

FALL 2013
Stat 8003: Statistical Methods I
Lecture 1

Jichun Xie

1 Syllabus

2 Examples of Statistics and Its Applications

2.1 Investigation of Salary Discrimination

When a group of workers believes that their employer is illegally discriminating against the group, legal remedies are often available. Usually such groups are minorities consisting of a racial, ethnic, gender, or age group. The discrimination may deal with salary, benefits, working conditions, mandatory retirement, etc. The statistical evidence is often crucial to the development of the legal case.

For example, if there is doubt about salary discrimination between male and female workers. What would statisticians do? Usually, they collect data from a subpoena by the legal team. The variables include salaries, years of experience, years of education, a measure of current job responsibility or complexity, a measure of the worker's current productivity, and, the last but not the least, gender. The statisticians consider linear regression model. First they put all the variables other than gender in the model, and then they put all the variables in the model. If the second model works much better for the data, or equivalently, gender is proven to be statistically significant, then this may be regarded as statistical evidence of salary discrimination between female and male employers. We will discuss in detail how to analyze data using linear regression model and how to judge whether a variable is statistically significant in this class.

2.2 Detection of Academic Fabrication

Five years ago, Duke University announced it had found the holy grail of cancer research. Dr. Anil Potti in the Dr. Joseph Nevins group discovered how to match a patient's tumor to the best chemotherapy drug. It was a breakthrough because every person's DNA is unique, so every tumor is different. A drug that kills a tumor in one person, for example, might not work in another. The research was published in the most prestigious medical journals. Duke is excited as well as the 112 patients who signed up for the revolutionary therapy. Doctors everywhere were eager to save lives with the new discovery. Later, however, two statisticians at MD Anderson Cancer Center began analyzing Dr. Potti's data to verify his results. However, they noticed that something really odd that they couldn't explain. They then emailed their questions to Duke. Dr. Potti admitted a few clerical errors, but he said that the new work confirmed his results. And Duke moved ahead. Dr. Nevins and Dr. Potti started a company to market the process. They made a fortune. Patients enrolled in the clinical trials are assigned with the treatment they would believe to be the best for them. However, at MD Anderson Cancer Center, the statisticians kept finding errors that they thought were alarming. The statisticians then wrote a statistical paper analyzing the errors they found in the revolutionary treatment. And they submitted the paper to *Annals of Applied Statistics*. They also contacted Duke, and Duke invited some external review committee to analyze Dr. Potti's investigation. After three months, the review committee concluded that Dr. Potti was not wrong. So the clinical trial went on. Things haven't been changed too much until later, the editor of a small independent newsletter, called "The Cancer Lette", got a tip from a confidential source: check Dr. Potti's Rhodes scholarship. Dr. Potti claimed he got the scholarship when he applied for federal grants. The trouble is that it wasn't true. Till then, Dr. Nevins realized that maybe Dr. Potti is faking the data. He then reviewed the original data and unfortunately his doubt has been confirmed. The data has been manipulated, and lots of the people, the patients, Duke including himself, have been deceived. It turned out that the therapy doesn't work at all. Their theory is wrong. But some of the patients have already died. Well, there were statistical evidences that the data might be manipulated when the clinical trials just started. However, the evidence was neglected or ignored. If these evidences could be treated with enough attention, maybe the fraud could be discovered earlier and fewer patients would die.

2.3 Statistical Tests for Drugs

Abuse of Diethylstilbestrol (DES)

Wikipedia: <http://en.wikipedia.org/wiki/Diethylstilbestrol>

Diethylstilbestrol (DES, former BAN stilboestrol) is a synthetic nonsteroidal estrogen that was first synthesized in 1938. It is also classified as an endocrine disruptor. Human exposure to DES occurred through diverse sources, such as dietary ingestion from supplemented cattle feed and medical treatment for certain conditions, including breast and prostate cancers. From about 1940 to 1970, DES was given to pregnant women in the mistaken belief it would reduce the risk of pregnancy complications and losses. In 1971, DES was shown to cause a rare vaginal tumor in girls and women who had been exposed to this drug in utero. The United States Food and Drug Administration subsequently withdrew DES from use in pregnant women. Follow-up studies have indicated DES also has the potential to cause a variety of significant adverse medical complications during the lifetimes of those exposed. The United States National Cancer Institute recommends women born to mothers who took DES undergo special medical exams on a regular basis to screen for complications as a result of the drug. Individuals who were exposed to DES during their mothers' pregnancies are commonly referred to as "DES daughters" and "DES sons".

