

## Announcements:

- You can turn in homework until 6pm, slot on wall across from 2202 Bren. Make sure you use the correct slot for Stat 7B (Utts).
- We are going to start using R Commander for homework, with today's assignment. See course website links under headings:
  - Installing R and R Commander and Getting Started
  - R Commander Handouts and Instruction Sheets
  - Pdf and Word versions; Word allows you to copy and paste
- Discussion this week is not for credit. Will go over using R Commander, and then general Q&A. If you plan to follow along with R Commander, install it *before* discussion and bring your computer.

## Homework (due Wed Jan 23):

Remember to see "How to use R Commander for assignment from Chapters 2 and 3" on course webpage. It includes instructions on how to save a graph. You can then print it or insert it in a document.

### Homework:

Ch. 2: # 104\*

Ch. 3: #24 (you do *not* need the original data; note typo in some printings of the book – intercept is 126, not 1.26)

Ch. 3: #98\* (Data on class website; worth double points)

\*Use R Commander for 2.104 and 3.98

Data for 3.98 link on website (in file called oldfaithful.txt)

TODAY: Chapter 3, Sections 3.1 and 3.2

## *Relationship between Two Quantitative Variables*

### Motivation

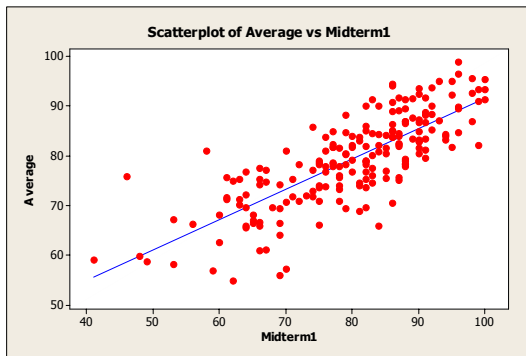
Measure 2 quantitative variables on the same units.

- How strongly related are they?
- In the future, if we know value of one, can we predict the other?

Example: After the 1<sup>st</sup> midterm, how well can we predict your final average for this class?

Data: Class of 200 students where both are known. Use it to create an equation to predict Final Average in future, when first Midterm score is known.

Scatter plot for the example (more later)



## Algebra Review for Linear relationship

Equation for a straight line:

$$y = b_0 + b_1x$$

$b_0$  = y-intercept, the value of  $y$  when  $x = 0$

$b_1$  = slope, the increase in  $y$  when  $x$  goes up by 1 unit

**Example** (deterministic = exact relationship): One pint of water weighs 1.04 pounds. ("A pint's a pound the world around.")

Suppose a bucket weighs 3 pounds. Fill it with  $x$  pints of water. Let  $y$  = weight of the filled bucket.

*How can we find  $y$ , when we know  $x$ ? Easy!*

### Example, continued

$b_0$  = y-intercept, the value of  $y$  when  $x = 0$

This is the weight of the empty bucket, so  $b_0 = 3$

$b_1$  = slope, the increase in  $y$  when  $x$  goes up by 1 unit; this is the added weight for adding 1 pint of water, i.e. 1.04 pounds.

The equation for the line:

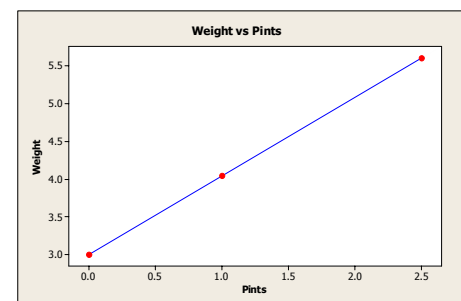
$$y = b_0 + b_1x$$

$$y = 3 + 1.04x$$

$x = 1$  pint  $\rightarrow y = 3 + 1.04(1) = 4.04$  pounds

$x = 2.5$  pints  $\rightarrow y = 3 + 1.04(2.5) = 5.6$  pounds

Plot of the line  $y = 3 + 1.04x$



You have just seen an example of a *deterministic relationship* – if you know  $x$ , you can calculate  $y$  exactly.

