# Simple Linear Regression

Lecture 1

August 28, 2007

Applied Correlation and Regression Methods in Education

# Today's Lecture

- Brief review of basic concepts:
  - ◆ Univariate statistical terminology.
  - ◆ Bivariate statistical terminology.
  - ◆ Simple linear regression.
- Simple linear regression.
- Partitioning the sum of squares.
- Tests of significance.
- Assumptions.
- Regression diagnostics (introduction).
- When $X$ is random.
- How to do today's topics in SPSS.

# Univariate statistics

A single random variable can be described by several characteristics of it's distribution:

■ Measures of central tendency (or location).
- ◆ Mean (a.k.a. expected value).
- ◆ Median (point where 50% of distribution falls above or below)
- ◆ Mode (most likely value)

■ Measures of spread.
- ◆ Variance and/or standard deviation.
- ◆ Range(s).

■ Measures that most normal people don't think about.
- ◆ Skewness - a measure of the symmetry of a distribution.
- ◆ Kurtosis - a measure of the peakedness of a distribution.

# Our Two Favorites

The mean:

$$\bar{X} = \sum_{i=1}^{N} \frac{X_i}{N}$$

The variance:

$$s_x^2 = \frac{\sum_{i=1}^{N}(X - \bar{X})^2}{N-1} = \frac{\sum_{i=1}^{N} X^2 - \frac{(\Sigma_{i=1}^{N} X)^2}{N}}{N-1}$$

# Additional Notes

■ Units associated with variance are squared (e.g., a variance of a distribution of height is given in $feet^2$).

■ More understandable is something where units aren't squared: standard deviation:

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

■ We divide the sum of squares by $N - 1$ to obtain an unbiased estimate of the variance, or by $N$ to obtain a "maximum likelihood" estimate of the variance (if a variable is normally distributed).

■ Variance/Standard Deviation are positive numbers $[0, \infty)$.

# Bivariate Distributions

Each of these images shows the distribution of a pair of variables:

# Bivariate Descriptive Statistics

■ Covariance: the joint covariation of two sets of variables from their respective means.

$$s_{xy} = \frac{\sum_{i=1}^{N}(X_i - \bar{X})(Y_i - \bar{Y})}{N-1} = \frac{\sum_{i=1}^{N}XY - \frac{(\Sigma_{i=1}^{N}X)(\Sigma_{i=1}^{N}Y)}{N}}{N-1}$$

◆ Units are products of units used to create measure (e.g., the covariance of a distribution of height and weight might be reported in *foot-pounds*).
◆ Covariance of a variable and itself is the variance.
◆ Covariance can take the value of any real number $(-\infty, \infty)$.

# Bivariate Descriptive Statistics

■ Because products of units can be hard to interpret, consider the correlation:

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

◆ Correlation is "unitless."
◆ Range of correlation is from $[-1, 1]$.

# A Single Hand Calculated Example

Consider some scores of students on Algebra and Geometry tests:

| Person | Algebra $(X)$ | Geometry $(Y)$ | $X^2$ | $Y^2$ | $XY$ |
|--------|---------------|----------------|-------|-------|------|
| A | 11 | 11 | 121 | 121 | 121 |
| B | 13 | 10 | 169 | 100 | 130 |
| C | 18 | 17 | 324 | 289 | 306 |
| D | 12 | 13 | 144 | 169 | 156 |
| E | 16 | 14 | 256 | 196 | 224 |
| N=5 | 70 | 65 | 1014 | 875 | 937 |

