

# Business Statistics

## Course notes

### Hedibert Freitas Lopes

Associate Professor of Econometrics and Statistics

The University of Chicago Booth School of Business

Email: [hlopes@ChicagoBooth.edu](mailto:hlopes@ChicagoBooth.edu)

<http://faculty.chicagobooth.edu/hedibert.lopes/research/>

1

The History of Science has suffered greatly from the use by teachers of second-hand material, and the consequent obliteration of the circumstances and the intellectual atmosphere in which the great discoveries of the past were made. A first-hand study is always instructive, and often...full of surprises.

Ronald A. Fisher

Our world, our life, our destiny, are dominated by uncertainty; this is perhaps the only statement we may assert without uncertainty.

Bruno de Finetti

If this [probability] calculus be condemned, then the whole of the sciences must also be condemned.

Henri Poincare

Those who ignore Statistics are condemned to reinvent it.

Bradley Efron

All models are wrong, but some are useful.

George E. P. Box

2

**TEXTBOOK**

Lind, Marchal and Wathen’s “Statistical Techniques in Business & Economics (12<sup>th</sup>, 13<sup>th</sup> or 14<sup>th</sup> editions)” plays a supporting role in this class, particularly for students who find handouts either too superficial or need additional examples/explanations to any given subject. The book contains several examples and solved problems.

**STATISTICAL PACKAGES**

Most of the computations in the classroom examples are simple enough to be performed by a scientific calculator and/or excel. Several of the computation and plots that appear in the lecture notes were obtained from MINITAB, R, Excel or MegaStat for Excel. MegaStat for Excel is a set of routines that can be easily “added-in” by Microsoft Excel. It comes with Lind, Marchal and Wathen’s textbook. However, excel by itself will be enough for most of our computations.

**HOMEWORK ASSIGNMENTS**

From 4 to 6 homework sets will be assigned, each one of which is invariably due one week after it has been handed out.

**GRADE POINT AVERAGE, FINAL NUMBER GRADE and LETTER GRADE**

The University of Chicago Graduate School of Business mandates a maximum (not minimum!) class grade point average (GPA) of 3.33. The overall class scores will be used to rank the class and grade cutoffs are chosen so that the highest class GPA is less than (or equal to) 3.33.

The final number grade (FNG) will be the weighted average of i) homework assignments average (HWA), ii) the midterm exam (MT) and iii) the final exam (FI). The weights are 20%, 30% and 50%, respectively. For example, suppose that your grades on HW1, HW2, HW3, HW4, MT and FI are 7.0, 8.0, 9.0, 10.0, 9.0 and 8.0, respectively, then the homework assignments average (HWA) is the average of HW1, HW2, HW3 and HW4, i.e.  $HWA=8.5$ . Therefore, your final number grade will be  $FNG = 0.2*HWA+0.3*MT+0.5*FI = 0.2*8.5+0.3*9.0+0.5*8.0 = 8.4$ . The letter grades I use are A, A-, B+, B, B-, C, D (lowest grading pass) and F (fail).

**CALCULATOR, CHEAT SHEET AND REQUESTS FOR RE-GRADING**

Bring your own calculator to all exams. For the midterm exam, a two-page (one sheet) “cheat sheet” is allowed. For the final exam, a four-page (two sheets) “cheat sheet” is allowed. All requests for re-grading of exams must be made in writing and must clearly state the basis of the request.

# Main topics

Exploratory data analysis

Probability

Statistical inference and hypothesis testing

Simple and multiple linear regression

#### UNIVARIATE EXPLORATORY DATA ANALYSIS

1. Graphical summaries of the data
2. Numerical descriptive measures
3. Boxplot

#### MULTIVARIATE EXPLORATORY DATA ANALYSIS

1. How to relate two things
2. Correlations and covariances
3. Linearly related variables
4. Portfolio example
5. Simple linear regression

#### BASIC PROBABILITY

1. Probability and random variables
2. Bivariate random variables
3. Marginal distribution
4. Conditional distribution
5. Independence
6. Computing joints from conditionals and marginals

#### MORE ON PROBABILITY

1. Continuous distributions
2. Normal distribution
3. Cumulative distribution function
4. Expectation as a long run average
5. Expected value and variance of continuous random variables
6. Random variables and formulas
7. Covariance/correlation for pairs of random variables
8. Independence and correlation

#### STATISTICAL INFERENCE

0. I.I.D. draws from the normal distribution
1. Binomial distribution
2. The central limit theorem
3. Estimating  $p$ , population and sample values
4. The sampling distribution of the estimator
5. Confidence interval for  $p$

#### HYPOTHESIS TESTING

1. Hypothesis testing
2. P-values.
3. Confidence intervals, tests, and p-values in general.

#### SIMPLE LINEAR REGRESSION

1. Simple linear regression model
2. Estimates and plug-in prediction
3. Confidence intervals and hypothesis testing
4. Fits, residuals, and R-squared

#### MULTIPLE LINEAR REGRESSION

1. Multiple linear regression model
2. Estimates and plug-in prediction
3. Confidence intervals and hypothesis testing
4. Fits, residuals, R-squared, and the overall F-test
5. Categorical explanatory variables: dummy variables

#### TOPICS IN REGRESSION

1. Residuals as diagnostics
2. Transformations as cures
3. Logistic regression
4. Understanding multicollinearity
5. Autoregressive models
6. Financial time series



5

## Univariate Exploratory Data Analysis

1. Graphical summaries of the data
  - 1.1 Dot plot
  - 1.2 Histogram
  - 1.3 Time series plot
2. Numerical descriptive measures
  - 2.1 Measures of central tendency
    - 2.1.1 The sample mean
    - 2.1.2 The median
  - 2.2 Measures of dispersion
    - 2.2.1 The sample variance
    - 2.2.2 The sample standard deviation
  - 2.3 Measure of asymmetric: skewness
  - 2.4 Measure of extremity: kurtosis
  - 2.5 Quantiles
  - 2.6 Empirical rule
3. Boxplot



6

## Summary of the lecture

- In this class you will learn how to graph  
small sets of quantitative observations: **dotplot**  
large sets of quantitative observations: **histogram**  
observations that are collected as time evolves: **time-series plot**
- You also will learn how to construct a **boxplot**, which can be prove useful when comparing observations from several samples
- Even though graphs are extremely useful and relatively simple to draw, in many situations numerical summaries are required, for instance as input into other systems.
- We will also talk about  
measures of central tendency (**mean and median**)  
measures of dispersion (**variance, standard deviation**)  
measure of asymmetry (**skewness**)  
measure of extremity (**kurtosis**)
- We will also discuss the **empirical rule** that says that roughly **68%** of the observations in any sample should fall within **one** sample standard deviation around the sample mean and **95%** should fall within **two** sample standard deviations around the sample mean.

7

## Book material

- **Chapter 1**  
Types of statistics (pages 6-7 (12 &13)\* ) and types of variables (pages 8-9 (12 & 13))
- **Chapter 2**  
Frequency distributions and Histogram (pages 25 -33 (12), 22-37 (13))
- **Chapter 3**  
Sample mean (page 58 (12 &13)) and sample median (page 62 (12& 13))  
Measures of dispersion (pages 71-77 (12), 71-80 (13))  
Empirical rule (page 80 (12), 82 (13))
- **Chapter 4**  
Dotplots (pages 97-98 (12), 99-100 (13))  
Boxplots (pages 108-111 (12), 110-113 (13))  
Skewness (pages 114-117 (12) , 113-117 (13))

\*Numbers in parentheses refer to the book edition

8

## 1. Graphical Summaries of the Data

### Two key ideas

#### **Exploratory (descriptive) issues:**

Look at the data (sample).  
Understand its structure without generalizing.

#### **Inference issues:**

Use data (sample) to generalize results to  
a larger population of interest.

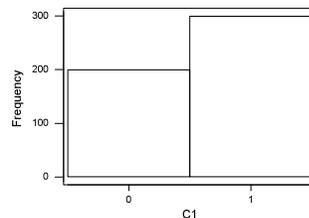
9

### Example

**Problem:** How many of 100,000 voters (population) prefer A over B? We can't ask them all!

**Solution:** Ask a sample of 500 voters.

Summarize, describe the data: 300 voters for A ( $A = 1$ ), 200 for B ( $B = 0$ ).  
We will learn how to generalize to the population. For now, we just learn how to analyze (describe) the data.



10

Let us look at some data. Data are the statistician's raw material, the numbers that we use to interpret reality.

All statistical problems involve either the collection, description and analysis of data, or thinking about the collection, description and analysis of data.

There are many aspects of data. Data may be: **univariate** (one variable per case) or **multivariate** (more than one variable per case).

There are also different types of data: **discrete** (transactions in a given day) and **continuous** (SP500)

11

### The Canadian Return Data

Here is a specific **data set** (or **sample**). We have 107 monthly returns on a broad based portfolio of Canadian assets (more on portfolios later).

```
canada
0.07  0.05  0.02 -0.04  0.08 -0.02 -0.05  0.02  0.03
0.00  0.03  0.08 -0.03  0.01  0.03  0.01  0.02  0.08
0.02 -0.02  0.00  0.01  0.02 -0.09  0.00  0.01 -0.07
0.07  0.00  0.02 -0.05 -0.04 -0.03  0.03  0.04  0.00
0.07  0.00  0.01  0.04 -0.02  0.02  0.01 -0.03  0.05
-0.02  0.00  0.01 -0.01 -0.05 -0.01  0.01  0.00  0.02
-0.02 -0.07  0.03 -0.04  0.03 -0.02  0.06  0.03  0.04
0.01 -0.01 -0.01  0.01 -0.05  0.09 -0.02  0.05  0.06
-0.05 -0.04 -0.01  0.01 -0.06  0.05  0.06  0.02 -0.01
-0.06  0.02 -0.05  0.06  0.04  0.02  0.04  0.02  0.02
0.00  0.00 -0.01  0.04  0.01  0.05 -0.01  0.02  0.04
0.02 -0.03 -0.03  0.05  0.04  0.08  0.07 -0.03
```

**Interpret:** Each number corresponds to a month. They are given in time order (go across columns first). Our first observation is .07. In the first month, the return was .07, in the 11th .03.

12

## 1.1 The dot plot

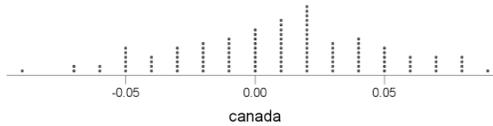
We are interested in ways to **summarize** or “**see**” the data.

The previous table was very unclear.

To display the returns we can use a simple graphical tool: **the dot plot**.

For each number simply place a dot above the corresponding point on the number line.

Dotplot for canada



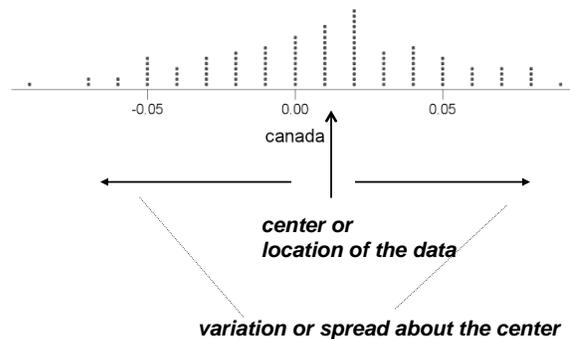
### Interpret:

The returns are *centered* or *located* at about .01.

The *spread* or *variation* in the returns is huge.

13

Dotplot for canada

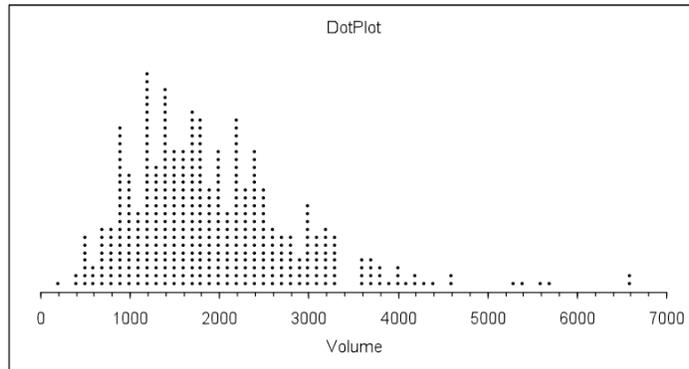


Notice that the data has a nice mound or bell shape. There is a central peak and right and left “tails” that die off roughly symmetrically.

14

Some data does not have the mound shape.

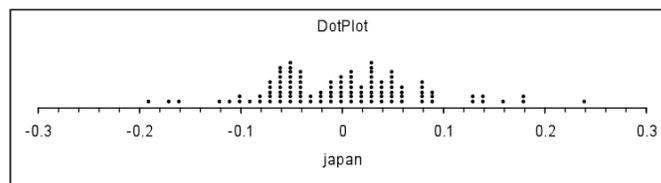
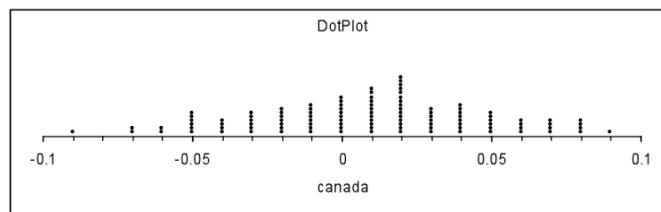
Daily volume of trades in the cattle pit.



It is skewed to the right or positively skewed.

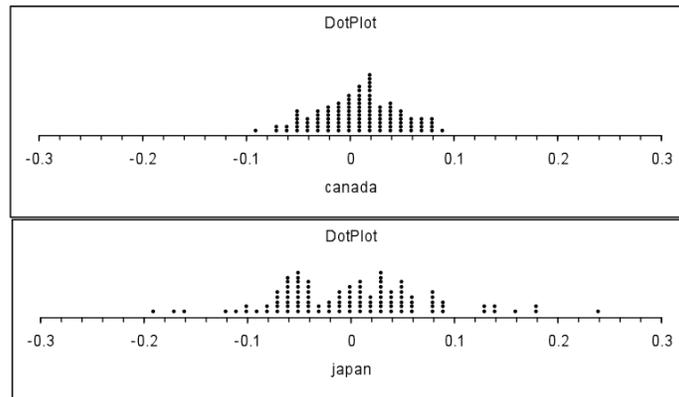
15

We also have data on countries other than Canada. Let us compare Canada with Japan.



16

It really helps to get things on the same scale.  
How is Japan different from Canada?



17

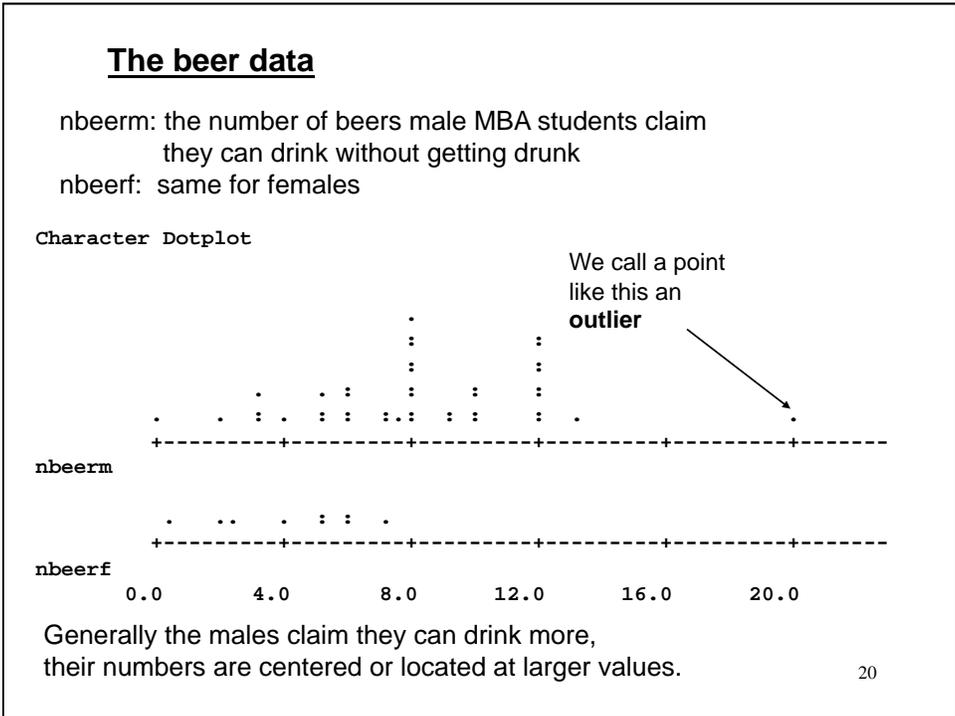
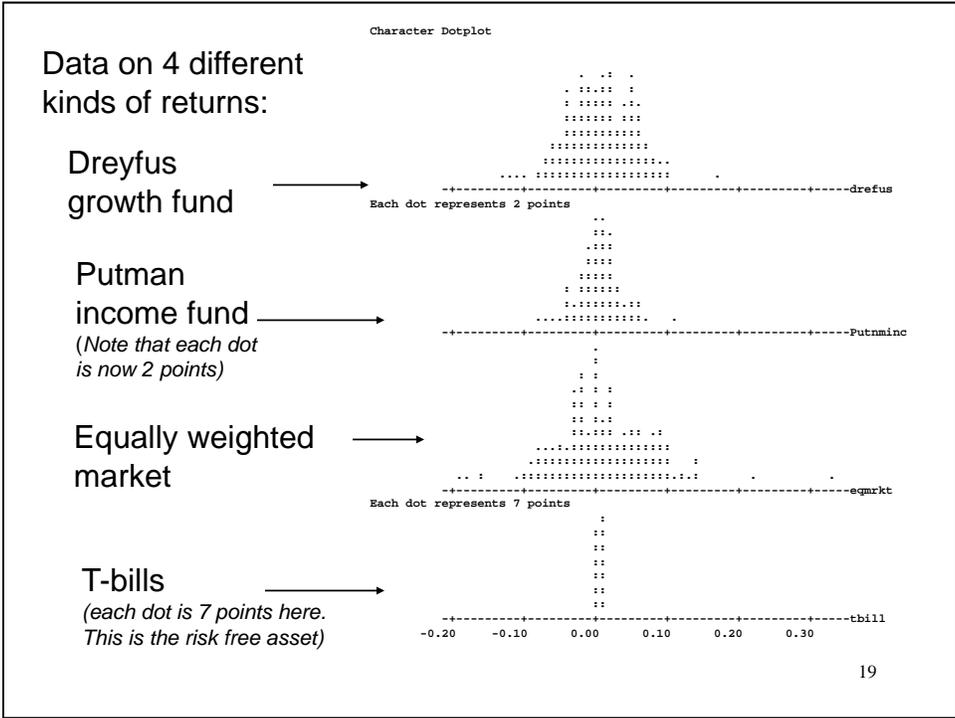
### Mutual fund data

Let us use the dot plot to compare returns on some other kinds of assets.

We will look at returns on different **mutual funds** such as the equally weighted market and T-bills.

The equally weighted market represents returns on a portfolio where you spread your money out equally over a wide variety of stocks.

18



## 1.2 The histogram

Sometimes the dot plot can look rather jumpy.

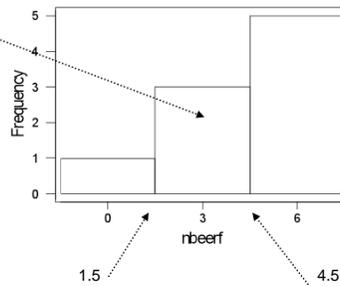
**The histogram** gives us a smoother picture of the data.

The height of each bar tells us how many observations are in the corresponding interval.

```
nbeerf
  4.0   2.0   5.0   6.0   0.5   7.0   6.0   2.5   5.0
```

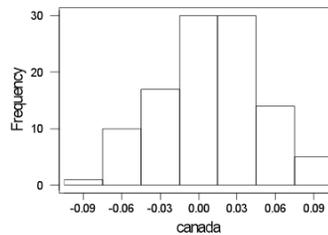
3 women have a number of beers between 1.5 and 4.5.

3 women have a number of beers in the interval (1.5, 4.5).

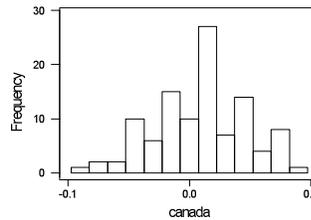


21

Here is the histogram of the Canadian returns.



The number of bars you use affects how “smooth” the picture looks.



22

### 1.3 The time series plot

We just looked at two kinds of data:

- 1) the return data
- 2) the number of beers

For the return data, each number corresponds to a month.  
For the beer data, each number corresponds to a person.

The return data has an important feature that the beer data does not have.

*It has an order!*

There is a first one, a second one, and ....

23

A sequence of observations taken over time is often called a **time series**.

We could have daily data (temperature),  
annual data (inflation),  
quarterly data (inflation, GDP)  
and so on.

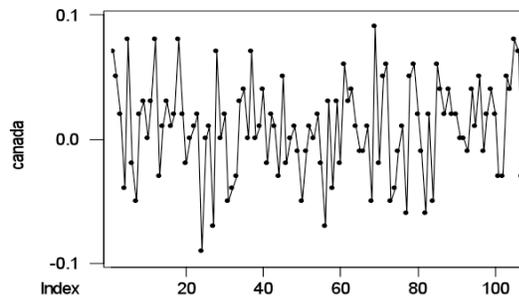
For time series data, the **time series plot** is an important way to look at the data.

24

Time series plot of the Canadian returns:

On the vertical axis we have returns.

On the horizontal axis we have "time".



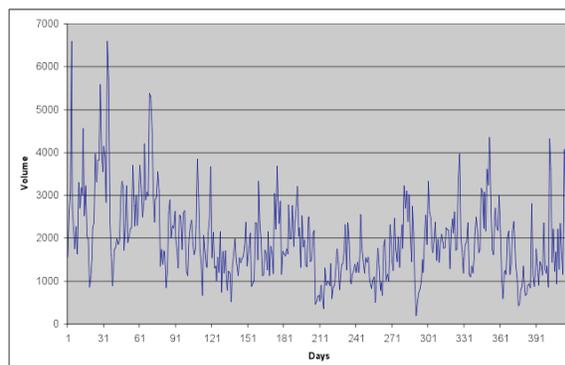
Do you see a pattern?

25

Time series plot of Daily volume of trades in the cattle pit:

On the vertical axis we have volumes.

On the horizontal axis we have days.

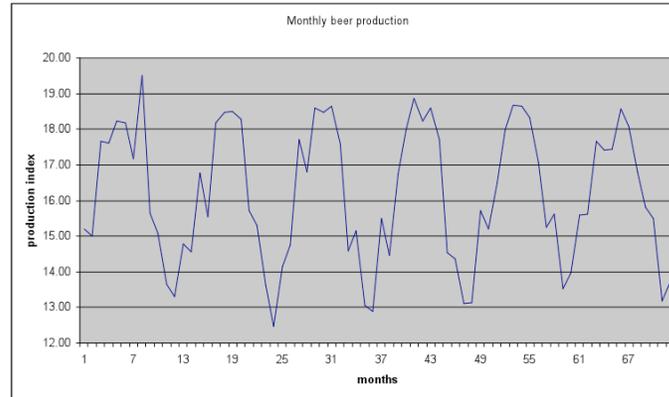


Do you see a pattern?

26

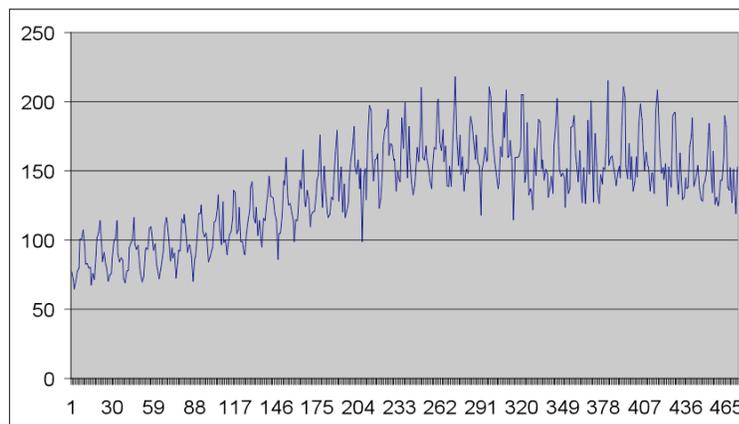
Monthly US beer production.

Now, do you see a pattern?



27

Australia: monthly production of beer.  
megalitres. April 1956 - Aug 1995



Two components: a seasonal (annual) cycle plus an increasing trend from 100 to 175, then a constant trend for the second half of the time series.

28

## 2. Numerical Descriptive Measures

We have looked at graphs.

Suppose we are now interested in having numerical summaries of the data rather than graphical representations.

We have seen that two important features of any data set are:

- 1) how spread out the data is, and
- 2) the central or typical value of the data set.

29

In this part of the notes we will describe methods to summarize a data set numerically.

First, we will introduce measures of central tendency to determine the “center” of a distribution of data values, or possibly the “most typical” data value.

Measures of central tendency include: **the mean** and **the median**.

Second, we will discuss measures of dispersion, such as **the sample standard deviation** and **the sample variance**.

30

## 2.1 Measures of Central Tendency

### 2.1.1 The sample mean

Suppose we collect  $n$  pieces of data. We need some way of describing the data. We write

$x_1, x_2, x_3, \dots, x_n$

the first number

the last number,  **$n$  is the number of numbers**, or the “number of observations.” You may also hear it referred to as the “sample size.”

They are the values that we observe.

31

Here,  $x$  is just a name for the set of numbers, we could just as easily use  $y$  (or Buddy).

$\underline{x}$

$x_1$	→	5	n=5
		2	
$x_3$	→	8	
		6	
		2	

Sometimes the order of the observations means something. In our return data the first observation corresponds to the first time period.

Sometimes it does not. In our beer data we just have a list of numbers, each of which corresponds to a student.

32

The **sample mean** is just the average of the numbers “x”:

$$\bar{x} = \frac{\text{sum}}{n} = \frac{x_1 + x_2 + \cdots + x_n}{n}$$

We often use the  $\bar{x}$  symbol to denote the mean of the numbers x.

We call it “x bar”.

33

**Here is a more compact way to write the same thing...**

Consider

$$x_1 + x_2 + \cdots + x_n$$

We use a shorthand for it (it is just **notation**):

$$\sum_{i=1}^n x_i = x_1 + x_2 + \cdots + x_n$$

This is **summation notation**

34

Using **summation notation** we have:

**The sample mean:**

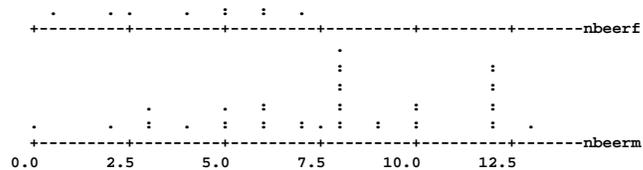
$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

35

### Graphical interpretation of the sample mean

Let us go back to our standard dot plots

Character Dotplot



In some sense, the men claim to drink more.  
To summarize this we can compute the average value  
for both men and women.  
(I deleted the outlier, I do not believe him!).

36

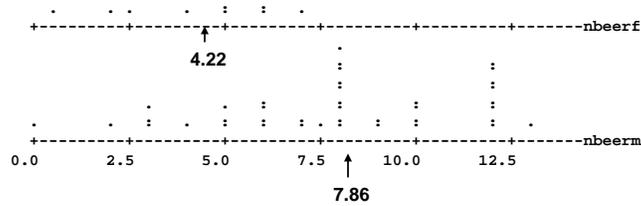
Mean of nbeerf = 4.2222

Mean of nbeerm = 7.8625

“On average women claim they can drink 4.2 beers. Men claim they can drink 7.8 beers”

In the picture, I think of the mean as the “center” of the data.

Character Dotplot



37

Let us compare the means of the Canadian and Japanese returns.

Mean of canada = 0.0090654

Mean of japan = 0.0023364

This is a big difference.

It was hard to see this difference in the dot plots (page 14) Because the difference is small compared to the variation.

38

### More on summation notation (take this as an aside)

Let us look at summation in more detail.

$\sum_{i=1}^n x_i$  means that for each value of  $i$ , from 1 to  $n$ , we add to the sum the value indicated, in this case  $x_i$ .

↑  
add in this value for each  $i$

39

To understand how it works let us consider some **examples**.

Think of each row as an observation on both  $x$  and  $y$ . To make things concrete, think of each row as corresponding to a year and let  $x$  and  $y$  be annual returns on two different assets.

$x$	$y$	year
0.07	0.11	1
0.06	0.05	2
0.04	0.09	3
0.03	0.03	4

In year 1 asset “ $x$ ” had return 7%.  
In year 4 asset “ $y$ ” had return 3%.

40

$$\begin{aligned}\sum_{i=1}^n x_i &= \sum_{i=1}^4 x_i = x_1 + x_2 + x_3 + x_4 \\ &= 0.07 + 0.06 + 0.04 + 0.03 \\ &= 0.2\end{aligned}$$

← compute x bar.

$$\bar{x} = \frac{0.2}{4} = 0.05$$

← compute y bar.

41

For each value of  $i$ , we can add in anything we want:

$$\sum_{i=1}^n (x_i - \bar{x}) =$$

$$\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) =$$

### 2.1.2 The median

After ordering the data, the median is the **middle value** of the data.

If there is an even number of data points, the median is the average of the two middle values.

#### **Example**

1,2,3,4,5                  Median = 3

1,1,2,3,4,5                Median =  $(2+3)/2 = 2.5$

43

### Mean versus median

Although both the mean and the median are good measures of the center of a distribution of measurements, the median is less sensitive to extreme values.

The median is not affected by extreme values since the numerical values of the measurements are not used in its computation.

#### **Example**

1,2,3,4,5	Mean: 3	Median: 3
1,2,3,4,100	Mean: 22	Median: 3

44

## 2.2 Measures of Dispersion

The mean and the median give us information about the central tendency of a set of observations, but they shed no light on the dispersion, or spread of the data.

**Example:** Which data set is more variable ?

5,5,5,5,5	Mean: 5
1,3,5,8,8	Mean: 5

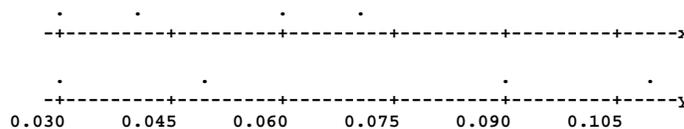
Do you only care about the average return on a mutual fund or you need a measure of risk, too?

Here is one ...

45

### 2.2.1 The Sample Variance

Character Dotplot



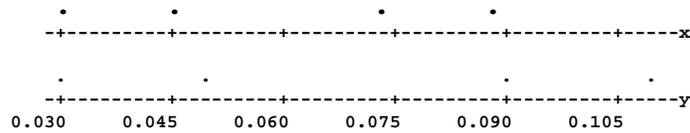
The y numbers are more *spread out* than the x numbers.  
We want a numerical measure of variation or spread.

The basic idea is to view variability in terms of distance between each measurement and the mean.

$$X_i - \bar{X}$$

46

Character Dotplot



$x$	$(x - \bar{x})$	$y$	$(y - \bar{y})$
0.07	0.02	0.11	0.04
0.06	0.01	0.05	-0.02
0.04	-0.01	0.09	0.02
0.03	-0.02	0.03	-0.04

47

We cannot just look at the distance between each measurement and the mean. **We need an overall measure of how big the differences are (i.e., just one number like in the case of the mean).**

Also, we cannot just sum the individual distances because the negative distances cancel out with the positive ones giving zero always (Why?).

**We average the squared distances and define**

$$\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$

48

So, the **sample variance** of the x data is defined to be:

**Sample variance:**

$$s_x^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

We use n-1 instead of n for technical reasons that will be discussed later.

**Think of it as the average squared distance of the observations from the mean.**

49

### Questions

- 1) What is the smallest value a variance can be?
- 2) What are the units of the variance?

It is helpful to have a measure of spread which is in the original units. The sample variance is **not** in the original units. We now introduce a measure of dispersion that solves this problem: **the sample standard deviation**

50

## 2.2.2 The sample standard deviation

It is defined as the square root of the sample variance (easy).

**The sample standard deviation:**

$$S_x = \sqrt{S_x^2}$$

The units of the standard deviation are the same as those of the original data.

51

### Example 1 (numerical)

Assume as before:  $Y - \bar{Y} = 0.04, -0.02, 0.02, -0.04$

$X - \bar{X} = 0.02, 0.01, 0.01, 0.02$

$$\begin{aligned} S_y^2 &= \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2 \\ &= \frac{1}{3} (0.04^2 + (-0.02)^2 + 0.02^2 + (-0.04)^2) \\ &= \frac{1}{3} (0.016 + 0.0004 + 0.0004 + 0.0016) \\ &= \frac{0.004}{3} = 0.00133 \\ S_y &= \sqrt{0.00133} \approx 0.0365 \end{aligned}$$

52

The sample standard deviation for the y data is bigger than that for the x data.

This numerically captures the fact that y has “more variation” about its mean than x.

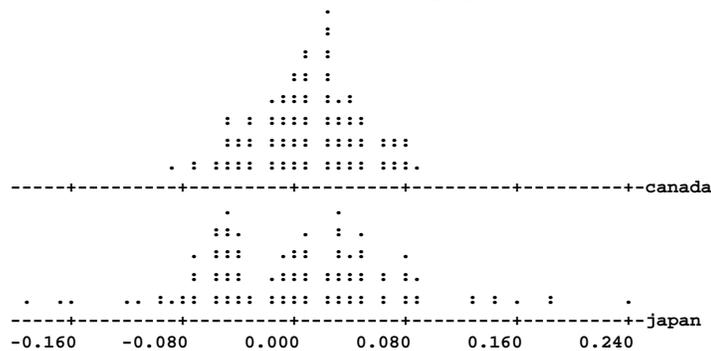
$$\begin{aligned}
 S_x^2 &= \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 \\
 &= \frac{1}{3} (0.02^2 + 0.01^2 + 0.01^2 + 0.02^2) \\
 &= \frac{1}{3} (0.004 + 0.0001 + 0.0001 + 0.0004) \\
 &= \frac{0.001}{3} \approx 0.000333 \\
 S_x &= \sqrt{0.000333} \approx 0.01826
 \end{aligned}$$

53

### Example 2 (graphical)

Character Dotplot

The standard deviations measure the fact that there is more spread in the Japanese returns



Variable	N	Mean	StDev
canada	107	0.00907	0.03833
japan	107	0.00234	0.07368

54

## 2.3 Measure of asymmetry: Skewness

Measures asymmetry of a distribution.

Symmetric data has zero skewness.

**Negatively skewness (the left tail is longer – mean < median)**

Occurs when the values to the left of (less than) the mean are fewer but farther from the mean than are values to the right of the mean.

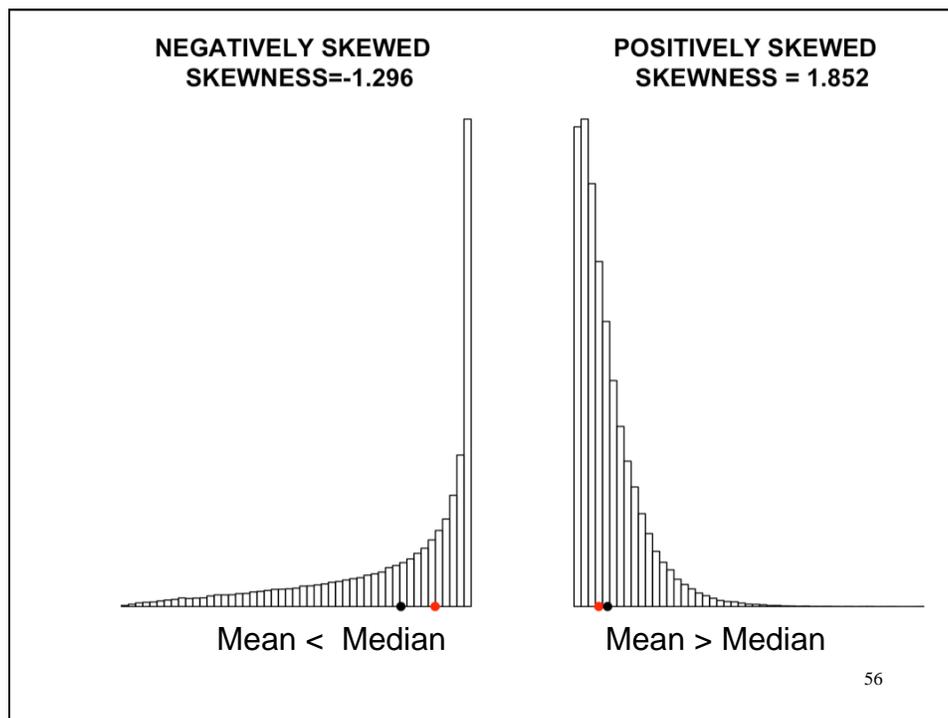
**Positively skewness (the right tail is longer – mean > median)**

Example: investment returns -5%, -10%, -15%, 30%

People like bets with positive skewness.

Willing to accept low, or even negative, expected returns when an asset exhibits positive skewness.

55



## 2.4 Measure of extremity: Kurtosis

Measures the degree to which exceptional values occur more frequently (high kurtosis) or less frequently (low kurtosis)

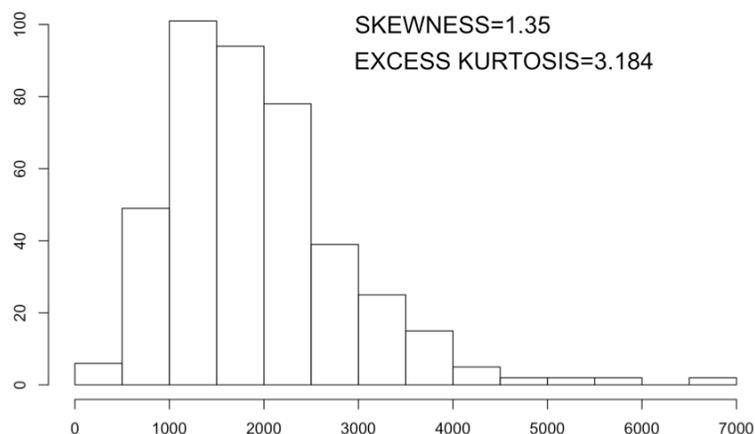
A reference distribution is the **normal distribution**, whose kurtosis is **three**.