**Definition:** In a **statistical relationship** there is *variation* in the possible values of  $y$  at each value of  $x$ .

If you know  $x$ , you can only find an *average* or *approximate* value for  $y$ .

We are interested in describing linear relationships between two quantitative variables. Usually we can identify one as the *explanatory variable* and one as the *response variable*. We always define:

$x$  = explanatory variable

$y$  = response variable

**Examples:**

	Explanatory Variable	Response Variable:
1. Son's height based on parents	$x$ = Average of mom's and dad's heights	$y$ = Son's height
2. Example 3.13	$x$ = Verbal SAT	$y$ = College GPA
3. Example 3.7	$x$ = Driver's age	$y$ = Distance (feet) they can read sign
4. Grades	$x$ = Midterm 1 score	$y$ = Final average

### Relating two quantitative variables

1. Graph – “Scatter plot” – to *visually see* relationship
2. Regression equation – to describe the “best” straight line through the data, and predict  $y$ , given  $x$  in the future.
3. Correlation coefficient – to *describe the strength and direction* of the linear relationship

**Example 1:** Can height of male student be predicted by knowing the average of his parents' heights?

**Example 2:** Can college GPA be predicted from Verbal SAT?

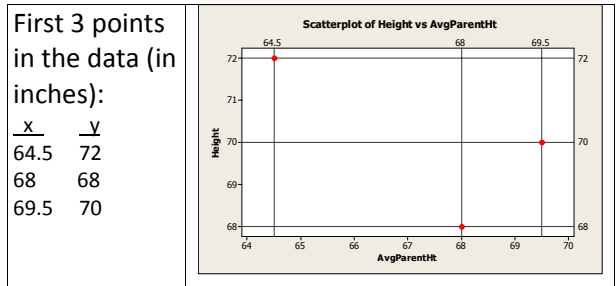
**Example 3:** Can the distance at which a driver can see a road sign be predicted from the driver's age?

**Example 4:** Can final average be predicted from midterm 1 score?

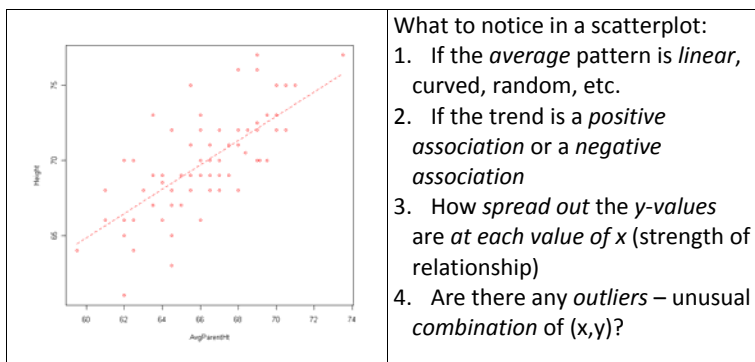
### Creating a scatter plot:

- Create axes with the appropriate ranges for  $x$  (horizontal axis) and  $y$  (vertical axis)
- Put in one “dot” for each  $(x,y)$  pair in the data set.

**Example 1:** Scatterplot of 3 points,  $x$  = avg parent ht,  $y$  = height



## Scatterplot of all 73 individuals, with a line through them



1. Average pattern looks *linear*
2. It's a *positive association* (as x goes up, y goes up, on average)
3. Student heights are quite spread out at each average parents' height
4. There are no obvious outliers in the combination of (x,y)

## REGRESSION LINE (REGRESSION EQUATION)

Basic idea: Find the "best" line to

1. *Estimate* the *average value of y* at a given value of x
2. *Predict* y in the future, when x is *known* but y is not

**Definition:** A **regression line** or **least squares line** is a straight line that best\* describes how values of a quantitative response variable (y) are related to a quantitative explanatory variable (x).

\*Best will be defined later.

