the reader. Hopefully, one of the following descriptions fits you:

- You've heard about big data and would like to learn more about it.
- You'd like to use big data in an application but don't have sufficient background in statistical modeling.
- You don't know how to implement statistical models in a software package.

Possibly all of these are true. This book should give you a good starting point for advancing your interest in this field. Clearly, you are already motivated.

This book does not assume any particularly advanced knowledge of mathematics and statistics. The ideas are developed from fairly mundane mathematical operations. But it may, in many places, require you to take a deep breath and not get intimidated by the formulas.

Icons Used in This Book

Throughout the book, we include several icons designed to point out specific kinds of information. Keep an eye out for them:

A Tip points out especially helpful or practical information about a topic. It may be hard-won advice on the best way to do something or a useful insight that may not have been obvious at first glance.

ARNING/

A Warning is used when information must be treated carefully. These icons point out potential problems or trouble you may encounter. They also highlight mistaken assumptions that could lead to difficulties.



Technical Stuff points out stuff that may be interesting if you're really curious about something, but which is not essential. You can safely skip these if you're in a hurry or just looking for the basics.



Remember is used to indicate stuff that may have been previously encountered in the book or that you will do well to stash somewhere in your memory for future benefit.

Beyond the Book

Besides the pages or pixels you're presently perusing, this book comes with even more goodies online. You can check out the Cheat Sheet at

www.dummies.com/cheatsheet/statisticsforbigdata.

We've also written some additional material that wouldn't quite fit in the book. If this book were a DVD, these would be on the Bonus Content disc. This handful of extra articles on various mini-topics related to big data is available at <u>www.dummies.com/extras/statisticsforbigdata</u>.

Where to Go From Here

You can approach this book from several different angles. You can, of course, start with <u>Chapter 1</u> and read straight through to the end. But you may not have time for that, or maybe you are already familiar with some of the basics. We suggest checking out the table of contents to see a map of what's covered in the book and then flipping to any particular chapter that catches your eye. Or if you've got a specific big data issue or topic you're burning to know more about, try looking it up in the index.

Once you're done with the book, you can further your big data adventure (where else?) on the Internet. Instructional videos are available on websites such as YouTube. Online courses, many of them free, are also becoming available. Some are produced by private companies such as Coursera; others are offered by major universities such as Yale and M.I.T. Of course, many new books are being written in the field of big data due to its increasing importance.

If you're even more ambitious, you will find specialized courses at the college undergraduate and graduate levels in subject areas such as statistics, computer science, information technology, and so forth. In order to satisfy the expected future demand for big data specialists, several schools are now offering a concentration or a full degree in Data Science.

The resources are there; you should be able to take yourself as far as you want to go in the field of big data. Good luck!

<u>Part I</u>

Introducing Big Data Statistics





Visit <u>www.dummies.com</u> for Great Dummies content online.

In this part ...

- Introducing big data and stuff it's used for
- Exploring the three Vs of big data
- Checking out the hot big data applications
- Discovering probabilities and other basic statistical idea

Chapter 1

What Is Big Data and What Do You Do with It?

In This Chapter

Understanding what big data is all about

Seeing how data may be analyzed using Exploratory Data Analysis (EDA)

Gaining insight into some of the key statistical techniques used to analyze big data

Big data refers to sets of data that are far too massive to be handled with traditional hardware. Big data is also problematic for software such as database systems, statistical packages, and so forth. In recent years, datagathering capabilities have experienced explosive growth, so that storing and analyzing the resulting data has become progressively more challenging.

Many fields have been affected by the increasing availability of data, including finance, marketing, and ecommerce. Big data has also revolutionized more traditional fields such as law and medicine. Of course, big data is gathered on a massive scale by search engines such as Google and social media sites such as Facebook. These developments have led to the evolution of an entirely new profession: the *data scientist*, someone who can combine the fields of statistics, math, computer science, and engineering with knowledge of a specific application. This chapter introduces several key concepts that are discussed throughout the book. These include the characteristics of big data, applications of big data, key statistical tools for analyzing big data, and forecasting techniques.

Characteristics of Big Data

The three factors that distinguish big data from other types of data are *volume, velocity,* and *variety.*

Clearly, with big data, the *volume* is massive. In fact, new terminology must be used to describe the size of these datasets. For example, one *petabyte* of data consists of 1.0×10^{15} bytes of data. That's 1,000 *trillion* bytes!