**High kurtosis** results in exceptional values that are called "fat tails." Fat tails indicate a higher percentage of very low and very high returns than would be expected with a normal distribution.

**Low kurtosis** results in "thin tails" and a wide middle with more values close to the average than there would be in a normal distribution, and tails are thinner than there would be in a normal distribution.

57

## Volume data



58

## Kurtosis: historical facts

- **KURTOSIS** was used by Karl Pearson in 1905 in "Das Fehlergesetz und seine Verallgemeinerungen durch Fechner und Pearson. A Rejoinder," *Biometrika*, **4**, 169-212, in the phrase "the degree of kurtosis." He states therein that he has used the term previously (*OED*). According to the *OED* and to Schwartzman the term is based on the Greek meaning a bulging, convexity.
- He introduced the terms *leptokurtic*, *platykurtic* and *mesokurtic*, writing in *Biometrika* (1905), **5**, 173: "Given two frequency distributions which have the same variability as measured by the standard deviation, they may be relatively more or less flat-topped than the normal curve. If more flat-topped I term them platykurtic, if less flat-topped leptokurtic, and if equally flat-topped mesokurtic" (*OED2*).
- In his "Errors of Routine Analysis" *Biometrika*, **19**, (1927), p. 160 Student provided a mnemonic:

"In case any of my readers may be unfamiliar with the term "kurtosis" we may define mesokurtic as "having  $\beta_2$  equal to 3," while platykurtic curves have  $\beta_2 < 3$  and leptokurtic  $> 3$ . The important property which follows from this is that platykurtic curves have shorter "tails" than the



normal curve of error and leptokurtic longer "tails." I myself bear in mind the meaning of the words by the above *mnemoria technica*, where the first figure represents platypus, and the second kangaroos, noted for "lepping," though, perhaps, with equal reason they should be hares!

59

Earliest Known Uses of Some of the Words of Mathematics: <http://jeff560.tripod.com/k.html>

## Computing skewness and excess kurtosis

Excess kurtosis is kurtosis minus 3.

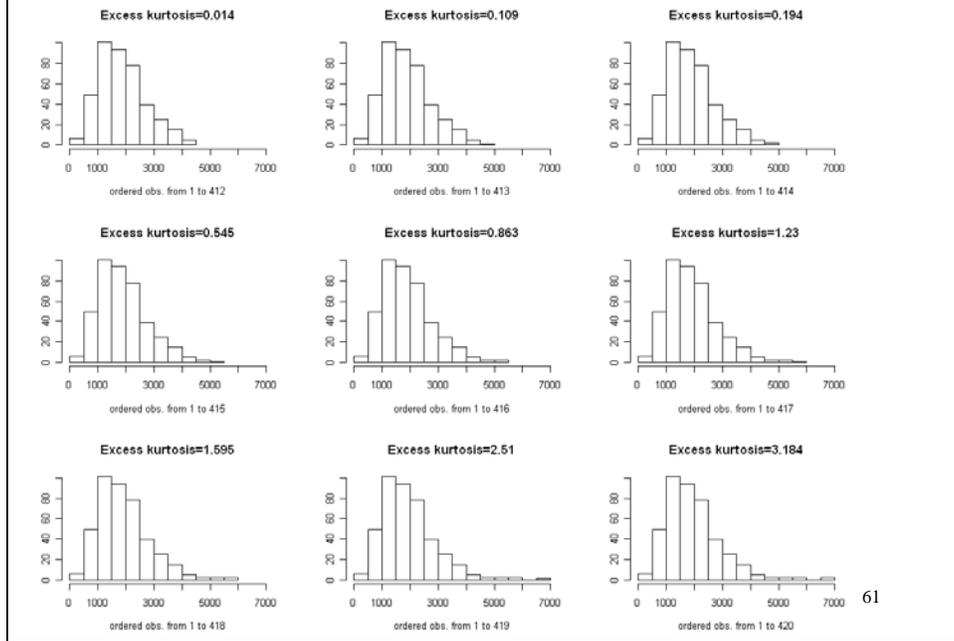
Excel computes excess kurtosis.

$$skewness = \frac{n}{(n-1)(n-2)} \sum_{i=1}^n \left( \frac{x_i - \bar{x}}{s} \right)^3$$

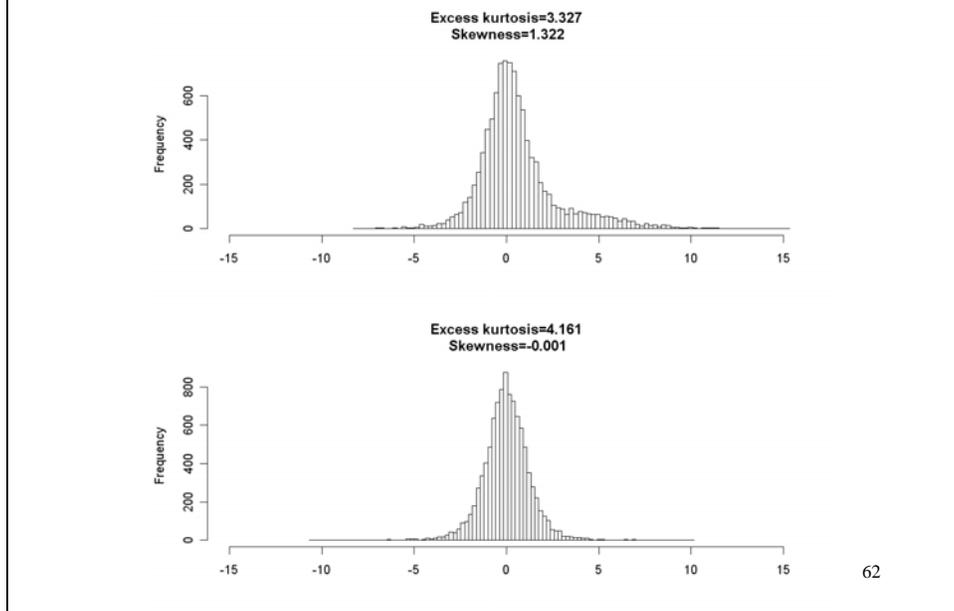
$$excess\ kurtosis = \frac{n(n+1)}{(n-1)(n-2)(n-3)} \sum_{i=1}^n \left( \frac{x_i - \bar{x}}{s} \right)^4 - \frac{3(n-1)^2}{(n-2)(n-3)}$$

60

## Volume data: kurtosis and outliers

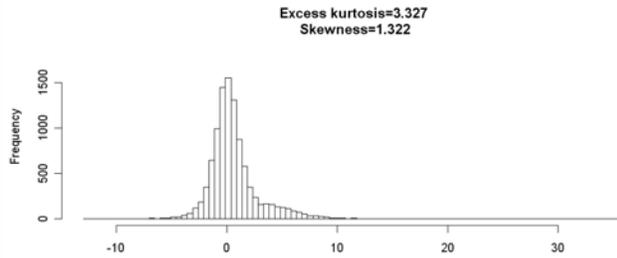


## Same kurtosis, different skewness

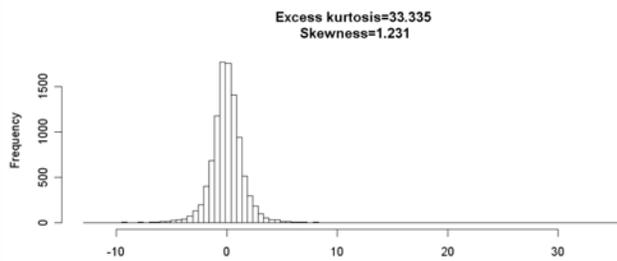


# Same skewness, different kurtosis

10 largest obs.



- 11.25794
- 11.26239
- 11.43341
- 11.48154
- 11.52330
- 11.94644
- 12.10322
- 12.33747
- 12.75935
- 15.32864



- 10.13302
- 10.98134
- 11.38262
- 11.73549
- 11.77891
- 12.84776
- 14.80519
- 15.38212
- 21.74778
- 35.23782

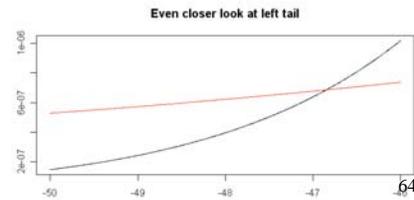
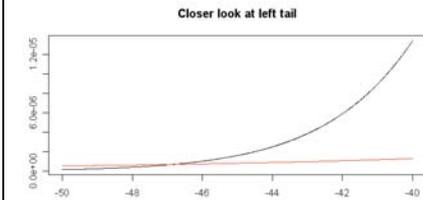
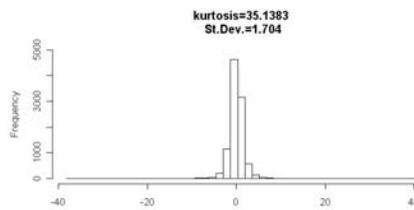
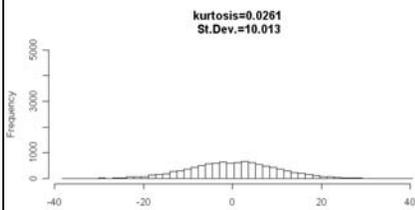
63

# Kurtosis and standard deviation

Left histogram: higher variability.

Left histogram: lower kurtosis or thinner tails.

Bottom curves: left tail behavior of both histograms.



64

## Same mean, variance, skewness Different kurtosis

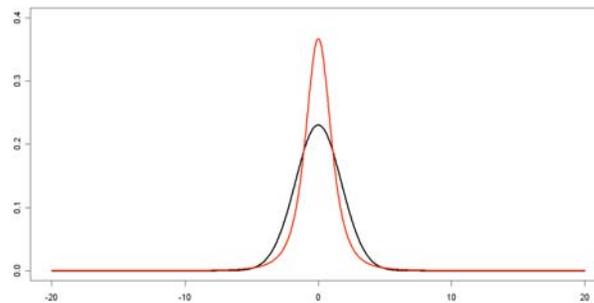
In both cases, mean=0, variance=3 and skewness=0.

**Excess kurtosis is 0.054 for the thin-tail distribution (black).**

**Excess kurtosis is 65.18 for the fat-tail distribution (red).**

Percentage of observations below cutoff

cutoff	Red	Black
-10	0.1064	0.0000
-9	0.1448	0.0000
-8	0.2038	0.0005
-7	0.2993	0.0065
-6	0.4636	0.0571
-5	0.7696	0.3571
-4	1.4004	1.6004
-3	2.8834	5.1393
-2	6.9663	11.8255
-1	19.5501	19.4970



65

## 2.5 Quantiles

**Quartiles:** divide the data into 4 equal parts.

Q1 = Median of the first half of the data

Q2 = Median

Q3 = Median of the second half of the data

**IQ = Interquartile range**

$$IQ = Q3 - Q1$$

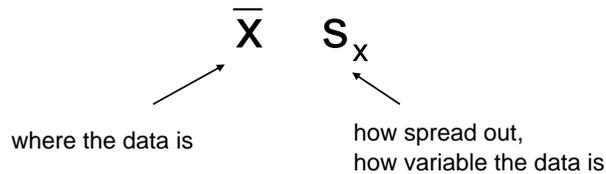
**Deciles:** divide the data into 10 equal parts.

**Percentiles:** divide the data into 100 equal parts.

66

## 2.6 The Empirical Rule

We now have **two numerical summaries** for the data



The mean is pretty easy to interpret (some sort of “center” of the data).

We know that the bigger  $s_x$  is, the more variable the data is, but how do we really interpret this number?

What is a big  $s_x$ , what is a small one ?

67

**The empirical rule will help us understand  $s_x$  and relate the summaries back to the dot plot (or the histogram).**

### Empirical Rule

For “mound shaped data”:

Approximately 68% of the data is in the interval

$$(\bar{X} - s_x, \bar{X} + s_x) = \bar{X} \pm s_x$$

Approximately 95% of the data is in the interval

$$(\bar{X} - 2s_x, \bar{X} + 2s_x) = \bar{X} \pm 2s_x$$

68

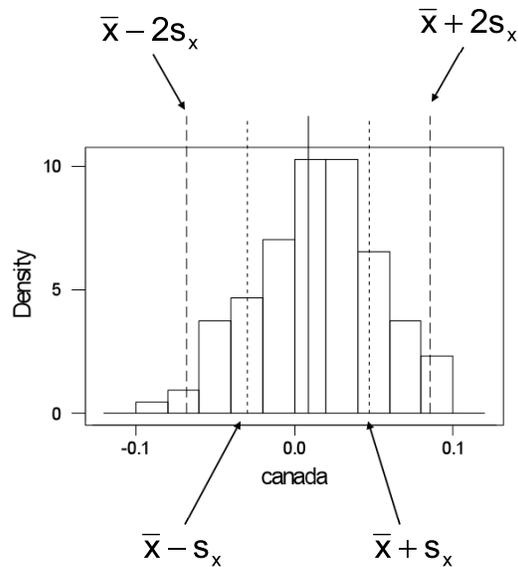
**Let us see this with the Canadian returns**

$\bar{x} = .00907$

$s_x = .03833$

The empirical rule says that roughly 95% of the observations are between the dashed lines and roughly 68% between the dotted lines.

Looks reasonable.

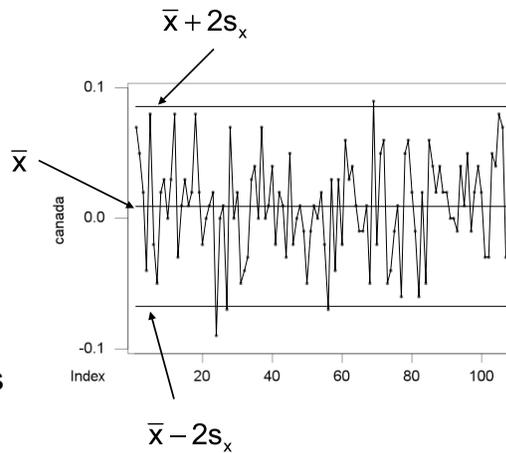


69

Same thing viewed from the perspective of the time series plot.

5% outside would be about 5 points.

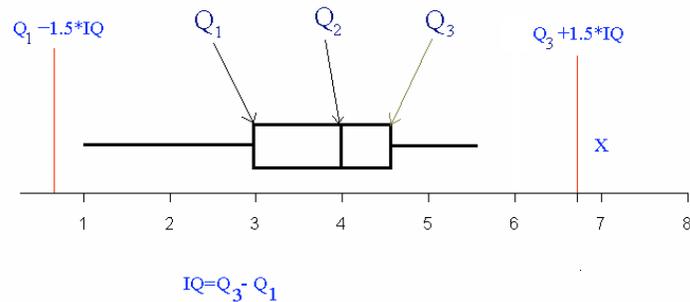
There are 4 points outside, which is pretty close.



70

### 3. BOXPLOT

1-2-2-3-3-4-4-4-4-4-4-5-5-5.5-7



1.0 is the smallest observation greater than  $Q_1 - 1.5 \cdot IQ$

5.5 is the largest observation lower than  $Q_3 + 1.5 \cdot IQ$

71

### Step by step illustration

Data: 65 69 70 63 63 72 63 60 69 66 71 73 70 65 74 69 69 87

Sort: 60 63 63 63 65 65 66 69 69 69 69 70 70 71 72 73 74 87

Q1 =

Q2 =

Q3 =

IQ =

1.5\*IQ =

$Q1 - 1.5 \cdot IQ =$

$Q3 + 1.5 \cdot IQ =$

72

## Solution

Sort: 60 63 63 63 65 65 66 69 69

69 69 70 70 71 72 73 74 87

$Q1 = 65$

$Q2 = 69$

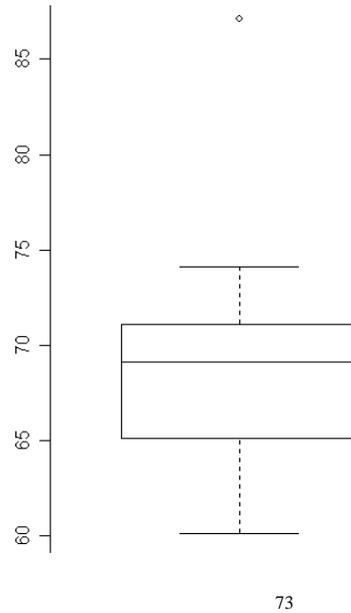
$Q3 = 71$

$IQ = Q3 - Q1 = 71 - 65 = 6$

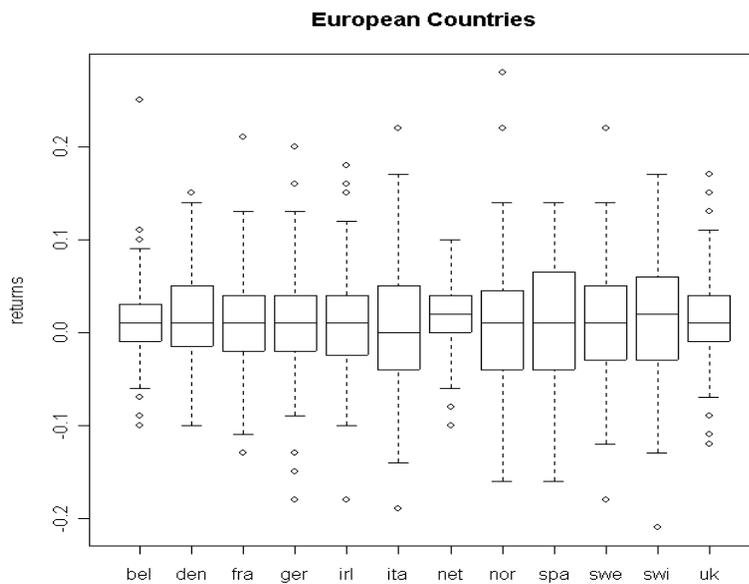
$1.5 * IQ = 9$

$Q1 - 1.5 * IQ = 65 - 9 = 56$

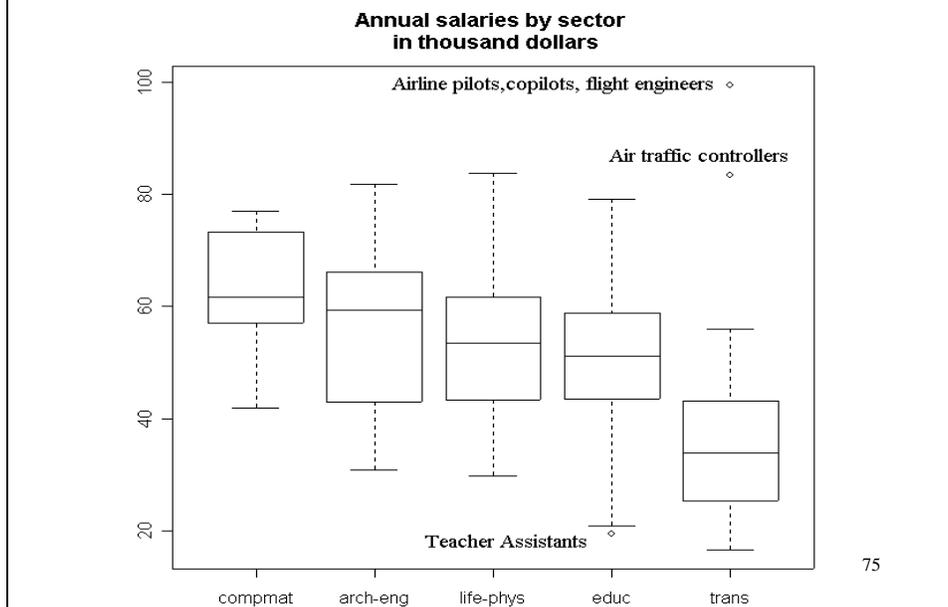
$Q3 + 1.5 * IQ = 71 + 9 = 80$



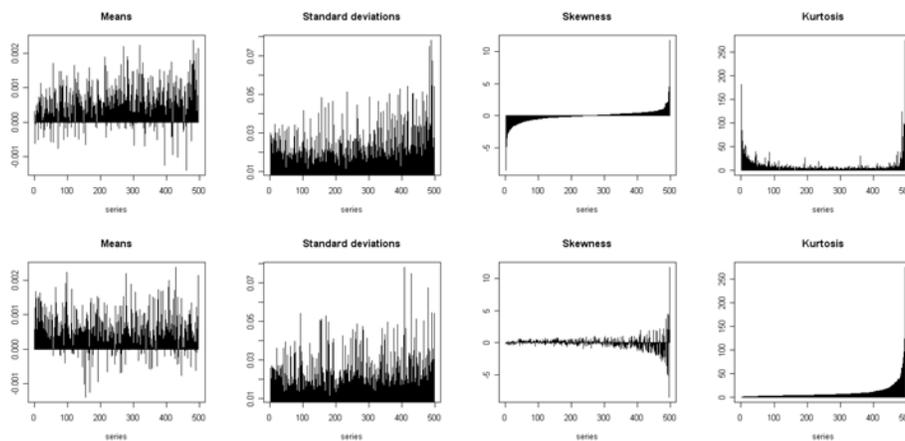
## Example: European returns



## Example: Annual salary (in thousands of dollars)



## Example: SP500 components



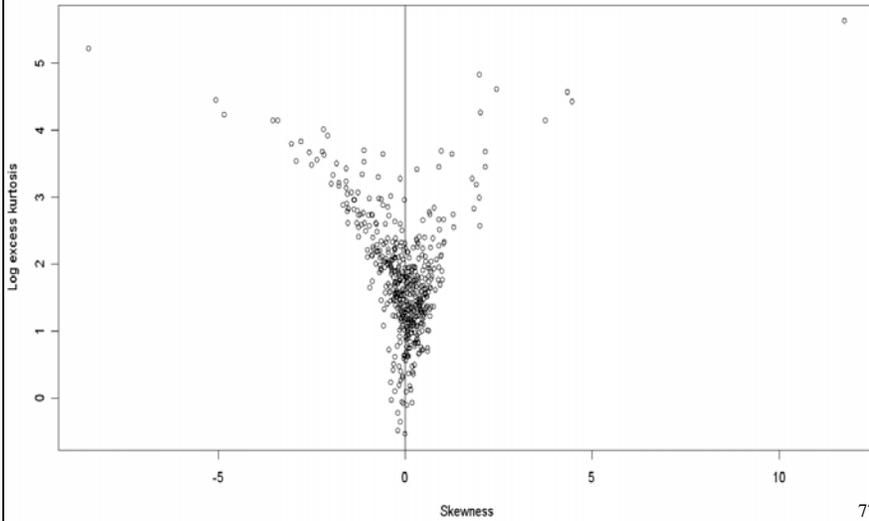
1<sup>st</sup> row: ordered by skewness

2<sup>nd</sup> row: ordered by kurtosis

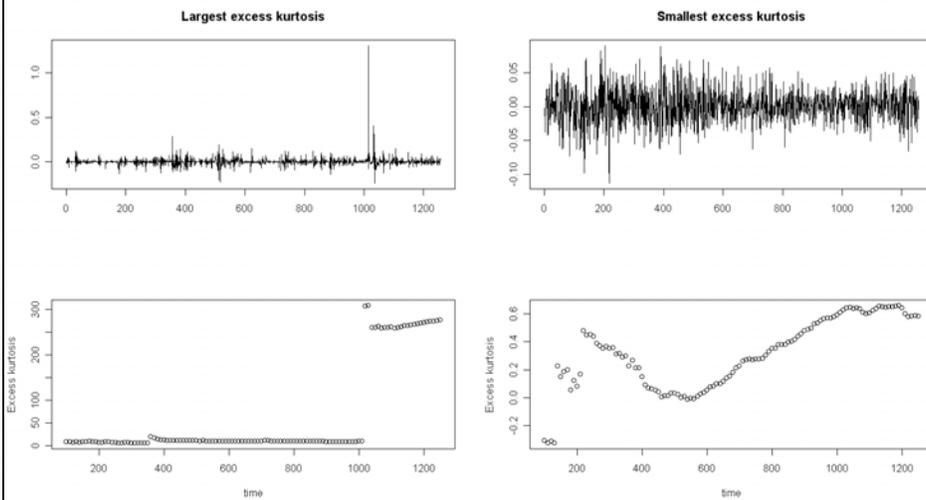
76

## S&P500: kurtosis and skewness

Skewness and logarithm of excess kurtosis for the S&P500 components.



## S&P500: Components with fattest and thinnest tails



The bottom graphs are excess kurtosis computed over time.

78

### Example: Number of siblings - MBA students

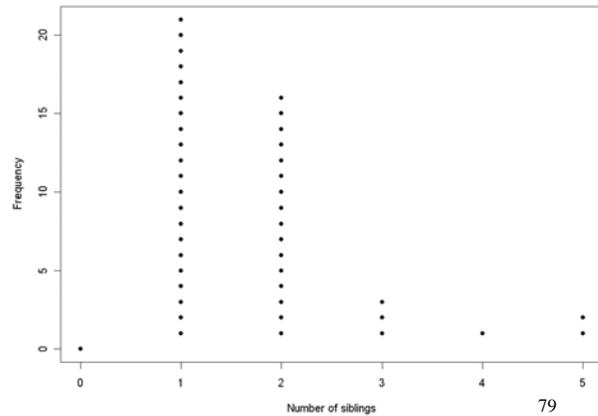
Data collected from Business Stats students on January 10<sup>th</sup> 2009 (41000-85):

0 1 1 1 1 1 1 1 1 1 1

1 1 1 1 1 1 1 1 1 1 1

2 2 2 2 2 2 2 2 2 2 2

2 2 2 2 2 3 3 3 4 5 5



79

$\bar{X} = 1.73$

Median = 1.50

Var = 1.087

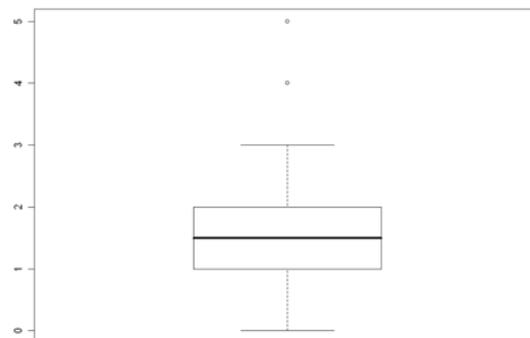
St.dev.=1.042

Q1 = 1.00

Q3 = 2.00

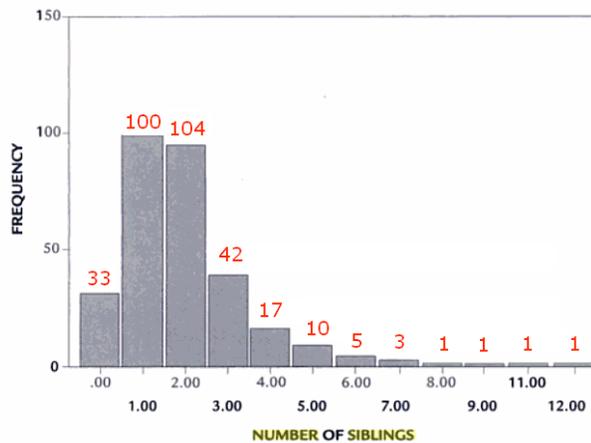
Skewness = 1.616

Excess kurtosis = 3.093



80

## Example: Number of siblings – Boston College

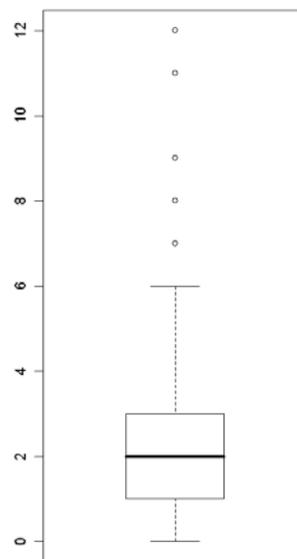


Source: Statistical Methods for Health Care Research (5<sup>th</sup> edition) by Barbara H. Munro. Publisher: Lippincott, Williams & Wilkins

**FIGURE 2-1.** Relative frequency distribution of number of siblings a child has. (Data collected with a grant funded by the National Institute of Nursing Research, R01 NR04838-01A2. P.I., Vessey, J. (2000). *Development of the CATS: Child-Adolescent Teasing Scale*. The William F. Connell School of Nursing, Boston College.)

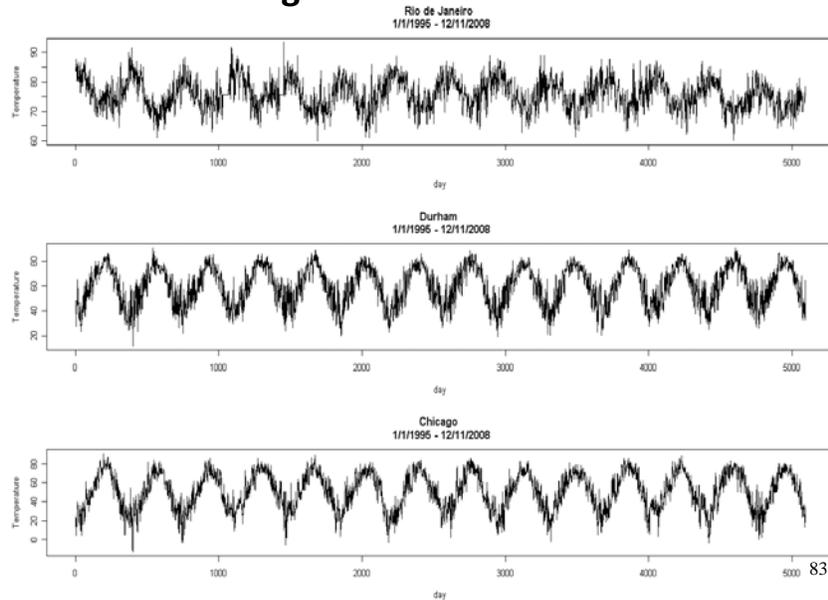
81

$\bar{X} = 2.022013$   
 $St.Dev. = 1.640233$   
 $Skewness = 2.165848$   
 $Excess\ kurtosis = 8.029811$   
 $Q1 = 1$   
 $Q2 = 2$   
 $Q3 = 3$



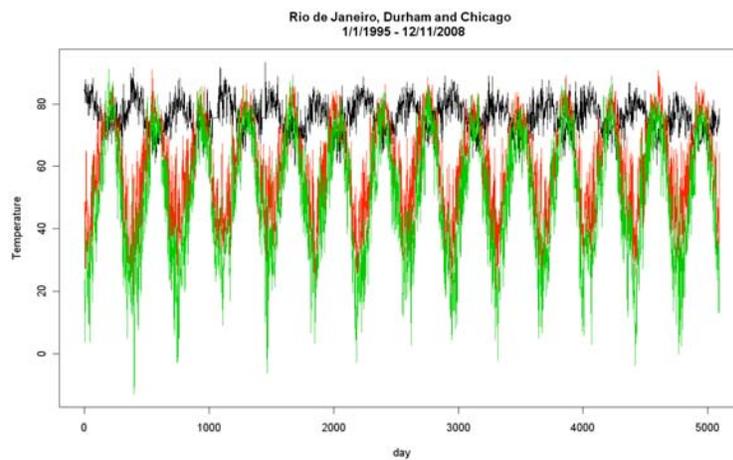
82

## Example: Average daily temperature in Rio de Janeiro, Durham and Chicago



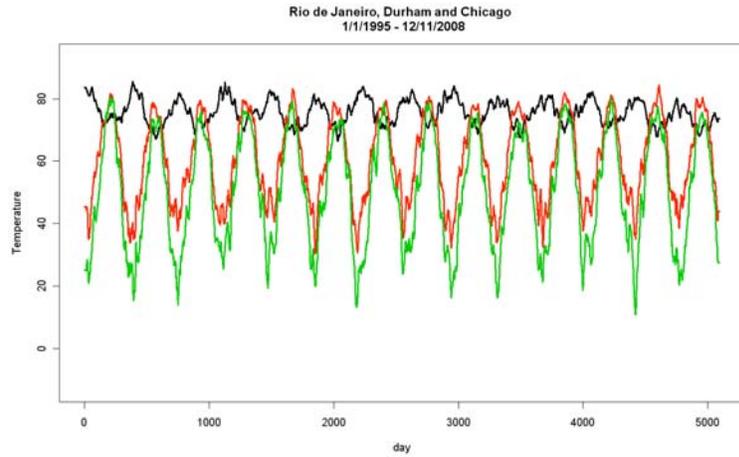
83

Seasonality is more pronounced in Durham and Chicago.  
Variability is also higher in Durham and Chicago.  
Longer winters in Chicago (really?!?)



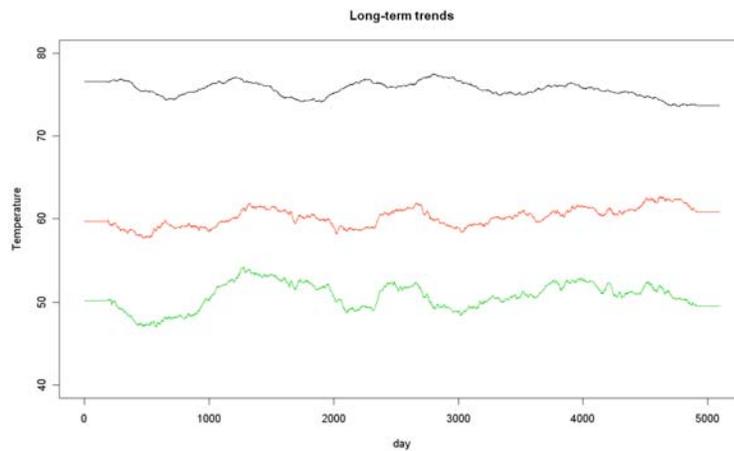
84

The time series were smoothed by replacing each observation by the average of 21 neighboring days, 10 to the left and 10 to the right of the observation.  
Smoothing the time series helps to highlight the short-term patterns.



85

The time series were smoothed by replacing each observation by the average of 364 neighboring days, 182 to the left and 182 to the right of the observation.  
Smoothing the time series helps to highlight the long-term patterns.



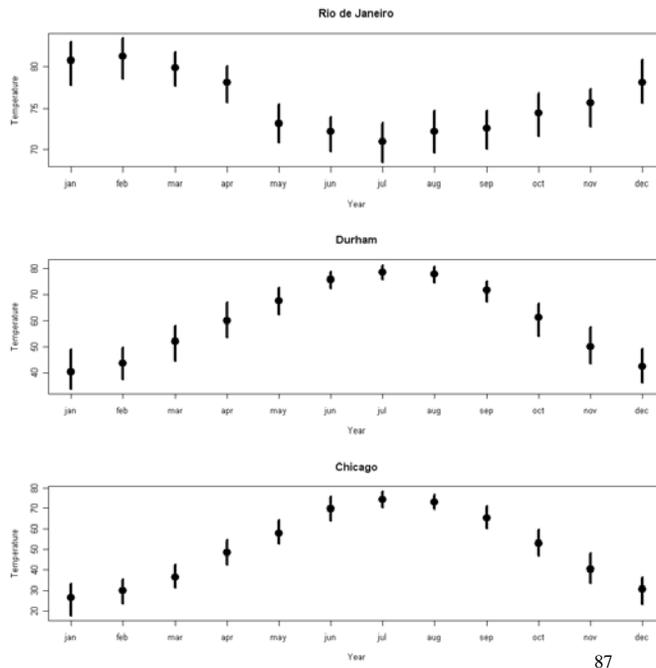
86

# Monthly behavior

**Rio:** variability seems to be constant throughout the year.

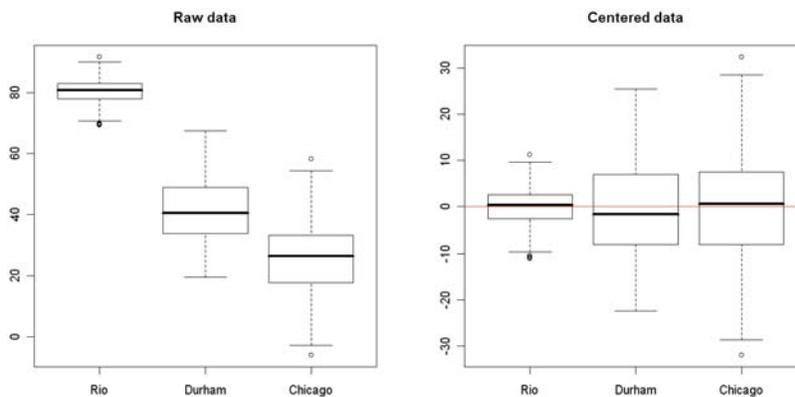
**Durham, Chicago:** variability seems to be higher during colder months than during warmer months.

Dot: medians  
vertical bar: Q1 to Q3



87

# January behavior



Rio is the warmest place in January (it is summer there!)  
Even Durham is much warmer than Chicago (what am I doing here?)  
Temperature in Chicago is the most variable.