2.4 Statistical Learning and Data Mining

Handwritten digit recognition

- Goal: identify single digits 0 – 9 based on images.
- Raw data: images that are scaled segments from five digit ZIP codes.
 - 16×16 eight-bit grayscale maps
 - Pixel intensities range from 0 (black) to 255 (white).
- Input data: a 256 dimension vector, or feature vectors with lower dimensions.

Foreground motion detection

- Goal: extract moving objects from a video sequence.
- Raw data: grayscale image sequence represented by matrices of size $m \times n \times t$, or color image sequence represented by 3 such arrays.
- Videos:
 - <http://www.youtube.com/watch?v=7pE-4eSMUs4>

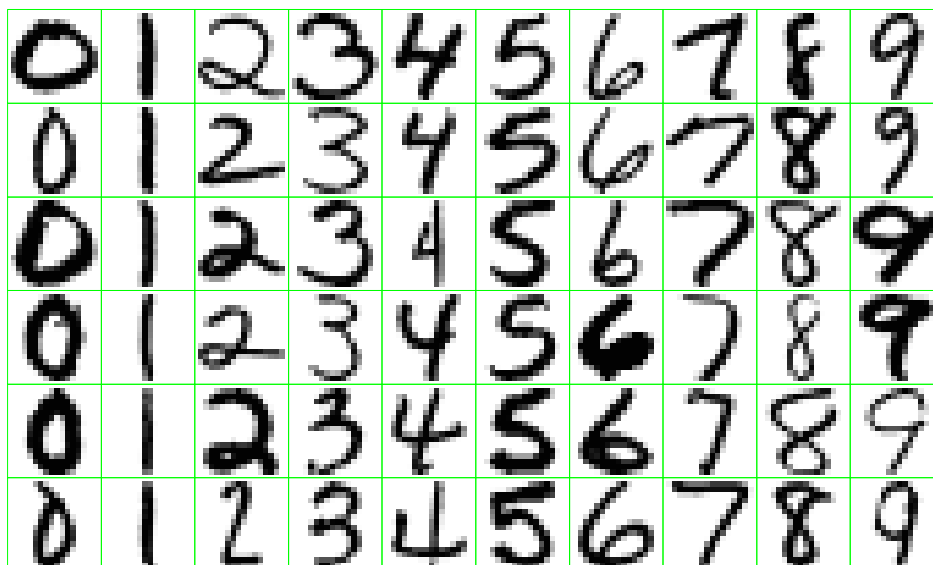


FIGURE 1.2. *Examples of handwritten digits from U.S. postal envelopes.*

Figure 1: Elements of Statistical Learning ©Hastie, Tibshirani & Friedman 2001 Chapter 1

- <http://www.youtube.com/watch?v=F6rxhlkgJkk>
- <http://www.youtube.com/watch?v=yplmDh0gNM8>

3 Introduction to LaTeX

References:

- LaTeX-project: <http://www.latex-project.org/>
- LaTeX wikibook: <http://en.wikibooks.org/wiki/LaTeX>

3.1 Basics

3.1.1 Global Structures

```
\documentclass{...}
\usepackage{...}
```

```
\begin{document}
...
\end{document}
```

- The preamble: the area between `\documentclass{...}` and `\begin{document}`.
- document text: the area between `\begin{document}` and `\end{document}`.

3.1.2 Text and Paragraph Formatting

```
\emph{emphasis}, \textbf{bold}, \textit{italic}
```

emphasis, **bold**, *italic*

```
\begin{verbatim}
The verbatim environment
  simply reproduces every
  character you input,
  including all  s p a c e s!
\end{verbatim}
```

The verbatim environment
simply reproduces every
character you input,
including all s p a c e s!

```
\begin{doublespace}
  This paragraph has \\ double \\ line spacing.
\end{doublespace}
```

This paragraph has
double
line spacing.

```
\begin{spacing}{2.5}
  This paragraph has \\ huge gaps \\ between lines.
\end{spacing}
```

This paragraph has

huge gaps

between lines.