$$\sum X = 70 \qquad \sum Y = 65 \qquad \sum XY = 937$$

$$\sum X^2 = 1014 \qquad \sum Y^2 = 875$$

$$\bar{X} = 14 \qquad \bar{Y} = 13$$

$$s_X^2 = 8.5 \qquad s_Y^2 = 7.5$$

$$s_{XY} = 6.75 \qquad r_{XY} = 0.845$$

# A Single Hand Calculated Example

$$r_{XY} = 0.845$$

# Previously...
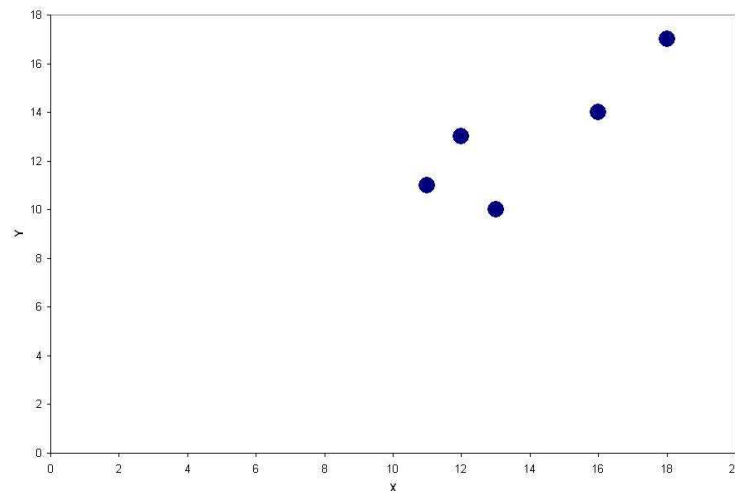
Recall the hand calculated sample problem with the following scatterplot:



$$r_{XY} = 0.845$$

- Correlation shows the degree of association between $X$ and $Y$...
  - ◆ ...but does so without scale.
- What if we wanted to predict $Y$ given a value of $X$?

# The Basics

■ Assume (for now) $X$ is fixed at pre-determined levels in an experiment - independent variable.

◆ For example, we have an experiment where subjects are given $X$ cups of coffee.

◆ Jon the Mad Scientist decides that subjects should be randomly assigned to a group drinking either 1, 2, 3, 4,or 5 cups of coffee.

■ The we want to estimate the *linear* effect of the independent variable $X$ on the dependent variable $Y$.

◆ For our example, we want to see how coffee drinking affects blood pressure.

◆ Blood pressure = $Y$ = dependent variable.

# The Basics

■ The linear regression model (for observation $i = 1, \ldots, N$):

$$Y_i = \alpha + \beta X_i + \epsilon_i$$

■ Don't be confused by the Greek alphabet, this is simply the equation for a line ($y = a + bx$).

■ $\alpha$ is the mean of the population when $X$ is zero...the $Y$ intercept.

■ $\beta$ is the slope of the line, the amount of increase in $Y$ brought about by a unit increase ($X' = X + 1$) in $X$.

■ $\epsilon_i$ is the random error, specific to each observation.

# Parameter Estimates

■ Chapter 2 parameterizes the linear regression model as:

$$Y = a + bX + e$$

■ To find estimates for $a$ and $b$ consider several possible choices:

◆ So that $\sum_{i}^{N} |e|$ is minimized.

◆ So that $\sum_{i}^{N} e^2$ is minimized.

◆ From some guy in the hallway.

# And The Winner Is...

■ Finding $a$ and $b$ that minimize:

$$\sum_i^N e^2$$

■ Using calculus, these happen to be:

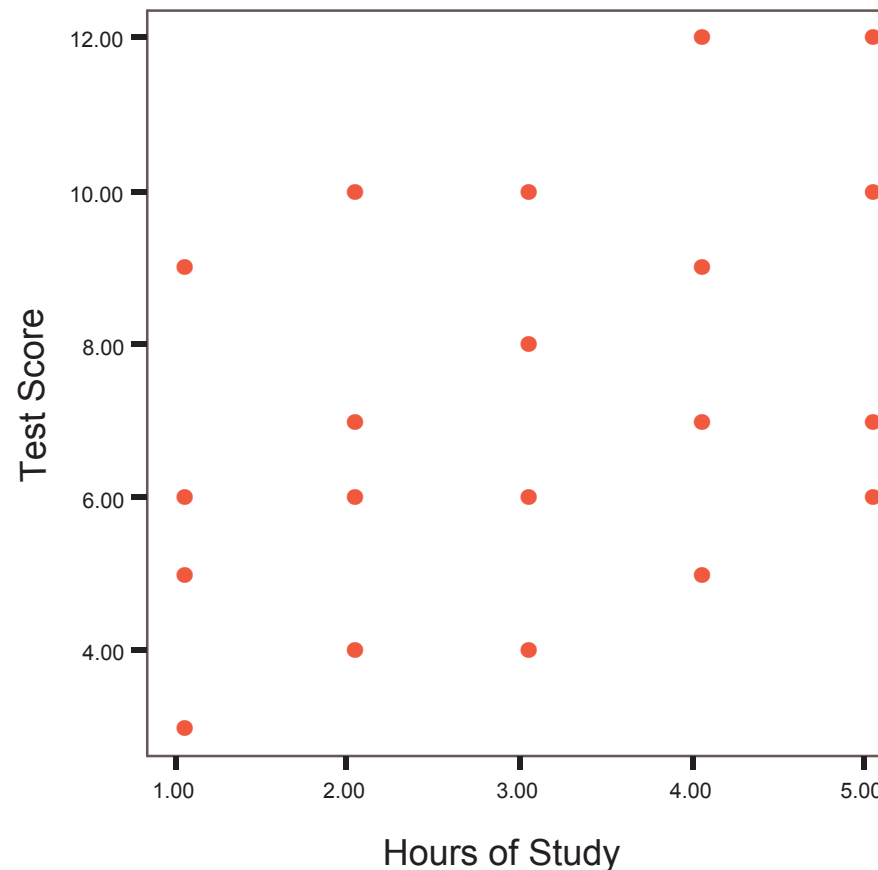$$b = \frac{\sum(X - \bar{X})(Y - \bar{Y})}{\sum(X - \bar{X})^2} = r_{xy}\frac{s_y}{s_x} = \frac{\sum xy}{\sum x^2}$$