Notation for the regression line is:

$$\hat{y} = b_0 + b_1x$$

"y-hat = b-zero + b-one times x"

**Example 1:**  $\hat{y} = 16.3 + 0.809x$

For instance, if parents' average height = 68 inches,

$$\hat{y} = 16.3 + 0.809x$$

$$16.3 + 0.809(68) = 71.3 \text{ inches}$$

Interpretation – the value 71.3 can be interpreted in two ways:

1. An *estimate* of the *average* height of all males whose parents' average height is 68 inches
2. A *prediction* for the height of a *one* male whose parents' average height is 68 inches

NOTE: It makes sense that we predict a male to be *taller* than the average of his parents. Presumably, a female would be predicted to be *shorter* than the average of her parents.

### Example 1, continued

Interpreting the y-intercept and the slope

*Intercept* = 16.3 is the estimated male height when parents' average height is 0. This makes no sense in this example!

*Slope* = +0.809 is the difference in estimated height for two males whose parents' average heights differ by 1 inch.

For instance, if parents' average height is 65 inches,

$$\hat{y} = 16.3 + 0.809(65) = 68.9 \text{ inches}$$

One inch higher parents' average height is 66 inches, and

$$\hat{y} = 16.3 + 0.809(66) = 69.7 \text{ inches}$$

(difference of .809 rounded to .8)

## Prediction Errors and Residuals

Individual  $y$  values can be written as:

$$y = \text{predicted value} + \text{prediction error}$$

or

$$y = \text{predicted value} + \text{residual}$$

or

$$y = \hat{y} + \text{residual}$$

For each individual,  $\text{residual} = y - \hat{y}$

$$= \text{"Observed value} - \text{predicted value"}$$

Example: Suppose the average of a guy's parents' heights is 66 inches, and he is 69 inches tall.

Observed data:  $x = 66$  inches,  $y = 69$  inches.

$$\text{Predicted height: } \hat{y} = 16.3 + 0.809(66) = 69.7 \text{ inches}$$

$$\text{Residual} = 69 - 69.7 = -0.7 \text{ inches}$$

The person is just 0.7 inches *shorter* than predicted.

$$y = \text{predicted value} + \text{residual}$$

$$69 = 69.7 + (-0.7)$$

Each  $y$  value in the original dataset can be written this way.

## DEFINING THE "BEST" LINE

**Basic idea:** Minimize how far off we are when we use the line to *predict*  $y$  by comparing to *actual*  $y$ .

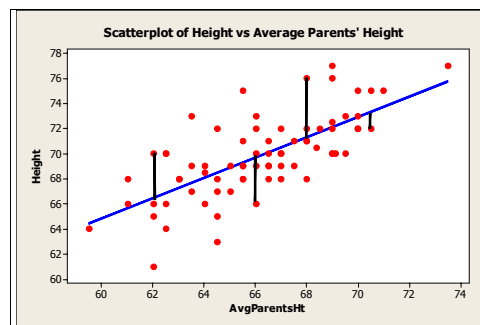
For each individual in the data

$$\text{"error"} = \text{"residual"} = y - \hat{y} = \text{observed } y - \text{predicted } y$$

**Definition:** The *least squares regression line* is the line that minimizes the sum of the squared residuals for all points in the dataset. The *sum of squared errors* = SSE is that minimum sum.

See picture on next page.

## ILLUSTRATING THE LEAST SQUARES LINE



SSE = 376.9 (average of about 5.16 per person, or about 2.25 inches when take square root)

**Example 1:**

This picture shows the residuals for 4 of the individuals. The blue line comes closer to *all of the points* than any other line, where "close" is defined by SSE =

$$\sum_{\text{all values}} \text{residual}^2$$

R Commander does the work for you!