A *byte* is a single unit of storage in a computer's memory. A byte is used to represent a single number, character, or symbol. A byte consists of eight *bits*, each consisting of either a 0 or a 1.

Velocity refers to the speed at which data is gathered. Big datasets consist of data that's continuously gathered at very high speeds. For example, it has been estimated that Twitter users generate more than a quarter of a million tweets *every minute*. This requires a massive amount of storage space as well as real-time processing of the data.

Variety refers to the fact that the contents of a big dataset may consist of a number of different formats, including spreadsheets, videos, music clips, email messages, and so on. Storing a huge quantity of these incompatible types is one of the major challenges of big data.

<u>Chapter 2</u> covers these characteristics in more detail.

Exploratory Data Analysis (EDA)

Before you apply statistical techniques to a dataset, it's important to examine the data to understand its basic properties. You can use a series of techniques that are collectively known as *Exploratory Data Analysis* (EDA) to analyze a dataset. EDA helps ensure that you choose the correct statistical techniques to analyze and forecast the data. The two basic types of EDA techniques are *graphical* techniques and *quantitative* techniques.

Graphical EDA techniques

Graphical EDA techniques show the key properties of a dataset in a convenient format. It's often easier to understand the properties of a variable and the relationships between variables by looking at graphs rather than looking at the raw data. You can use several graphical techniques, depending on the type of data being analyzed. <u>Chapters 11</u> and <u>12</u> explain how to create and use the following:

- Box plots
- 🛩 Histograms
- Normal probability plots
- Scatter plots

Quantitative EDA techniques

Quantitative EDA techniques provide a more rigorous method of determining the key properties of a dataset. Two of the most important of these techniques are Interval estimation (discussed in <u>Chapter 11</u>).
Hypothesis testing (introduced in <u>Chapter 5</u>).

Interval estimates are used to create a *range* of values within which a variable is likely to fall. *Hypothesis* testing is used to test various propositions about a dataset, such as

- The mean value of the dataset.
- The standard deviation of the dataset.
- The probability distribution the dataset follows.

Hypothesis testing is a core technique in statistics and is used throughout the chapters in $\underline{Part \ III}$ of this book.

Statistical Analysis of Big Data

Gathering and storing massive quantities of data is a major challenge, but ultimately the biggest and most important challenge of big data is putting it to good use.

For example, a massive quantity of data can be helpful to a company's marketing research department only if it can identify the key drivers of the demand for the company's products. Political polling firms have access to massive amounts of demographic data about voters; this information must be analyzed intensively to find the key factors that can lead to a successful political campaign. A hedge fund can develop trading strategies from massive quantities of financial data by finding obscure patterns in the data that can be turned into profitable strategies. Many statistical techniques can be used to analyze data to find useful patterns:

- Probability distributions are introduced in <u>Chapter 4</u> and explored at greater length in <u>Chapter 13</u>.
- Regression analysis is the main topic of <u>Chapter 15</u>.
- Time series analysis is the primary focus of <u>Chapter</u> <u>16</u>.
- Forecasting techniques are discussed in <u>Chapter 17</u>.

Probability distributions

You use a *probability distribution* to compute the probabilities associated with the elements of a dataset. The following distributions are described and applied in this book:

- Binomial distribution: You would use the binomial distribution to analyze variables that can assume only one of two values. For example, you could determine the probability that a given percentage of members at a sports club are left-handed. See <u>Chapter 4</u> for details.
- Poisson distribution: You would use the Poisson distribution to describe the likelihood of a given number of events occurring over an interval of time. For example, it could be used to describe the probability of a specified number of hits on a website over the coming hour. See <u>Chapter 13</u> for details.
- Normal distribution: The normal distribution is the most widely used probability distribution in most disciplines, including economics, finance, marketing, biology, psychology, and many others. One of the characteristic features of the normal distribution is symmetry — the probability of a variable being a

given distance below the mean of the distribution equals the probability of it being the same distance above the mean. For example, if the mean height of all men in the United States is 70 inches, and heights are normally distributed, a randomly chosen man is equally likely to be between 68 and 70 inches tall as he is to be between 70 and 72 inches tall. See <u>Chapter</u> <u>4</u> and the chapters in <u>Parts III</u> and <u>IV</u> for details. The normal distribution works well with many applications. For example, it's often used in the field of finance to describe the returns to financial assets. Due to its ease of interpretation and implementation, the normal distribution is sometimes used even when the assumption of normality is only approximately correct.