88

## Example: Highest temperatures in the USA

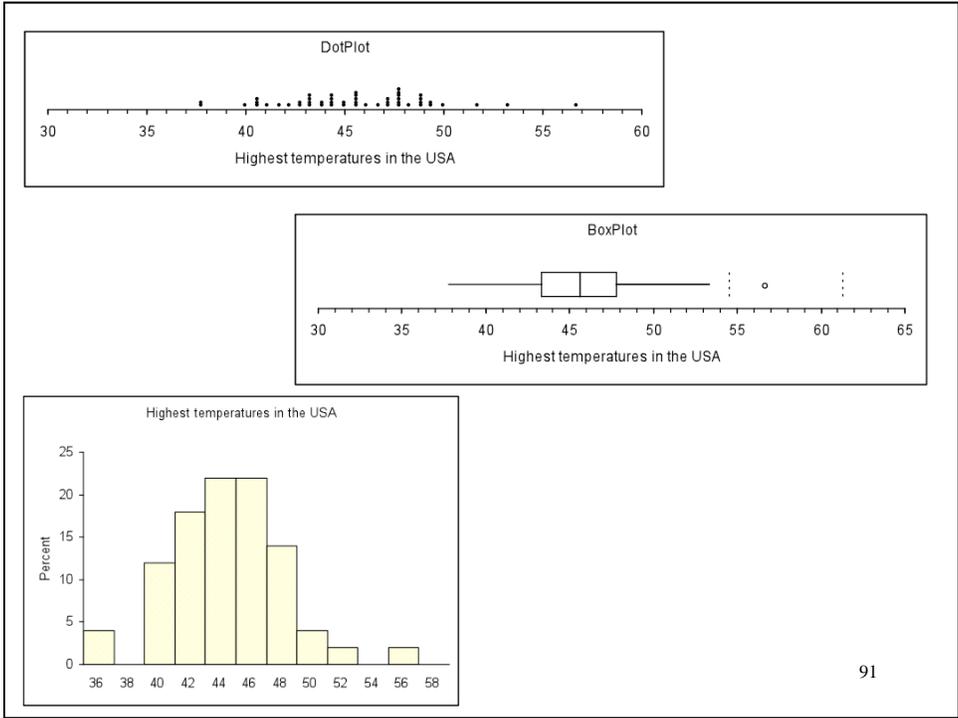
HAWAII	37.8	GEORGIA	44.4	IOWA	47.8
ALASKA	37.8	ALABAMA	44.4	NEBRASKA	47.8
RHODE-ISLAND	40	WEST-VIRGINIA	44.4	WASHINGTON	47.8
CONNECTICUT	40.6	MICHIGAN	44.4	IDAHO	47.8
MAINE	40.6	TENNESSEE	45	COLORADO	47.8
VERMONT	40.6	OHIO	45	OREGON	48.3
NEW-HAMPSHIRE	41.1	LOUISIANA	45.6	TEXAS	48.9
MASSACHUSETTS	41.7	KENTUCKY	45.6	OKLAHOMA	48.9
NEW-YORK	42.2	WISCONSIN	45.6	ARKANSAS	48.9
FLORIDA	42.8	MINNESOTA	45.6	SOUTH-DAKOTA	48.9
MARYLAND	42.8	WYOMING	45.6	KANSAS	49.4
DELAWARE	43.3	MISSISSIPPI	46.1	NORTH-DAKOTA	49.4
VIRGINIA	43.3	INDIANA	46.7	NEW-MEXICO	50
NEW-JERSEY	43.3	ILLINOIS	47.2	NEVADA	51.7
NORTH-CAROLINA	43.3	UTAH	47.2	ARIZONA	53.3
SOUTH-CAROLINA	43.9	MONTANA	47.2	CALIFORNIA	56.7
PENNSYLVANIA	43.9	MISSOURI	47.8		

89

### Highest temperatures

Count	50
Mean	45.604
sample variance	13.901
sample standard deviation	3.728
Minimum	37.8
Maximum	56.7
Range	18.9
mean - 2s	38.147
mean + 2s	53.061
percent in interval (95.44%)	92.0%
mean - 3s	34.419
mean + 3s	56.789
percent in interval (99.73%)	100.0%
Skewness	0.279
Kurtosis	0.728
1st quartile	43.300
Median	45.600
3rd quartile	47.800
interquartile range	4.500

90



91

## Example: US 2004 unemployment rates

### US 2004 unemployment rates | (as percentage of the labor force)

Index	State	Rate	Index	State	Rate
1	HAWAII	3.3	27	UTAH	5.2
2	NORTH DAKOTA	3.4	28	KENTUCKY	5.3
3	SOUTH DAKOTA	3.5	29	WEST VIRGINIA	5.3
4	VERMONT	3.7	30	TENNESSEE	5.4
5	VIRGINIA	3.7	31	COLORADO	5.5
6	NEBRASKA	3.8	32	KANSAS	5.5
7	NEW HAMPSHIRE	3.8	33	NORTH CAROLINA	5.5
8	WYOMING	3.9	34	PENNSYLVANIA	5.5
9	DELAWARE	4.1	35	ALABAMA	5.6
10	MARYLAND	4.2	36	ARKANSAS	5.7
11	NEVADA	4.3	37	LOUISIANA	5.7
12	MONTANA	4.4	38	MISSOURI	5.7
13	GEORGIA	4.6	39	NEW MEXICO	5.7
14	MAINE	4.6	40	NEW YORK	5.8
15	IDAHO	4.7	41	OHIO	6.1
16	MINNESOTA	4.7	42	TEXAS	6.1
17	FLORIDA	4.8	43	CALIFORNIA	6.2
18	IOWA	4.8	44	ILLINOIS	6.2
19	NEW JERSEY	4.8	45	MISSISSIPPI	6.2
20	OKLAHOMA	4.8	46	WASHINGTON	6.2
21	CONNECTICUT	4.9	47	SOUTH CAROLINA	6.8
22	WISCONSIN	4.9	48	MICHIGAN	7.1
23	ARIZONA	5.0	49	OREGON	7.4
24	MASSACHUSETTS	5.1	50	ALASKA	7.5
25	INDIANA	5.2	51	DISTRICT OF COLUMBIA	8.2
26	RHODE ISLAND	5.2			

Mean ( $\bar{x}$ ) = 5.2078431  
variance = 1.1691373  
standard deviation ( $s$ ) = 1.0812665

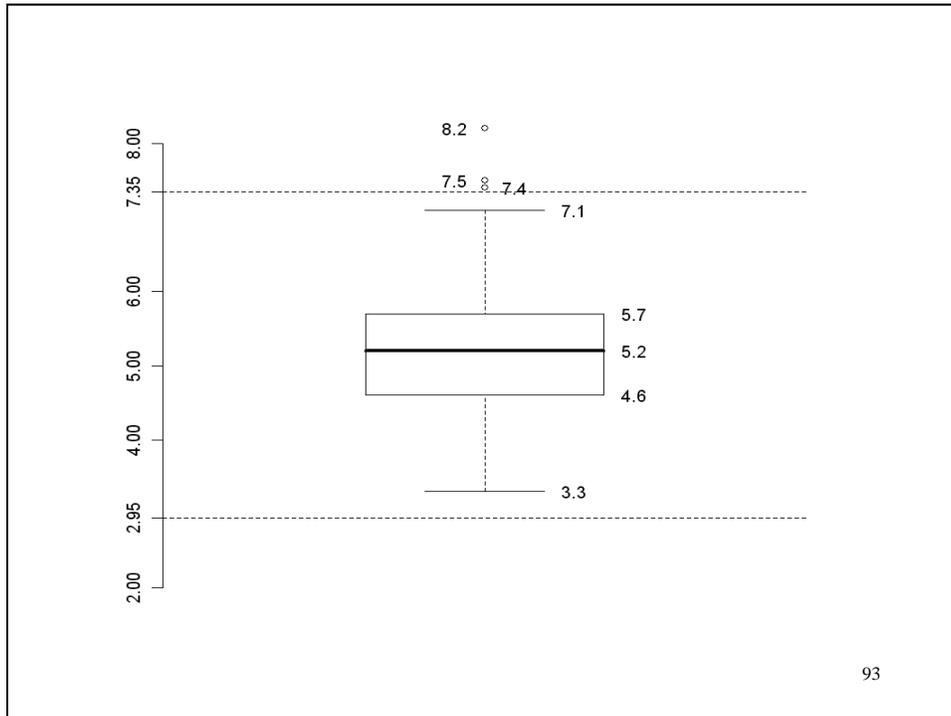
Q1 = 4.6 (Georgia)  
Q2 = 5.2 (Rhode Island)  
Q3 = 5.7 (New Mexico)

skewness = 0.4798145  
kurtosis = 0.3317919

Empirical rule actual coverage

[ $\bar{x}-1*s$ ;  $\bar{x}+1*s$ ] = [4.13; 6.289110] 72.55%  
[ $\bar{x}-2*s$ ;  $\bar{x}+2*s$ ] = [ 3.05; 7.370376] 94.12%  
[ $\bar{x}-3*s$ ;  $\bar{x}+3*s$ ] = [ 1.96; 8.451643] 100.00%

92



## Multivariate Exploratory Data Analysis



1. How to relate two things
2. Correlations and covariances
3. Linearly related variables
  - 3.1 Mean and variance of a linear function
  - 3.2 Linear combinations
  - 3.3 Mean and variance of a linear combination: 2 inputs
  - 3.4 Mean and variance of a linear combination: 3 inputs
  - 3.5 Mean and variance of a linear combination: k inputs
4. Portfolio example
5. Simple linear regression

# Summary of the lecture

In this class you will learn how to

- Relate two sets of variables: **sample linear correlation coefficient**
- Compute sample mean, variance and standard deviation of **linear combinations** of variables
- Study the practical example of **portfolio allocation**

## Book

Skewness (pages 114-117 (12)\*, 113-117 (13))

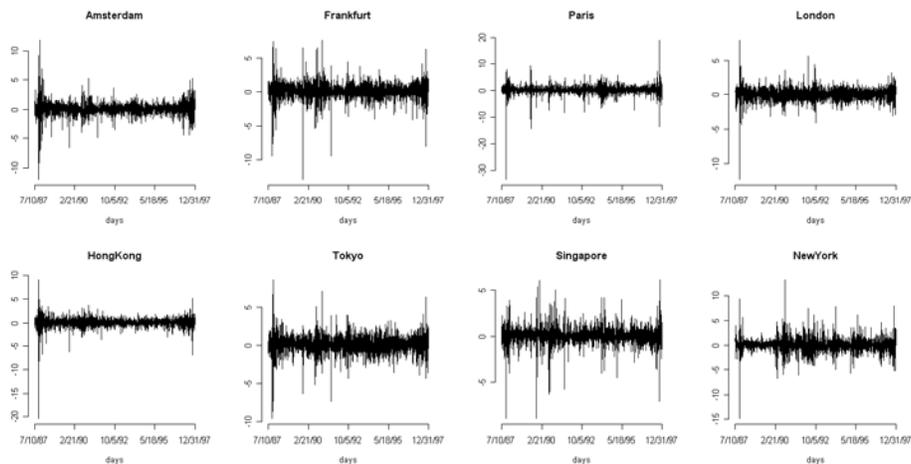
What is correlation analysis? (pages 429-435 (12), 458-465 (13))

\*Number in parenthesis refers to the book edition

95

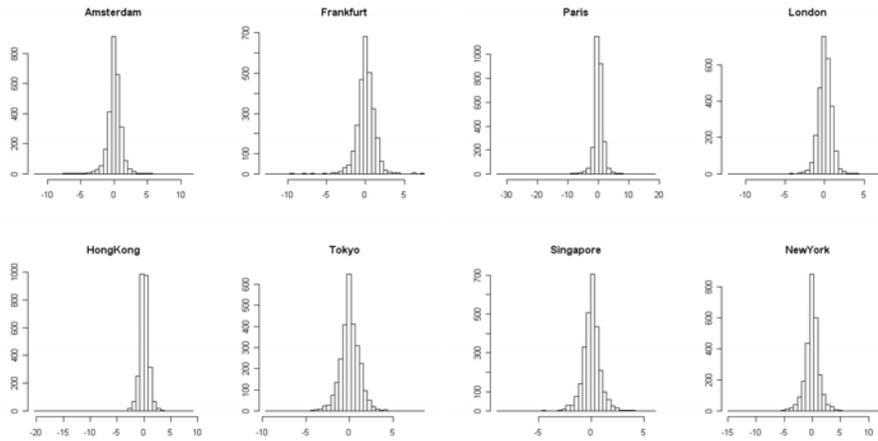
## Example: Comparing international stock returns

July 10, 1987 until December 31, 1997 (2733 days) - Amsterdam (EOE) , Frankfurt (DAX), Paris (CAC40), London (FTSE100), Hong Kong (Hang Seng) Tokyo (Nikkei), Singapore (Singapore All Shares), New York (S&P500).



96

# Histograms



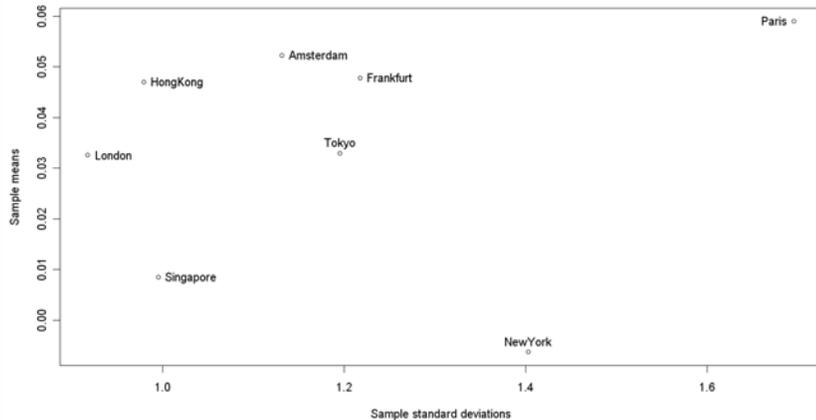
97

# Statistical summary

Country	mean	stdev	skewness	kurtosis
<b>Amsterdam</b>	0.0522	1.1320	-0.3902	18.0457
<b>Frankfurt</b>	0.0478	1.2178	-0.8355	12.5641
<b>Paris</b>	0.0590	1.6956	-3.1012	67.3491
<b>London</b>	0.0325	0.9181	-1.4047	23.1069
<b>HongKong</b>	0.0470	0.9798	-3.5694	77.4448
<b>Tokyo</b>	0.0329	1.1956	-0.3647	7.0778
<b>Singapore</b>	0.0085	0.9956	-1.1182	13.5078
<b>NewYork</b>	-0.0064	1.4030	0.1065	10.8264

98

It is considered good to have  
a large mean return  
and  
a small standard deviation.



99

## 1. How to Relate Two Things

The mean and standard deviation help us summarize a bunch of numbers which are measurements of just one thing (one variable)

A fundamental and totally different question is **how one thing relates to another.**

In this section of the notes we look at **scatter plots** and how **covariance** and **correlation** can be used to summarize them.

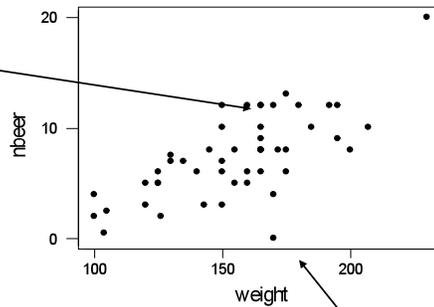
When examining two things (variables) at the time, the **scatter plot** will be our main graphical tool whereas **covariance** and **correlation** will be our main numerical summaries.

100

## Example

Is the number of beers you can drink related to your weight?

<u>nbeer</u>	<u>weight</u>	
12.0	192	1
12.0	160	2
5.0	155	3
5.0	120	4
7.0	150	5
13.0	175	6
4.0	100	7
12.0	165	8
12.0	165	9
12.0	150	10
.	.	.
.	.	.



Scatter plot

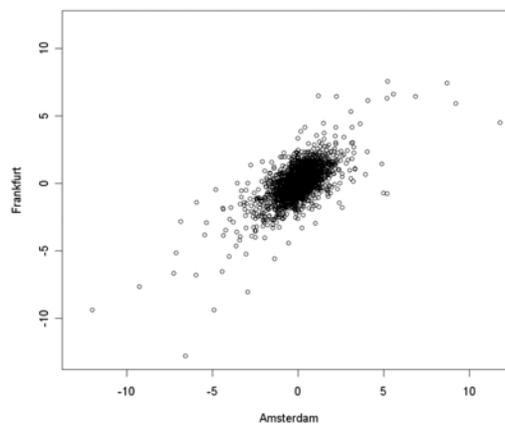
Now we think of each pair of numbers as an observation.  
Each pair corresponds to a person.  
Each person has two numbers associated with him/her,  
# beers and weight.  
Each pair corresponds to a point on the plot.

101

## Example

International stock returns: Amsterdam and Frankfurt.

Each point corresponds to a day.



102

In general we have observations

$(x_i, y_i)$  ← the  $i$ th observation is a pair of numbers

and each point on the plot corresponds to an observation.

Our data looks like:

x	y	i
12.0	192	1
12.0	160	2
5.0	155	3
5.0	120	4
7.0	150	5
13.0	175	6
4.0	100	7
12.0	165	8

The plot enables us to see the relationship between x and y

.....

103

## 2. Covariance and Correlation

In both examples it does look like there is a relationship.

Even more, the relationship looks linear in that it looks like we could draw a line through the plot to capture the pattern.

**Covariance** and **correlation** summarize how strong a **linear** relationship there is between two variables.

In our first example weight and # beers were two variables. In our second example our two variables were two kinds of returns.

In general, we think of the two variables as x and y.

104

## Historical note

**1885:** Sir Francis Galton: studying the heights of children versus the heights of parents.

There's a *regression-back-to-the-mean effect*. If your parents are on average higher than the average, you'll regress back to the average.

**1888:** Co-relation: slope of the least-squares regression line for data in standardized (by median and quartile range) form

**1896:** Karl Pearson, product moment definition  
The misuse of correlation has multiplied faster than the proper use of it !

105

**The sample covariance** between x and y:

$$s_{xy} = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

**The sample correlation** between x and y:

$$r_{xy} = \frac{s_{xy}}{s_x s_y}$$

So, the correlation is just the covariance divided by the two standard deviations.

106

We will get some intuition about these formulae, but first let us see them in action. How do they summarize data for us? Let us start with the correlation.

**Correlation**, the facts of life:

$$-1 \leq r_{xy} \leq 1$$

The **closer r is to 1** the stronger the linear relationship is with a **positive slope**.  
When one goes up, the other tends to go up.

The **closer r is to -1** the stronger the linear relationship is with a **negative slope**.  
When one goes up, the other tends to go down.

107

The correlations corresponding to the two **scatter plots** we looked at are:

Correlation of amsterdam and frankfurt = 0.677

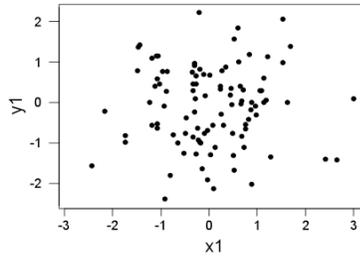
Correlation of nbeer and weight = 0.692

The **larger correlation between nbeer and weight** indicates that the linear relationship is stronger.

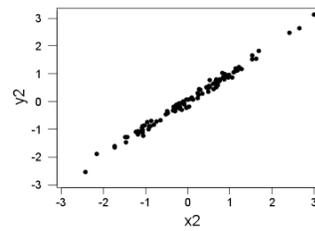
Let us look at some more examples.

108

Correlation of  
y1 and x1 = 0.019

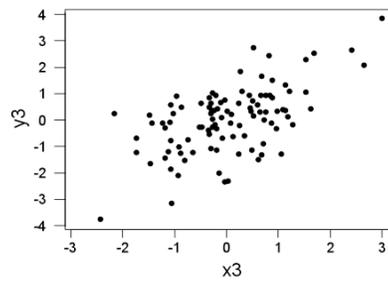


Correlation of  
y2 and x2 = 0.995

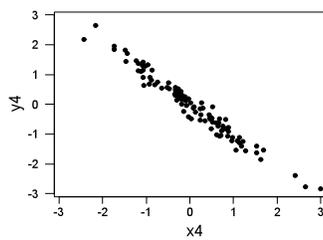


109

Correlation of  
y3 and x3 = 0.586

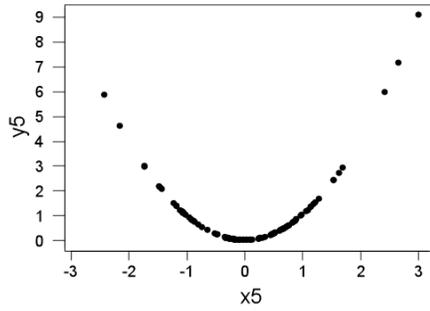


Correlation of  
y4 and x4 = -0.982



110

Correlation of  $y_5$  and  $x_5 = 0.210$

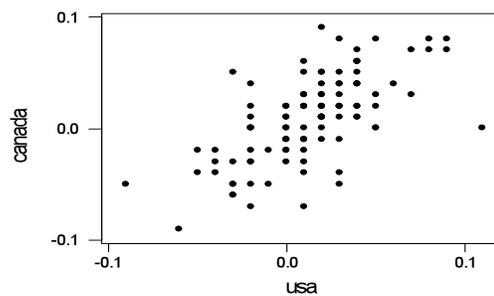


The correlation only measures **linear** relationships (here the value is small but there is a strong nonlinear relationship between  $y_5$  and  $x_5$ .)

111

### Example: The country data

Which countries go up and down together?  
I have data on 23 countries.  
That would be a lot of plots!



112

## Example: International stock returns

To summarize, we can compute all pair wise correlations:

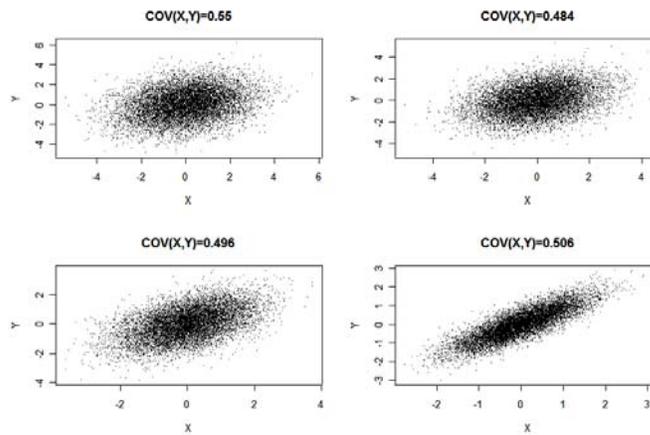
	Amsterdam	Frankfurt	Paris	London	HongKong	Tokyo	Singapore	NewYork
Amsterdam	1.000							
Frankfurt	0.678	1.000						
Paris	0.345	0.393	1.000					
London	0.657	0.481	0.280	1.000				
HongKong	0.408	0.284	0.177	0.419	1.000			
Tokyo	0.653	0.607	0.298	0.565	0.340	1.000		
Singapore	0.307	0.371	0.462	0.248	0.174	0.292	1.000	
NewYork	0.284	0.295	0.267	0.302	0.118	0.243	0.298	1.000

Why is this blank?

113

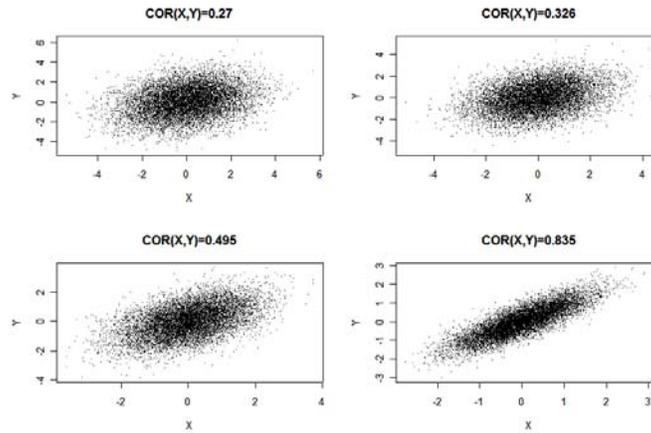
## Why compute both covariance and correlation?

The four covariances are around 0.5.....



114

....but the four correlations are rather different.



**Note: Correlations are unit free. They are between -1 and 1. Covariance, on the other hand, carries the units of X and of Y.**

115

### 3 Linearly Related Variables

We have studied data sets that display some kind of relation with each other (the mutual fund returns and the market returns, for instance).

Sometimes there is an exact linear relation between variables:

$$y = c_0 + c_1 x$$

Can we say something about the **sample mean** of y if all we know is the sample mean of x (and vice versa)?

Can we say something about the **sample standard deviation** of y if all we know is the sample standard deviation of x (and vice versa)?

We will answer these questions in the sequel.

116

### Example

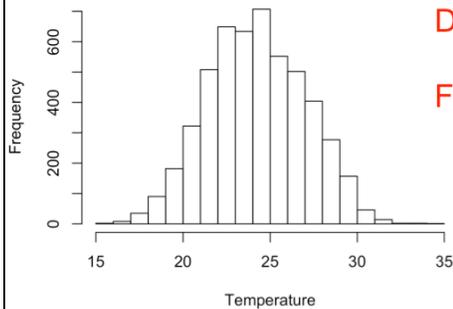
Suppose we have daily temperatures (in Celsius degree) in Rio de Janeiro from January 1<sup>st</sup>, 1995 to December, 11<sup>th</sup> 2008.

We also know that the sample mean and the sample variance for the daily temperature for this period are 24.24C and 2.78C.

What in the hell are Celsius degree?

Don't panic!!!!

$F = 32 + 1.8C$



117

In general, we like to use the symbols  $y$  and  $x$  for the two variables

The variable  $y$  is a linear function of the variable  $x$  if:

$$y = c_0 + c_1x$$

$c_0$  : the intercept

$c_1$  : the slope

We think of the  $c$ 's as constants (fixed numbers) while  $x$  and  $y$  vary.

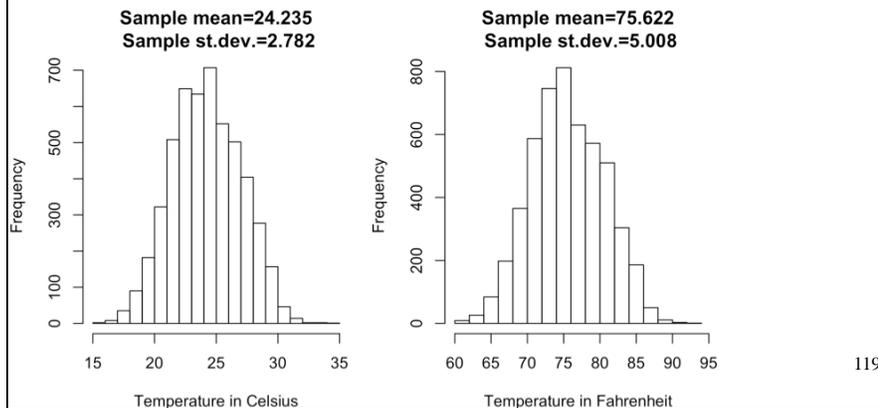
118

### 3.1 Mean and variance of a linear function

How are the mean and variance of  $y$  related to those of  $x$ ?

Let us look at our temperature example.

It is not a coincidence that  $32+1.8*24.235=75.622$  and that  $1.8*2.782=5.008$



Suppose

$$y = c_0 + c_1 X$$

Then,

$$\bar{y} = c_0 + c_1 \bar{X}$$

$$s_y^2 = c_1^2 s_x^2$$

$$s_y = |c_1| s_x$$

Recall that  $|x|$  is the absolute value of  $x$ . For instance,  $|-5|=5$  and  $|10|=10$

120

### **3.2. Linear combinations**

We may want a variable to be related to several others instead of just one. We will assume that Y is a function of X,Z,...rather than just a function of X.

When a variable y is linearly related to several others, we call it a **linear combination**.

$$y = C_0 + C_1X_1 + C_2X_2 + \dots + C_kX_k$$

y is a linear combination of the x's.  
c<sub>i</sub> is the coefficient of x<sub>i</sub>.

121

### **Example: house pricing**

Home	Nbhd	Offers	SqFt	Brick	Bed	Bath	Price
1	2	2	1790	No	2	2	114300
2	2	3	2030	No	4	2	114200
3	2	1	1740	No	3	2	114800
4	2	3	1980	No	3	2	94700
5	2	3	2130	No	3	3	119800
6	1	2	1780	No	3	2	114600

We will see later, when studying multiple linear regression, that the price can be modeled as a linear combination of the other variables.

The following formula relates the expected sales price of a house (Price) to its size (SqFt), number of bedrooms (Bed) and number of bathrooms (Bath):

$$\text{Price} = -5640.83 + 35.64*\text{SqFt} + 10459.93*\text{Bed} + 13546.13*\text{Bath}$$

122

**Example: Portfolio allocation**

Let us use country returns and suppose that we had put 0.5 into USA and 0.5 into Hong Kong, ie.

$$\text{port} = 0.5 \cdot \text{hongkong} + 0.5 \cdot \text{usa}$$

What would our returns have been?

hongkong	usa	port
0.02	0.04	0.030
0.06	-0.03	0.015
0.02	0.01	0.015
-0.03	0.01	-0.010
0.08	0.05	0.065
.....		

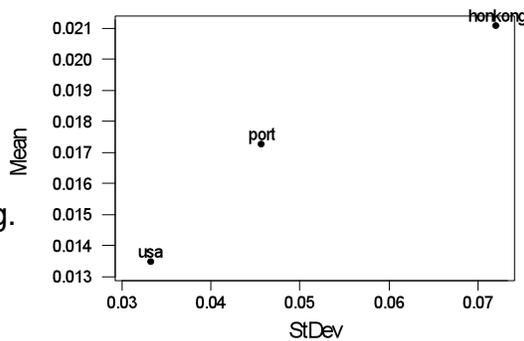
For each month, we get the portfolio return as  $\frac{1}{2} \cdot \text{hongkong} + \frac{1}{2} \cdot \text{usa}$ .

123

How do the returns on this portfolio compare with those of Hong Kong and USA?

It looks like the mean for my portfolio is right in between the means of USA and Hong Kong.

What about the standard deviation?



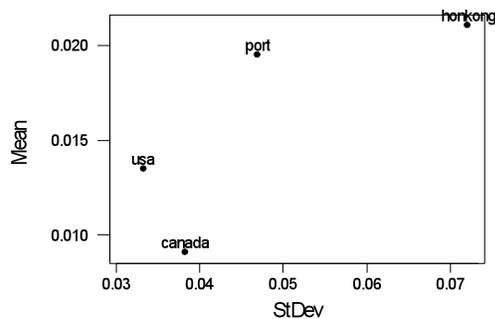
124

Let us try a portfolio with three stocks.

Let us go short on Canada (i.e., we borrow Canada to invest in the other stocks), ie.

$$\text{port} = -0.5 \cdot \text{canada} + 1.0 \cdot \text{usa} + 0.5 \cdot \text{hongkong}$$

Clearly,  
forming  
portfolios  
is an interesting  
thing to do!



125

Basic question: **why would we form portfolios?**

Maybe the portfolio has a nice mean and variance (i.e. nice “average return” and nice “risk”)

There are some basic formulae that relate the mean and standard deviation of a linear combination to the means, variances and covariances of the input variables.

We can apply these formulae to understand how the mean and variance of a portfolio depend on the input assets. These formulae constitute the basic part of the tool-kit of those who really understand finance.

126

### 3.3. Mean and variance of a linear combination: 2 inputs

First, we consider the case where we have only two inputs.

#### **2 inputs:**

$$\text{Suppose } y = c_0 + c_1x_1 + c_2x_2$$

Then,

$$\bar{y} = c_0 + c_1\bar{x}_1 + c_2\bar{x}_2$$

$$s_y^2 = c_1^2s_{x_1}^2 + c_2^2s_{x_2}^2 + 2c_1c_2s_{x_1x_2}$$

127

#### Example: Portfolio means

$$\text{Port} = 0.5 \cdot \text{hongkong} + 0.5 \cdot \text{usa}$$

hongkong	usa	port
0.02	0.04	0.030
0.06	-0.03	0.015
0.02	0.01	0.015
-0.03	0.01	-0.010
0.08	0.05	0.065
.....		

For each month, we get the portfolio return as  $\frac{1}{2} \cdot \text{hongkong} + \frac{1}{2} \cdot \text{usa}$ .

The mean returns on USA, and Hong Kong are 0.01346, and 0.02103

The mean return on Port is  $0.5 \cdot 0.01346 + 0.5 \cdot 0.02103 = 0.01724$

128

Let us do the same exercise for the variance:

Covariance matrix

	hongkong	usa	
hongkong	0.00521		The diagonals are variances, The off diagonals are Covariances.
usa	0.00103	0.00111	

As before, we apply the formula:

$$\begin{aligned}\text{Var}(\text{Port}) &= (0.5)^2 \cdot 0.00521 + (0.5)^2 \cdot 0.00111 + 2 \cdot (0.5) \cdot (0.5) \cdot 0.001 \\ &= 0.25 \cdot 0.00521 + 0.25 \cdot 0.00111 + 0.5 \cdot 0.001 = 0.0021.\end{aligned}$$

129

Let us do it one more time:

$$\text{Port} = 0.25 \cdot \text{usa} + 0.75 \cdot \text{hongkong}$$

$$\begin{aligned}\text{Var}(\text{Port}) &= \\ &= (0.25)^2 \cdot 0.00111 + \\ &+ (0.75)^2 \cdot 0.00521 + \\ &+ (2) \cdot (0.25) \cdot (0.75) \cdot (0.00103) \\ &= 0.0033\end{aligned}$$

130

### 3.4. Mean and variance of a linear combination: 3 inputs

Second, we consider the case where we have three inputs.

#### **3 inputs:**

Suppose

$$y = c_0 + c_1x_1 + c_2x_2 + c_3x_3$$

Then,

$$\bar{y} = c_0 + c_1\bar{x}_1 + c_2\bar{x}_2 + c_3\bar{x}_3$$

$$s_y^2 = c_1^2s_{x_1}^2 + c_2^2s_{x_2}^2 + c_3^2s_{x_3}^2 + 2(c_1c_2s_{x_1x_2} + c_1c_3s_{x_1x_3} + c_2c_3s_{x_2x_3})$$

131

### Example: Portfolio based on fidel, eqmrkt and windsor funds.

```
port = 0.1*fidel+0.4*eqmrkt+0.5*windsor
```

Covariance matrix

	fidel	eqmrkt	windsor
fidel	0.003202		
eqmrkt	0.003190	0.004700	
windsor	0.002410	0.002990	0.0023658

$$\begin{aligned}\text{Var}(\text{port}) &= (0.1)*(0.1)*0.003202 + \\ & (0.4)*(0.4)*0.0047 + \\ & (0.5)*(0.5)*0.0023658 + \\ & 2*((0.1)*(0.4)*0.00319+(0.1)*(0.5)*0.00241+(0.4)*(0.5)*0.00299) = \\ & 0.0030676\end{aligned}$$

132

### 3.5. Mean and variance of a linear combination: k inputs

**K inputs:** Suppose

$$y = c_0 + c_1x_1 + c_2x_2 + \dots + c_kx_k$$

then,

$$\bar{y} = c_0 + c_1\bar{x}_1 + c_2\bar{x}_2 + \dots + c_k\bar{x}_k$$

$$s_y^2 = c_1^2s_{x_1}^2 + c_2^2s_{x_2}^2 + \dots + c_k^2s_{x_k}^2 + 2\sum_{i<j} c_i c_j s_{x_i x_j}$$

133

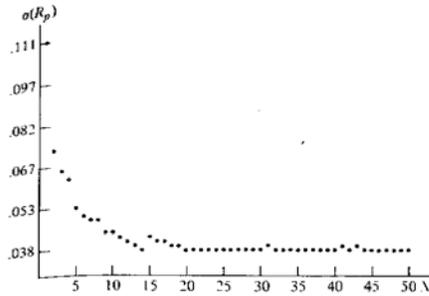
### 4. Portfolio example

Cut from a Finance Textbook:

Fama [1976] has illustrated this result empirically.<sup>11</sup> His results are shown in Fig. 6.18. He randomly selected 50 securities listed on the New York Stock Exchange and calculated their standard deviations using monthly data from July 1963 to June 1968. Then a single security was selected randomly. Its standard deviation of return was around 11%. Next, this security was combined with another (also randomly selected) to form an equally weighted portfolio of two securities. The standard deviation fell to around 7.2%. Step by step more securities were randomly added to the portfolio until all 50 securities were included. Almost all of the diversification was obtained after the first 10-15 securities were randomly selected. In addition the portfolio stan-

<sup>11</sup> See Fama [1976], *Foundations of Finance*, pp. 253-254.

134



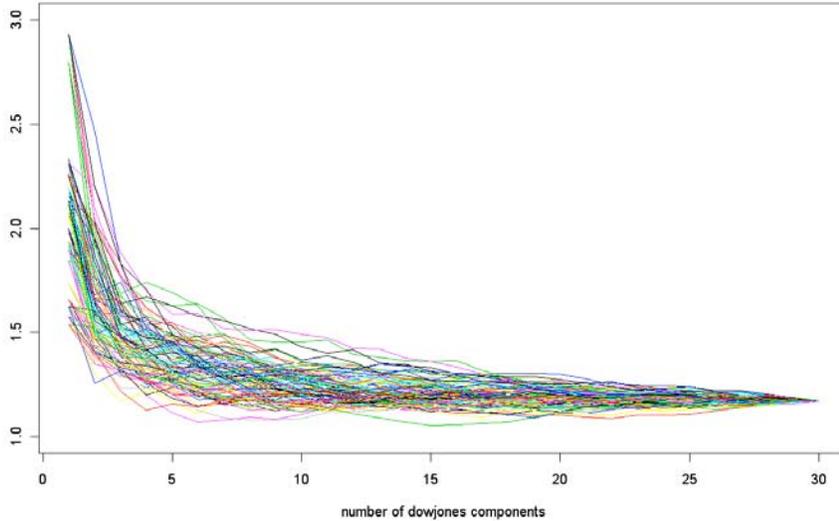
**Figure 6.18**  
The standard deviation of portfolio return as a function of the number of securities in the portfolio. (From Fama, E. F., *Foundations of Finance*, reprinted with permission of the author.)

standard deviation quickly approached a limit which is roughly equal to the average covariance of all securities. One of the practical implications is that most of the benefits of diversification (given a random portfolio selection strategy) can be achieved with fewer than 15 stocks.