■ And:

$$a = \bar{Y} - b\bar{X}$$

■ LS estimates are considered BLUE: Best Linear Unbiased Estimators.

# An Example of Simple Linear Regression

■ Table 2.1 in Pedhazur lists data from an experiment where $X$ was the number of hours given for study, and $Y$ is the score on a test.

# Example (continued)

■ We can tell that:
  ◆ $\sum(X - \bar{X})^2 = \sum x^2 = 40$
  ◆ $\sum(X - \bar{X})(Y - \bar{Y}) = \sum xy = 30$
  ◆ $\bar{X} = 3.0$
  ◆ $\bar{Y} = 7.3$

■ So:

$$b = \frac{\sum xy}{\sum x^2} = \frac{30}{40} = 0.75$$

$$a = \bar{Y} - b\bar{X} = 7.3 - (0.75 \times 3.0) = 5.05$$

■ Given these estimates, the linear regression line is given by:

$$Y' = 5.05 + 0.75X$$

# Example (continued)

Test Score = 5.05 + 0.75 * X
R−Square = 0.17

# Properties of the LS Estimates

■ Notice:

$$Y' = a + bX = (\bar{Y} - b\bar{X}) + bX = \bar{Y} + b(X - \bar{X})$$

■ This means when $X = \bar{X}$, $Y' = \bar{Y}$.

■ Furthermore, this implies that for every simple linear regression, the point $(\bar{X}, \bar{Y})$ will always fall on the regression line.

■ Furthermore, if $b = 0$ (no relation between $X$ and $Y$), the "best" guess (in a LS sense) for a value of $Y$ would be $\bar{Y}$.

# How "Good" Is Our Line?

- ■ You've fit the model...you have your estimates...now what?
- ■ We need to determine how <u>well</u> $Y$ is predicted by $X$.
- ■ Specifically, how much variability in $Y$ can be accounted for by the linear regression where $X$ is a predictor?
- ■ A math adventure:

$$Y = \bar{Y} + (Y' - \bar{Y}) + (Y - Y')$$

$$Y - \bar{Y} = (Y' - \bar{Y}) + (Y - Y')$$

$$\sum (Y - \bar{Y})^2 = \sum (Y' - \bar{Y})^2 + \sum (Y - Y')^2$$

$$\sum (Y - \bar{Y})^2 = ss_{reg} + ss_{res}$$

# Sums of Squares

■ The Sums of Squares for our dependent variable, $Y$, is denoted in Pedhazur by $\sum y^2$.

■ As shown before, in a regression, these can be partitioned into two components:

◆ Part due to the regression line: $ss_{reg}$

◆ Part due to error: $ss_{res}$

■ Dividing both components by $\sum y^2$ gives the proportion of sum of squares due accounted for by the regression and the proportion accounted for by error.

# Example Computation

■ In our studying v. test score example from before, we know that:

◆ $\sum y^2 = 130.2$

◆ $b = 0.75$

◆ $\sum xy = 30$

■ To compute $ss_{reg}$, one of several formulas can be used:

$$ss_{reg} = \frac{(\sum xy)^2}{\sum x^2} = b \sum xy = b^2 \sum x^2 = 0.75 \times 30 = 22.5$$

■ $ss_{res} = \sum y^2 - ss_{reg} = 130.2 - 22.5 = 107.7$

# Example Continued

■ Dividing both by $\sum y^2$ gives:

◆ Proportion explained by regression: 0.173

◆ Proportion explained by error: 0.827

# Variance Accounted For

■ Instead of approaching fit via Sums of Squares, one could look at it as a spinoff of the Pearson correlation:

$$r_{xy}^2 = \frac{\left(\sum xy\right)^2}{\sum x^2 \sum y^2}$$

$$ss_{reg} = r_{xy}^2 \sum y^2$$

# Variance Accounted For - Example

■ From our studying v. test score example from before, we know that:

$$r_{xy} = 0.416$$

■ So, $r_{xy}^2 = 0.173$...