The Student's t-distribution: The Student's tdistribution is similar to the normal distribution, but with the Student's t-distribution, extremely small or extremely large values are much more likely to occur. This distribution is often used in situations where a variable exhibits too much variation to be consistent with the normal distribution. This is true when the properties of small samples are being analyzed. With small samples, the variation among samples is likely to be quite considerable, so the normal distribution shouldn't be used to describe their properties. See <u>Chapter 13</u> for details.

Note: The Student's t-distribution was developed by W.S. Gosset while employed at the Guinness brewing company. He was attempting to describe the properties of small sample means.

The chi-square distribution: The chi-square distribution is appropriate for several types of applications. For example, you can use it to determine whether a population follows a particular probability distribution. You can also use it to test whether the variance of a population equals a specified value, and to test for the independence of two datasets. See <u>Chapter 13</u> for details.

The F-distribution: The F-distribution is derived from the chi-square distribution. You use it to test whether the variances of two populations equal each other. The F-distribution is also useful in applications such as regression analysis (covered next). See <u>Chapter 14</u> for details.

Regression analysis

Regression analysis is used to estimate the strength and direction of the relationship between variables that are *linearly* related to each other. <u>Chapter 15</u> discusses this topic at length.



Two variables X and Y are said to be *linearly* related if the relationship between them can be written in the form

Y = mX + b

where

m is the *slope*, or the change in Y due to a given change in X

b is the *intercept*, or the value of *Y* when X = 0

As an example of regression analysis, suppose a corporation wants to determine whether its advertising expenditures are actually increasing profits, and if so, by how much. The corporation gathers data on advertising and profits for the past 20 years and uses this data to estimate the following equation:

Y = 50 + 0.25X

where

Y represents the annual profits of the corporation (in millions of dollars).

X represents the annual advertising expenditures of the corporation (in millions of dollars).

In this equation, the slope equals 0.25, and the intercept equals 50. Because the slope of the regression line is 0.25, this indicates that on average, for every \$1 million increase in advertising expenditures, profits rise by \$.25 million, or \$250,000. Because the intercept is 50, this indicates that with no advertising, profits would still be \$50 million.

This equation, therefore, can be used to forecast future profits based on planned advertising expenditures. For example, if the corporation plans on spending \$10 million on advertising next year, its expected profits will be as follows:

Y = 50 + 0.25X

Y = 50 + 0.25(10) = 50 + 2.5 = 52.5

Hence, with an advertising budget of \$10 million next year, profits are expected to be \$52.5 million.

Time series analysis

A *time series* is a set of observations of a single variable collected over time. This topic is talked about at length

in <u>Chapter 16</u>. The following are examples of time series:

- The daily price of Apple stock over the past ten years.
- The value of the Dow Jones Industrial Average at the end of each year for the past 20 years.
- The daily price of gold over the past six months.

With time series analysis, you can use the statistical properties of a time series to predict the future values of a variable. There are many types of models that may be developed to explain and predict the behavior of a time series.

One place where time series analysis is used frequently is on Wall Street. Some analysts attempt to forecast the future value of an asset price, such as a stock, based entirely on the history of that stock's price. This is known as *technical analysis*. Technical analysts do not attempt to use other variables to forecast a stock's price — the only information they use is the stock's own history.



Technical analysis can work only if there are inefficiencies in the market. Otherwise, all information about a stock's history should already be reflected in its price, making technical trading strategies unprofitable.

Forecasting techniques

Many different techniques have been designed to forecast the future value of a variable. Two of these are time series regression models (<u>Chapter 16</u>) and simulation models (<u>Chapter 17</u>).

Time series regression models

A *time series regression model* is used to estimate the trend followed by a variable over time, using regression techniques. A *trend line* shows the direction in which a variable is moving as time elapses.

As an example, <u>Figure 1-1</u> shows a time series that represents the annual output of a gold mine (measured in thousands of ounces per year) since the mine opened ten years ago.



© John Wiley & Sons, Inc.

Figure 1-1: A time series showing gold output per year for the past ten years.

The equation of the trend line is estimated to be

Y = 0.9212X + 1.3333

where

X is the year.

Y is the annual production of gold (measured in thousands of ounces).

This trend line is estimated using regression analysis. The trend line shows that on average, the output of the mine grows by 0.9212 thousand (921.2 ounces) each year.