3.1.3 Labels and Cross Referencing

Another good point of LaTeX is that you can easily reference almost anything that is numbered (sections, figures, formulas), and LaTeX will take care of numbering, updating it whenever necessary. The commands to be used do not depend on what you are referencing, and they are:

```
\label{marker}
```

You give the object you want to reference a marker, you can see it like a name.

```
\ref{marker}
```

You can reference the object you have marked before. This prints the number that was assigned to the object.

```
\pageref{marker}
```

It will print the number of the page where the object is.

3.2 Mathematics

One of the greatest motivating forces for Donald Knuth when he began developing the original TeX system was to create something that allowed simple construction of mathematical formulas, whilst looking professional when printed. The fact that he succeeded was most probably why TeX (and later on, LaTeX) became so popular within the scientific community. Typesetting mathematics is one of LaTeX's greatest strengths. It is also a large topic due to the existence of so much mathematical notation.

If your document requires only a few simple mathematical formulas, plain LaTeX has most of the tools that you will need. If you are writing a scientific document that

contains numerous complicated formulas, you might need the `amsmath` package and the `mathtools` packages. Include

```
\usepackage{amsmath,mathtools}
```

in the preamble.

3.2.1 Mathematics Environments

Text:

```
x in the math text mode:  $x$ ,  $\backslash(x\backslash)$ 
```

x in the math text mode: x , x

Displayed:

```
There are Three ways to display a math equation:

$$x + y = z$$


$$[x + y = z]$$


$$\begin{equation*} x + y = z \end{equation*}$$

```

There are Three ways to display a math equation:

$$x + y = z$$

$$x + y = z$$

$$x + y = z$$

Displayed with equation number:

```
\begin{equation}\label{eq:xyz}
  x + y = z
\end{equation}
```

$$x + y = z \tag{1}$$

Equation (`\ref{eq:xyz}`) on page `\pageref{eq:xyz}` describes a relationship between x , y and z .

Equation (1) on page 7 describes a relationship between x , y and z .

Multi-equation display:

```
\begin{align*}
x + y + z &= 6 \\
x - y &= 2 \\
x - z &= 4
\end{align*}
```

$$\begin{aligned}x + y + z &= 6 \\ x - y &= 2 \\ x - z &= 4\end{aligned}$$

3.2.2 Symbols

Mathematics has many symbols! One of the most difficult aspects of learning LaTeX is remembering how to produce symbols. There are of course a set of symbols that can be accessed directly from the keyboard:

```
+ - = ! / ( ) [ ] < > | ' :
```

Beyond those listed above, distinct commands must be issued in order to display the desired symbols. There are a great deal of examples such as Greek letters, set and relations symbols, arrows, binary operators, etc. For example,

```
\forall x \in X, \quad \exists y \leq \epsilon
```

$$\forall x \in X, \quad \exists y \leq \epsilon$$

Fortunately, there's a tool that can greatly simplify the search for the command for a specific symbol. Search Detexify and you will find a website that allows you search for the command by inputting handwriting symbols. There are also apps for iPhone and Android phones. Some softwares such as WinEdt incorporates a list of basic symbols in the toolbox. Another option would be to look in the "The Comprehensive LaTeX Symbol List" available at:

<http://www.ctan.org/tex-archive/info/symbols/comprehensive>

Now we introduce some basic symbols that are commonly used.

Greek letters. Greek letters are commonly used in mathematics, and they are very easy to type in *math mode*. You just have to type the name of the letter after a backslash: if the first letter is lowercase, you will get a lowercase Greek letter, if the first letter is uppercase (and only the first letter), then you will get an uppercase letter.