**Dowjones components: January 1997 to December 2006 - 2516 observations**

Company	tick	mean	stdev	skew	kurt
1 ALCOA	AA	0.058	2.258	0.354	2.909
2 American Intl Group	AIG	0.060	1.890	0.229	3.206
3 AMERICAN EXPRESS	AXP	0.078	2.157	0.106	3.181
4 BOEING CO	BA	0.049	2.132	-0.348	6.105
5 CITIGROUP	C	0.087	2.191	0.292	5.832
6 CATERPILLAR	CAT	0.079	2.111	-0.079	3.107
7 DU PONT (EI)	DD	0.031	1.898	0.164	2.835
8 DISNEY (WALT) CO	DIS	0.043	2.190	0.057	6.740
9 GENERAL ELECTRIC	GE	0.057	1.840	0.218	3.777
10 GENERAL MOTORS CORP	GM	0.027	2.249	0.273	4.270
11 HOME DEPOT	HD	0.080	2.298	-0.672	12.629
12 HONEYWELL INTL	HON	0.047	2.334	0.196	14.129
13 HEWLETT-PACKARD	HPQ	0.072	2.796	0.219	5.505
14 IBM	IBM	0.062	2.076	0.202	6.607
15 INTEL CORP	INTC	0.054	2.930	-0.086	4.904
16 JOHNSON&JOHNSON	JNJ	0.057	1.539	-0.270	6.966
17 JP MORGAN CHASE	JPM	0.059	2.313	0.351	5.494
18 COCA-COLA CO	KO	0.017	1.659	0.030	4.232
19 MCDONALDS CORP	MCD	0.048	1.844	0.102	4.180
20 3M CO	MMM	0.047	1.625	0.254	3.687
21 Altria Group	MO	0.074	2.057	0.156	7.074
22 MERCK & CO	MRK	0.035	1.915	-1.067	18.368
23 MICROSOFT CORP	MSFT	0.073	2.249	0.122	6.019
24 PFIZER	PFE	0.051	1.990	-0.106	2.661
25 PROCTER & GAMBLE	PG	0.057	1.736	-2.455	45.731
26 AT&T	T	0.047	1.998	0.058	2.766
27 UNITED TECH CORP	UTX	0.078	1.933	-1.216	19.988
28 VERIZON COMMUNICATIONS	VZ	0.039	1.894	0.249	3.844
29 WAL-MART STORES	WMT	0.078	1.974	0.310	2.713
30 EXXON MOBIL CORP	XOM	0.067	1.575	0.150	2.634

## 100 replications of Fama's exercise



137

### Weights for the minimum variance portfolio

Companies	tick	weight
1 ALCOA	AA	0.0111
2 American Intl Group	AIG	0.0038
3 AMERICAN EXPRESS	AXP	-0.0439
4 BOEING CO	BA	0.0379
5 CITIGROUP	C	-0.0453
6 CATERPILLAR	CAT	0.0138
7 DU PONT (EI)	DD	0.0055
8 DISNEY (WALT) CO	DIS	0.0387
9 GENERAL ELECTRIC	GE	-0.0574
10 GENERAL MOTORS CORP	GM	0.0367
11 HOME DEPOT	HD	-0.0151
12 HONEYWELL INTL	HON	-0.0271
13 HEWLETT-PACKARD	HPQ	0.0170
14 IBM	IBM	0.0732
15 INTEL CORP	INTC	-0.0313
16 JOHNSON&JOHNSON	JNJ	0.1427
17 JP MORGAN CHASE	JPM	-0.0062
18 COCA-COLA CO	KO	0.0845
19 MCDONALDS CORP	MCD	0.1005
20 3M CO	MMM	0.1287
21 Altria Group	MO	0.0836
22 MERCK & CO	MRK	0.0332
23 MICROSOFT CORP	MSFT	0.0555
24 PFIZER	PFE	-0.0153
25 PROCTER & GAMBLE	PG	0.0910
26 AT&T	T	0.0085
27 UNITED TECH CORP	UTX	0.0139
28 VERIZON COMMUNICATIONS	VZ	0.0891
29 WAL-MART STORES	WMT	0.0315
30 EXXON MOBIL CORP	XOM	0.1410

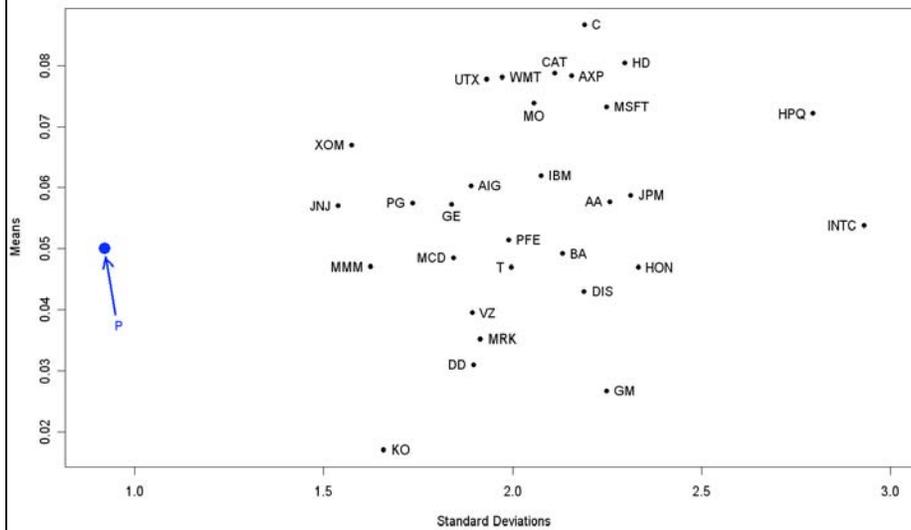
You will learn everything about the minimum variance portfolio in an Investments course.

For now, just keep in mind it is a portfolio whose variance is smaller than other portfolios.

Mean = 0.050  
 Variance = 0.921  
 Stdev = 0.960  
 Kurtosis = 3.056  
 Skewness = -0.161

138

## Mean-standard deviation plot



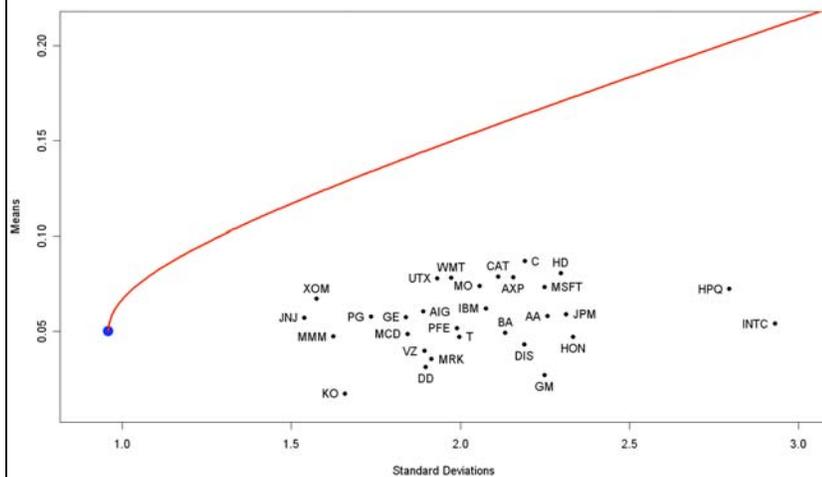
139

## Portfolios based on the 30 Dowjones components

**Blue dot:** Minimum variance portfolio.

**Red line:** Minimum variance portfolio for a given mean return target.

Positive and negative weights are allowed, as long as they add up to 1.



140

# Portfolios based on 8 Dowjones components

**Blue dot:** Minimum variance portfolio.

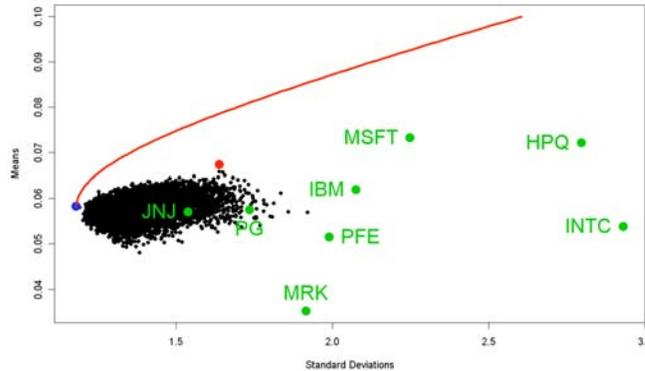
**Red line:** Minimum variance portfolio for a given mean return target.

Positive and negative weights are allowed, as long as they add up to 1.

**Black dots:** Several randomly selected portfolios with weights between 0 and 1 and adding up to 1.

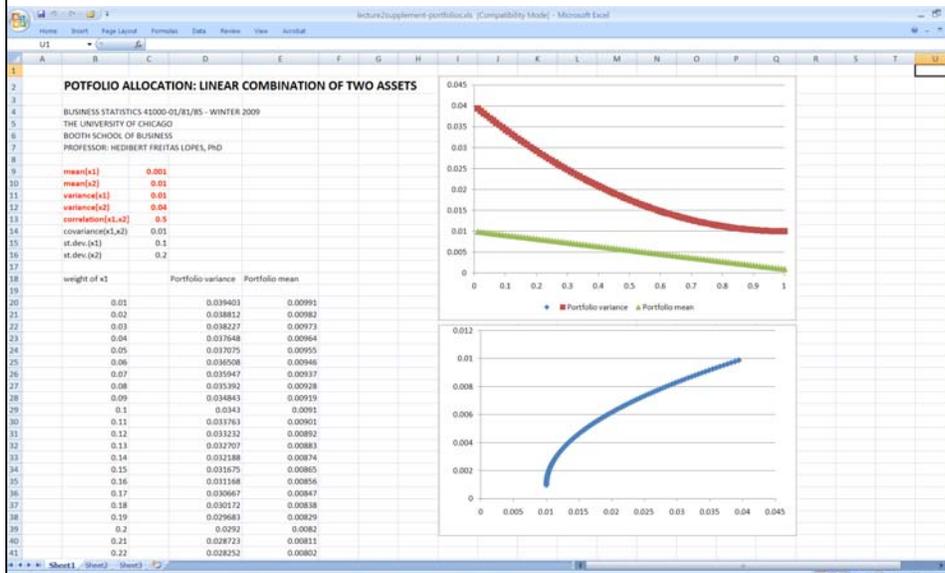
## Red dot portfolio

COMPANY	weight
HEWLETT-PACKARD	0.03
IBM	0.12
INTEL CORP	0.17
MICROSOFT CORP	0.09
JOHNSON&JOHNSON	0.09
MERCK & CO	0.27
PFIZER	0.06
PROCTER & GAMBLE	0.18



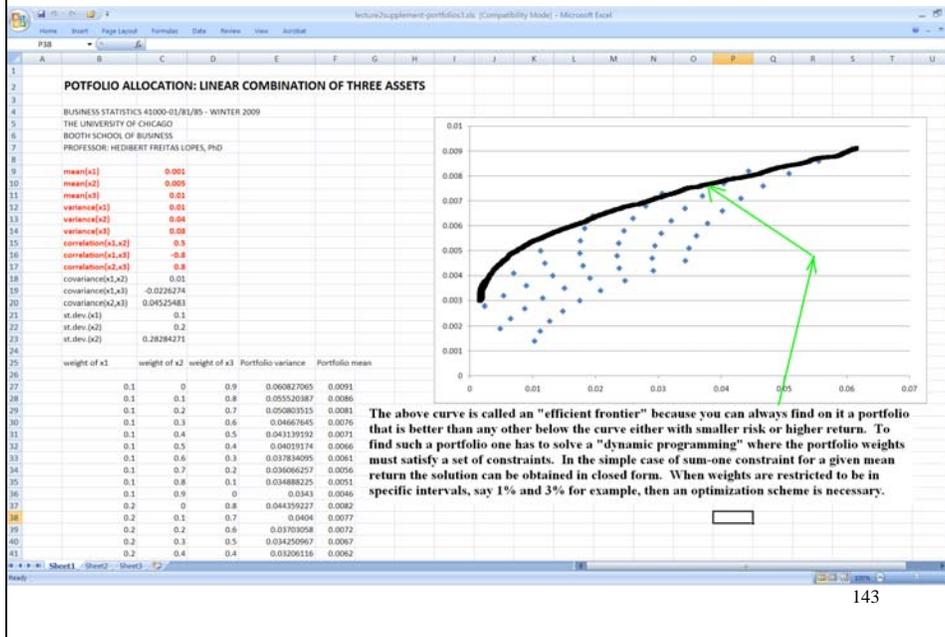
141

# Excel: Constrained portfolios with 2 assets



142

## Excel: Constrained portfolios with 3 assets

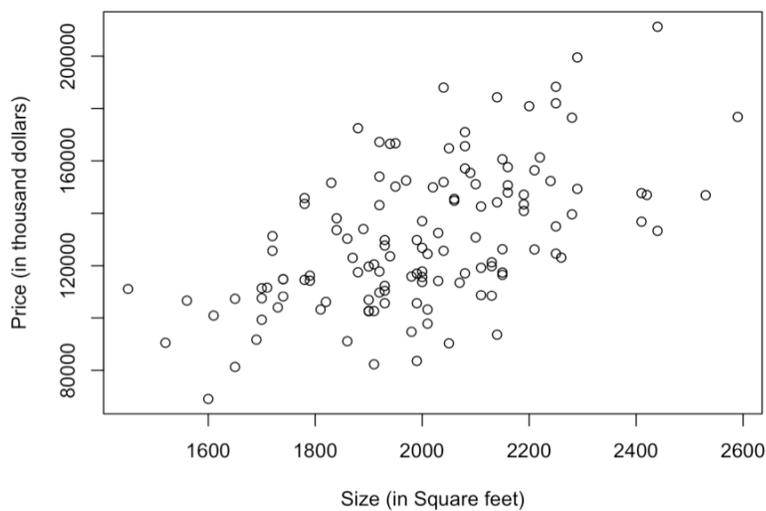


143

## 5. Simple Linear Regression

This is data on 128 homes.

$x$ =size (square feet)  $y$  = price (dollars)



144

Covariance matrix

	SqFt	Price
SqFt	44762.89	3143533
Price	3143533.22	721930821

Hard to say what "721930821" means.

Correlation matrix

	SqFt	Price
SqFt	1.0000000	0.5529822
Price	0.5529822	1.0000000

That is better!

Size and Price are clearly linearly correlated!

145

But what is the equation of the line you would draw through the data?

Linear regression fits a line to the plot.

When I "run a regression" I get values for the intercept and the slope

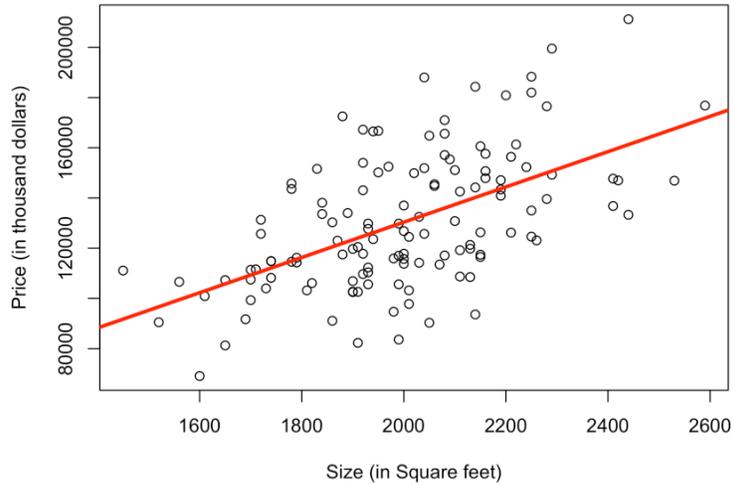
PRICE = intercept + slope\*SIZE

PRICE = -10091.13 + 70.23\*SIZE

146

Here is the scatter plot with the line drawn through it.

Looks reasonable!



### 5.1. Regression and Prediction

Suppose you had a house and you knew the size = 2000 but you do not know the price.

How could you use regression to guess or "predict" the price?

Just plug the size into the equation of the line:

$$\begin{aligned}\text{estimated price} &= -10091.13 + 70.23 \cdot 2000 \\ &= 130368.9\end{aligned}$$

148

Correlation and covariance are "symmetric".

The covariance between  $y$  and  $x$  is the same thing as the covariance between  $x$  and  $y$ .

**Regression is not symmetric.**

We regress  $y$  on  $x$ .

$y$ : dependent variable

$x$ : independent variable.

We say that " $y$  depends on  $x$ ".

In our example  $y$ =price depends on  $x$ =size.

149

## Basic Probability



1. Probability and Random Variables
2. Bivariate Random Variables
3. The Marginal Distribution
4. The Conditional Distribution
5. Independence
6. Computing Joints from Conditionals and Marginals

150

## Summary of the lecture

In this lecture we will enter the realm of **statistical modeling**. However, in order to set the stage for more complex scenarios, such as estimation, hypothesis testing and linear regression, we must introduce the notation, the jargon of **probability**. We begin by

- Defining probability and presenting properties;
- **Discrete random variables**: where the outcomes are countable, such as number of votes for candidate A per county, number of children per family, and number of collisions monthly claimed in a certain insurance company;
- **Bivariate random variables by contingency tables**: For instance, should salary level have 4 categories (low,medium,high,extreme) and happiness have 3 categories (unhappy, indifferent, happy), then one could argue that there are 8 joint levels of salary by happiness in a 4 by 3 contingency table;
- **Marginal distributions**: Looking at the margins of a table;
- **Conditional distributions**: looking at a column/row of a table.

151

## Book material

- Chapter 5:
  - Probability, experiment, outcome and event (141-142 (12), 140-141 (13))
  - Events mutually exclusive (143 (12), 142 (13))
  - Events collectively exhaustive (page 144 (12), 143 (13))
  - Classical probability (143 (12), 142 (13))
  - Empirical probability (144 (12), 143 (13))
  - Subjective probability (145 (12), 144 (13))
  - Rules for computing probabilities (147-154 (12), 174-155 (13))
  - Contingency tables (155-157 (12), 156-158 (13))
- Chapter 6
  - Discrete random variable (184 (12 &13))

152

In this section of the course we learn about **random variables** and **probability**.

This is a very important topic that gets used in a variety of situations.



In order to think about many real world problems we have to face the fact that we are **uncertain** about some important aspects of the situation.

153

## Monty Hall Problem

<http://www.youtube.com/watch?v=mhlc7peGIGg>

## Birthday Problem

154

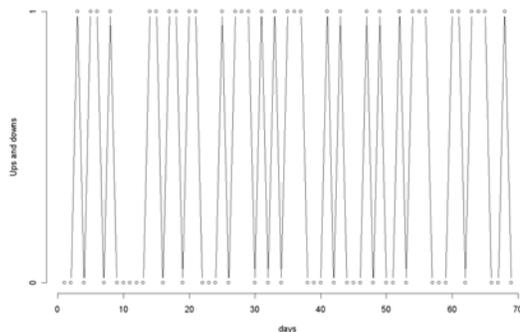
# 1. Probability and Random Variables

## Example 1: S7P500 ups and downs in 2008

69 days  
 33 ups (33 1's)  
 36 downs (36 0's)

Date	SP500 x(t)	Diff x(t)-x(t-1)	Up=1,Down=0
1/2/2008	1447.16		
1/3/2008	1447.16	0	0
1/4/2008	1411.63	-35.53	0
1/7/2008	1416.18	4.55	1
1/8/2008	1390.19	-25.99	0
1/9/2008	1409.13	18.94	1
.	.	.	.
.	.	.	.
.	.	.	.
4/3/2008	1369.31	1.78	1
4/4/2008	1370.4	1.09	1
4/7/2008	1372.54	2.14	1
4/8/2008	1365.54	-7	0
4/9/2008	1354.49	-11.05	0
4/10/2008	1360.55	6.06	1
4/11/2008	1332.83	-27.72	0

155



$$\frac{0+0+1+0+1+\dots+0+0+1+0}{69} = \frac{33(1)+36(0)}{69} = 0.478$$

The average tells us the percentage of days that resulted in a positive SP500 return.

48% of the days resulted in a positive return.

156

### What will happen the next day?

•Let  $X$  denote the outcome. Then  $X$  is either 0 or 1.

• $X$  is a numerical quantity about which we are uncertain.

•**Random Variable:** We do not know what  $X$  will be, but we **do** know that it will be either 1 or 0 with certain probabilities.

### What are these probabilities?

Tough questions! 47.8% is simply a rough estimate of the actual chance that SP500 is up in a given day. It is a rough estimate because it is based only on a very recent past, which may or may not represent the TRUE process driving the SP500 movement.

157

### Example 2: Tossing a “fair” coin

Let us see a (much simpler) example where we are more comfortable assessing these probabilities

They are  $\Pr(X=1)=0.5$  and  $\Pr(X=0)=0.5$ .

The probability of a 1 is 0.5.

The probability of a 0 is 0.5.

### What does it mean?

The two possible outcomes are equally likely  
**(by the very nature of a coin).**

Over the long run, if we tossed the coin over and over again, we expect a 1 (or, equivalently, a zero) 50% of the time.

158

## Probability as the long-run frequency

How often it happens

That is, if we toss the coin  $n$  times with  $n$  **really big** and

$n_1$  is the number of 1's

$n_0$  is the number of 0's

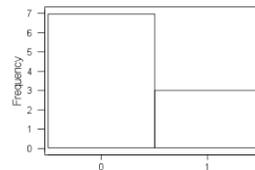
then,

$$\frac{n_1}{n} \approx .5 \quad \frac{n_0}{n} \approx .5$$

159

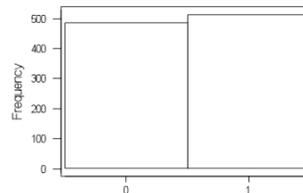
### 10 tosses

Of course, if we toss a coin 10 times we do not necessarily expect to get exactly 5 heads and 5 tails.



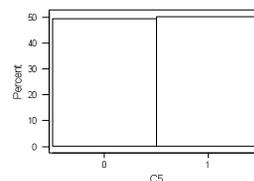
### 1,000 tosses

If we toss it 1000 times we expect the proportion to work out in the long run.



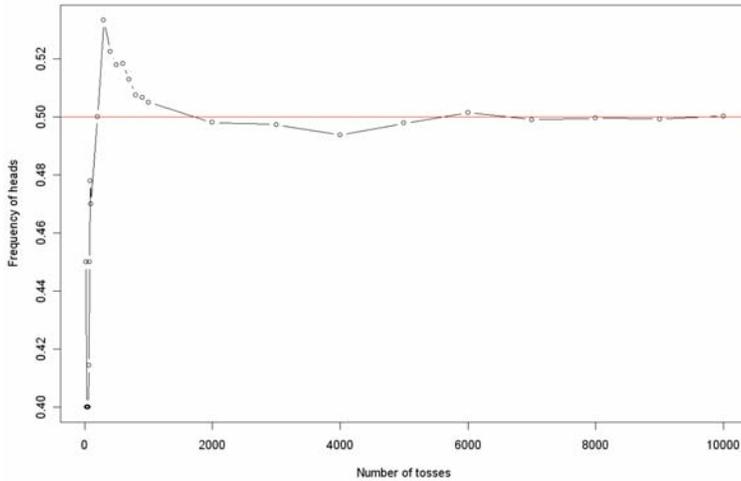
### 10,000 tosses

For "all" the tosses we expect to get 50% heads. For some, we could get something different. The closer to "all" we get, the more likely it is that the observed fraction will be close to .5.



160

The larger the “sample size” the closer the observed frequency of heads is to true probability of 50%.

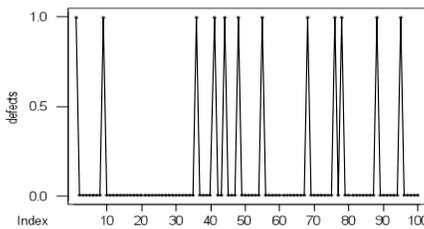


161

### Example 3: defects

Suppose we are making computer chips.

We record 1 if defective 0 if good.



Mean of defects = 0.12000

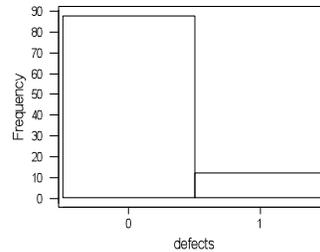
12% are defective

162

Suppose we are about to make the next chip.

**What will happen?**

We will get either a one (a defective part) or a zero (a good part) with some probabilities.



Again, we think of  $Y$  as an uncertain quantity (a random variable) with two possible outcomes, 1 and 0 (defective and good) having probabilities:

$$\Pr(Y = 1) = ? \quad \Pr(Y = 0) = 1 - ?$$

163

**Important**

Unlike the coin example, it is not obvious what to use for the probabilities here (why?).

In our sample we have 12% defectives.

Does that mean that the probability of a defective = .12?

164

Of course, **NOT!**

Later in the course we will think of the sample frequency as an **estimate** of the true probability.

So, we might estimate probabilities:

$$\Pr(Y = 1) = .12 \quad \Pr(Y = 0) = .88$$

**But, we could be wrong!**

165

**Example 4: tossing 3 coins simultaneously**

Suppose we toss three coins.

Let H:head and T: tail.

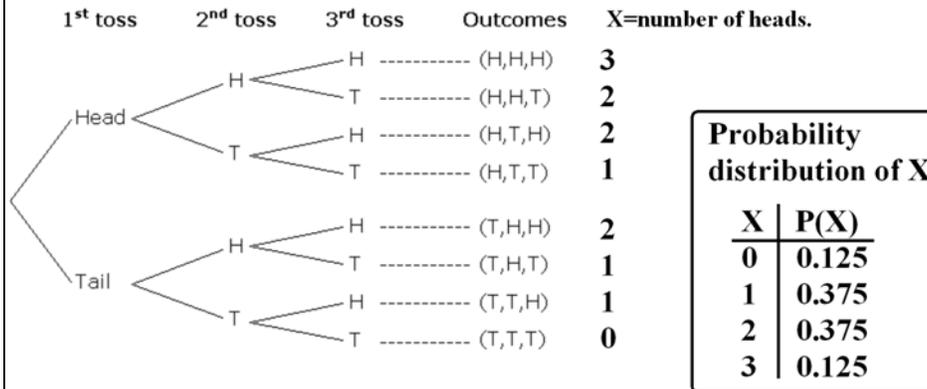
Then, the eight possible outcomes are

HHH, HHT, HTH, HTT, THH, THT, TTH, TTT.

Let X denote the number of heads (it is a random variable).  
X has three possible outcomes: 0, 1, 2 or 3.

166

## Tree diagram



167

## Definition of discrete random variable

A **discrete Random Variable** is a numerical quantity we are unsure about. We quantify our uncertainty by:

1. Listing the numbers it could turn out to be, i.e., the possible outcomes.
2. Assigning to each number a probability.  
Probabilities are numbers between 0 and 1 and sum up to 1.

The word “discrete” refers to the fact that we just have a list of outcomes. Later we will study continuous random variables where “any” outcome is possible.

168

For the random variable denoted by  $X$ , we often use  $x$  to denote a possible outcome.

**Example**

	$\Pr ( X=x )$	$x$	This table gives the <b>probability distribution</b> of the random variable $X$ .
$X:$	0 . 25	0	
	0 . 50	1	
	0 . 25	2	

Each probability tells us **how often** the corresponding outcome happens.

**Interpret.** 25% of the time we get 2 heads.

**Important:** a probability distribution is a list of probabilities, one for each outcome.

169

**Notation**

We use various notations for the probability that the random variable  $X$  takes on the value (outcome)  $x$ :

$$\Pr(X = x), \Pr(x), p_x(x), p(x)$$

These all mean the same thing.

With  $p(x)$  it must be understood from the context that you are talking about the outcome  $x$  of the random variable  $X$ .

170

**Example 5:**

Suppose we toss a die, let  $z$  denote the outcome:

$z$	$p(z)$
1	$\frac{1}{6}$
2	$\frac{1}{6}$
3	$\frac{1}{6}$
4	$\frac{1}{6}$
5	$\frac{1}{6}$
6	$\frac{1}{6}$

171

**Note:**

To get the probability that any one of a bunch of outcomes occurs we sum up their probabilities.

$$P(a < X < b) = \sum_{a < x < b} p(x)$$

**Example 5 (cont.)**

Suppose you roll a die.

Let  $X$  be the number.

$$P(2 < X < 5) = P(X=3) + P(X=4) = 1/6 + 1/6 = 2/6 = 1/3.$$

172

**Example 6:**

Suppose we toss two dice.  
Let Y denote the sum.

y	p(y)
2	$\frac{1}{36}$
3	$\frac{2}{36}$
4	$\frac{3}{36}$
5	$\frac{4}{36}$
6	$\frac{5}{36}$
7	$\frac{6}{36}$
8	$\frac{5}{36}$
9	$\frac{4}{36}$
10	$\frac{3}{36}$
11	$\frac{2}{36}$
12	$\frac{1}{36}$

What is the probability of getting more than 8?

$$\Pr(Y > 8) = \Pr(Y=9) + \Pr(Y=10) + \Pr(Y=11) + \Pr(Y=12)$$

173

**Example 7: Investing in an asset**

Suppose you are considering investing in an asset.

Let R denote the return next month.  
We think of R as a random variable.  
We do not know what the return will be (it is random) but we assume we know what the possible outcomes and probabilities are.  
In other words, we are truly modeling a future event.

r	0.05	0.10	0.15
Pr(R=r)	0.1	0.5	0.4

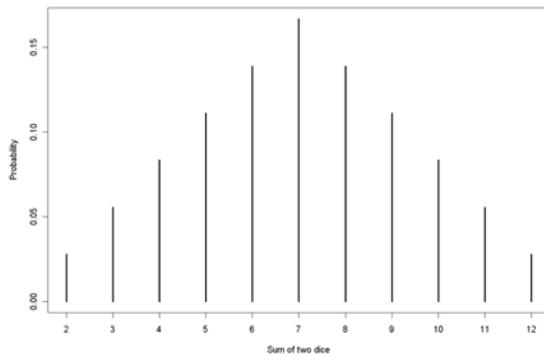
The probability that the return will be greater than 0.05 is 0.9.

174

## Graphing discrete random variables

We can use a graph to see the probability distribution of a random variable. Simply plot  $p(y)$  versus  $y$ :

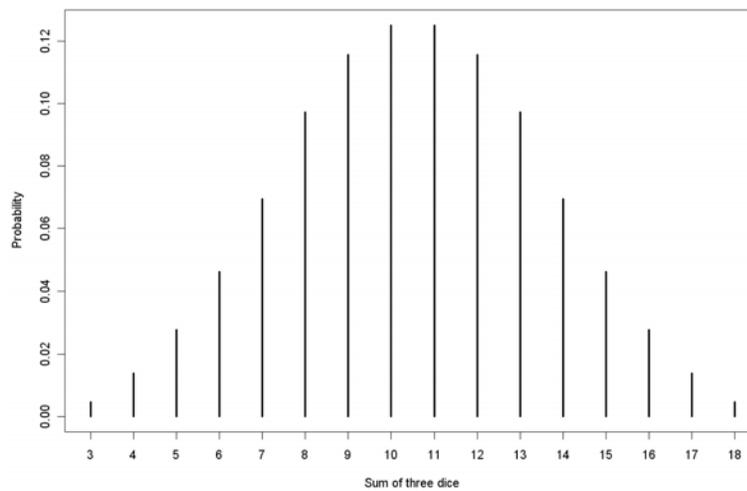
**Example 8:**  $Y =$  the sum of two dice.



$y$	$p(y)$
2	$\frac{1}{36}$
3	$\frac{2}{36}$
4	$\frac{3}{36}$
5	$\frac{4}{36}$
6	$\frac{5}{36}$
7	$\frac{6}{36}$
8	$\frac{5}{36}$
9	$\frac{4}{36}$
10	$\frac{3}{36}$
11	$\frac{2}{36}$
12	$\frac{1}{36}$

175

**Example 9:**  $Y =$  the sum of three dice



176

### The Bernoulli distribution

One of the most famous discrete random variable.

The situation where something happens or not and we want to talk about the probability of it happening is our most basic scenario.

To describe this situation we use a random variable which is 1 if something happens and 0 otherwise and probability ("it happens") =  $p$ .

Such a random variable is said to have the **Bernoulli distribution**.

Notation:  $Y \sim \text{Bernoulli}(p)$  means  $P(Y=1)=p$ ,  $P(Y=0)=1-p$

**Example 10: Toss a coin.**  $X=1$  if head, 0 else.

Then,

$$X \sim \text{Bernoulli}(0.5).$$

177

The random variable  $X$  has the Bernoulli distribution with **parameter  $p$**  (between 0 and 1) if

$$\Pr(X = 1) = p$$

$$\Pr(X = 0) = 1 - p$$

In general, we think of  $X=1$  as the thing happens and  $X=0$  as the thing does not happen.

178

## Something to think about

The word random variable refers to the outcome before it happens.

A random variable describes what we think will happen.

After we have an outcome (say, after we toss a coin), the obtained value is sometimes called a *draw* from the common distribution (it is a **data point** or an observation from **the sample**).

179

The Bernoulli distribution is named after Jakob Bernoulli, who was born in Basel, Switzerland on December 27, 1654 and lived until August 16, 1705. He is one of the eight prominent mathematicians in the Bernoulli family.

Erhard Weigel	1650	Universitat Leipzig
Gottfried <b>Leibniz</b>	1666	Universitat Altdorf
<b>Jakob Bernoulli</b>	????	
<b>Johann Bernoulli</b>	1694	
Leonhard <b>Euler</b>	1726	Universitat Basel
Joseph Louis <b>Lagrange</b>		Ecole Polytechnique
Simeon Denis <b>Poisson</b>		Ecole Polytechnique
Michel Chasles	1814	Ecole Polytechnique
Hubert Anson Newton	1850	Yale University
Eliakim Hastings Moore	1885	Yale University
Robert Lee Moore	1905	The University of Chicago
John Kline	1916	University of Pennsylvania
Donald Flanders	1927	University of Pennsylvania
Jacob Wolfowitz	1942	New York
Jack Kiefer	1952	Columbia
Lawrence Brown	1964	Cornell
James Berger	1974	Cornell
Peter Müller	1991	Purdue
<b>Hedibert Freitas Lopes</b>	<b>2000</b>	<b>Duke</b>

180

## 2. Bivariate Discrete Random Variables

Let  $X$  be the return on the nasdaq.

Let  $Y$  be the return on the djia.

We can think of both as random variables

We need probability to describe what both turn out to be

Could there be a relationship? If one "turns out big," will the other tend to be big as well?

	djia < -4	-4 <= djia < -3	-3 <= djia < -2	-2 <= djia < -1	-1 <= djia < 0	0 <= djia < 1	1 <= djia < 2	2 <= djia < 3	3 <= djia < 4	djia >=4	TOTAL
nasdaq < -4	0.7	0.2	0.3	0.6	0.6	0.1	0.0	0.0	0.0	0.0	2.5
-4 <= nasdaq < -3	0.1	0.2	0.6	1.0	0.5	0.5	0.0	0.0	0.0	0.0	3.0
-3 <= nasdaq < -2	0.0	0.2	1.6	3.1	1.5	0.5	0.0	0.0	0.0	0.0	7.1
-2 <= nasdaq < -1	0.0	0.1	0.5	4.4	5.3	1.2	0.2	0.0	0.0	0.0	11.8
-1 <= nasdaq < 0	0.0	0.0	0.1	1.4	15.3	6.7	0.5	0.0	0.0	0.0	24.1
0 <= nasdaq < 1	0.0	0.0	0.3	7.8	19.1	1.6	0.1	0.0	0.0	0.0	29.0
1 <= nasdaq < 2	0.0	0.0	0.0	0.1	1.1	6.9	4.0	0.4	0.0	0.0	12.4
2 <= nasdaq < 3	0.0	0.0	0.0	0.0	0.5	1.4	2.3	0.9	0.0	0.0	5.3
3 <= nasdaq < 4	0.0	0.0	0.0	0.0	0.2	0.3	0.8	0.6	0.4	0.1	2.4
nasdaq >=4	0.0	0.0	0.0	0.0	0.0	0.3	0.5	0.5	0.5	0.6	2.5
TOTAL	0.8	0.7	3.2	10.9	32.8	37.0	10.2	2.6	0.9	0.8	100.0

Source: Jan/2004 to Dec/2008 - <http://finance.yahoo.com>

181

We give the **bivariate** probability distribution of a **pair of random variables** by:

1. Listing out all the possible **pairs of values** that they could take on.
2. For each pair we give a probability.  
The sum of the probabilities over all pairs = 1.

182

**Example 10:**

**SP&500 and Dowjones ups and downs in 2008**

Let  $X=1$  if SP&500 is up and  $X=0$  if it is down

Let  $Y=1$  if DOW is up and  $Y=0$  if it is down

Then, the joint distribution of  $X$  and  $Y$  is given by this table

$(x,y)$	$p(x,y)$
$(0,0)$	0.478
$(0,1)$	0.072
$(1,0)$	0.044
$(1,1)$	0.406

We simply list out all possibilities for the pairs and give each one a probability.

183

**Example 11: Tossing two coins**

Let  $X$  be the result of tossing a coin ( $1=H$ ,  $0=T$ ).

Let  $Y$  be the result from a second coin toss.

Then, the joint distribution of  $X$  and  $Y$  is given by this table

$(x,y)$	$p(x,y)$
$(0,0)$	0.25
$(0,1)$	0.25
$(1,0)$	0.25
$(1,1)$	0.25

We simply list out all possibilities for the pairs and give each one a probability.