You could use this trend line to predict the output next year (the 11th year of operation) by substituting 11 for X, as follows:

Y = 0.9212X + 1.3333

Y = 0.9212(11) + 1.3333 = 11.4665

Based on the trend line equation, the mine would be expected to produce 11,466.5 ounces of gold next year.

Simulation models

You can use *simulation* models to forecast a time series. Simulation models are extremely flexible but can be extremely time-consuming to implement. Their accuracy also depends on assumptions being made about the time series data's statistical properties.

Two standard approaches to forecasting financial time series with simulation models are historical simulation and Monte Carlo simulation.

Historical simulation

Historical simulation is a technique used to generate a probability distribution for a variable as it evolves over time, based on its past values. If the properties of the variable being simulated remain stable over time, this technique can be highly accurate. One drawback to this approach is that in order to get an accurate prediction, you need to have a lot of data. It also depends on the assumption that a variable's past behavior will continue into the future.

As an example, <u>Figure 1-2</u> shows a histogram that represents the returns to a stock over the past 100 days.



Distribution of Returns

This histogram shows the probability distribution of returns on the stock based on the past 100 trading days. The graph shows that the most frequent return over the past 100 days was a loss of 2 percent, the second most frequent was a loss of 3 percent, and so on. You can use the information contained within this graph to create a probability distribution for the most likely return on this stock over the coming trading day.

Monte Carlo simulation

Monte Carlo simulation is a technique in which random numbers are substituted into a statistical model in order to forecast the future values of a variable. This methodology is used in many different disciplines, including finance, economics, and the hard sciences, such as physics. Monte Carlo simulation can work very well but can also be extremely time-consuming to implement. Also, its accuracy depends on the statistical model being used to describe the behavior of the time series.

As you can see, we've got a lot to cover in this book. But don't worry, we take it step by step. In <u>Part I</u>, we look at what big data is. We also build a statistical toolkit that we carry with us throughout the rest of the book. <u>Part II</u> focuses on the (extremely important) process of preparing data for the application of the techniques just described. Then we get to the good stuff in <u>Parts III</u> and <u>IV</u>. Though the equations can appear a little intimidating at times, we have labored to include examples in every chapter that make the ideas a little more accessible. So, take a deep breath and get ready to begin your exploration of big data!

Chapter 2

Characteristics of Big Data: The Three Vs

In This Chapter

Understanding the characteristics of big data and how it can be classified

Checking out the features of the latest methods for storing and analyzing big data

The phrase *big data* refers to *datasets* (collections of data) that are too massive for traditional database management systems (DBMS) to handle properly. The rise of big data has occurred for several reasons, such as the massive increase in e-commerce, the explosion of social media usage, the advent of video and music websites, and so forth.

Big data requires more sophisticated approaches than those used in the past to handle surges of information. This chapter explores the characteristics of big data and introduces the newer approaches that have been developed to handle it.

Characteristics of Big Data

The three main characteristics that define big data are generally considered to be volume, velocity, and variety. These are the three Vs. *Volume* is easy to understand. There's *a lot* of data. *Velocity* suggests that the data comes in faster than ever and must be stored faster than ever. *Variety* refers to the wide variety of data structures that may need to be stored. The mixture of incompatible data formats provides another challenge that couldn't be easily managed by DBMS.

Volume

Volume refers, as you might expect, to the quantity of data being generated. A proliferation of new sources generates massive amounts of data on a continuous basis. The sources include, but are certainly not limited to, the following:

- 🖊 Internet forums
- 🛩 YouTube
- 🖊 Facebook
- 🖊 Twitter
- Cellphones (videos, photos, texts)
- Internet search engines
- Political polling

The volume of data being created is accelerating rapidly, requiring new terminology to describe these massive quantities. This terminology includes names that describe progressively larger amounts of storage. These names can sound quite strange in a world where people are familiar with only megabytes (MB) and gigabytes (GB), and maybe terabytes (TB). Some examples are the *petabyte* (PB), the *zettabyte* (ZB), and the *yottabyte* (YB).

You are likely familiar with the megabyte: one thousand kilobytes, or one million bytes of storage. A gigabyte refers to one *billion* bytes of storage. Until recently, the storage capacity of hard drives and other storage devices was in the range of hundreds of gigabytes, but in 2015

1TB, 2TB, and 4TB internal and external hard drives are now common.

The next step up is the terabyte, which refers to one *trillion* bytes. One trillion is a *large* number, expressed as a one followed by *twelve* zeros:

1,000,000,000,000

You can write this number using *scientific notation* as 1.0×10^{12} .