`\alpha, \beta, \gamma, \Gamma, \pi, \Pi, \phi, \varphi, \Phi`

$\alpha, \beta, \gamma, \Gamma, \pi, \Pi, \phi, \varphi, \Phi$

Operators. An operator is a function that is written as a word: e.g. trigonometric functions (sin, cos, tan), logarithms and exponentials (log, exp). LaTeX has many of these defined as commands:

`\cos (2\theta) = \cos^2 \theta - \sin^2 \theta`

$\cos(2\theta) = \cos^2 \theta - \sin^2 \theta$

For certain operators such as limits, the subscript is placed underneath the operator:

`\lim_{x \to \infty} \exp(-x) = 0`

$\lim_{x \rightarrow \infty} \exp(-x) = 0$

For the modular operator there are two commands: `\bmod` and `\pmod`

`a \bmod b, \quad x \equiv a \pmod b`

$a \bmod b, \quad x \equiv a \pmod b$

Powers and indices. Powers and indices are equivalent to superscripts and subscripts in normal text mode. The caret (^) character is used to raise something, and the underscore (_) is for lowering. If more than one expression is raised or lowered, they should be grouped using curly braces ({ and }).

`k_{n+1} = n^2 + k_n^2 - k_{n-1}`

$k_{n+1} = n^2 + k_n^2 - k_{n-1}$

An underscore (_) can be used with a vertical bar (|) to denote evaluation using subscript notation in mathematics:

`f(n) = n^5 + 4n^2 + 2 |_{n=17}`

$f(n) = n^5 + 4n^2 + 2|_{n=17}$

Fractions and binomials. A fraction is created using the `\frac{numerator}{denominator}` command. (For those who need their memories refreshed, that's the top and bottom respectively!). Likewise, the binomial coefficient (aka the Choose function) may be written using the `\binom` command:

`\frac{n!}{k!(n-k)!} = \binom{n}{k}`

$\frac{n!}{k!(n-k)!} = \binom{n}{k}$

For relatively simple fractions, it may be more aesthetically pleasing to use powers and indices:

`^3/_7, \quad 3/7`

$$^3/_7, \quad 3/7$$

Roots. The `\sqrt` command creates a square root surrounding an expression. It accepts an optional argument specified in square brackets ([and]) to change magnitude:

`\sqrt{\frac{a}{b}}, \quad \sqrt[n]{1+x+x^2+x^3+\ldots}`

$$\sqrt{\frac{a}{b}}, \quad \sqrt[n]{1+x+x^2+x^3+\ldots}$$

Sums and integrals. The `\sum` and `\int` commands insert the sum and integral symbols respectively, with limits specified using the caret (^) and underscore (_). The typical notation for sums is:

`\sum_{i=1}^{10} t_i`

$$\sum_{i=1}^{10} t_i$$

The limits for the integrals follow the same notation. It's also important to represent the integration variables with an upright d, which in math mode is obtained through the `\mathrm{}` command, and with a small space separating it from the integrand, which is attained with the `\,` command.

`\int_0^{\infty} \mathrm{e}^{-x}\,,\mathrm{d}x`

$$\int_0^{\infty} \mathrm{e}^{-x} \, \mathrm{d}x$$

Automatic sizing. Very often mathematical features will differ in size, in which case the delimiters surrounding the expression should vary accordingly. This can be done automatically using the `\left`, `\right`, and `\middle` commands. Any of the previous delimiters may be used in combination with these:

`\left(\frac{x^2}{y^3}\right)`

$$\left(\frac{x^2}{y^3}\right)$$

`P\left(A=2\middle|\frac{A^2}{B}>4\right)`

$$P\left(A=2\left|\frac{A^2}{B}>4\right.\right)$$

Curly braces are defined differently by using `\left\{` and `\right\}`,

```
\left\{\frac{x^2}{y^3}\right\}
```

$$\left\{\frac{x^2}{y^3}\right\}$$

If a delimiter on only one side of an expression is required, then an invisible delimiter on the other side may be denoted using a period (`.`).

```
\left.\frac{x^3}{3}\right|_0^1
```

$$\left.\frac{x^3}{3}\right|_0^1$$

Matrices and arrays. A basic matrix may be created using the matrix environment[3]: in common with other table-like structures, entries are specified by row, with columns separated using an ampersand (&) and a new rows separated with a double backslash (`\\`)

```
\begin{matrix}
a & b & c \\
d & e & f \\
g & h & i
\end{matrix}
```