184

**Notation:**

$$p(x,y) = \Pr(X = x \text{ and } Y = y)$$

As before, we might also write

$$p_{XY}(x,y)$$

The **joint bivariate** distribution of X and Y is specified by the numbers

$$p(x,y)$$

for all possible x and y (for all possible pairs).

The distribution is discrete in that there is just a list (a finite number) of possible (x,y) pairs.

185

**Note:** An alternative way to display the probabilities is:

		X		(x,y)	p(x,y)
		0	1		
Y	0	0.478	0.044	(0,0)	0.478
	1	0.072	0.406	(0,1)	0.072
				(1,0)	0.044
				(1,1)	0.406

We have a two way table where each spot in the table corresponds to a possible (x,y) pair. At each spot we give the probability of the corresponding pair.

186

**Example 12: Investing in 2 assets**

Let X and Y be returns on two different assets.			X		
			5%	10%	15%
		5%	0.10	0.07	0.07
What does this table say about the relationship between X and Y?	Y	10%	0.03	0.30	0.03
		15%	0.05	0.05	0.30
What is the probability that they are equal?					

187

Probability means the same thing as in the univariate case

We expect to see the pair (x,y)=(10%,10%) 0.30 of the time.			X		
			5%	10%	15%
		5%	0.10	0.07	0.07
	Y	10%	0.03	0.30	0.03
		15%	0.05	0.05	0.30

188

### 3. The Marginal Distribution

The joint distribution of X and Y tells us what we expect to happen for **both of them**.

From this, we should be able to figure out what happens for **one of them**.

That is, we should be able to get

$$p_X(x) \quad \text{and} \quad p_Y(y)$$

from

$$p_{XY}(x, y)$$

189

#### Example 12 (cont.)

		X		
		5%	10%	15%
Y	5%	0.10	0.07	0.07
	10%	0.03	0.30	0.03
	15%	0.05	0.05	0.30

What is  $p_X(5\%)$  ?

$$\begin{aligned} p_X(5\%) &= p_{XY}(5\%, 5\%) + p_{XY}(5\%, 10\%) + p_{XY}(5\%, 15\%) \\ &= 0.10 + 0.03 + 0.05 = 0.18 \end{aligned}$$

190

### The marginal distributions

Given the joint distribution of X and Y defined by

$$p_{XY}(x, y)$$

the marginal (individual) distributions of X and Y are given by,

$$p_X(x) = \sum_{\text{all } y} p_{XY}(x, y)$$

$$p_Y(y) = \sum_{\text{all } x} p_{XY}(x, y)$$

191

### Example 12 (cont.)

Let us write out the marginal distributions (or, based on a common jargon, **the marginals**) using our standard two way table.

		X			$p_Y(y)$
		5%	10%	15%	
Y	5%	0.10	0.07	0.07	0.24
	10%	0.03	0.30	0.03	0.36
	15%	0.05	0.05	0.30	0.40
$p_X(x)$		0.18	0.42	0.40	1.00

192



















































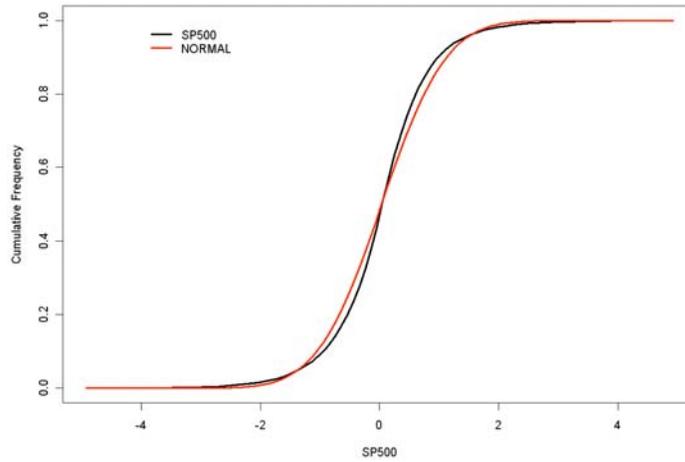




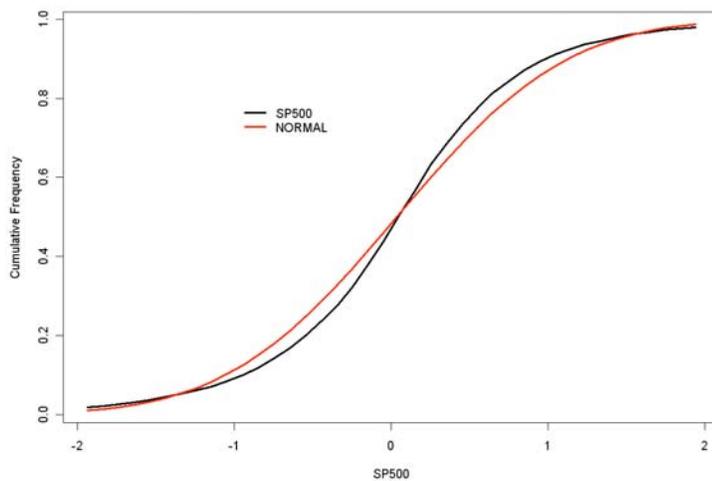




Comparing the empirical c.d.f. of S&P500 returns with the normal model with mean 0.0368 and variance 0.7291.



**A closer look between -2 and 2:**  
The normal model IS NOT a good model for the SP500 returns.











































The **correlation** between random variables (discrete or continuous) is

$$\rho_{XY} = \frac{\sigma_{XY}}{\sigma_X \sigma_Y}$$

293

$\rho$ : the basic facts

$$-1 \leq \rho \leq 1$$

If  $\rho$  is close to 1, then it means there is a line, with positive slope, such that  $(X, Y)$  is likely to fall close to it.

If  $\rho$  is close to -1, same thing, but the line has a negative slope.

294







## 0. I.I.D Draws from the Normal Distribution

We want to use the normal distribution to model data in the real world.

Surprisingly often, data **looks like** i.i.d. draws from a normal distribution.

301

**Note:** We can have i.i.d. draws from any distribution.

By writing

$$X_1, X_2, \dots, X_n \sim N(\mu, \sigma^2) \text{ i.i.d.}$$

we mean that each random variable  $X$  will be an independent draw from the same normal distribution.

We have not formally defined independence for continuous distributions, but our intuition is the same as before!

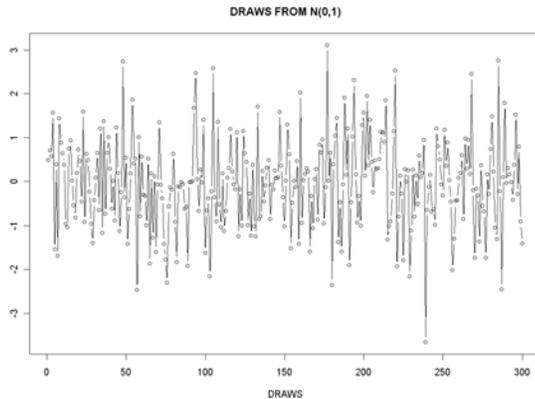
Each draw is has no effect on the others, and the same normal distribution describes what we think each variable will turn out to be.

302

### What do i.i.d. normal draws look like?

We can simulate i.i.d draws from the normal distribution.

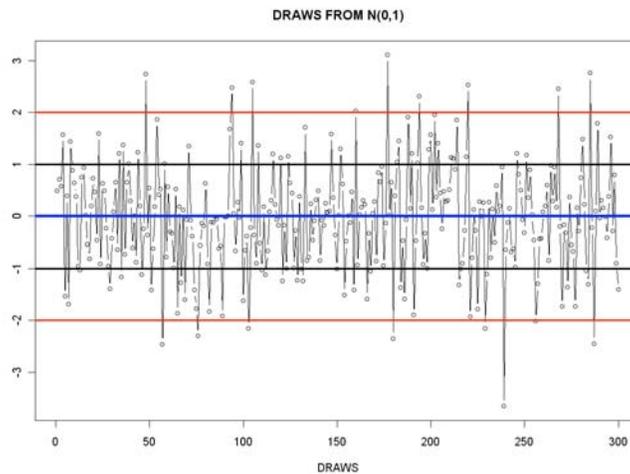
Here are 300 "draws" simulated from the standard normal distribution.



There is no pattern, they look "random"

303

Same with lines drawn in at  $\mu=0$ ,  $\pm 1s$  and  $\pm 2s$



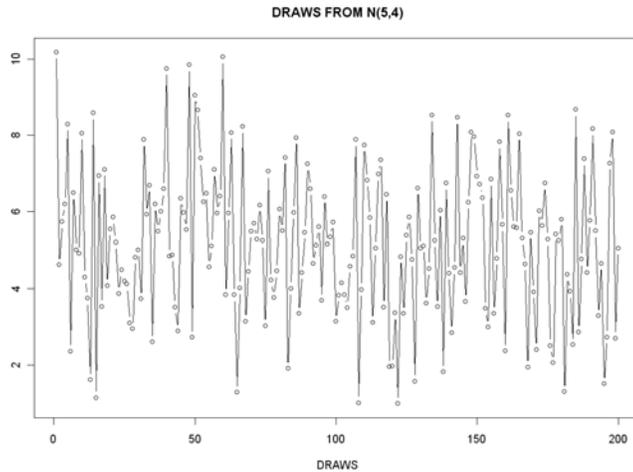
In the long run, 95% will be between +2 and -2.

*Do you remember the empirical rule?*

304

**Draws from a normal other than the standard one.**

These are 200 i.i.d. draws from  $N(5,4)$ , ie. a normal distribution with mean 5, variance 4 and, therefore, standard deviation 2.

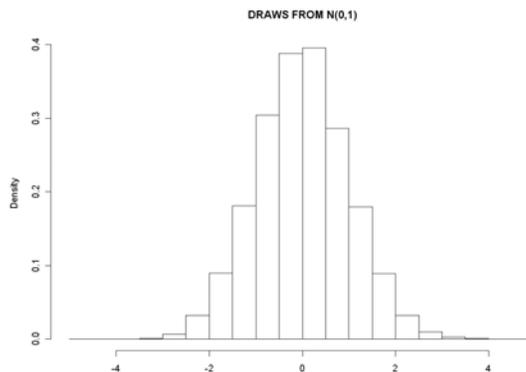


305

Here is the histogram of 5000 draws from the standard normal.

The height of each bar tells us the number of observations in the interval.

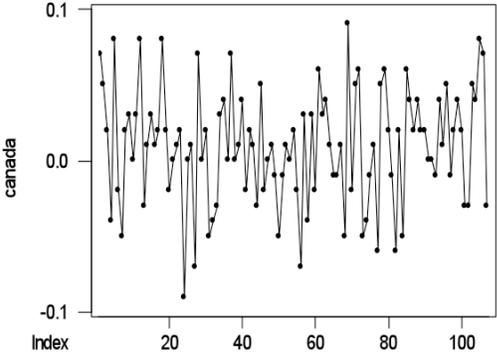
All the intervals have the same width.



306

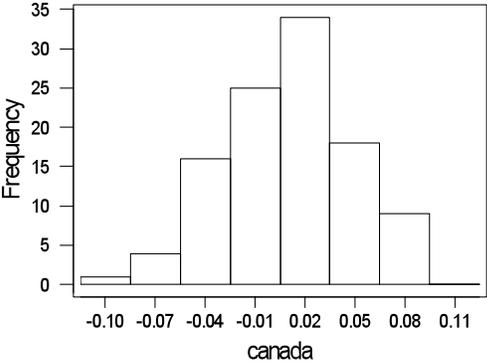


We look, once again, at the Canadian returns data.  
We have monthly returns from Feb '88 to Dec '96.



No  
apparent  
pattern!

309



Normality  
seems  
reasonable!

**Conclude: The returns look like i.i.d. normal draws!**

310



















*Before we take the sample*  $\hat{p}$  is a random variable.

We wonder how close  $\hat{p}$  will be to  $p$ .

*After we take the sample*, the resulting sample proportion  $\hat{p}$  is just a number, it is just our estimate of  $p$ .

329

#### 4. The Sampling Distribution of the Estimator

Well, we have our plan.

*What are our chances?*

After we have our sample we are either close or not.

Before we have the sample we can think about what the properties of our estimator are.

*How wrong could we be ?*

330



Don't confuse the *probability distribution* of  $\hat{p}$  with how the 1's and 0's are "distributed" in the population.

The distribution of 1's and 0's in the population is summarized by the unknown proportion  $p$ .

Notice that the probability distribution of  $\hat{p}$  when  $n=100$ , for instance, is not the same as the probability distribution of  $\hat{p}$  when  $n=1000$ .

333

We can compute the mean and variance of our estimator to summarize its properties:

$$E(\hat{p}) = E\left(\frac{Y}{n}\right) = \frac{1}{n}E(Y) = \frac{np}{n} = p$$

The estimator is **unbiased**.

Our estimate can turn out to be too big or too small, but it has no tendency to be wrong.

334

### Question

Suppose instead of asking 700 randomly chosen people, you asked 700 friends.

Would the proportion of democratic voters in that sample be an unbiased estimate of the population proportion?

335

the variance:

$$\begin{aligned}\text{Var}(\hat{p}) &= \text{Var}\left(\frac{Y}{n}\right) = \frac{1}{n^2} \text{Var}(Y) \\ &= \frac{np(1-p)}{n^2} = \frac{p(1-p)}{n}\end{aligned}$$

Not too useful by itself.

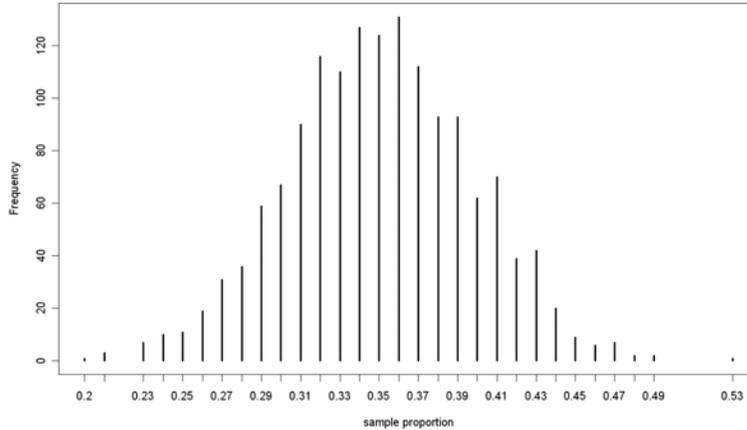
But we can combine it with the central limit theorem to get:

336





Let us suppose that now the same 1500 persons toss the coin 100 times each.



**Information accumulation:** None of the 1500 persons obtained less than 20 or more than 53 heads when tossing the coin 100 times.

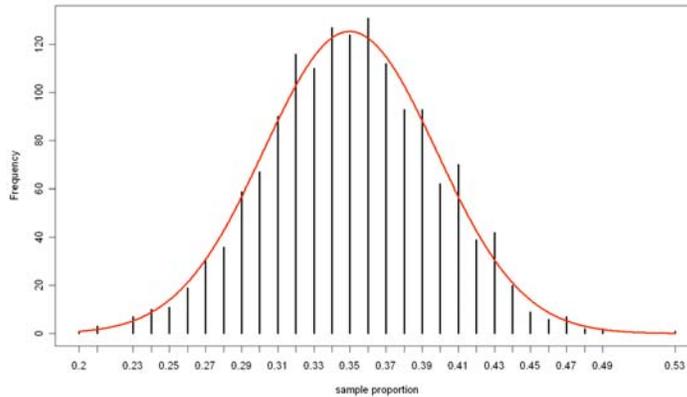
341

The  
true  
proportion  
of  
heads  
is  
**35%!**

342

Since the true proportion of heads is  $p=0.35$ , we can check how good the normal approximation is.

$$\hat{p} \sim N\left(0.35, \frac{0.35(1-0.35)}{100}\right) = N(0.35, 0.047697^2)$$



343

The **approximate** probabilities (under normality) are

$$\Pr(20 \text{ heads or less}) = 0.08308472\%$$

and

$$\Pr(53 \text{ heads or more}) = 0.008038164\%$$

The **true** probabilities are

$$\Pr(20 \text{ heads or less}) = 0.07836153\%$$

and

$$\Pr(53 \text{ heads or more}) = 0.007757356\%$$

344

**Example:**

$$\mu = p$$

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

suppose

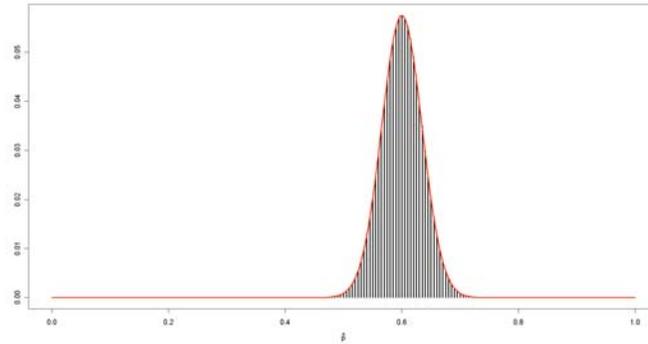
$$p = 0.6$$

$$n = 200$$

then

$$m = 0.6$$

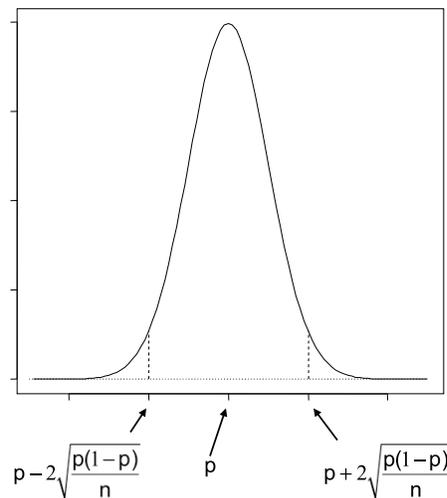
$$s = 0.0346$$



The normal curve tells us what kinds of estimates we could get if we about to take a sample of size  $n=200$  and the true population  $p = 0.6$ .

345

In general  
this is what  
we expect  
 $\hat{p}$  to be like:



Notice that the bigger  $n$  is, the better our chances are!!

346

**Example (cont.)**

Sample size:  $n=100$

True proportion:  $p=0.35$

Estimated proportion:  $\hat{p} \sim N(0.35, 0.002275)$

The approximate **95% probability interval** for  $\hat{p}$  is  $(0.35 - 2 \cdot 0.047697 ; 0.35 + 2 \cdot 0.047697) = (0.255; 0.445)$ .

**Example (cont.)**

Sample size:  $n=200$

True proportion:  $p=0.6$

Estimated proportion:  $\hat{p} \sim N(0.6, 0.0012)$

The approximate **95% probability interval** for  $\hat{p}$  is  $(0.6 - 2 \cdot 0.0346 ; 0.6 + 2 \cdot 0.0346) = (0.531; 0.669)$ .

347

**5. Confidence Interval for  $p$**

Well, that's all very well, but we still don't have an answer to our real question:

Given the data, how do we feel about  $p$ ?

The **confidence interval** is the classic solution.

It builds directly on all that we have done.

348

Confidence Interval for p:

*How different is our estimate from p?*

$$\Pr\left(p - 2\sqrt{\frac{p(1-p)}{n}} < \hat{p} < p + 2\sqrt{\frac{p(1-p)}{n}}\right) \approx 0.95$$

$$\hat{p} \approx p \pm 2\sqrt{\frac{p(1-p)}{n}}$$

$$\hat{p} \approx p \pm 2\sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$$

Since we don't know p, we just plug in the estimate for the standard deviation. *This is wrong, but we hope not too wrong!*

349

The difference between the sample and population proportions is approximately:

$$2\sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$$

350

**Example (cont.):**

Front page of Chicago trib, 1/14/2004:

"700 likely Illinois voters in the November general election were polled".

$$\hat{p} = 0.48$$

(abuse of notation!)

"48% would not like to see Bush re-elected."

"The survey has an error margin of four percentage points among general election voters.."

$$2\sqrt{\frac{0.48 * (1 - 0.48)}{700}} = 0.038$$

351

So the difference between our estimate of 0.48 and the **unknown true value** is about 0.038.

The **95% confidence interval** for the true p is

$$0.48 \pm 0.038$$

"estimate +/- error"

**Interval: (0.442 ; 0.518)**

352

Is that a big interval ?

If the election is tomorrow and we want to know the winner it is big.

If the election is three months away and last month Bush was at 70% approval then the interval is small enough to tell us things have really changed.

353

Do our estimates of  $p$  always pan out?

**Example:** Leading up to a democratic primary in Wisconsin, a poll of 600 showed Kerry with  $53\% \pm 4\%$  and Edwards with  $16\% \pm 4\%$ . The actual results a few days later were Kerry 40% and Edwards 34%.

**Example:** Results are based on telephone interviews with 1,002 national adults, aged 18 and older, conducted Feb. 9-12, 2004. For results based on the total sample of national adults, one can say with 95% confidence that the margin of sampling error is  $\pm 3$  percentage points.

**In addition to sampling error, question wording and practical difficulties in conducting surveys can introduce error or bias into the findings of public opinion polls.**

In practice, getting a random sample, or, more generally, a sample that is not biased towards some particular subset, can be tough !!

354



### Note

We use the term **standard error** to denote the estimate of a standard deviation.

**Before you get the sample**, you have an (approximate) 95% chance the true value will be in the confidence interval. After you get the data and compute the interval it is either in there or not.

We call the interval a "confidence interval" rather than a probability interval to emphasize this difference.

The "root n" in the formula precisely captures the fact that with larger samples we know more !!

357

### Question:

How much do I know about the parameter?

### Answer:

*Confidence interval small:* I know a lot.

*Confidence interval big:* I know little.

358

**Example:**

Suppose  $\hat{p} = 0.2$  and  $n=100$ .

Standard error:  $s.e. = 0.04$

suppose  $\hat{p} = 0.2$  and  $n=10,000$ . (n went up by a factor of 100)

Standard error:  $s.e.= 0.004$  (s.e. went down by 1/10)

If I want to half the s.e., I have to increase the sample size by a factor of 4!

**This is the “the tragedy of root n”.**

359

**Example:** How many observations should you collect to guarantee that, on average, the different between the true  $p$  and the estimated  $\hat{p}$ , namely  $\hat{p}$ , is less than 0.01?

What you want is to find  $n$  such that

$$2 \cdot \sqrt{p(1-p)/n} < 0.01$$

or

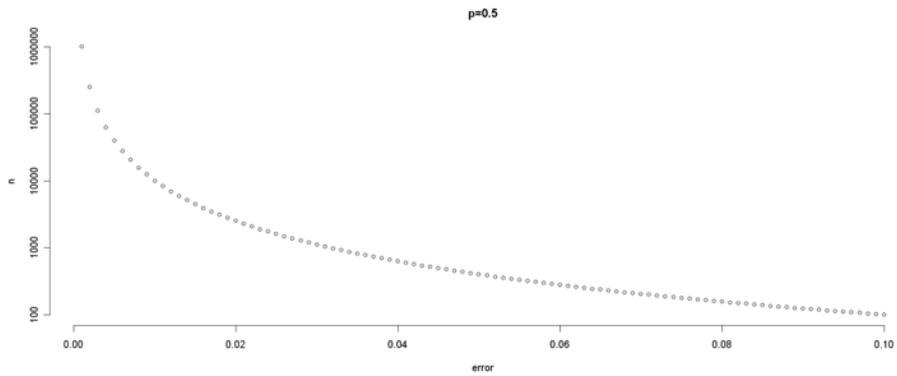
$$n > 40000 \cdot p(1-p).$$

$p$	$n$
0.1	3600
0.3	8400
0.5	<b>10000</b> $\Leftarrow$ A conservative decision maker would probably choose $n$ around 10000
0.6	9600
0.8	6400

360

If now you wanted the different between  $p$  and  $\hat{p}$  to be, on average, less than 0.04 (like in example 1)? Again, you want to find  $n$  such that  $2 \cdot \sqrt{p(1-p)/n} < 0.04$  or  $n > 2500 \cdot p(1-p)$ .

$p$	0.1	0.3	0.5	0.6	0.8
$n$	225	525	625	600	400



361

## Hypothesis testing



1. Hypothesis testing
2. P-values.
4. Confidence intervals, tests, and p-values in general.

362

## 1. Hypothesis testing for p

**Example:** Suppose we have an important manufacturing process. The manager **claims** that the defect rate is 10%.

What does this mean?

If defects are i.i.d. Bernoulli with  $p = 0.1$ , then *in the long run* we will have 10% defective.

We want to **test** the claim or **hypothesis** that  $p=0.1$ .

363

### **Experiment 1:**

Suppose we make 5 parts and 1 of the parts is defective.  
The estimated defect rate is 0.2.

**What does that tell us about  $p=0.1$ ?**

### **Experiment 2:**

Suppose we make 20 parts and 4 of the parts is defective.  
The estimated defect rate is 0.2.

**What does that tell us about  $p=0.1$ ?**

### **Experiment 3:**

Suppose we make 1000 parts and 200 parts are defective  
The estimated defect rate is 0.2.

**What does that tell us about  $p=0.1$ ?**

364



Under the hypothesis that  $p=0.1$ , the data is

Experiment 1: Highly probable => 32.80%

Experiment 2: Somewhat likely => 8.98%

Experiment 3: Very unlikely => 0.00%

### Basic Intuition (and strategy)

If the outcome of an experiment is very unlikely under the tested hypothesis, then the data provides evidence to reject the hypothesis.

367

**Clearly,  
we  
have  
to  
trust  
the  
data!**

368



**The flip side of the coin:**

**If  $p=0.1$** , what kind of value can we expect for  $\hat{p}$ ?

Recall that, under the hypothesis that  $p=0.1$ , it follows that

$$\begin{aligned}\hat{p} &\approx N\left(p, \frac{p(1-p)}{n}\right) \\ &\approx N\left(0.1, \frac{0.1(1-0.1)}{100}\right) \\ &\approx N(0.1, 0.03^2)\end{aligned}$$

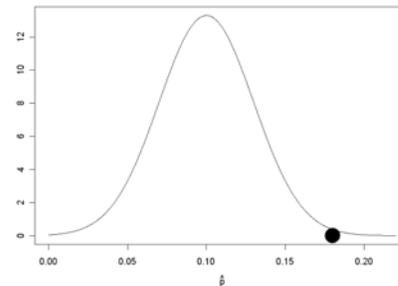
**If  $p=0.1$** , then the possible values of  $\hat{p}$  will be (approximately) normal with mean 0.1 and variance  $0.03^2$ .

371

If  $p=0.1$ , then

$$\hat{p} \approx N(0.1, 0.03^2)$$

There is a very small probability of getting a value as big as 0.18 (which is what we obtain from our specific sample).



It is very unlikely to obtain a value that big given that  $p=0.1$ .

Since we trust what we see (the estimated value from the data) we **infer** that a distribution with  $p=0.1$  is **not likely** to be the generating one.

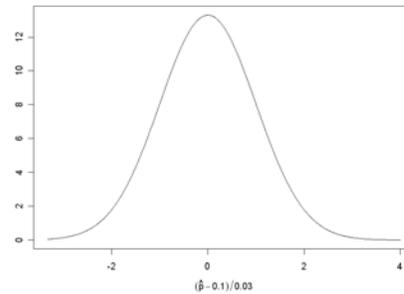
**We should probably reject the claim.**

372

It is easy to see that 0.18 is roughly 2.7 standard deviations to the right of 0.1:

$$\frac{0.18 - 0.1}{0.03} = 2.67$$

In other words,  
obtaining 0.18 from a normal distribution with mean 0.1 and variance  $0.03^2$  is the same as obtaining 2.67 from a normal distribution with mean 0 and variance 1 (the standard normal).



**2.67 is pretty unlikely.**  
**It is reasonable to reject the claim.**

373

Basic Logic:

*If the null hypothesis  $p=p^0$  is true then,*

$$\frac{\hat{p} - p^0}{\sqrt{\frac{p^0(1-p^0)}{n}}}$$

should look like a draw from the standard normal distribution !!

374



### Note (2)

If we do not reject, we **do not say** that we accept.

We say that we **fail to reject**.

This is because if we do not reject we have not proven that the null is true, we just **do not have enough evidence to reject it**.

377

### **Note (3)**

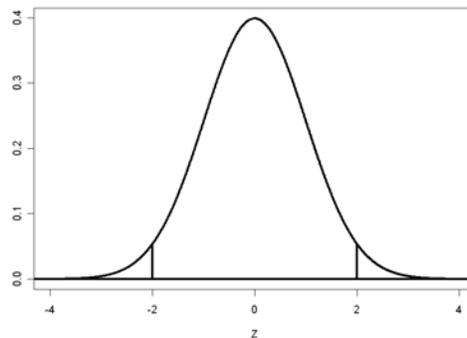
The **level** has the interpretation:

$$\Pr(\text{reject } H_0 \mid H_0 \text{ is true}) = 0.05$$

#### **Decision rule:**

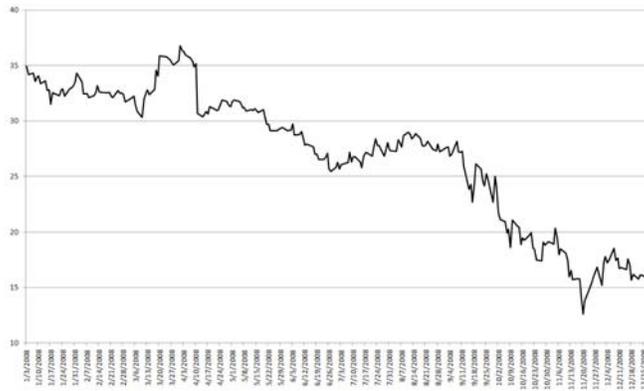
Reject  $H_0$  whenever the *test statistics* is bigger than 2 or smaller than  $-2$ .

**If  $H_0$  is true**, then 5% of the time, on average, the above decision rule will be a mistake.



378

**Example:** Let us check the claim that  $H_0$ : the daily closing price of GE in 2008 is just as likely to go up as down.



**Model:** Assume that, day to day, it is i.i.d Bernoulli ( $p$ ) whether the price of GE goes up or not. Record a 1 if it goes down and a 0 if it goes up. Then,  $p$  is the probability that the stock goes down. **We want to test  $H_0: p=0.5$ .**

379

**Data summary:**

It went down 133 days out of 252 days.

It went up 119 days out of 252 days.

The estimated  $p$  is  $133/252 = 0.52778$

The *test statistic* is

$$\frac{\hat{p} - p^0}{\sqrt{\frac{p^0(1-p^0)}{n}}} = \frac{0.52778 - 0.5}{\sqrt{\frac{0.5(1-0.5)}{252}}} = \frac{0.02778}{0.03149704} = 0.8819876$$

Since 0.8819876 is in the interval (-2,2), we DO NOT have strong evidence to reject  $H_0$ . **We fail to reject  $H_0$ .**

380

## 2. p-values

**Example:** Suppose that an i.i.d. sample of size  $n=100$  is taken from a **Bernoulli( $p$ ) model**, for some unknown value  $p$  (just like with the previous GE example). **We want to test  $H_0: p = 0.2$ .**

**Case I:** Suppose the data produces  $\hat{p} = 0.278$ .  
*Test statistic:*  $(0.278-0.2)/\sqrt{0.2*0.8/100} = 1.95$ .

**Case II:** Suppose the data produces  $\hat{p} = 0.282$ .  
*Test statistic:*  $(0.282-0.2)/\sqrt{0.2*0.8/100} = 2.05$ .

### **Not very interesting decision rule:**

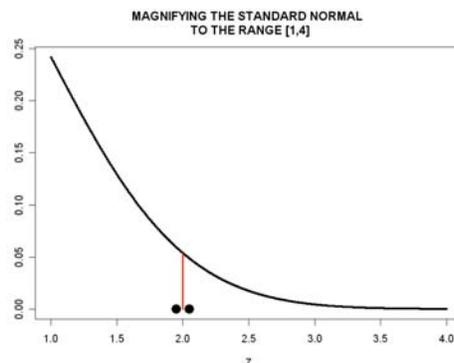
Failing to reject  $H_0$  in Case I and Rejecting  $H_0$  in Case II.  
The **evidence** is only a little different,  
but we **act** totally differently !!

381

**Remember our basic idea:** Reject if what we see is unlikely given the hypothesis.

The standard normal tells us what kind of *test statistic* we should get if the null hypothesis is true.

The farther out in the tail the *test statistics* is, the more we want to reject !!



**Rather than picking a cutoff, the p-value measures how far out in the tail the test stat is.**

382

Null hypothesis  $H_0: p = p^0$

$$\text{test statistic} = \frac{\hat{p} - p^0}{\sqrt{\frac{p^0(1-p^0)}{n}}}$$

The p-value for  $H_0$  is defined as

$$\text{p-value} = 1 - P(Z < |\text{test statistic}|)$$

where  $Z \sim N(0,1)$ .

**p-value is the probability of getting a *test statistic* as far out or farther than the one we got.**

383

**Example:**

Suppose the *test statistic* = 1.  
What is the p-value?

Suppose the *test statistic* = 2.  
What is the p-value?

Suppose the *test statistic* = 3.  
What is the p-value?

Suppose the *test statistic* = 4.  
What is the p-value?

384

Suppose the *test statistic* = 1.  
What is the p-value?

**0.3173105**

Suppose the *test statistic* = 2.  
What is the p-value?

**0.04550026**

Suppose the *test statistic* = 3.  
What is the p-value?

**0.002699796**

Suppose the *test statistic* = 4.  
What is the p-value?

**0.00006334248**

385

Here is a table of *test statistics* and p-values.

test-statistics	pvalue
0.0	1.000000
0.5	0.617075
1.0	0.317311
1.5	0.133614
2.0	0.045500
2.5	0.012419
3.0	0.002700
3.5	0.000465
4.0	0.000063
4.5	0.000007
5.0	0.000001
5.5	0.000000
6.0	0.000000
6.5	0.000000
7.0	0.000000
7.5	0.000000
8.0	0.000000
8.5	0.000000
9.0	0.000000
9.5	0.000000
10.0	0.000000

The p-value is just a  
measure of how "far out"  
the *test statistic* is.

386

**Example (cont.):**

Null hypothesis:  $p=0.1$ .

Sample size:  $n=100$  parts.

Sample proportion of defective: 0.18.

*Test statistic:*  $(0.18-0.1)/0.03 = 2.666667$ .

The p-value is 0.007660761.

Strong data evidence against the null hypothesis.

**Example (cont.):**

Null hypothesis:  $p=0.5$ .

Sample size:  $n=252$  days.

Sample proportion of downs: 0.52778.

*Test statistic:*  $(0.52778-0.5)/0.03149704 = 0.8819876$ .

The p-value is 0.3777835.

Lack of data evidence against the null hypothesis.

387

**Rejection and the p-value**

If the *test statistic* is less than 2 (in absolute value) then the p-value is greater than 0.05.

If the *test statistic* is greater than 2 (in absolute value) then the p-value is less than 0.05.

If you want to accept/reject you can just look at the p-value.

But the p-value tells you much more.

**The p-value tells you about the strength of the data evidence against a particular hypothesis.**

388

To test the null hypothesis at level 0.05,  
we reject if the p-value is less than 0.05.

To test the null hypothesis at level  $\alpha$ ,  
we reject if the p-value is less than  $\alpha$ .

***SMALL P-VALUE  
BIG TEST STATISTIC  
REJECT***

389

### 3. Confidence Intervals, Tests, and p-values in General

We have discussed confidence intervals for two **parameters**:

**NORMAL**

$m$ , the mean of i.i.d. normal observations

**BERNOULLI**

$p$ , the probability of 1, for i.i.d. Bernoulli observations

390

**More generally**, we could have a parameter which we could call  $q$ .

$q$  represents a true feature of the process or **population** under study.

Given a **sample** we obtain an estimate of  $q$ , say  $\hat{\Theta}$ .

Here are some examples:

391

Let  $\hat{\theta}$  denote an estimate of  $\theta$ .

	$\theta$	$\hat{\theta}$	
The expected value	$\mu = E(X)$	$\bar{x}$	The sample mean
The probability of success	$p$	$\hat{p}$	The ratio of successes over number of trials.
The standard deviation	$\sigma$	$s_x$	The sample standard deviation

We think of each sample quantity as an estimate of the corresponding "population" quantity (assuming our observations are i.i.d)

392

## Confidence Intervals

Because of the variation inherent in our data, we know our estimates **could be wrong**.

How wrong can we be?

The **standard error** tells us.

In general, we have (at least approximately, by the central limit theorem, for a sufficient number of observations) a 95% chance that the true value will be within 2 standard errors of the estimate.

393

In general:  $\hat{\theta} \pm 2se(\hat{\theta})$

$$m: \quad \bar{X} \pm 2se(\bar{X}) \quad se(\bar{X}) = \frac{s_x}{\sqrt{n}}$$

$$\text{Bernoulli } p: \quad \hat{p} \pm 2se(\hat{p}) \quad se(\hat{p}) = \frac{\sqrt{\hat{p}(1-\hat{p})}}{\sqrt{n}}$$

Now that we have the basic idea, we can look at confidence intervals for any quantity without necessarily knowing the details (i.e., the formula per se).

394

**Example:** We can get a confidence interval for  $s$  in the i.i.d. normal model!!

*Results for one-sample analysis for canada*

*Summary measures*

Sample size	107
Sample mean	0.009
Sample standard deviation	0.038

*Confidence interval for mean*

Confidence level	95.0%
Sample mean	0.009
Std error of mean	0.004
Degrees of freedom	106
Lower limit	0.002
Upper limit	0.016

We don't know how the confidence interval for  $s$  is computed!!

*Confidence interval for standard deviation*

Confidence level	95.0%
Sample standard deviation	0.038
Degrees of freedom	106
Lower limit	0.034
Upper limit	0.044

We're not going into the details anymore !!<sup>395</sup>

**Hypothesis Tests**

Here someone has some hypothesis about the real world.

Given the data we ask:

Could this data have arisen **if the hypothesis is true?**

**The p-value provides an answer for us.**

A small p-value means something weird happened if the hypothesis were true. We reject the hypothesis!