With scientific notation, a number is expressed as a constant multiplied by a power of ten. For example, 3,122 would be expressed as 3.122×10^3 , because 10^3 equals 1,000. The constant always has one digit before the decimal point, and the remaining digits come after the decimal point.

For larger units of storage, the notation goes like this:

 1.0×10^{15} bytes = one petabyte 1.0×10^{18} bytes = one exabyte 1.0×10^{21} bytes = one zettabyte 1.0×10^{24} bytes = one yottabyte

Here's an interesting name for a *very* large number: 1.0×10^{100} is called a *googol*. The name of the search engine Google is derived from this word. Speaking of Google, the company is currently processing over 20 petabytes of information each day, which is more than the estimated amount of information currently stored at the Library of Congress.

Velocity

As the amount of available data has surged in recent years, the speed with which it becomes available has also accelerated dramatically. Rapidly received data can be classified as the following:

- 🖊 Streaming data
- Complex event processing

Streaming data is data transferred to an application at an extremely high speed. The classic example would be the movies you download and watch from sources such as Netflix and Amazon. In these cases, the data is being downloaded while the movie is playing. If your Internet connection isn't very fast, you've probably noticed annoying interruptions or glitches as the data downloads. In those cases, you need more *velocity*.

Streaming is useful when you need to make decisions in real time. For example, traders must make split-second decisions as new market information becomes available. An entire branch of finance known as *market microstructure* analyzes how prices are generated based on real-time trading activity. *High-frequency trading* (HFT) uses computer algorithms to generate trades based on incoming market data. The data arrives at a high speed, and the assets are held for only fractions of a second before being resold.

Complex event processing (CEP) refers to the use of data to predict the occurrence of events based on a specific set of factors. With this type of processing, data is examined for patterns that couldn't be found with more traditional approaches, so that better decisions may be made in real time. An example is your GPS device's ability to reroute you based on traffic and accident data.

Variety

In addition to traditional data types (numeric and character fields in a file), data can assume a large number of different forms. Here are just a few:

- Spreadsheets
- Word-processing documents
- 🖊 Videos
- 🖊 Photos
- 🖊 Music
- 🖊 Emails
- 🖊 Text messages

With such a variety of formats, storing and analyzing these kinds of data are extremely challenging. The formats are incompatible with each other, so combining them into one large database is problematic.



This is one of the major challenges of big data: finding ways to extract useful information from multiple types of disparate files.

Traditional Database Management Systems (DBMS)

A traditional DBMS stores data and enables it to be easily retrieved. There are several types of database management systems, which can be classified according to the way data is organized and cross-referenced. This section focuses on three of the most important types: relational model, hierarchical model, and network model databases.

Relational model databases

With a relational database, the data is organized into a series of *tables*. Data is accessed by the row and column in which it's located. This model is very flexible and is easy to expand to include new information. You simply add more records to the bottom of an existing table, and you can create new categories by simply adding new rows or columns.

Table 2-1 shows a simple example of a table in a relational database.

Name	Title	Years with Company	Annual Salary
Smith, John	Senior Accountant	8	\$144,000
Jones, Mary	VP, Research and Development	24	\$250,000
Williams, Tony	CFO	13	\$210,000

Table 2-1 Employee Data Organized as a Relational Database

The data in <u>Table 2-1</u> is organized as a series of *records in a table.* Each record contains information about one employee. Each record contains four *fields:* Name, Title, Years with Company, and Annual Salary.

Using this setup, you can find information about employees very quickly and easily. For example, if the human resources department wants to determine which employees have been with the company for at least ten years, a new table — with information drawn from this table — could be generated to list the employees. Table 2-2 shows the new table.

Table 2-2 Employees Who Have Been with the Company at Least Ten Years

Name	Years with Company
Jones, Mary	24
Williams, Tony	13

The relational database user accesses the information with a special type of software known as a *query language.* One of the most widely used query languages is SQL (Structured Query Language).



✓ The "structure" of Structured Query Language is quite simple and is basically the same for all relational database systems. Syntax differs slightly from system to system. But in all cases, queries follow the same format (though not all elements need always be present).

Select (list of data fields you want to see)

From (list of tables containing the data)

Where (list of filtering and other conditions you want to use)

Group by (instructions for summarizing the data)

Having (list of conditions on the summarized data)

Order by (sorting instructions).