**Statistics -> Fit models -> Linear regression**

Then highlight the variables you want (response = y and explanatory = x) in the popup box. The results look like this:

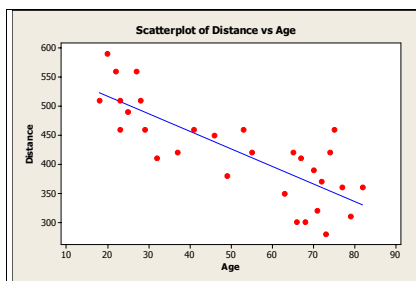
```
Call:
lm(formula = Height ~ AvgHt, data = UCDavisM)

Residuals:
    Min       1Q   Median       3Q      Max
-5.4768 -1.3305 -0.2858  1.2427  5.7142

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept) 16.3001     6.3188   2.580  0.0120 *
AvgHt        0.8089     0.0954   8.479 2.16e-12 ***
```

**EXAMPLE OF A NEGATIVE ASSOCIATION**

- A study was done to see if the distance at which drivers could read a highway sign at night changes with age.
- Data consist of n = 30 (x, y) pairs where x = Age and y = Distance at which the sign could first be read (in feet).



The regression equation is

$$\hat{y} = 577 - 3x$$

Notice *negative slope*

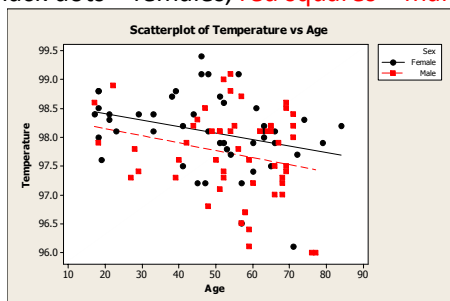
Ex:  $577 - 3(20) = 577 - 60 = 517$

Age	Pred. distance
20 years	517 feet
50 years	427 feet
80 years	337 feet

**Interpretation of slope and intercept?**

**Separating Groups in Regression and Correlation**

Example: Body temperature for 100 adults aged 17 to 84  
 Black dots = females, red squares = males



Note females slightly higher at all ages. Regression equations:

Males:  $\hat{y} = 98.4 - .0126(age)$

Females:  $\hat{y} = 98.6 - .0112(age)$

**Not easy to find the best line by eye!**

Applets:

[http://onlinestatbook.com/stat\\_sim/reg\\_by\\_eye/index.html](http://onlinestatbook.com/stat_sim/reg_by_eye/index.html)

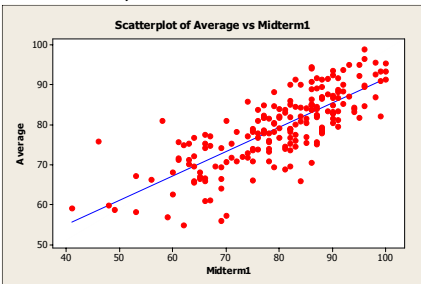
<http://www.rossmanchance.com/applets/Reg/index.html>

<http://illuminations.nctm.org/LessonDetail.aspx?ID=L455>

(More with this last one next time – influence of outliers.)

#### Example 4: Predicting final average from midterm

- Relationship is linear, positive association
- Regression equation:  $\hat{y} = 30.6 + 0.6x$
- For instance: midterm = 80, predicted avg = 78.6  
 $x = 100, \hat{y} = 90.6$   
 $x = 50, \hat{y} = 60.6$



#### SUMMARY OF WHAT YOU SHOULD KNOW

1. How to read a scatterplot to look for
  - a. Linear trend or not (curved, etc.)
  - b. positive or negative association (or neither)
  - c. strength of relationship (how close points are to line; more on this next time)
  - d. outliers
2. Given a regression equation,
  - a. Use it to *predict*  $y$  and *estimate*  $y$  for *given*  $x$  (useful when using the equation in the future,  $x$  known,  $y$  not)
  - b. Interpret slope and intercept
  - c. Find residual for a given individual, when given  $x$  and  $y$  for that individual.

#### Homework (due **Wed**, Jan 23):

Ch. 2: # 104\*

Ch. 3: #24 (you do *not* need the original data; note typo in some printings of the book – intercept is 126, not 1.26)

Ch. 3: #98\* (Data on class website; worth double points)

\*Use R Commander for 2.104 and 3.98

Data for 3.98 link on website (in file called oldfaithful.txt)