In particular, if the p-value  $< \alpha$ , we reject at level  $\alpha$ !

**Example:** Assuming Canadian returns are i.i.d. normal, we can test the null hypothesis that  $H_0: m = m^0$ .

**Results for one-sample analysis for canada**

**Summary measures**

Sample size	107
Sample mean	0.009
Sample standard deviation	0.038

**Test of mean=0 versus two-tailed alternative**

Hypothesized mean	0.000
Sample mean	0.009
Std error of mean	0.004
Degrees of freedom	106
t-test statistic	2.447
p-value	0.016

Here is the p-value for  $H_0: m = 0$ . We reject at level 5%.<sup>397</sup>

Again, even though we don't know the details of the test, we have some sense of how to interpret it.

But,

it only means something if we understand what hypothesis is being tested!!!

The calculation of the p-value assumes iid returns!!

If the returns are not iid, it is garbage!!!

**You don't have to understand the details of the test, you do have to understand the modeling assumptions that underlie it !!**

398

**Example:** There is a test for whether a sequence looks like it is i.i.d.!!

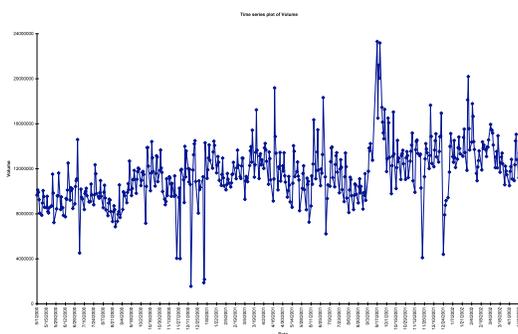
**Runs Test Results for canada**

Number of obs	107	<b>Null hypothesis:</b> Ho: data are i.i.d.
Number above cutoff	61	
Number below cutoff	46	
Number of runs	60	
E(R)	53.449	
Stdev(R)	5.045	
Z-value	1.298	The p-value is 0.2
p-value (2-tailed)	0.194	Fail to reject !!

399

**Example:** Daily volume of shares traded.

**Null hypothesis:**  
Ho: data are i.i.d.



**Runs Test Results for Volume**

Number of obs	498
Number above cutoff	213
Number below cutoff	285
Number of runs	74
E(R)	244.795
Stdev(R)	10.913
Z-value	-15.650
p-value (2-tailed)	0.000

400

## **Summary**

In general, given a model we compute a confidence interval as estimate  $\pm 2$  standard errors.

In general we can assess a hypothesis by the p-value.  
Small p-value  $\Rightarrow$  reject.

The standard errors and p-values are computed given the basic assumptions of the model. To use them properly, you must understand what these are !!

401

## **Warning: Tests are not infallible.**

Inevitably, for complex hypotheses, the tests will be more sensitive to some alternatives than others.

***The best test is the intra-ocular test:*** look at your data, it should hit you right between the eyes !!

402

## Simple Linear Regression

1. The Simple Linear Regression Model
2. Estimates and Plug-in Prediction
3. Confidence Intervals and Hypothesis Tests
4. Fits, residuals, and R-squared

403

## Book material

- What is correlation analysis and drawing the line of regression (pages 429-445 (12), 458-477 (13))
- Assumptions underlying linear regression (pages 449-450 (12), 480-482 (13))
- The standard error of estimate Confidence and prediction intervals (pages 446-448 and 451-454 (12), 477-480 and 482-486 (13))
- The relationships among the coefficient of correlation, the coefficient of determination, and the standard error of estimate (pages 457-459 (12), 489-491 (13))

404

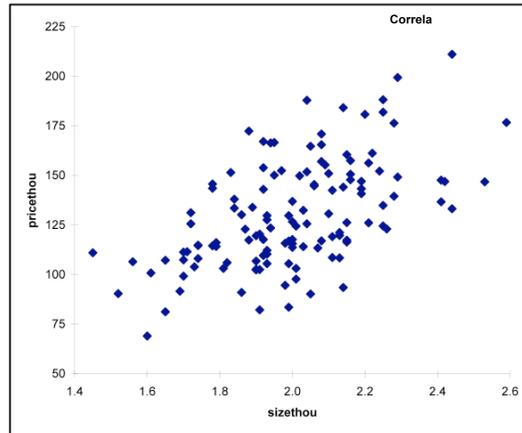
## 1. The Simple Linear Regression Model

price vs size  
from the housing  
data we looked  
at before.

Two numeric  
variables.

We want to  
build a formal  
probability model  
for the variables.

price: thousands of dollars  
size: thousands of square feet



405

Do you remember conditional probabilities?

Regression looks at the conditional distribution of  $Y$  given  $X$ .

Instead of coming up with a story for the joint  $p(x,y)$ , regression just talks about  $p(y|x)$ :

Given that I know  $x$ , what will  $y$  be?

### Example 1:

Given I know that  $x = 6'5''$  (height), what will  $y$  (weight) be?

406

### Why regression is so popular?

Lots of reasons but two would be:

- (i) Sometimes you know  $x$  and just need to predict  $y$ , as in the house price data;
- (ii) As we discussed before, the conditional distribution is an excellent way to think about the relationship between two variables.

407

### What kind of model should we use?

In the housing data, the "overall linear relationship" is striking.

Given  $x$ ,  $y$  is approximately a linear function of  $x$ .

**$y = \text{linear function of } x + \text{error}$**

408

### The Simple Linear Regression Model

$$Y_i = \alpha + \beta x_i + \varepsilon_i$$

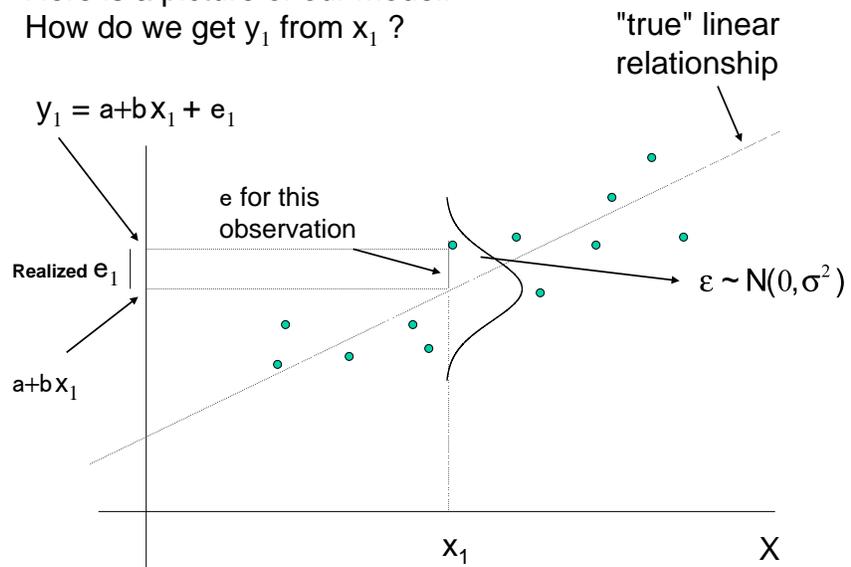
$$\varepsilon_i \sim N(0, \sigma^2) \quad \text{iid}$$

We need the normal distribution to describe what kinds of errors we might get !!!

How far  $y_i$  is from the line  $\alpha + \beta x_i$  ?

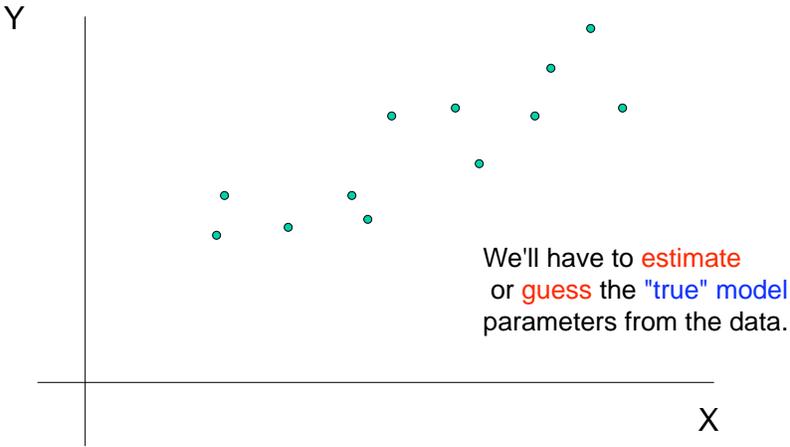
409

Here is a picture of our model.  
How do we get  $y_1$  from  $x_1$  ?



410

Of course, the model is "behind the curtain",  
all we see are the data.

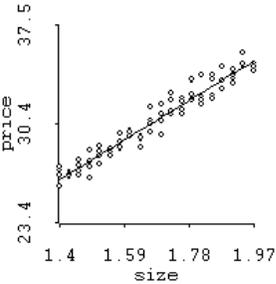
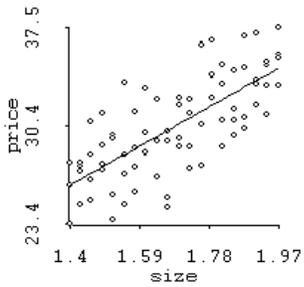


411

### The role of $s$

$s$  large

$s$  small



We need  $s$  in the model to describe how close  
the relationship is to linear, how big the errors are.

412

Another way to think about the model

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

$$\varepsilon \sim N(0, \sigma^2)$$

$\varepsilon$  independent of  $X$

is,

$$Y \mid x \sim N(\alpha + \beta x, \sigma^2)$$

since given  $x$ ,  $Y$  is just the normal  $\varepsilon$  plus the constant  $\alpha + \beta x$ .

Note that we dropped the subscripts.  
( $Y$  instead of  $Y_i$ ).  
Here we just write  $Y$  and  $x$ .  
We must assume that the model applies to all  $(x, y)$  pairs we have seen (the data) and those we wish to think about in the future.

413

Given the model, and  $x$ ,  
what do you think  $Y$  will be?

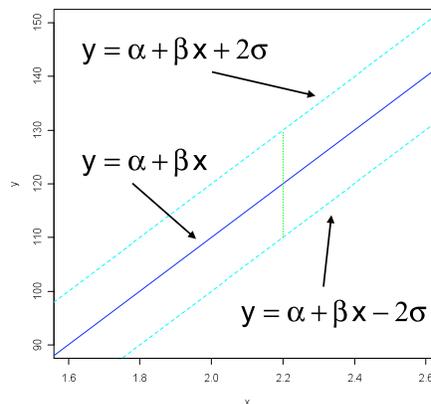
Your guess:

$$\alpha + \beta x$$

How wrong could you be?

$$\pm 2\sigma$$

$$Y = \alpha + \beta x \pm 2\sigma$$



Of course we don't know the  $\alpha$ 's and  $\sigma$ 's so  
we have to estimate them !!

414

## 2. Estimates and Plug-in Prediction

### Example 2:

Here is the output from the regression of price on size

#### Results of multiple regression for pricethou

#### Summary measures

Multiple R	0.5530
R-Square	0.3058
Adj R-Square	0.3003
SE of Est	22.4755

#### ANOVA Table

Source	df	SS	MS	F	p-value
Explained	1	28036.3627	28036.3627	55.5011	0.0000
Unexplained	126	63648.8516	505.1496		

#### Regression coefficients

	Coefficient	Std Err	t-value	p-value	Lower limit	Upper limit
Constant	-10.0911	18.9661	-0.5321	0.5956	-47.6245	27.4422
sizethou	70.2263	9.4265	7.4499	0.0000	51.5716	88.8810

**a**

**b**

$s_e$

a is our estimate of a  
b is our estimate of b  
 $s_e$  is our estimate of  $s$ .

415

Now we think of the fitted regression line as an estimate of the true line.

If the **fitted line** is

$$y = a + bx$$

then **a** is our estimate of **a** and **b** is our estimate of **b**.

"SE of Est" is our estimate of  $s$ .

We'll denote this by  $s_e$ .

We may give the formulas for the estimators later!

416

If we plug in our estimates for the true values then a "plug-in" predictive interval given  $x$  is:

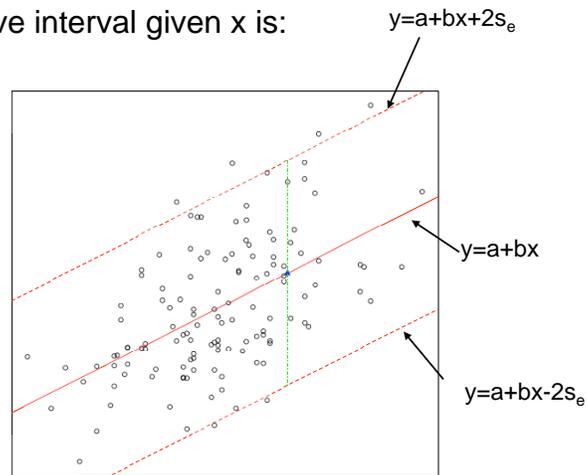
$$y = a+bx \pm 2s_e$$

Suppose we know  $x=2.2$ .

$$a+bx = 144.41$$

$$2s_e = 44.95$$

interval for  $y =$   
 $144.41 \pm 44.95$



417

summary:

<u>parameter</u>	<u>estimate</u>
$a$	$a$
$b$	$b$
$s$	$s_e$

plug-in predictive interval given a value for  $x$ :

$$a+bx \pm 2s_e$$

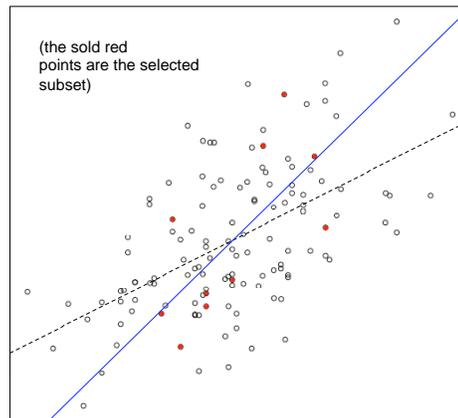
418

### 3. Confidence Intervals and Hypothesis Tests

I randomly picked  
10 of the houses  
out of our data set.

With just those  
10 observations,  
I get the solid line  
as my estimated  
line.

The dashed line  
uses all the data.



Which line would you rather use to predict?

419

With more data we expect we have a better  
chance that our estimates will be close to the  
true (or "population" values).

The "true line" is the one that "generalizes"  
to the size and price of future houses,  
not just the ones in our current data.

How big is our error?

We have standard errors and confidence intervals  
for our estimates of the true slope and intercept.

420

Let  $s_a$  denote the standard error associated with the estimate a.  
 Let  $s_b$  denote the standard error associated with the estimate b.

*Results of multiple regression for pricethou*

**Summary measures**

Multiple R	0.5530
R-Square	0.3058
Adj R-Square	0.3003
StErr of Est	22.4755

**ANOVA Table**

Source	df	SS	MS	F	p-value
Explained	1	28036.3627	28036.3627	55.5011	0.0000
Unexplained	126	63648.8516	505.1496		

**Regression coefficients**

	Coefficient	Std Err	t-value	p-value	Lower limit	Upper limit
Constant	-10.0911	18.9661	-0.5321	0.5956	-47.6245	27.4422
sizethou	70.2263	9.4265	7.4499	0.0000	51.5716	88.8810



Notation: it might make more sense to use  $se(b)$  instead of  $s_b$ , but I am following the book.

421

95% confidence interval for a:

$$a \pm t_{val} * s_a$$

$$t_{val} = \text{TINV}(.05, n - 2) \quad (\text{in excel})$$

95% confidence interval for b:

$$b \pm t_{val} * s_b$$

$$t_{val} = \text{TINV}(.05, n - 2) \quad (\text{in excel})$$

*estimate  
+/-  
2 standard errors  
!!!!!!*

If  $n$  is bigger than 30 or so,  $t_{val}$  is about 2.

422

**Example 2 (cont.)**

For the housing data the 95% confidence interval for the slope is:

$$70.23 \pm 2(9.43) = 70.23 \pm 18.86 = (51.4, 89.1)$$

big !! (what are the units?)

423

With only 10 observations  $b=135.50$  and  $s_b = 49.77$ .

Note how much bigger the standard error is than with all 128 observations!!

=tinv(.05,8)	
2.306006	

$$135.5 \pm 2.3 * (50) = (20.5, 250.5)$$

really big !!

424

Note:

It the confidence interval for slope and intercept are big the plug-in predictive interval can be misleading!!

There are ways to correct for plugging in estimates but we won't cover them.

The predictive interval just gets bigger!!

425

### **Example 2 (cont.)**

#### *Results of multiple regression for pricethou*

##### **Summary measures**

Multiple R	0.5530
R-Square	0.3058
Adj R-Square	0.3003
StErr of Est	22.4755

##### **ANOVA Table**

Source	df	SS	MS	F	p-value
Explained	1	28036.3627	28036.3627	55.5011	0.0000
Unexplained	126	63648.8516	505.1496		

##### **Regression coefficients**

	Coefficient	Std Err	t-value	p-value	Lower limit	Upper limit
Constant	-10.0911	18.9661	-0.5321	0.5956	-47.6245	27.4422
sizehou	70.2263	9.4265	7.4499	0.0000	51.5716	88.8810

$$b \pm 2 * s_b$$

426

### Hypothesis tests on coefficients:

To test the null hypothesis

$$H_0 : \alpha = \alpha^0 \text{ vs. } H_a : \alpha \neq \alpha^0$$

**We reject** at level .05 if

$$|t| = \left| \frac{a - \alpha^0}{s_a} \right| > t_{val}$$

$$t_{val} = \text{TINV}(.05, n - 2)$$

Otherwise, **we fail to reject**.

*t is the  
"t statistic"*

*reject if  
the t statistic  
is bigger  
than 2 !!*

Intuitively, we reject if estimate is more than 2 se's away from proposed value.

427

### Same for slope:

To test the null hypothesis

$$H_0 : \beta = \beta^0 \text{ vs. } H_a : \beta \neq \beta^0$$

**We reject** at level .05 if

$$|t| = \left| \frac{b - \beta^0}{s_b} \right| > t_{val}$$

$$t_{val} = \text{TINV}(.05, n - 2)$$

Otherwise, **we fail to reject**.

Intuitively, we reject if estimate is more than 2 se's away from proposed value.

428

Note:

the hypothesis:  $H_0: b = 0$

is often tested.

Why?

$$Y \mid x \sim N(\alpha + \beta x, \sigma^2)$$

If the slope = 0, then the conditional distribution of Y does not depend on x => they are independent !  
(under the assumptions of our model)

429

### Example 2 (cont.)

Stats packages automatically print out the t-statistics for testing whether the **intercept=0** and whether the **slope=0**.

Results of multiple regression for pricethou

Summary measures

Multiple R	0.5530
R-Square	0.3058
Adj R-Square	0.3003
SE of Est	22.4755

ANOVA Table

Source	df	SS	MS	F	p-value
Explained	1	28036.3627	28036.3627	55.5011	0.0000
Unexplained	126	63648.8516	505.1496		

Regression coefficients

	Coefficient	Std Err	t-value	p-value	Lower limit	Upper limit
Constant	-10.0911	18.9661	-0.5321	0.5956	-47.6245	27.4422
sizethou	70.2263	9.4265	7.4499	0.0000	51.5716	88.8810

To test  $b=0$ , the t-statistic is  $(b-0)/s_b = 70.2263/9.4265 = 7.45$

We reject the null at level 5% because the t-stat is bigger than 2 (in absolute value).

430



### **Example 3: The market model**

In finance, a popular model is to regress stock returns against returns on some market index, such as the S&P 500.

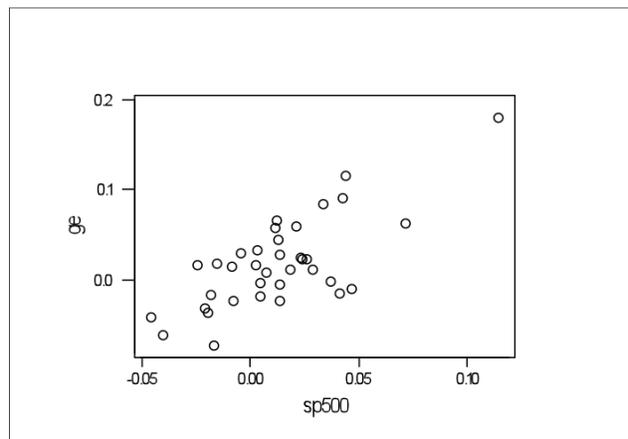
The slope of the regression line, referred to as “beta”, is a measure of how sensitive a stock is to movements in the market.

Usually, a beta less than 1 means the stock is less risky than the market, equal to 1 same risk as the market and greater than 1, riskier than the market.

433

We will examine the market model for the stock General Electric, using the S&P 500 as a proxy for the market.

Three years of monthly data give 36 observations.



434

Regression output:

The regression equation is  
 $ge = 0.00301 + 1.20 \text{ sp500}$

Predictor	Coef	Stdev	t-ratio	p
Constant	0.003013	0.006229	0.48	0.632
sp500	1.1995	0.1895	6.33	0.000

$s = 0.03454$        $R\text{-sq} = 54.1\%$        $R\text{-sq}(\text{adj}) = 52.7\%$

We can test the hypothesis that the slope is zero:  
that is, **are GE returns related to the market?**

435

The test statistic is

$$t = \frac{b - 0}{s_b} = \frac{1.2}{.1895} = 6.33$$

and

$$tval = 2.03$$

so we reject the null hypothesis at level .05. We could have looked at the p-value (which is smaller than .05) and said the same thing right away.

436

We now test the hypothesis that GE has the same risk as the market: that is, the slope equals 1.

The t statistic is:

$$t = \frac{1.1995 - 1}{.1895} = 1.055$$

Now, 1.055 is less than 2.03 so **we fail to reject.**

What is the p-value ??

437

What is the 95% confidence interval for the GE beta?

$$1.2 \pm 2(.2) = [.8, 1.6]$$

**Question:** what does this interval tell us about our level of certainty about the beta for GE?

438

#### 4. Fits, resids, and R-squared

Our model is:

$$Y = \alpha + \beta X + \varepsilon$$

We think of each  $(x_i, y_i)$  as having been generated by

$$Y_i = \alpha + \beta x_i + \varepsilon_i$$

part of y that depends on x

part of y that has nothing to do with x

439

It turns out to be useful to estimate these two parts for each observation in our sample.

For each  $(x_i, y_i)$  in the data:

$$\alpha + \beta x_i \approx a + bx_i$$

$$\varepsilon_i = y_i - (\alpha + \beta x_i) \approx y_i - (a + bx_i) = e_i$$

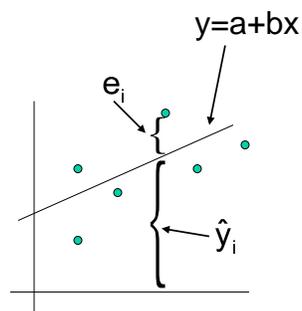
have,

$$\hat{y}_i = a + bx_i, \quad e_i = y_i - \hat{y}_i$$

$$y_i = \hat{y}_i + e_i$$

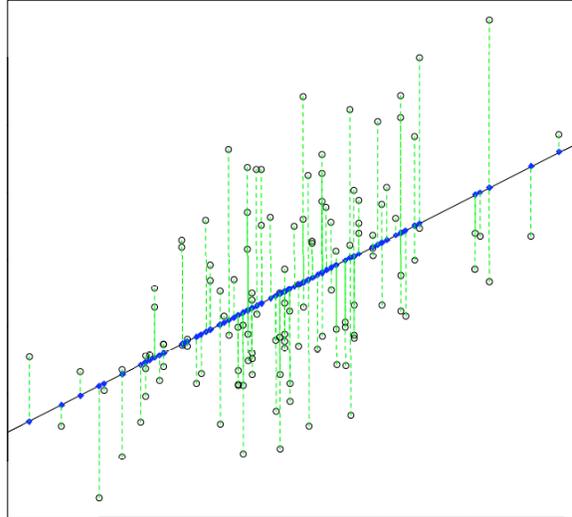
$\hat{y}_i$  : fitted value for  $i^{\text{th}}$  observation.

$e_i$  : residual for  $i^{\text{th}}$  observation.



440

Fits and  
resids  
for the  
housing data.



441

Regression chooses  
a,b so that:

$$\bar{e} = 0$$

$$\text{cor}(e,x) = 0$$

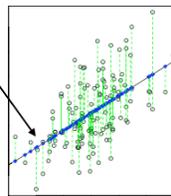
Intuition:

model:  $E(e)=0, \text{cor}(x,e)=0$   
=> make sample quantities  
exactly so:

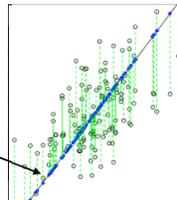
y vs x

resid off line vs x

reg  
line



slope  
too  
big



442

Note:

$$\begin{aligned}\text{cor}(e, x) = 0 &\Rightarrow \text{cor}(e, a + bx) = 0 \\ &\Rightarrow \text{cor}(e, \hat{y}) = 0\end{aligned}$$

Have:

$$\begin{aligned}y_i &= \hat{y}_i + e_i \\ \text{cor}(e, \hat{y}) &= 0 \quad \bar{e} = 0\end{aligned}$$

443

$$y_i = \hat{y}_i + e_i$$

$\Rightarrow$

$$\bar{y} = \bar{\hat{y}} + \bar{e} = \bar{\hat{y}} \quad \text{because residuals have 0 sample average}$$

$$s_y^2 = s_{\hat{y}}^2 + s_e^2 \quad \text{because residuals and fits have 0 sample correlation.}$$

$\Rightarrow$

$$\sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 + \sum_{i=1}^n e_i^2$$

total variation in y =

variation explained by x + unexplained variation

444

## R-squared

$$R^2 = \frac{\text{explained}}{\text{total}} = \frac{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

$$= 1 - \frac{\sum_{i=1}^n e_i^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

$0 \leq R^2 \leq 1$  the closer R-squared is to 1, the better the fit.

445

### Results of multiple regression for pricethou

#### Summary measures

Multiple R	0.5530
R-Square	0.3058
Adj R-Square	0.3003
StErr of Est	22.4755

$R^2$

$$\sum_{i=1}^n (\hat{y}_i - \bar{y})^2$$

#### ANOVA Table

Source	df	SS	MS	F	p-value
Explained	1	28036.3627	28036.3627	55.5011	0.0000
Unexplained	126	63648.8516	505.1496		

#### Regression coefficients

	Coefficient	Std Err	t-value	p-value	Lower limit	Upper limit
Constant	-10.0911	18.9661	-0.5321	0.5956	-47.6245	27.4422
sizethou	70.2263	9.4265	7.4499	0.0000	51.5716	88.8810

$$\sum_{i=1}^n e_i^2$$

R-squared =  $28036.3627 / (28036.3627 + 63648.8516) = 0.3057894$

446

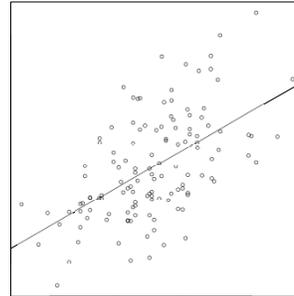
Note:

$R^2$  is also equal to the square of the correlation between  $y$  and  $x$ .

**Table of correlations**

	SqFt	Price	Fitted Values	Residuals
SqFt	1.000			
Price	0.553	1.000		
Fitted Values	1.000	0.553	1.000	
Residuals	0.000	0.833	0.000	1.000

line has intercept 0 and slope 1



$$.553^2 = 0.305809$$

$y$

Note:  $\text{cor}(y, x) = \text{cor}(y, \hat{y})$

$R^2$  = the square of the correlation between  $y$  and the fits !!

$\hat{y}$

447

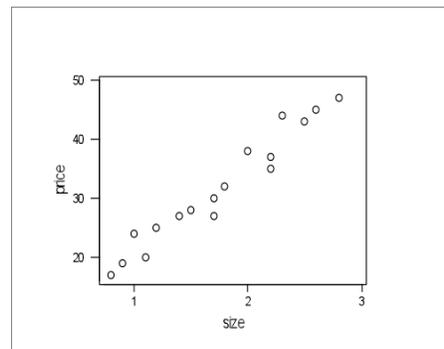
### Example 4

Housing data from a different neighborhood.

price: thousands of dollars  
size: thousands of square feet

The **correlation** is .974.

$$R^2 = .974^2 = 0.948676$$



448

## Regression output:

The regression equation is  
price = 5.76 + 14.8 size

Predictor	Coef	Stdev	t-ratio	p
Constant	5.763	1.633	3.53	0.003
size	14.8159	0.8829	16.78	0.000

$s = 2.210$        $R\text{-sq} = 94.9\%$        $R\text{-sq}(\text{adj}) = 94.6\%$

### Analysis of Variance

SOURCE	DF	SS	MS	F	p
Regression	1	1374.7	1374.7	281.58	0.000
Error	15	73.2	4.9		
Total	16	1447.9			

Fit	Stdev.Fit	95% C.I.	95% P.I.
38.358	0.669	( 36.932, 39.783)	( 33.436, 43.279)

449

For any  $x$ , the plug-in predictive interval has error

$\pm 2s_e = \pm 4.4$  thousands of dollars: **big!!!**

Even though  $R^2$  is big, we still have a lot of predictive uncertainty !!!

I think people over-emphasize  $R^2$ .  
I like  $s_e$  !!

450

## Multiple Linear Regression

1. The Multiple Linear Regression Model
2. Estimates and Plug-in Prediction
3. Confidence Intervals and Hypothesis Tests
4. Fits, residuals, R-squared, and the overall F-test
5. Categorical Explanatory Variables: Dummy Variables

451

## Book material

- What is correlation analysis and drawing the line of regression (pages 429-445 (12), 458-477 (13))
- Assumptions underlying linear regression (pages 449-450 (12), 480-482 (13))
- The standard error of estimate Confidence and prediction intervals (pages 446-448 and 451-454 (12), 477-480 and 482-486 (13))
- The relationships among the coefficient of correlation, the coefficient of determination, and the standard error of estimate (pages 457-459 (12), 489-491 (13))
- Multiple regression analysis (pages 475-483 (12), 512-519 (13))

452

## 1. The Multiple Linear Regression Model

The plug-in predictive interval for the price of a house given its size is quite large.

How can we improve this?

If we know more about a house, we should have a better idea of its price !!

453

Our data has more variables than just size and price:

The first 7 rows are:

(price and size /1000)

Home	Nbhd	Offers	SqFt	Brick	Bedrooms	Bathrooms	Price	pricethou	sizethou
1	2	2	1790	No	2	2	114300	114.3	1.79
2	2	3	2030	No	4	2	114200	114.2	2.03
3	2	1	1740	No	3	2	114800	114.8	1.74
4	2	3	1980	No	3	2	94700	94.7	1.98
5	2	3	2130	No	3	3	119800	119.8	2.13
6	1	2	1780	No	3	2	114600	114.6	1.78
7	3	3	1830	Yes	3	3	151600	151.6	1.83

Suppose we know the number of bedrooms and bathrooms a house has as well as its size, then what would our prediction for price be ?

454

### The Multiple Linear Regression Model

$$Y_i = \alpha + \beta_1 x_{i1} + \beta_2 x_{i2} + \cdots + \beta_k x_{ik} + \varepsilon_i$$

$$\varepsilon_i \sim N(0, \sigma^2) \quad \text{iid}$$

y is a linear combination of the x variables + error.

The error works exactly the same way as in simple linear reg!!  
We assume the e are independent of all the x's.

455

Another way to think about the model

$$Y | \mathbf{x} = (x_1, x_2, \dots, x_k) \sim N(\mu_x, \sigma^2)$$

$$\mu_x = \alpha + \beta_1 x_1 + \beta_2 x_2 + \cdots + \beta_k x_k$$

Y is normal with the mean depending on the x's through a linear combination.

456

If we model price as depending on size, nbed, nbath, then we have:

$$\text{Price}_i = \alpha + \beta_1 \text{nbed}_i + \beta_2 \text{nbath}_i + \beta_3 \text{size}_i + \varepsilon_i$$

Given data, we have estimates of  $a$ ,  $b_i$ , and  $s$ .

- $a$  is our estimate of  $a$ .
- $b_i$  is our estimate of  $b_i$ .
- $s_e$  is our estimate of  $s$ .

457

## 2. Estimates and Plug-in Prediction

Here is the output from the regression of price on size (SqFt), nbed (Bedrooms) and nbath (Bathrooms):

*Results of multiple regression for pricethou*

**Summary measures**

Multiple R	0.6630	
R-Square	0.4396	
Adj R-Square	0.4260	
StErr of Est	20.3565	$s_e$

**ANOVA Table**

Source	df	SS	MS	F	p-value
Explained	3	40300.9877	13433.6626	32.4180	0.0000
Unexplained	124	51384.2266	414.3889		

**Regression coefficients**

	Coefficient	Std Err	t-value	p-value	Lower limit	Upper limit
Constant	-5.6408	17.2004	-0.3279	0.7435	-39.6852	28.4035
Bedrooms	10.4599	2.9123	3.5916	0.0005	4.6956	16.2242
Bathrooms	13.5461	4.2187	3.2110	0.0017	5.1962	21.8961
sizethou	35.6427	10.6673	3.3413	0.0011	14.5292	56.7561

So, for example,  $b_2 = 13.5461$

458

Our estimated relationship is:

$$\text{Price} = -5.64 + 10.46 \cdot \text{nbed} + 13.55 \cdot \text{nbath} + 35.64 \cdot \text{size} \\ \pm 2(20.36)$$

Interpret:

With size, and nbath **held fixed**, adding one bedroom adds 10.460 thousands of dollars.

With nbed and nbath held fixed, 1 square foot increases the price \$36.

459

Suppose a house had size = 2.2, 3 bedrooms and 2 bathrooms.

What is your (estimated) idea of the price?

$$-5.64 + 10.46 \cdot 3 + 13.55 \cdot 2 + 35.64 \cdot 2.2 = 131.248$$

$$2s_e = 40.72$$

$$131.248 \pm 40.72$$

This is our multiple regression plug-in predictive interval.

The error is still estimated to be  $\pm 2s_e$  !

460

Note:

When we regressed price on size the coefficient was about 70.

Now the coefficient for size is about 36.

Without nbath and nbed in the regression, an increase in size can be associated with an increase in nbath and nbed *in the background*.

If all I know is that one house is a lot bigger than another I might expect the bigger house to have more beds and baths!

With nbath and nbed held fixed, the effect of size is smaller.

461

Note:

With just size, our predictive +/- was

$$2 * 22.467 = 44.934$$

With nbath and nbed added to the model the +/- is

$$2 * 20.36 = 40.72$$

The additional information makes our prediction more precise (but not a whole lot in the case, we still need some "better x's").

462

### 3. Confidence Intervals and Hypothesis Tests

95% confidence interval for a:

$$a \pm tval * s_a$$

tval = TINV(.05,n - k - 1) (in excel)

95% confidence interval for  $b_i$ :

$$b_i \pm tval * s_{b_i}$$

tval = TINV(.05,n - k - 1) (in excel)

*estimate*  
+/-  
2 *standard errors*

!!!!!!

(recall the k is the number of x's)

463

#### Results of multiple regression for *pricethou*

##### Summary measures

Multiple R	0.6630
R-Square	0.4396
Adj R-Square	0.4260
StErr of Est	20.3565

##### ANOVA Table

Source	df	SS	MS	F	p-value
Explained	3	40300.9877	13433.6626	32.4180	0.0000
Unexplained	124	51384.2266	414.3889		

##### Regression coefficients

	Coefficient	Std Err	t-value	p-value	Lower limit	Upper limit
Constant	-5.6408	17.2004	-0.3279	0.7435	-39.6852	28.4035
Bedrooms	10.4599	2.9123	3.5916	0.0005	4.6956	16.2242
Bathrooms	13.5461	4.2187	3.2110	0.0017	5.1962	21.8961
size <sub>thou</sub>	35.6427	10.6673	3.3413	0.0011	14.5292	56.7561

eg  $s_{b_2} = 4.22$

the interval for  $b_2$  is 13.57 +/- 2(4.22)

StatPro prints out all the confidence intervals.

464

Hypothesis tests on coefficients:

To test the null hypothesis

$$H_0 : \alpha = \alpha^0 \text{ vs. } H_a : \alpha \neq \alpha^0$$

**We reject** at level .05 if

$$|t| > t_{val} \text{ where, } t = \frac{a - \alpha^0}{s_a}$$

$$t_{val} = \text{TINV}(.05, n - k - 1)$$

Otherwise, **we fail to reject.**

*t is the  
"t statistic"*

*reject if  
the t statistic  
is bigger  
than 2 !!*

Intuitively, we reject if estimate is more than 2 se's away from proposed value.

465

Same for slope:

To test the null hypothesis

$$H_0 : \beta_i = \beta_i^0 \text{ vs. } H_a : \beta_i \neq \beta_i^0$$

**We reject** at level .05 if

$$|t| > t_{val} \text{ where, } t = \frac{b_i - \beta_i^0}{s_{b_i}}$$

$$t_{val} = \text{TINV}(.05, n - k - 1)$$

Otherwise, **we fail to reject.**

Intuitively, we reject if estimate is more than 2 se's away from proposed value.

466

## Example

Packages automatically print out the t-statistics for testing whether the intercept=0 and whether each slope=0 as well as the associated p-values.