$$\begin{matrix} a & b & c \\ d & e & f \\ g & h & i \end{matrix}$$

When writing down arbitrary sized matrices, it is common to use horizontal, vertical and diagonal triplets of dots (known as ellipses) to fill in certain columns and rows. These can be specified using the `\cdots`, `\vdots` and `\ddots` respectively:

```
A_{m,n} =
\begin{pmatrix}
a_{1,1} & a_{1,2} & \cdots & a_{1,n} \\
a_{2,1} & a_{2,2} & \cdots & a_{2,n} \\
\vdots & \vdots & \ddots & \vdots \\
a_{m,1} & a_{m,2} & \cdots & a_{m,n}
\end{pmatrix}
```

$$A_{m,n} = \begin{pmatrix} a_{1,1} & a_{1,2} & \cdots & a_{1,n} \\ a_{2,1} & a_{2,2} & \cdots & a_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m,1} & a_{m,2} & \cdots & a_{m,n} \end{pmatrix}$$

3.3 Tables

3.3.1 The Tabular Environment

The tabular environment can be used to typeset tables with optional horizontal and vertical lines. LaTeX determines the width of the columns automatically. The first line of the environment has the form:

```
\begin{tabular}[pos]{table spec}
```

The *table spec* argument tells LaTeX the alignment to be used in each column and the vertical lines to insert. The number of columns does not need to be specified as it is inferred by looking at the number of arguments provided. It is also possible to add vertical lines between the columns here. The following symbols are available to describe the table columns (some of them require that the package array has been loaded):

l	left-justified column
c	centered column
r	right-justified column
p'width'	paragraph column with text vertically aligned at the top
m'width'	paragraph column with text vertically aligned in the middle (requires array package)
b'width'	paragraph column with text vertically aligned at the bottom (requires array package)
	vertical line
	double vertical line

The optional parameter pos can be used to specify the vertical position of the table relative to the baseline of the surrounding text. In most cases, you will not need this option. It becomes relevant only if your table is not in a paragraph of its own. You can use the following letters:

b	bottom
c	center (default)
p	top

In the first line you have pointed out how many columns you want, their alignment and the vertical lines to separate them. Once in the environment, you have to introduce the text you want, separating between cells and introducing new lines. The commands you have to use are the following:

<code>&</code>	column separator
<code>\\</code>	start new row (additional space may be specified after <code>\\</code> using square brackets, such as <code>\\[6pt]</code>)
<code>\hline</code>	horizontal line

3.3.2 Basic Examples

This example shows how to create a simple table in LaTeX. It is a three-by-three table, but without any lines.

```
\begin{tabular}{ l c r }
  1 & 2 & 3 \\
  4 & 5 & 6 \\
  7 & 8 & 9 \\
\end{tabular}
```

```
1 2 3
4 5 6
7 8 9
```

Expanding upon that by including some vertical lines:

```
\begin{tabular}{ l | c || r }
  1 & 2 & 3 \\
  4 & 5 & 6 \\
  7 & 8 & 9 \\
\end{tabular}
```

```
1 | 2 || 3
4 | 5 || 6
7 | 8 || 9
```

To add horizontal lines to the very top and bottom edges of the table:

```
\begin{tabular}{ l | c || r }
\hline
  1 & 2 & 3 \\
  4 & 5 & 6 \\
  7 & 8 & 9 \\
\end{tabular}
```

```
\hline
\end{tabular}
```

1	2	3
4	5	6
7	8	9

And finally, to add lines between all rows, as well as centering (notice the use of the center environment – of course, the result of this is not obvious from the preview on this web page):

```
\begin{center}
\begin{tabular}{l | c || r }
\hline
1 & 2 & 3 \\ \hline
4 & 5 & 6 \\ \hline
7 & 8 & 9 \\ \hline
\end{tabular}
\end{center}
```

1	2	3
4	5	6
7	8	9

3.3.3 Floating with Table

To tell LaTeX we want to use our table as a float, we need to place a tabular environment in a table environment, which is able to float and add a label and caption.

