

LINEAR REGRESSION CHAP 7

The object [of statistics] is to discover methods of condensing information concerning large groups of allied facts into brief and compendious expressions suitable for discussion.

Francis Galton (1822-1911)

#### Relationship between 2 Quantitative Variables

In the mid-20<sup>th</sup> century Dr. Mildred Trotter, a forensic anthropologist, determined relationships between dimensions of various bones and a person's height. The relationships she found are still in use today in the effort to identify missing persons based on skeletal remains (think C.S.I).

One relationship compares the length of the femur to the person's height.

Measure the length of your femur along the outside of your leg (in inches). Also measure your height (also in inches) if you don't already know it. Record both at the front of the room.



#### Scatterplot of Height vs. Femur



Femur length is the **predictor** variable and height is the **response** variable.

Association between variables: **Direction** – positive **Form** – linear **Strength** – moderate (correlation = 0.591)



A **linear model** is a line that passes "through the middle" of the plotted points.



The line we will create is called the **"Least Squares Regression Line."** It minimizes something called the **squared residuals** (**residuals** are also called **errors**).

#### Linear Model

A **residual** is the difference between the observed response and what the model predicts for the response:



residual = observed y – predicted y

residual =  $y - \hat{y}$  "*y*-hat"

Example: The indicated point is for a femur length of 15 in. The **observed height** is 62 in., and the model **predicts a height** of 64.4 in. The residual is:

residual = 62 - 64.4

$$= -2.4$$

The model <u>over-predicts</u> height for this data point by 2.4 in.

# Least Squares Regression Equation

#### The form of the **linear model** is:



Note 1 – These formulas are on your AP Formula Sheet Note 2 – See optional Math Box on p. 180 to "look under the hood"

## Least Squares Regression Equation

If we standardized the data, the standard deviation *s* would be 1:



Since r can never be larger than 1, the numerator indicates that the predicted y tends to be closer to the mean (in standard deviations) the its corresponding x. This property is called **regression to the mean**.



So, when *x* increases by 1SD, *y* increases by only *r* times 1SD. This is called *regression to the mean*.



### Back to Our Data ...

8



Variable	Mean	Std. dev.
femur	16.10	1.71
height	65.87	4.10
	<i>r</i> = 0.591	

Compute the Least Squares Regression Line Model:

$$b_1 = \frac{rs_y}{s_x} \qquad b_0 = \overline{y} - b_1 \overline{x}$$
$$\hat{y} = b_0 + b_1 x$$

#### Back to Our Data ...



Calculate the Least Squares Regression Model using your calculator (see p. 193 for TI Tips).

## **Regression on the Computer**

10

#### Simple linear regression results:

Dependent Variable: height Independent Variable: femur height = 43.06999 + 1.4157449 femur Sample size: 34 R (correlation coefficient) = 0.59071902 R-sq = 0.34894896 Estimate of error standard deviation: 3.3562315

Parameter estimates:

Parameter	Estimate	Std. Err.	Alternative	DF	T-Stat	P-value
Intercept	43.06999	5.5348135	≠ 0	32	7.7816516	<0.0001
Slope	1.4157449	0.3418507	≠ <b>0</b>	32	4.1414119	0.0002

#### Interpret the Slope and Intercept

#### height = 43.07 + 1.42 (femur)

Predicted height if femur length is 0 inches! For every 1-inch increase in femur length we <u>expect</u> <u>an increase of about</u> 1.42 inches in height.



12

#### height = 43.07 + 1.42 (femur)

#### Predict a person's height if the femur length is: 1.18.5 inches 2.25 inches



A prediction made outside the range of the predictor (explanatory) data is called **extrapolation**.

Proceed with caution. The model may not hold true beyond the observed data.



Make a **residual plot**. If there is <u>no apparent pattern</u> to the residuals, then the Least Squares Regression (LSR) Model is reasonable.



No pattern means the model captured all of the "interesting" relationship between the variables.

The mean of the residuals will always be O, and the smaller the residuals from the model, the closer the model fits the data.

Software usually plots the residuals against the predicted values (heights in this case). When there is only one explanatory variable, the result is the same when the residuals are plotted against the explanatory variable (other than scale).



Residuals vs. Predicted Values



Residuals vs. Explanatory Variable

See TI Tips on p. 193 for instructions on creating a residual plot on your calculator.

Note – Every time you have your calculator create a LSR Model, a residuals list (RESID) is created. **Every <u>time</u>**. The most common mistake made by AP Stat students with regards to this is to forget to rerun the LSR command when working with a new set of data.



No left-over underlying pattern, so the linear model is reasonable.

Curved pattern left-over after regression, so linear model did not capture entire variable relationship.

Residuals increase in size with an increase in *x*, so linear model did not capture entire variable relationship.

#### <u>After</u> Deciding On The Model ...

A statistic we will call "R-squared\*" tells us something about how well the model captures the relationship in the <u>variability</u> present in the association between the predictor and response variables.

It goes something like this ...

\*It has a technical name, but no one really ever uses it.



What if I told you a new student was joining the class next week, and I want you to predict his/her height? Further, the only model you have is based on the mean height of students in the AP Stats class (y-bar = 65.9 in.).

What height is your best guess?

That's right, the mean height of 65.9 in.







The mean height is not a satisfactory guess ... we need more information.



#### **R-squared**

The ratios of these two errors (summed for all data points) is *r*-squared.





R-squared is the fraction of variation in height accounted for by the regression model with femur length. Specifically ...

... About 34.9% of the <u>variation</u> in <u>height</u> can be <u>attributed</u> to the linear regression model with <u>femur</u> length.



R-squared will always be between 0% and 100%.

R-squared describes the <u>percent of variation in the</u> <u>response variable</u> that can be explained (or attributed to) the regression with the explanatory variable.

In the case of heights predicted from femur lengths, if only 34.9% of the observed variation of heights can be attributed to differences in femur length, there must be other factors that contribute to differences in heights.

What might some of these factors be?

r vs. R-squared

Correlation *r* describes the strength and direction of a linear association between two quantitative variables.

R-squared describes the amount of variation observed in the response variable accounted for through the linear model with the predictor variable.

### What if the Variables are Reversed?

Because the Least Squares Regression Line <u>minimizes the vertical</u> <u>distances to the observed data</u>, you can't ever (no, not ever ... don't even think it!) use a regression equation to predict the explanatory variable from the response variable.



height = 43.07 + 1.42 (femur)



femur = -0.13 + 0.25 (height)



- •If the scatterplot doesn't look straight, don't fit a line!
- •Watch out for outliers (in all directions).
- •Do use a residual plot to see if the linear model looks reasonable.
- •Don't use R-squared to determine if a linear model is reasonable. R-squared only tells you something about the amount of variation in the response variable that can be traced back to the <u>a particular predictor</u> <u>variable</u>.
- •Association does not imply causation. A regression model with high r and R-squared in no way means the predictor variable causes the response variable.
- •Extrapolate with caution.



#### "I think the Government is behind this plot.

www.causeweb.org John Landers