Results of multiple regression for pricethou

### Summary measures

Multiple R	0.6630
R-Square	0.4396
Adj R-Square	0.4260
StErr of Est	20.3565

### ANOVA Table

Source	df	SS	MS	F	p-value
Explained	3	40300.9877	13433.6626	32.4180	0.0000
Unexplained	124	51384.2266	414.3869		

### Regression coefficients

	Coefficient	Std Err	t-value	p-value	Lower limit	Upper limit
Constant	-5.6408	17.2004	-0.3279	0.7435	-39.6852	28.4035
Bedrooms	10.4599	2.9123	3.5916	0.0005	4.6956	16.2242
Bathrooms	13.5461	4.2187	3.2110	0.0017	5.1962	21.8961
sizethou	35.6427	10.6673	3.3413	0.0011	14.5292	56.7561

eg.  $\frac{b_3 - 0}{s_{b_3}} = 35.64/10.67=3.34 \Rightarrow \text{reject}$

467

## 4. Fits, resids, and R-squared

In multiple regression the fit is:

$$\hat{y}_i = a + b_1 x_{i1} + b_2 x_{i2} + \dots + b_k x_{ik}$$

"the part of y related to the x's "

as before, the residual is the part left over:

$$e_i = y_i - \hat{y}_i$$

468

In multiple regression, the residuals have sample mean 0 and are uncorrelated with each of the x's and the fitted values:

Table of correlations

	SqFt	Bedrooms	Bathrooms	Price	Fitted Values	Residuals
SqFt	1.000					
Bedrooms	0.484	1.000				
Bathrooms	0.523	0.415	1.000			
Price	0.553	0.526	0.523	1.000		
Fitted Values	0.834	0.793	0.789	0.663	1.000	
Residuals	0.000	0.000	0.000	0.749	0.000	1.000

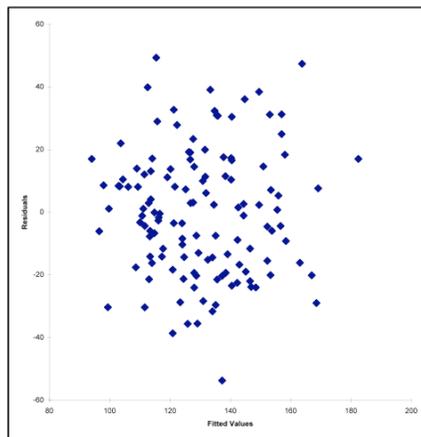
$$y_i = \hat{y}_i + e_i$$

estimated x part of y

estimated part of y  
that has nothing to do with x's

469

This is the plot of the residuals from the multiple regression of price on size, nbath, nbed vs the fitted values. We see the 0 correlation.



The correlation is also 0, for each of the x's.

470

$$y_i = \hat{y}_i + e_i$$

$$\text{cor}(\hat{y}, e) = 0, \quad \bar{e} = 0$$

So, just as with one x we have:

$$\sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 + \sum_{i=1}^n e_i^2$$

total variation in y =

variation explained by x + unexplained variation

471

### R-squared

$$R^2 = \frac{\text{explained}}{\text{total}} = \frac{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

$$= 1 - \frac{\sum_{i=1}^n e_i^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

$0 \leq R^2 \leq 1$  the closer R-squared is to 1, the better the fit.

472

In our housing example:

*Results of multiple regression for pricethou*

**Summary measures**

Multiple R	0.6630
R-Square	0.4396
Adj R-Square	0.4260
SErr of Est	20.3565

**ANOVA Table**

Source	df	SS	MS	F	p-value
Explained	3	40300.9877	13433.6626	32.4180	0.0000
Unexplained	124	51384.2266	414.3689		

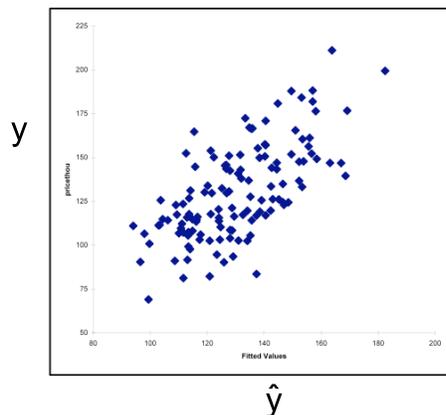
**Regression coefficients**

	Coefficient	Std Err	t-value	p-value	Lower limit	Upper limit
Constant	-5.6408	17.2004	-0.3279	0.7435	-39.6852	28.4035
Bedrooms	10.4599	2.9123	3.5916	0.0005	4.6956	16.2242
Bathrooms	13.5461	4.2187	3.2110	0.0017	5.1962	21.8961
sizethou	35.6427	10.6673	3.3413	0.0011	14.5292	56.7561

$$R^2 = \frac{40301}{40301+51384} = .439$$

473

$R^2$  is also the square of the correlation between the fitted values and  $y$ :



Regression finds the linear combination of the  $x$ 's which is most correlated with  $y$ .

$$\text{cor}(\hat{y}, y) = .663$$

$$.663^2 = 0.439569$$

(with just size, the correlation between fits and  $y$  was .553)

474

The "Multiple R" is the correlation between y and the fits

Results of multiple regression for pricethou

$$\text{cor}(\hat{y}, y) = .663$$

Summary measures

Multiple R	0.6630
R-Square	0.4396
Adj R-Square	0.4260
StErr of Est	20.3565

ANOVA Table

Source	df	SS	MS	F	p-value
Explained	3	40300.9877	13433.6626	32.4180	0.0000
Unexplained	124	51384.2266	414.3889		

Regression coefficients

	Coefficient	Std Err	t-value	p-value	Lower limit	Upper limit
Constant	-5.6408	17.2004	-0.3279	0.7435	-39.6852	28.4035
Bedrooms	10.4599	2.9123	3.5916	0.0005	4.6956	16.2242
Bathrooms	13.5461	4.2187	3.2110	0.0017	5.1962	21.8961
sizethou	35.6427	10.6673	3.3413	0.0011	14.5292	56.7561

475

The overall F-test

The p-value beside "F" if testing the null hypothesis:

$$H_0 : \beta_1 = \beta_2 = \dots = \beta_k = 0 \text{ (all the slopes are 0)}$$

Results of multiple regression for pricethou

Summary measures

Multiple R	0.6630
R-Square	0.4396
Adj R-Square	0.4260
StErr of Est	20.3565

ANOVA Table

Source	df	SS	MS	F	p-value
Explained	3	40300.9877	13433.6626	32.4180	0.0000
Unexplained	124	51384.2266	414.3889		

Regression coefficients

	Coefficient	Std Err	t-value	p-value	Lower limit	Upper limit
Constant	-5.6408	17.2004	-0.3279	0.7435	-39.6852	28.4035
Bedrooms	10.4599	2.9123	3.5916	0.0005	4.6956	16.2242
Bathrooms	13.5461	4.2187	3.2110	0.0017	5.1962	21.8961
sizethou	35.6427	10.6673	3.3413	0.0011	14.5292	56.7561

We reject the null, at least some of the slopes are not 0.

476

## 5. Categorical Explanatory Variables: Dummy Variables

Here, again, is the first 7 rows of our housing data:

Home	Nbhd	Offers	SqFt	Brick	Bedrooms	Bathrooms	Price	pricethou	sizethou
1	2	2	1790	No	2	2	114300	114.3	1.79
2	2	3	2030	No	4	2	114200	114.2	2.03
3	2	1	1740	No	3	2	114800	114.8	1.74
4	2	3	1980	No	3	2	94700	94.7	1.98
5	2	3	2130	No	3	3	119800	119.8	2.13
6	1	2	1780	No	3	2	114600	114.6	1.78
7	3	3	1830	Yes	3	3	151600	151.6	1.83

Does whether a house is brick or not affect the price of the house?

This is a categorical variable.

How can we use multiple regression with categorical x's ???!

What about the neighborhood? (location, location, location !!)

477

### Adding a Binary Categorical x

To add "brick" as an explanatory variable in our regression we create the dummy variable which is 1 if the house is brick and 0 otherwise:

Home	Nbhd	Offers	SqFt	Brick	Bedrooms	Bathrooms	Price	sizethou	pricethou	brickdum
1	2	2	1790	No	2	2	114300	1.79	114.3	0
2	2	3	2030	No	4	2	114200	2.03	114.2	0
3	2	1	1740	No	3	2	114800	1.74	114.8	0
4	2	3	1980	No	3	2	94700	1.98	94.7	0
5	2	3	2130	No	3	3	119800	2.13	119.8	0
6	1	2	1780	No	3	2	114600	1.78	114.6	0
7	3	3	1830	Yes	3	3	151600	1.83	151.6	1
8	3	2	2160	No	4	2	150700	2.16	150.7	0
9	2	3	2110	No	4	2	119200	2.11	119.2	0
10	2	3	1730	No	3	3	104000	1.73	104	0
11	2	3	2030	Yes	3	2	132500	2.03	132.5	1
12	2	2	1870	Yes	2	2	123000	1.87	123	1
13	1	4	1910	No	3	2	102600	1.91	102.6	0
14	1	5	2150	Yes	3	3	126300	2.15	126.3	1

the  
"brick dummy"

•  
•  
•

478

Note:

I created the dummy by using the excel formula:

=IF(Brick="Yes",1,0)

but we'll see that StatPro has a nice utility for creating dummies.

479

As a simple first example, let's regress price on size and brick.

Here is our model:

$$\text{Price}_i = \alpha + \beta_1 \text{size}_i + \beta_2 \text{brickdum}_i + \varepsilon_i$$

How do you interpret  $b_2$  ?

480

What is the expected price of a brick house given the size?

$$E(\text{Price} \mid \text{size} = s, \text{brick}) = \alpha + \beta_1 s + \beta_2$$

What is the expected price of a non-brick house given the size?

$$E(\text{Price} \mid \text{size} = s, \text{nonbrick}) = \alpha + \beta_1 s$$

$\beta_2$  is the expected difference in price between a brick and non-brick house.

481

Note:

You could also create a dummy which was 1 if a house was non brick and 0 if brick.

That would be fine, but the meaning of  $\beta_2$  which change.

You can't put both dummies in though because given one, the information in the other is redundant.

482

Let's try it !!

*Results of multiple regression for pricethou*

**Summary measures**

Multiple R	0.6884
R-Square	0.4739
Adj R-Square	0.4655
SErr of Est	19.6441

**ANOVA Table**

Source	df	SS	MS	F	p-value
Explained	2	43448.6791	21724.3396	56.2964	0.0000
Unexplained	125	48236.5352	385.8923		

**Regression coefficients**

	Coefficient	Std Err	t-value	p-value	Lower limit	Upper limit
Constant	-9.4443	16.5771	-0.5697	0.5699	-42.2525	23.3639
sizehou	66.0584	8.2653	7.9922	0.0000	49.7003	82.4165
brickdum	23.4451	3.7098	6.3198	0.0000	16.1029	30.7873

$\pm 2se = 39.3$ , this is the best we've done !

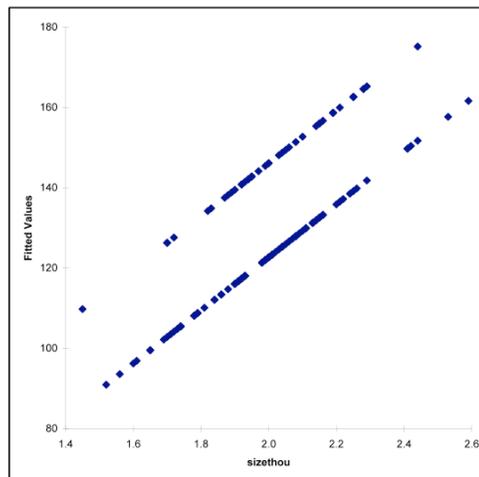
what is the brick effect:

$$23.4 \pm 2(3.7) = 23.4 \pm 7.4$$

483

We can see the effect of the dummy by plotting the fitted values vs size.

The upper line is for the brick houses and the lower line is for the non-brick houses.



484

We can interpret  $b_2$  as a shift in the intercept.

Notice that our model assumes that the price difference between a brick and non-brick house does not depend on the size!

The two variables do not "interact".

Sometimes we expect variables to interact.

485

Now let's add brick to the regression of price on size, nbath, and nbed:

*Results of multiple regression for pricethou*

**Summary measures**

Multiple R	0.7634
R-Square	0.5828
Adj R-Square	0.5692
StErr of Est	17.6345

**ANOVA Table**

Source	df	SS	MS	F	p-value
Explained	4	53435.3823	13358.8456	42.9580	0.0000
Unexplained	123	38249.8320	310.9742		

**Regression coefficients**

	Coefficient	Std Err	t-value	p-value	Lower limit	Upper limit
Constant	-5.2794	14.9004	-0.3543	0.7237	-34.7739	24.2151
Bedrooms	10.8731	2.5237	4.3084	0.0000	5.8776	15.8686
Bathrooms	9.8184	3.6993	2.6541	0.0090	2.4959	17.1409
sizethou	35.8006	9.2409	3.8742	0.0002	17.5088	54.0923
brickdum	21.9091	3.3712	6.4989	0.0000	15.2361	28.5821

$$\pm 2se = 35.2$$

Adding brick seems to be a good idea !!

486

I created one dummy for each the neighborhoods.

Home	Nbhd	Offers	SqFt	Brick	Bedrooms	Bathrooms	Price	Nbhd_1	Nbhd_2	Nbhd_3
1	2	2	1790	No	2	2	114300	0	1	0
2	2	3	2030	No	4	2	114200	0	1	0
3	2	1	1740	No	3	2	114800	0	1	0
4	2	3	1980	No	3	2	94700	0	1	0
5	2	3	2130	No	3	3	119800	0	1	0
6	1	2	1780	No	3	2	114600	1	0	0
7	3	3	1830	Yes	3	3	151600	0	0	1
8	3	2	2160	No	4	2	150700	0	0	1

.  
.
  
.

eg. Nbhd\_1 indicates if the house is in neighborhood 1 or not

487

Now we add any two of the three dummies.  
Given any two, the information in the third is redundant.

Let's first do price on size and neighborhood:

$$\text{Price}_i = \alpha + \beta_1 \text{size}_i + \beta_2 N1_i + \beta_3 N2_i + \varepsilon_i$$

where now I've use N1 to denote the dummy for neighborhood 1 and same for 2.

488

$$\text{Price}_i = \alpha + \beta_1 \text{size}_i + \beta_2 \text{N1}_i + \beta_3 \text{N2}_i + \varepsilon_i$$

$$E(\text{Price} \mid \text{size} = s, \text{neighborhood3}) = \alpha + \beta_1 s$$

$$E(\text{Price} \mid \text{size} = s, \text{neighborhood2}) = \alpha + \beta_1 s + \beta_3$$

$$E(\text{Price} \mid \text{size} = s, \text{neighborhood1}) = \alpha + \beta_1 s + \beta_2$$

$\beta_3$ : difference between hood 2 and hood 3

$\beta_2$ : difference between hood 1 and hood 3

The neighborhood corresponding to the dummy we leave out becomes the "base case" we compare to.

489

Let's try it!

*Results of multiple regression for pricethou*

**Summary measures**

Multiple R	0.8277
R-Square	0.6851
Adj R-Square	0.6774
SErr of Est	15.2601

**ANOVA Table**

Source	df	SS	MS	F	p-value
Explained	3	62809.1498	20936.3833	89.9053	0.0000
Unexplained	124	28876.0645	232.8715		

**Regression coefficients**

	Coefficient	Std Err	t-value	p-value	Lower limit	Upper limit
Constant	62.7765	14.2477	4.4061	0.0000	34.5763	90.9766
Nbhd_1	-41.5353	3.5337	-11.7542	0.0000	-48.5294	-34.5412
Nbhd_2	-30.9666	3.3688	-9.1922	0.0000	-37.6344	-24.2988
sizethou	46.3859	6.7459	6.8762	0.0000	33.0340	59.7379

+/- 2se = 30.52 !!!

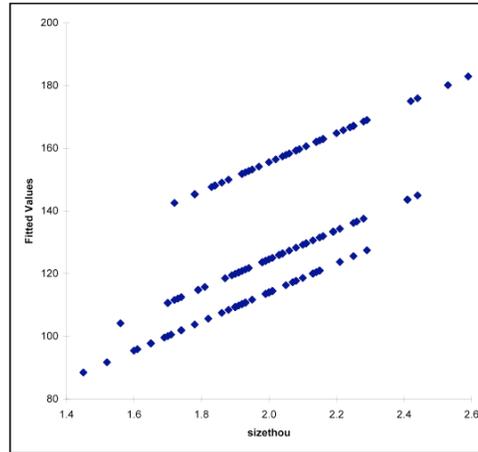
490

Here is  
fits vs size.

Which line  
corresponds  
to which  
neighborhood ?

Where do you  
want to live ?

Again we  
are assuming  
that size and  
neighborhood do not  
interact.



491

ok, let's try price on size, nbed, nbath, brick, and  
neighborhood.

**Results of multiple regression for pricethou**

**Summary measures**

Multiple R	0.8972
R-Square	0.8050
Adj R-Square	0.7954
SErr of Est	12.1547

**ANOVA Table**

Source	df	SS	MS	F	p-value
Explained	6	73809.1440	12301.5240	83.2669	0.0000
Unexplained	121	17876.0703	147.7361		

**Regression coefficients**

Coefficient	Std Err	t-value	p-value	Lower limit	Upper limit
Constant	52.0032	11.5181	4.5149	0.0000	29.2000 74.8063
Bedrooms	1.9022	1.9023	0.9999	0.3193	-1.8639 5.6682
Bathrooms	6.8269	2.5628	2.6638	0.0088	1.7532 11.9007
Nbhd_1	-34.0837	3.1690	-10.7554	0.0000	-40.3576 -27.8099
Nbhd_2	-29.2180	2.8637	-10.2030	0.0000	-34.8874 -23.5486
size2hou	35.9304	6.4044	5.6102	0.0000	23.2511 48.6097
Brick_Yes	18.5078	2.3963	7.7235	0.0000	13.7637 23.2519

$\pm 2s_e = 24 !!$

492

## Maybe we don't need bedrooms:

### Results of multiple regression for pricethou

#### Summary measures

Multiple R	0.8963
R-Square	0.8034
Adj R-Square	0.7954
SErr of Est	12.1547

#### ANOVA Table

Source	df	SS	MS	F	p-value
Explained	5	73661.4233	14732.2847	99.7203	0.0000
Unexplained	122	18023.7910	147.7360		

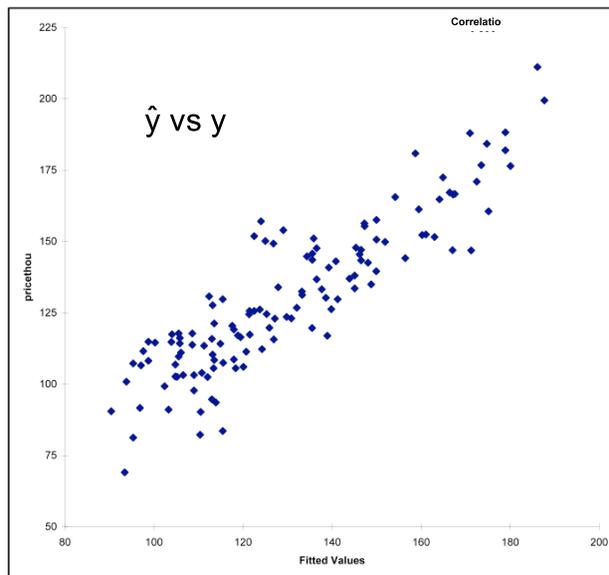
#### Regression coefficients

	Coefficient	Std Err	t-value	p-value	Lower limit	Upper limit
Constant	53.6295	11.4027	4.7032	0.0000	31.0567	76.2023
Bathrooms	7.2304	2.5308	2.8569	0.0050	2.2204	12.2405
Nbhd_1	-35.3137	2.9205	-12.0916	0.0000	-41.0952	-29.5322
Nbhd_2	-30.1452	2.7094	-11.1262	0.0000	-35.5087	-24.7817
sizethou	37.9050	6.0924	6.2217	0.0000	25.8445	49.9656
Brick_Yes	18.3121	2.3883	7.6674	0.0000	13.5843	23.0400

Dropping bedrooms did not increase  $s_e$  or decrease R-Square so no need to bother with it.

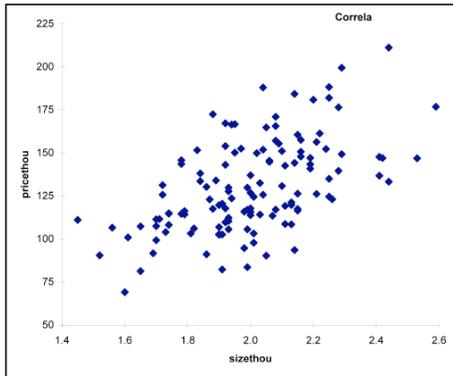
493

Regression finds a linear combination of the variables that is like  $y$ .

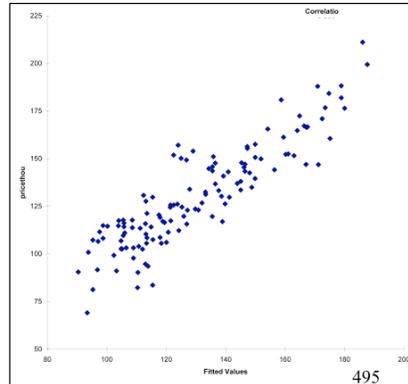


494

price vs size:

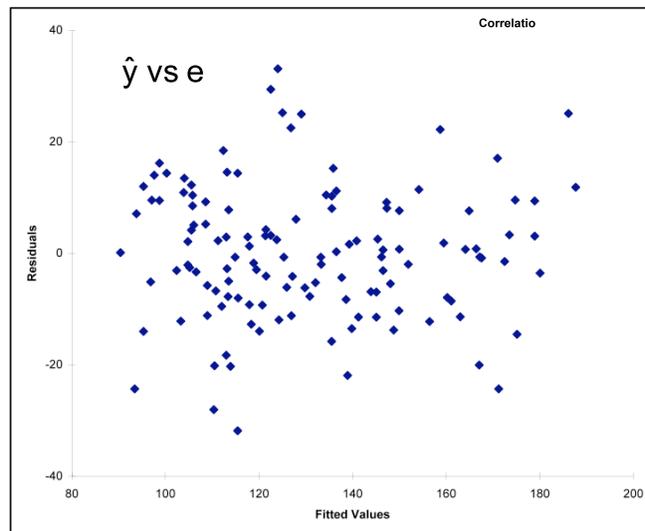


price vs combination of size, nbath, brick, nbhd



495

The residuals are the part of  $y$  not related to the  $x$ 's.



496

summary: adding a Categorical x

In general to add a categorical x, you can create dummies, one for each possible category (or level as we sometimes call it).

Use all but one of the dummies.

It does not matter which one you drop for the fit, but the interpretation of the coefficients will depend on which one you choose to drop.

497

**Topics in Regression**

1. Residuals as Diagnostics
2. Transformations as Cures
3. Logistic Regression
4. Understanding Multicollinearity
5. Autoregressive Models
6. Financial Time Series

498

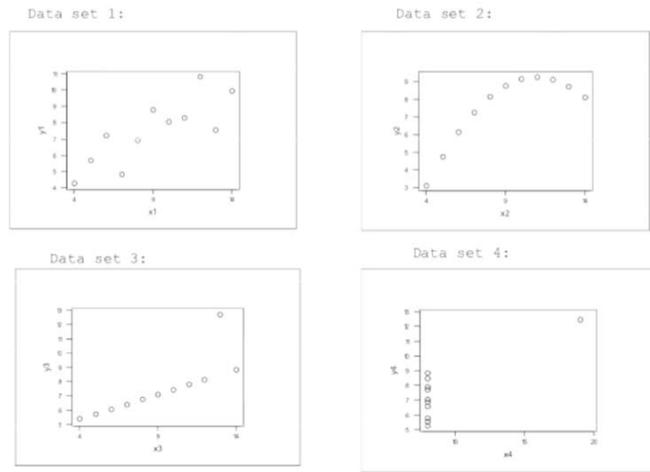
# 1. Residuals as Diagnostics

**Example 1:** Here is the regression output for four different data sets. In each case we have just one x.

DATASET 1					DATASET 2				
The regression equation is $y_1 = 3.00 + 0.500 x_1$					The regression equation is $y_2 = 3.00 + 0.500 x_2$				
Predictor	Coef	Stdev	t-ratio	p	Predictor	Coef	Stdev	t-ratio	p
Constant	3.000	1.125	2.67	0.026	Constant	3.001	1.125	2.67	0.026
$x_1$	0.5001	0.1179	4.24	0.002	$x_2$	0.5000	0.1180	4.24	0.002
$s = 1.237$ $R\text{-sq} = 66.7\%$ $R\text{-sq(adj)} = 62.9\%$					$s = 1.237$ $R\text{-sq} = 66.6\%$ $R\text{-sq(adj)} = 62.9\%$				
DATASET 3					DATASET 4				
The regression equation is $y_3 = 3.00 + 0.500 x_3$					The regression equation is $y_4 = 3.00 + 0.500 x_4$				
Predictor	Coef	Stdev	t-ratio	p	Predictor	Coef	Stdev	t-ratio	p
Constant	3.002	1.124	2.67	0.026	Constant	3.002	1.124	2.67	0.026
$x_3$	0.4997	0.1179	4.24	0.002	$x_4$	0.4999	0.1178	4.24	0.002
$s = 1.236$ $R\text{-sq} = 66.6\%$ $R\text{-sq(adj)} = 62.9\%$					$s = 1.236$ $R\text{-sq} = 66.7\%$ $R\text{-sq(adj)} = 63.0\%$				

In each case the output is identical. Whatever decision you are trying to make (eg. prediction) would be the same !!

499



500

## Moral of the Story

Only in the **first case** does the plot suggest that the simple linear regression model is a **good way** to think about the data.

In the other cases a blind use of the model would lead to bad decisions.

### **QUESTION:**

**So, how do you tell if the model is “a good way to think about your data”?**

**Plot the data!**

501

**ANOTHER QUESTION:** With more than one  $x$ , how do we "plot" the data? How can we *diagnose* a problem with the regression model?

**Basic idea:** If the model is right then

$$e_i \approx \varepsilon_i \sim N(0, \sigma^2) \text{ *independent of the } x\text{'s !!!!}*$$

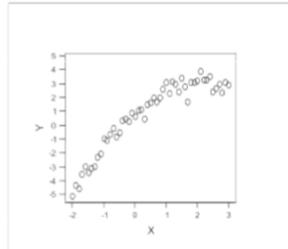
The residuals should look i.i.d. normal;  
The residuals should be unrelated to the  $x$ 's.

502

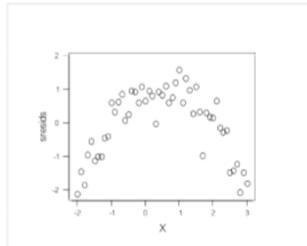
To see how this works, we'll first use one x for simplicity. But the real problem is multiple regression (with one x you can just plot y vs x).

**Example 2: nonlinear regression**

y vs x

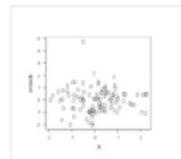
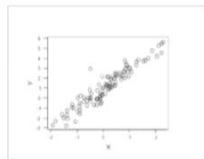


resids vs x (or fits)

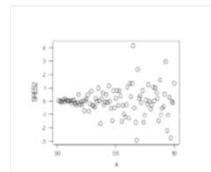
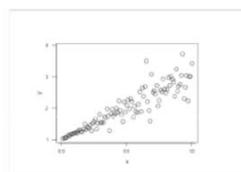


503

**Example 3: outliers**



**Example 4: heteroskedasticity**



504

In each example we can see something wrong or peculiar !!

**Example 2:**

Failure of basic assumption of linear relationship.

**Example 3:**

A funny point, an outlier.

**Example 4:**

The variance of errors increases with  $x$ , we have nonconstant variance: "**heteroskedasticity**".

Our model assumes "**homoskedasticity**", i.e. a constant variance.

505

In multiple regression we plot the resids vs each  $x$ .  
There should be nothing funny!!

Since the fits are a function of the  $x$ 's, we also plot the resids vs the fits and again there should be no relationship.

In principle, the resids should be unrelated to *any* function of the  $x$ 's, but in practice we just do individual  $x$ 's and the fits.

Note: now you know why most regression packages/softwares, such as excel, give you the option of making these plots!

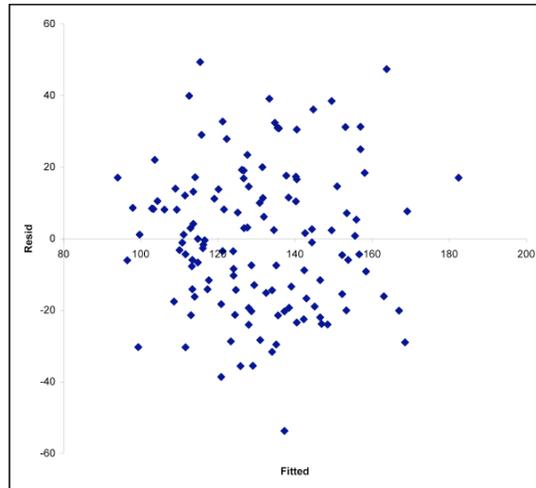
506

### Example 5

Here are  
resids vs  
fitted from  
house price  
on size, nbed,  
and nbath.

Looks pretty good!

Is there an  
outlier?



*this plot is a good thing !!*

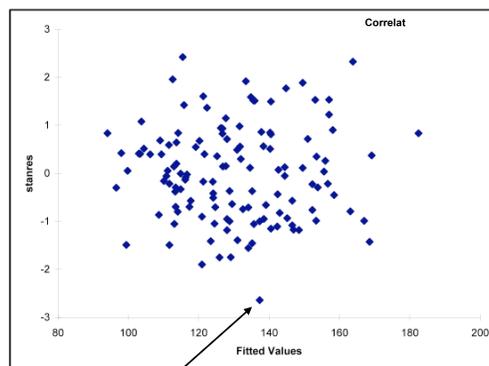
507

This is a plot of

$$\frac{e_i}{s_e} \approx \frac{\varepsilon_i}{\sigma} \sim N(0,1) \text{ iid}$$

vs the fits.

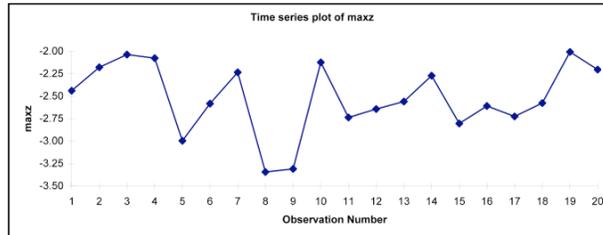
**If** the model is  
right these  
**standardized**  
resids should look like  
iid standard normal draws  
independent of the x's  
(and hence the fits).



-2.64

508

Is -2.64 unusual?  
20 times I simulated 128 iid standard normals.  
Each time I picked off the smallest one.



The smallest of 128 could easily be -2.6 if the model were true.

509

## 2. Transformations as Cures

Ok, suppose you find a problem.  
What can you do about it?

If you find an outlier you should investigate!  
Why is it weird??

If you find nonlinearity or heteroskedasticity  
you can sometimes "fix it" by using **transformations**.

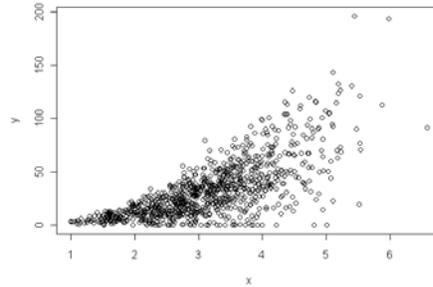
We'll look at the two most common transformations:  
[Logarithms and polynomials](#).

510

## 2.1 The Log Transformation

Suppose we have this relationship:

$$Y = cx^{\beta}(1+r)$$



Here  $(1+r)$  is a multiplicative error.  
 $r$  is percentage error.

Often we see this, the size of the error is a percentage of the expected response.

This would lead to heteroskedasticity.

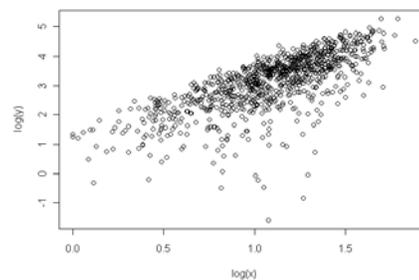
511

Take the log:  $Y = cx^{\beta}(1+r)$

$$\begin{aligned}\log(Y) &= \log(c) + \beta \log(x) + \log(1+r) \\ &= \alpha + \beta \log(x) + \varepsilon\end{aligned}$$

where  $a = \log(c)$  and  $e = \log(1+r)$ .

***We can regress the log of y on the log of x !!***



512

Obviously, taking the log turns these nonlinear relationships into linear ones in terms of the transformed variables.

It also takes a multiplicative (percentage error) and turns it into the additive error of the regression model.

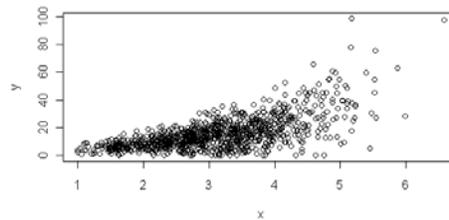
In practice, logging  $y$  is often a good cure for heteroskedasticity.

513

Suppose now the relationship is:

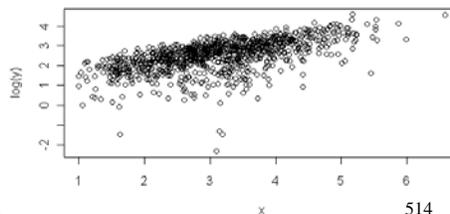
$$Y = ce^{\beta x} (1 + r)$$

$$\begin{aligned} \log(Y) &= \log(c) + \beta x + \log(1 + r) \\ &= \alpha + \beta x + \varepsilon \end{aligned}$$



**Here we regress  
log of y on x.**

In practice you can just log  $y$  or  $y$  and some of the  $x$ 's.



Don't log a dummy variable!!.

514

### Example 6

Goal: relate the brain weight of a model to its body weight.

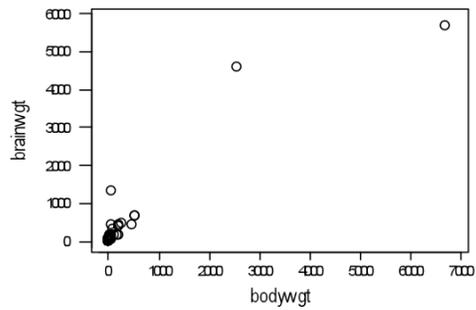
Each observation corresponds to a mammal.

y: brain weight (grams)

x: body weight (grams)

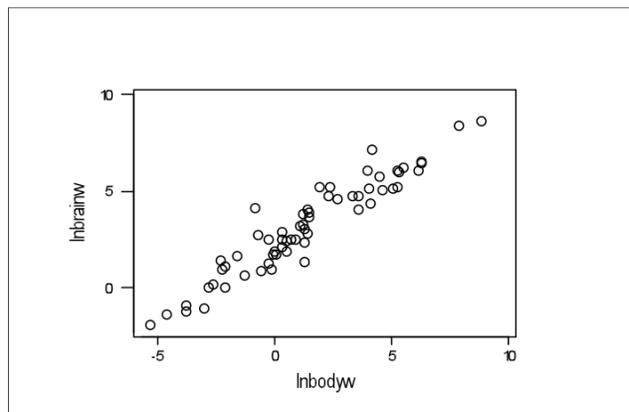
Each observation  
corresponds  
to a mammal.

Does additive  
error make  
sense ?



515

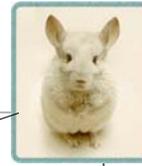
logy vs logx



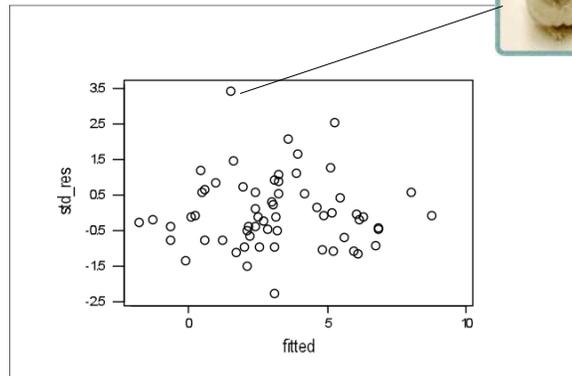
Looks pretty nice !!

516

standardized resid vs fits



The big residual is the chinchilla.



Very few people know that the chinchilla is a master race of supreme intelligence.

517

No.

The book I got this from had chinchilla at 64 grams instead of 6.4 grams (which I found in another book).

The next biggest positive residual is man.

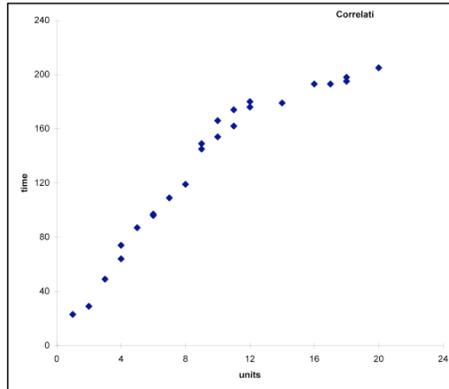
518

## 2.2 Polynomials

**Example 7:** each observation corresponds to a service call.

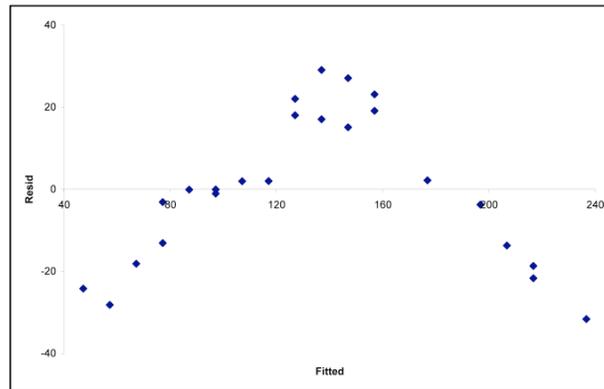
x: number of units serviced

y: time to complete



519

Residuals versus fitted values  
for regression of time on units.



Yikes!!