The table environment is a type of floats just as figure is. In fact, they bear a lot of similarities (positioning, captions, *etc.*).

```
\begin{table}[position specifier]
\centering
\begin{tabular}{table spec}
... your table ...
\end{tabular}
\caption{This table shows some data}
\label{tab:myfirsttable}
\end{table}
```

You can set the optional parameter position specifier to define the position of the table,

where it should be placed. The following characters are all possible placements. Using sequences of it define your "wishlist" to LaTeX.

h	where the table is declared (here)
t	at the top of the page
b	at the bottom of the page
p	on a dedicated page of floats
!	override the default float restrictions. E.g., the maximum size allowed of a "b" float is normally quite small; if you want a large one, you need this ! parameter as well.

Table 1 is an example of floating table.

```
\begin{table}[h]
  \centering
  \begin{tabular}{ l | c || r }
    \hline
    1 & 2 & 3 \\ \hline
    4 & 5 & 6 \\ \hline
    7 & 8 & 9 \\ \hline
  \end{tabular}
  \caption{This table shows some data}
  \label{tab:myfirsttable}
\end{table}
```

1	2	3
4	5	6
7	8	9

Table 1: This table shows some data

3.4 Floats, Figures and Captions

3.4.1 Figures

To create a figure that floats, use the figure environment.

```
\begin{figure}[placement specifier]
... figure contents ...
\end{figure}
```

The placement specifier parameter exists as a compromise, and its purpose is to give the author a greater degree of control over where certain floats are placed. Specifier Permission

h	Place the float here, i.e., approximately at the same point it occurs in the source text (however, not exactly at the spot)
t	Position at the top of the page.
b	Position at the bottom of the page.
p	Put on a special page for floats only.
!	Override internal parameters LaTeX uses for determining "good" float positions.
H	Places the float at precisely the location in the LaTeX code. Requires the float package, e.g., <code>\usepackage{float}</code> . This is somewhat equivalent to h!.

3.4.2 Captions

It is always good practice to add a caption to any figure or table. Fortunately, this is very simple in LaTeX. All you need to do is use the `\caption{'text'}` command within the float environment. Because of how LaTeX deals sensibly with logical structure, it will automatically keep track of the numbering of figures, so you do not need to include this within the caption text.

The location of the caption is traditionally underneath the float. However, it is up to you to therefore insert the caption command after the actual contents of the float (but still within the environment). If you place it before, then the caption will appear above the float.

Figure 2 and Figure 3 are two examples.

```
\begin{figure}[h!]
  \caption{A picture of a dog.}
  \centering
  \includegraphics[width=0.5\textwidth,
    natwidth=1024,natheight=768]{dog.jpg}\label{fig:dog}
\end{figure}
```

```
\begin{figure}[h!]
  \centering
  \reflectbox{\includegraphics[width=0.5\textwidth,natwidth=1920,
    natheight=1080]{dog.jpg} }
  \caption{A picture of the same dog
    looking the other way!}
```


Figure 2: A picture of a dog.



```
\end{figure}
```



Figure 3: A picture of the same dog looking the other way!

3.4.3 Multiple Figures in One Float

A useful extension is the subcaption package, which uses subfloats within a single float. This gives the author the ability to have subfigures within figures, or subtables within table floats. Subfloats have their own caption, and an optional global caption. An example will best illustrate the usage of this package:

```
\begin{figure}[t]  
  \centering  
  \begin{subfigure}[b]{0.45\textwidth}
```



(a) A dog



(b) A cat

Figure 4: Pictures of animals

```

\centering
\includegraphics[width=\textwidth,natwidth=1024,
    natheight=768]{dog.jpg}
\caption{A dog}
\label{fig:dog2}
\end{subfigure}%
\quad %add desired spacing between images, e. g. ~, \quad, \
quad etc.
%(or a blank line to force the subfigure onto a new line)
\begin{subfigure}[b]{0.45\textwidth}
\centering
\includegraphics[width=\textwidth,natwidth=1000,
    natheight=781]{cat.jpg}
\caption{A cat}
\label{fig:cat}
\end{subfigure}
\caption{Pictures of animals}\label{fig:animals}
\end{figure}

```