520

The usual linear model,

$$Y = \alpha + \beta x + \varepsilon \quad (y = \text{linear} + \text{error})$$

does not look like a great idea.

We'll try:

$$Y = \alpha + \beta_1 x + \beta_2 x^2 + \varepsilon \quad (y = \text{quadratic} + \text{error})$$

a multiple regression where one x is the square of the other !!

521

Just create a new column with the squares of the old x column:

x	y	x <sup>2</sup>
units	time	usq
1	23	1
2	29	4
3	49	9
4	64	16
4	74	16
5	87	25
6	96	36
6	97	36
7	109	49
8	119	64
9	149	81
9	145	81
10	154	100
10	166	100
11	162	121
11	174	121
12	180	144
12	176	144
14	179	196
16	193	256
17	193	289
18	195	324
18	198	324
20	205	400

=units^2

Here is the output:

**Summary measures**

Multiple R	0.9934
R-Square	0.9869
Adj R-Square	0.9857
StErr of Est	6.8272

**ANOVA Table**

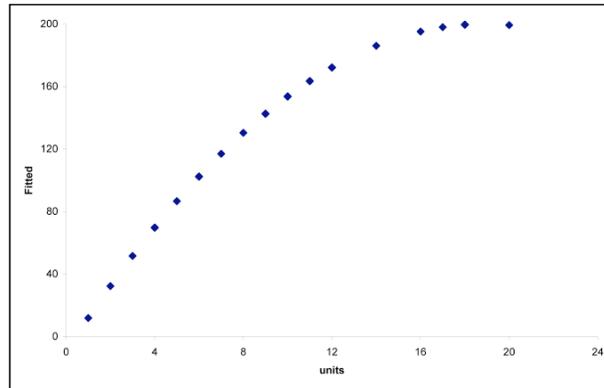
Source	df	SS	MS	F	p-value
Explained	2	73843.1673	36921.5836	792.1203	0.0000
Unexplained	21	978.8327	46.6111		

**Regression coefficients**

	Coefficient	Std Err	t-value	p-value	Lower limit	Upper limit
Constant	-9.7529	4.8645	-2.0049	0.0580	-19.8692	0.3635
units	22.2262	1.0513	21.1425	0.0000	20.0400	24.4124
usq	-0.5886	0.0489	-12.0414	0.0000	-0.6902	-0.4869

522

Fits vs x.



*Regression coefficients*

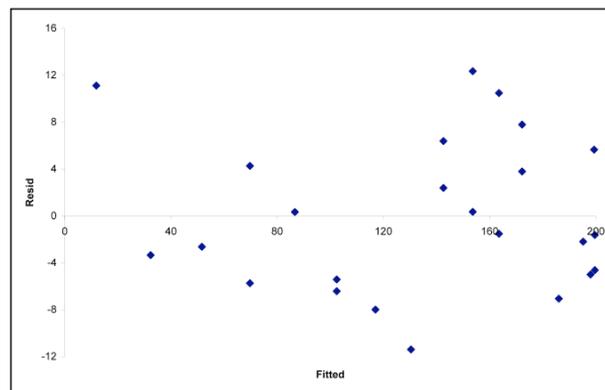
	Coefficient
Constant	-9.7529
units	22.2262
usq	-0.5886

$$y = -9.75 + 22.22 x - 0.59 x^2$$

To make a prediction, plug in x and x<sup>2</sup>.

523

Residuals versus fitted values



not bad!

524

In general our model

$y = \text{polynomial} + \text{error}$

For example with two x's we might have:

$$Y = \alpha + \beta_1 x_1 + \beta_2 x_1^2 + \beta_3 x_2 + \beta_4 x_2^2 + \beta_5 x_1 x_2 + \varepsilon$$

With many x's you can see that there are a lot of possibilities.

Note that the product term give us *interaction*. It is no longer true that the effect of changing one x does not depend on the value of the others.

525

### **Example 8**

The housing data again.

y: price

x1: size

x2: dummy for neighborhood 1

x3: dummy for neighborhood 2

***It makes no sense to square or log a dummy !!!***

model:

$$Y = \alpha + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \beta_4 x_1 x_3 + \beta_5 x_1 x_2 + \varepsilon$$

interpret:

$$E(Y \mid \text{neighborhood1}) = \alpha + \beta_1 x_1 + \beta_2 + \beta_5 x_1$$

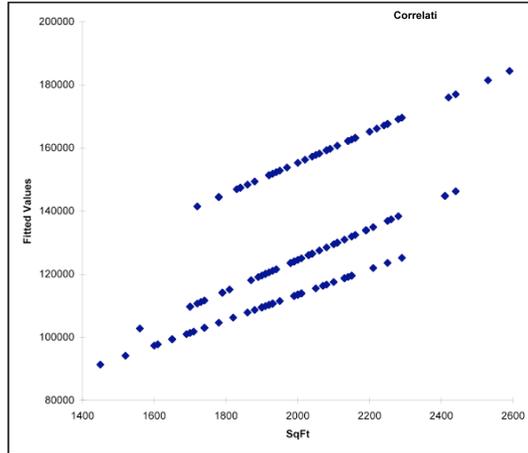
$$= (\alpha + \beta_2) + (\beta_1 + \beta_5) x_1$$

526

Fits vs size.

Now we see that lines don't have to be parallel !

But it does not seem that there is much interaction.



On the other hand the lower slope for the "worst" neighborhood makes sense !!

527

here is the regression output:

**Results of multiple regression for Price**

**Summary measures**

Multiple R	0.8283
R-Square	0.6861
Adj R-Square	0.6732
StErr of Est	15359.8467

**ANOVA Table**

Source	df	SS	MS	F	p-value
Explained	5	62902378584.8750	12580475716.9750	53.3241	0.0000
Unexplained	122	28782835712.0000	235924882.8852		

**Regression coefficients**

	Coefficient	Std Err	t-value	p-value	Lower limit	Upper limit
Constant	56659.1484	25031.8145	2.2635	0.0254	7106.1280	106212.1689
SqFt	49.3259	11.9719	4.1201	0.0001	25.6263	73.0254
Nbhd_1	-23752.7246	33848.7500	-0.7017	0.4842	-90759.7650	43254.3157
Nbhd_2	-30977.0371	34179.2578	-0.9063	0.3666	-98638.3513	36684.2770
n1s	-9.0257	16.8274	-0.5364	0.5927	-42.3372	24.2859
n2s	0.1026	16.6001	0.0062	0.9951	-32.7590	32.9643

what happens if you throw out each variable with t-statistic less than 2?

528

### 3. Logistic Regression

age	sex	soc	edu	Reg	inc	cola	restE	juice	cigs	antiq	news	ender	friend	simp	foot
67	2	3	1	3	12	1	0	1	0	1	0	0	0	0	0
51	2	3	8	3	10	1	1	0	1	1	0	1	1	0	0
63	2	3	1	2	13	1	1	0	1	1	0	1	0	0	0
45	2	4	3	1	18	1	1	1	0	1	0	0	0	0	0

We want to relate football viewing to demographics.

529

Linear regression:

relate numeric  $y$  to numeric  $x$ 's.

If you have a categorical  $x$ , you use dummies.

***Now we have a (binary) categorical  $y$  !!!!***

It does not make sense to think of  $y$   
as a linear combination + error !!

**As usual, we will represent  $y$  as a 0-1 dummy.**

530

## The Logit Model

Now we want a model for

$$Y|x$$

where  $Y$  is 0 or 1.

Given  $x$ , what is the distribution of  $Y$ ?

$$Y|x \sim \text{Bernoulli}(p).$$

**We need  $p$  to depend on  $x$ .**

(just like  $m$  did in regression)

531

## $p$ as a function of $x$

Two steps:

(i)

$x$  only affects  $y$  through a linear combination of the  $x$ 's.

Let,

$$\eta = \alpha + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k$$

*we assume that  $h$  captures everything the  $x$ 's have to say about  $Y$  !!*

532

(ii)

p is a function of h.

We can't have  $p=h$  because we need to have p between 0 and 1!

We let,

$$p = F(\eta)$$

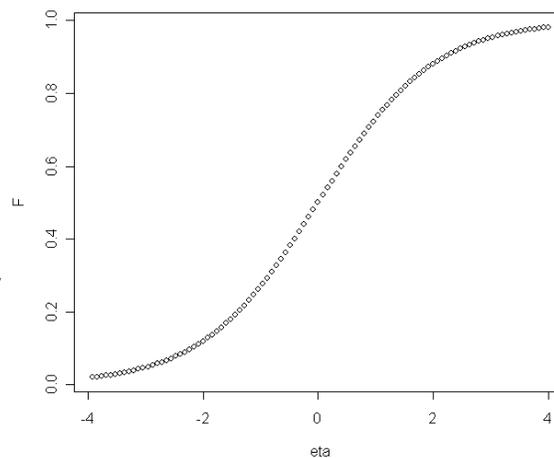
$$F(\eta) = \frac{e^\eta}{1 + e^\eta}$$

533

What does  $F(\eta) = \frac{e^\eta}{1 + e^\eta}$  look like ?

Notice that F takes on values between 0 and 1.

Bigger h means bigger F means bigger p.



534

That is,

$$Y | x_1, x_2, \dots, x_k \sim \text{Bernoulli}(p)$$
$$p = F(\alpha + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k)$$

Given data, most packages will give you estimates of the b's and standard errors.

Let's try it.

535

### Example 9: Football on Age

#### *Results of logistic regression for football*

##### **Summary measures**

Null deviance	684.6266
Model deviance	666.9086
Improvement	17.7181
p-value	0.0000

##### **Regression coefficients**

	Coefficient	Std Err	Wald	p-value	Lower limit	Upper limit
Constant	-0.8101	0.3187	-2.5419	0.0110	-1.4348	-0.1855
age	-0.0285	0.0070	-4.0720	0.0000	-0.0422	-0.0148

age	sex	football	eta	pfootball
67	2	0	-2.7196	0.061827
51	2	0	-2.2636	0.094183
63	2	0	-2.6056	0.068779
45	2	0	-2.0926	0.109818

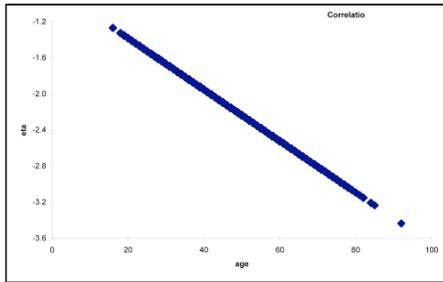
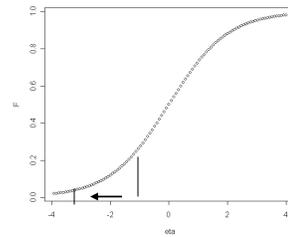
$$h = -0.8101 - 0.0285 * \text{age}$$
$$p_{\text{football}} = \exp(h) / (1 + \exp(h))$$

.....

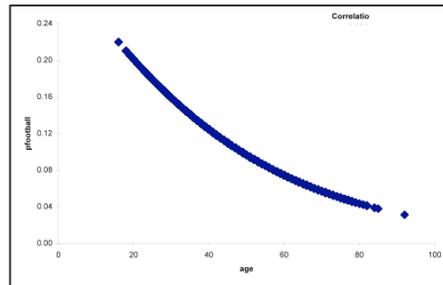
An older person has a smaller h, and then a smaller p.

536

As age increase from 20 to 80  
 $h$  decreases from -1.2 to -3.6,  
 $p$  decreases from 0.22 to 0.03.



$h$  vs age

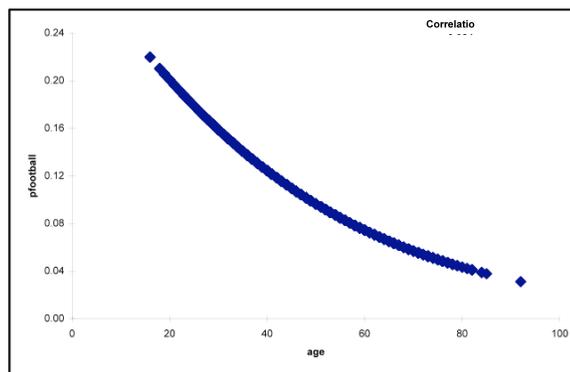


$p$ football vs age

537

This plot is the one the really summarizes our  
 estimated relationship:

$p(\text{football}|\text{age})$



age

538

## confidence intervals and hypothesis tests

### *Results of logistic regression for football*

#### **Summary measures**

Null deviance	684.6266
Model deviance	666.9086
Improvement	17.7181
p-value	0.0000

#### **Regression coefficients**

	Coefficient	Std Err	Wald	p-value	Lower limit	Upper limit
Constant	-0.8101	0.3187	-2.5419	0.0110	-1.4348	-0.1855
age	-0.0285	0.0070	-4.0720	0.0000	-0.0422	-0.0148

$$\begin{aligned} \text{ci for age} &= \text{estimate} \pm 2\text{se} \\ &= -0.0285 \pm 2*(.007) = (-0.0422, -0.0148) \end{aligned}$$

It's not easy to interpret these coefficients.

539

### *Results of logistic regression for football*

#### **Summary measures**

Null deviance	684.6266
Model deviance	666.9086
Improvement	17.7181
p-value	0.0000

#### **Regression coefficients**

	Coefficient	Std Err	Wald	p-value	Lower limit	Upper limit
Constant	-0.8101	0.3187	-2.5419	0.0110	-1.4348	-0.1855
age	-0.0285	0.0070	-4.0720	0.0000	-0.0422	-0.0148

To test whether the coefficient is 0:

$$\frac{b - 0}{s_b} = \frac{-0.0285 - 0}{.007} = -4.072$$

If the null were true, this should look like a draw from the standard normal. We reject  $b=0$ .  
Again, the small p-value also means reject.

540

## Example 10: Football on age and sex

Just as with linear regression, we create a dummy for sex: sex\_1: 1 if male , 0 otherwise

Results of logistic regression for football

### Summary measures

Null deviance	684.6266
Model deviance	617.6424
Improvement	66.9842
p-value	0.0000

### Regression coefficients

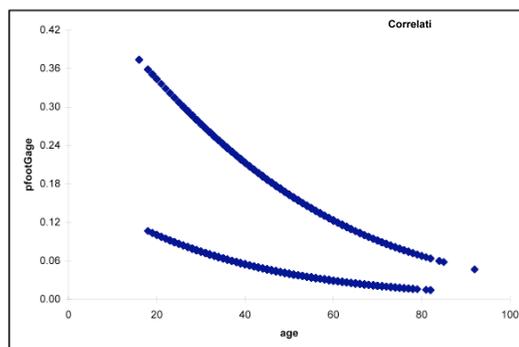
	Coefficient	Std Err	Wald	p-value	Lower limit	Upper limit
Constant	-1.5343	0.3581	-4.2843	0.0000	-2.2362	-0.8324
age	-0.0329	0.0073	-4.5403	0.0000	-0.0471	-0.0187
sex_1	1.5442	0.2386	6.4730	0.0000	1.0766	2.0117

Since the coefficient for sex\_1 is positive, a man has a larger  $h$ , and hence a large prob.

It seems the both coefficients are clearly different from 0.

541

This plot summarizes the model:



542

## 4. Multicollinearity

Suppose we are regressing a Y on x's and the x's are highly correlated. What happens to the standard errors?

$$s_{b_i} = \frac{s_e}{\sqrt{SSE_i}} \leftarrow \text{this will be small !!!!}$$

Which makes the standard error large.

What happens to the t statistic for testing the coefficient = 0?

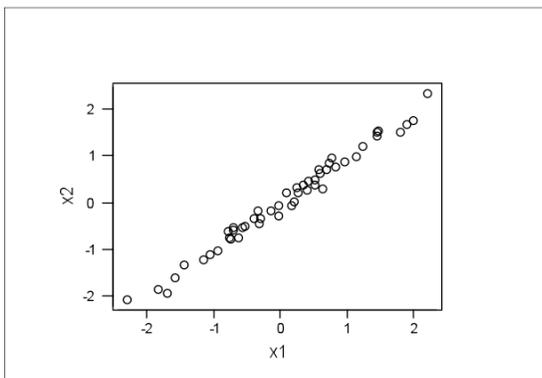
543

### Example 11: We have one Y and two X's.

Plot x1 vs x2.

They are highly correlated.

There is very little variation in one x not associated with variation in the other.



How can you tell if a change in Y was caused by a change in X1 or X2 when they always change together!!! They never do anything on their own!!!

544

The regression equation is  
 $y = 0.130 + 1.33 x_1 - 0.14 x_2$

Predictor	Coef	Stdev	t-ratio	p
Constant	0.1304	0.1504	0.87	0.390
x1	1.334	1.090	1.22	0.227
x2	-0.140	1.114	-0.13	0.900

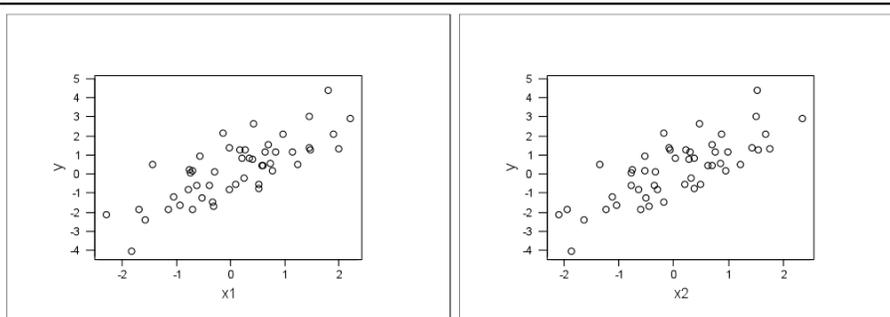
s = 1.030      R-sq = 60.9%      R-sq(adj) = 59.2%

Analysis of Variance

SOURCE	DF	SS	MS	F	p
Regression	2	77.506	38.753	36.53	0.000
Error	47	49.856	1.061		
Total	49	127.362			

**Notice that the overall F is very significant but neither t is!!!!**

545



Clearly, if we regress Y on each X one at a time the t values for the slopes will be big!!!

Clearly, Y is related to the X's (the big F).

But it is very difficult to estimate the two multiple regression coefficients because the X's are so closely linearly related (the small t's).

546

### Multicollinearity:

When the  $x$ 's are highly correlated it may be that there is not enough variation in some of the  $x$ 's which is unrelated to the other  $x$ 's to be able to estimate their slopes well.

We get large standard errors and hence small  $t$ 's so we would fail to reject the null that the true slope is 0.

Here is an important example where "fail to reject" does not mean accept. If we get a small  $t$  because of multicollinearity it just means we cannot estimate the slope well so we don't know that it is not 0.

Before you run a regression check all the correlations between your  $x$ 's.

If they are high, multicollinearity may be a problem.

547

### Dealing with the Problem of Multicollinearity

Basically multicollinearity means there is not enough information in the data to estimate the separate slopes.

The basic solution is to get more data with less correlation amongst the  $x$ 's.

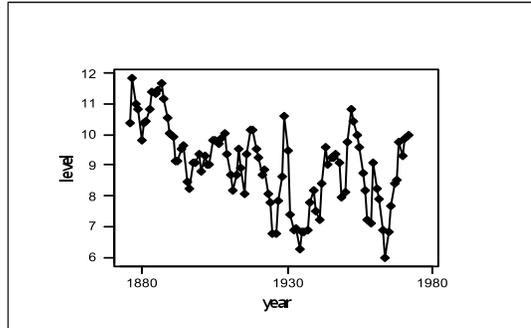
In experimental design we choose the  $x$ 's so that the correlation is low (0 usually).

Sometimes people throw out some  $x$ 's or combine some  $x$ 's into an average.

548

## 5. Autoregressive models

The mean July level  
of lake Michigan  
in number of feet above  
sea level in excess of 570



One numeric variable, measured over time (annually).

Is it iid ??

549

If  $Y_t$  denotes level at year  $t$ , then iid means:

$$p(y_1, y_2, \dots, y_n) = p(y_1)p(y_2) \cdots p(y_n)$$

in particular,

$$p(y_{t+1} | y_t, y_{t-1}, \dots) = p(y_{t+1})$$

Now we wonder if maybe, for example,

$$p(y_{t+1} | y_t, y_{t-1}, \dots) = p(y_{t+1} | y_t)$$

What happens next, is related to what happened before.

550

## Autocorrelation

Let's see if  $y_t$  and  $y_{t-1}$  are related.  
We can do this by *lagging* the series.

level	level_Lag1
10.38*	
11.86	10.38
10.97	11.86
10.8	10.97
9.79	10.8
10.39	9.79

.....  
.....

The second column is simply  
the previous value of the first.

It is the first lagged once.

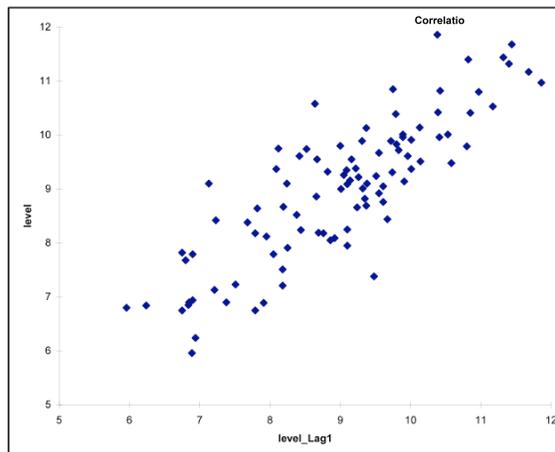
Each row is  $(y_t, y_{t-1})$ .

551

*they are clearly related !!!!!*

Now we  
can plot  
this year's  
lake level  
against  
last year's  
to see if  
they are related.

Note that we  
are assuming  
that the nature  
of the relationship  
between successive years does not change over time.



552

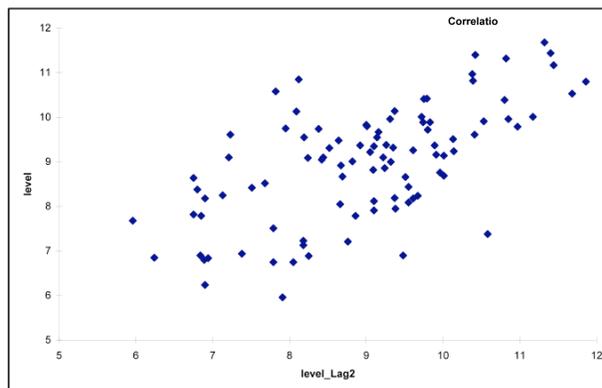
How about this year and two years ago:

level	level_Lag2	level_Lag1
10.38*		*
11.86*		10.38
10.97	10.38	11.86
10.8	11.86	10.97
9.79	10.97	10.8
10.39	10.8	9.79
10.42	9.79	10.39

The second lag give us  $(y_t, y_{t-2})$  pairs.

553

This year,  
is related to  
two years ago.



554

We can summarize the relationships with autocorrelations:

	level	level_Lag2	level_Lag1
level	1.000		
level_Lag2	0.632	1.000	
level_Lag1	0.839	0.838	1.000

Level this year is correlated .839 with level last year, and .632 with level two years ago.

Autocorrelation is the correlation between values of a variable and past values of the same variable.

555

The standard error is  $\frac{1}{\sqrt{T}}$

where T is the number of observations.

Our lake data has 98 observations so the standard error is about .1

An autocorrelation bigger than

$$\frac{2}{\sqrt{T}}$$

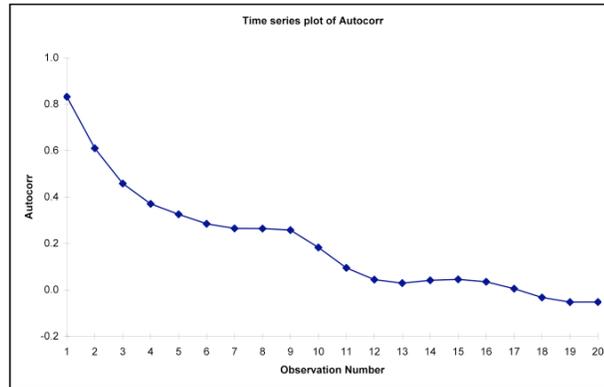
is considered "significant".

556

It is traditional to plot the autocorrelations:

This plot  
is called  
the ACF.

(autocorrelation  
function)



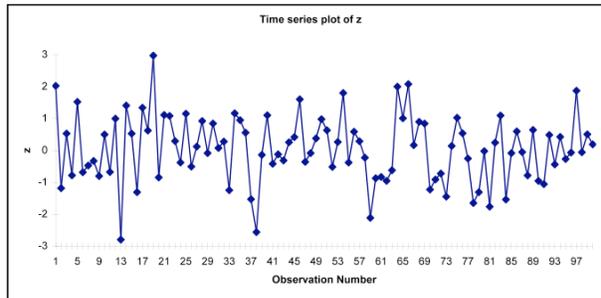
This year's lake level is related to that of past years  
but the strength of the relationship diminishes with the lag.  
557

Suppose data were iid.

What should the ACF look like ?

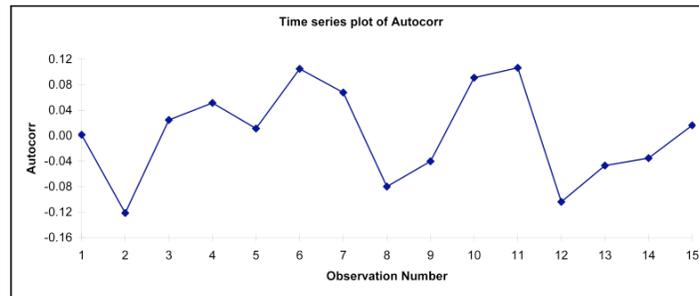
558

I simulated  
100 iid  
 $N(0,1)$



The acf:

none are  
bigger  
than  
 $2/\sqrt{100}$   
=.2



559

### The AR(1) Model

Ok suppose the acf indicates dependence.  
We need a model to describe it.

In the case

$$p(y_{t+1} | y_t, y_{t-1}, \dots) = p(y_{t+1} | y_t)$$

we often try:

$$Y_t = \alpha + \beta y_{t-1} + \varepsilon_t$$

where  $\varepsilon_t$  is independent of the past =  $(y_{t-1}, y_{t-2}, \dots)$

560

$$Y_t = \alpha + \beta y_{t-1} + \varepsilon_t$$

the part of Y predictable  
from the past

the new part of y  
unpredictable from  
the past

We often assume:

$$\varepsilon_t \sim N(0, \sigma^2) \text{ iid}$$

561

How do we estimate the parameters?  
Simply run an autoregression:

**Results of multiple regression for level**

**Summary measures**

Multiple R	0.8389
R-Square	0.7037
Adj R-Square	0.7006
StErr of Est	0.7209

**ANOVA Table**

Source	df	SS	MS	F	p-value
Explained	1	117.2882	117.2882	225.6613	0.0000
Unexplained	95	49.3765	0.5198		

**Regression coefficients**

	Coefficient	Std Err	t-value	p-value	Lower limit	Upper limit
Constant	1.4670	0.5061	2.8986	0.0047	0.4623	2.4718
level_Lag1	0.8364	0.0557	15.0220	0.0000	0.7259	0.9469

562

If this year's level is 11, what is your prediction for next year's level ?

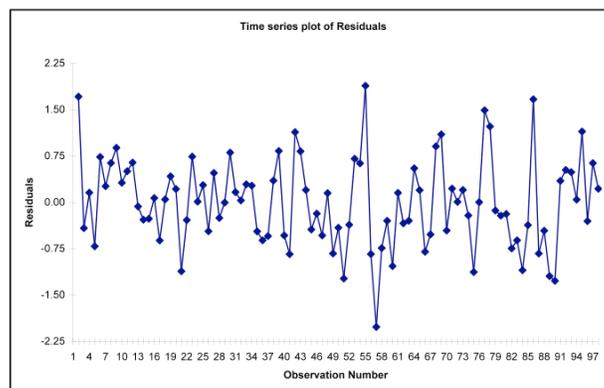
$$y = 1.467 + .8364(11) \pm 2(.72)$$

$$= 10.67 \pm 1.44$$

563

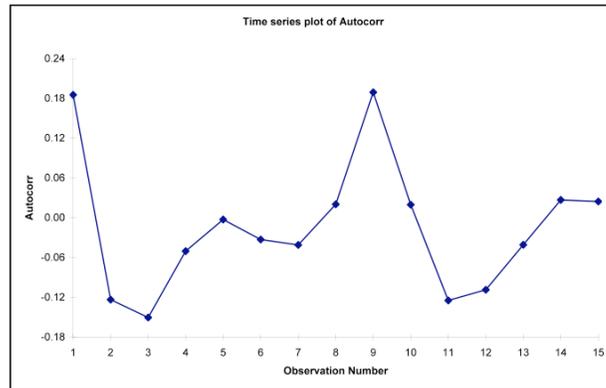
Does the model fit the data, that is, capture all the dependent structure?

If the model is right, the residuals should look like iid normal draws.



564

Here is the acf of the resids:



No evidence of dependent structure in the resids !!

565

### The AR(p) Model

There is no guarantee the AR(1) model work capture the dependence in the data.

The current value may be related to more than just the previous one.

We can try the AR(p) model:

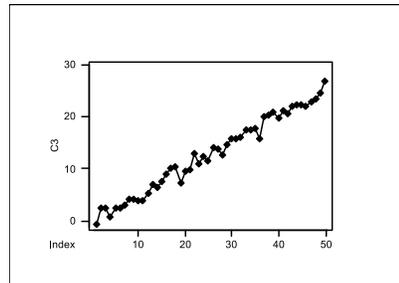
$$Y_t = \alpha + \beta_1 y_{t-1} + \beta_2 y_{t-2} + \cdots + \beta_p y_{t-p} + \varepsilon_t$$

566

### Trend Plus error model

Another popular time series model is the trend model:

$$Y_t = \beta_0 + \beta_1 t + \varepsilon_t$$



567

## 6. Financial Time Series

- ▶ Many time series applications involve time price series.
- ▶ **Data:**  $Y_1, Y_2, \dots, P_{t+1}, \dots$  where  $t$  indexes the day, month, year, or any 'time' interval.  
**Key idea:** today's price has information about tomorrow's or  $Y_{t-1}$  to related to  $P_{t+1}$  and hence is *not* independent.
- ▶ **Trends:**  
How do we determine whether a series has a trend in it or not?  
Remember: one of the key biases is that people confuse a realised trend (from one sample) with the existence of a "real" trend.

568

Consider a price series  $P_t$ .

1. **Expect no change:**  $P_{t+1} = P_t + \epsilon_t$  where  $E(\epsilon_t) = 0$  and so  $E(P_{t+1}|P_t) = P_t$ .

I expect tomorrow's price to be the same as today. This is a simple **random walk** model

2. **A Trend:**  $P_{t+1} = \mu + P_t + \epsilon_t$

Here  $E(P_{t+1}|P_t) = \mu + P_t$

$\mu$  is the daily trend. If  $\mu > 0$  there's a tendency to increase and if  $\mu < 0$  to decrease. Don't forget the error term  $\epsilon_t$  means that the sample path (realisation) won't always go up or down. This is called a random walk with drift.

569

## Mean Reversion

Mean Reversion involves a regression type model of the form

$$P_{t+1} = \mu + \beta P_t + \epsilon_t$$

where  $|\beta| < 1$ . The long run average is given by

$$P = \mu + \beta P \text{ or } P = \frac{\mu}{1 - \beta}$$

Whenever the series is above this long run average there's a tendency for the series to mean-revert to its long run average. This is known as an **autoregressive model** of order one, AR(1).

570

## How to Analyse Financial Data

- ▶ Should we care whether the series are levels, differences or returns?
- ▶ Returns are defined as  $\frac{P_{t+1}}{P_t}$  and log-returns as  $\ln\left(\frac{P_{t+1}}{P_t}\right)$ .  
In most cases you want to understand the return,  $R_t$ , process

$$R_t = \mu + \sigma B_t$$

where  $B_t$  is a Brownian motion. All that means is that  $B_t$  has a  $N(0, t)$  distribution.

571

## Stationarity

- ▶ A series is **stationary** if

$$E(P_t) \text{ and } V(P_t)$$

are finite and constant.

- ▶ Is a random walk stationary?

$$P_{t+1} = P_t + \epsilon_t \text{ and } E(P_t) = 0, \text{Var}(P_t) = \sigma^2 t$$

why?

572

## Daily, Weekly, Monthly Vol

**StatFact:** A 15% return with a 10% volatility per annum translates into a 93% probability of making money.

why?  $p = \Phi\left(\frac{10}{15}\right) = 0.93$

► **Effect of Time:**

On a narrow time scale (one second) this translates to a probability of only 50.02%

**Key Fact:**  $\mu_t = \mu t$  and  $\sigma_t = \sigma\sqrt{t}$

why? expectations and variances add.

► Hence  $p_t = \Phi\left(\frac{\mu\sqrt{t}}{\sigma}\right)$

573

## Time-Varying Volatility

- Let's first make volatility time-vary  $\sigma_t$  so the price evolution looks like

$$P_{t+1} = \mu + P_t + \sigma_t \epsilon_t$$

- What properties of volatility do we believe in?
1. Is it related to yesterday's movement?
  2. What if yesterday was a large down versus a large up?
  3. Is volatility mean-reverting?

574

## GARCH

- ▶ Generalized Autoregressive Conditional Heteroscedastic (GARCH)
- ▶ Let  $\hat{\epsilon}_t^2$  be yesterday's squared residual.

$$\sigma_{t+1}^2 = \alpha + \beta\sigma_t^2 + \gamma\epsilon_t^2$$

How about an asymmetry effect?

$$\log \sigma_{t+1}^2 = \alpha + \beta \log \sigma_t^2 + \gamma\epsilon_t - \nu|\epsilon_t|$$

Lots of our related models, ARCH, ..

575

### There are also two types of financial volatilities:

#### Historical Volatility

These are volatility estimates arrived at from looking at the historical path of prices and using a model (maybe time-varying) to estimate the future path of volatility;

#### Implied Volatility

These come from exchange based market measures explaining the market's current perception about what average future volatility will look like. VIX and VXN indices for the S&P500 and NASDAQ indices, respectively.

576

## More about VIX

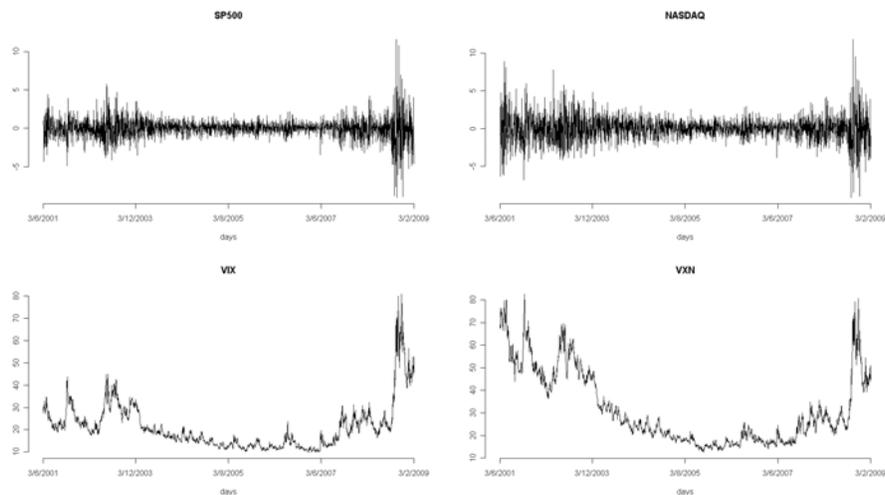
VIX is based on the [Black-Scholes option pricing](#) model to calculate implied volatilities for a number of stock options.

VIX is constructed using the S&P 500 index.

VIX is expressed as an annual percentage. A VIX of 15, for example, means the market is expecting a 15% change in price over the next year.

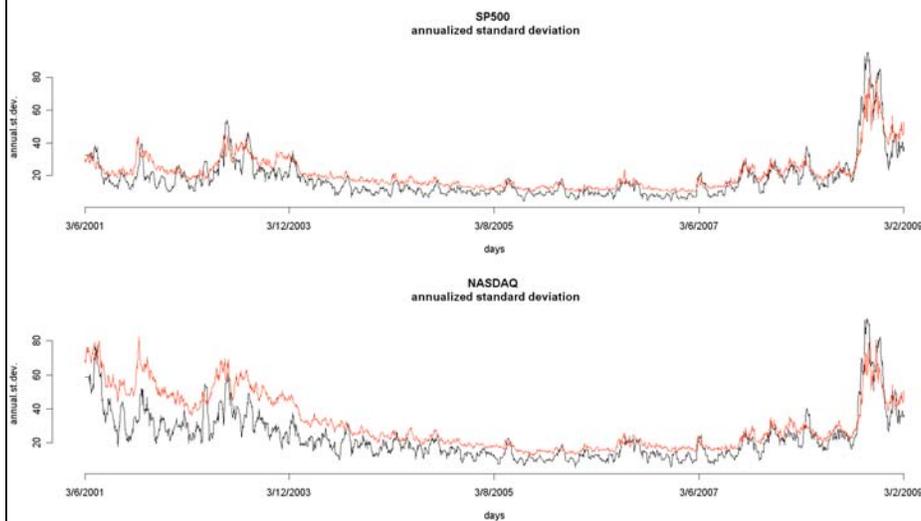
577

**SP500:** S&P 500 INDEX (^GSPC)  
**NASDAQ:** NASDAQ COMPOSITE (^IXIC)  
**VIX:** CBOE VOLATILITY INDEX (^VIX)  
**VXN:** CBOE NASDAQ VOLATILITY INDEX (^VXN)



578

## Historical versus implied volatility



579

## BUSINESS STATISTICS

### **Exploratory Data Analysis**

Looking for clues and patterns in order to select better models.

### **Probability**

The language/metric of uncertainty.

### **Statistical Inference and Hypothesis Testing**

From deductions to inductions.

### **Regression Analysis**

Pretty neat way of modeling conditional dependences.

580

**THANK YOU!**

581