■ Which is the percent of variance accounted for by the regression.

■ Therefore, $1 - r_{xy}^2 = 0.827$; the percent of variance due to error, or unexplained variance.

# Statistical Significance

■ Statistical significance <span style="color:red">is</span> literally the likelihood of the null hypothesis in a hypothesis test.

$$H_0: \quad b = 0$$
$$H_1: \quad b \neq 0$$

■ Statistical significance <span style="color:red">is not</span> related to the meaningfulness of the result.

◆ Take $b = 0.0001$.

◆ With $se(b) = 0.0000001$.

◆ $b$ is "statistically significant" but in reality, there is a negligible relationship between $X$ and $Y$.

# Hypothesis Tests in Simple Linear Regression

■ There are several types of hypothesis tests that can be used in simple linear regression.

■ Most of these tests are all related: a test for a non-zero relationship between two variables, $X$ and $Y$.

■ Today, we have covered three different measures of association:

◆ $b$ - $ss_{reg}$

◆ $r^2$

◆ $r_{xy}$

# Hypothesis Tests in Simple Linear Regression

- Each of the following tests evolve statistically from placing distributional assumptions on the disturbance term.

$$Y = a + bX + e$$

- Typically, we say that errors follow a normal distribution:

$$e \sim N(0, \sigma_e^2)$$

- By placing distributional assumptions
  - ◆ Sampling distributions are defined.
  - ◆ Hypothesis tests are enabled.
  - ◆ Assumptions must be recognized and checked for validity.

# Testing the Entire Regression

■ A test of all meaningful regression parameters is as follows:

$$H_0 : \quad b = 0$$
$$H_1 : \quad b \neq 0$$

■ Because we only have one predictor variable, this is very straight forward.

■ The test statistic is $F$-distributed, and is computed by:

$$F = \frac{\frac{ss_{reg}}{df_{reg}}}{\frac{ss_{res}}{df_{res}}} = \frac{\frac{ss_{reg}}{k}}{\frac{ss_{res}}{N-k-1}}$$

■ F-table gives significance for a given Type I error rate $\alpha$, OR USE THE COMPUTER.

# Testing the Entire Regression - Example

- From our example, the test of the entire regression is:

$$H_0: \quad b = 0$$
$$H_1: \quad b \neq 0$$

$$F = \frac{\frac{ss_{reg}}{df_{reg}}}{\frac{ss_{res}}{df_{res}}} = \frac{\frac{ss_{reg}}{k}}{\frac{ss_{res}}{N-k-1}} = \frac{\frac{22.5}{1}}{\frac{107.7}{18}} = 3.76$$

- Typing "=fdist(3.76,1,18)" into MS Excel gives a p-value of 0.0683.

- We could infer that there is no statistical evidence to suggest $b$ is different from $0$.

- This means there is no statistical relationship between hours of study and the score on a test.

# Testing the VAC

■ Because of the relationship of $r_{xy}^2$ with $ss_{reg}$ the previous test could be phrased as:

$$F = \frac{\frac{r_{xy}^2}{df_{reg}}}{\frac{1-r_{xy}^2}{df_{res}}} = \frac{\frac{r_{xy}^2}{k}}{\frac{1-r_{xy}^2}{N-k-1}} = \frac{\frac{0.1728}{1}}{\frac{1-0.1728}{18}} = 3.76$$

■ Notice this is identical to the $F$ we received before. Therefore our conclusions are the same...

# Testing the Regression Coefficient

■ Finally, the regression coefficient between $X$ and $Y$ can tested for difference from zero:

$$H_0 : \quad b = 0$$

$$H_1 : \quad b \neq 0$$

■ The variance of $Y$ conditional on $X$ is given by:

$$s_{y.x}^2 = \frac{\sum (Y - Y')}{N - k - 1} = \frac{ss_{res}}{N - k - 1}$$

■ The standard error of the regression coefficient, $s_b$, is based on the conditional variance of $Y|X$:

$$s_b = \sqrt{\frac{s_{y.x}^2}{\sum x^2}}$$

# Testing the Regression Coefficient

■ Finally, the test statistic for the following hypothesis test is:

$$H_0: \quad b = 0$$
$$H_1: \quad b \neq 0$$

$$t = \frac{b}{s_b}$$

■ This statistic is $t$ distributed, with N-k-1 degrees of freedom (which is the denominator of $s_{y.x}^2$).

■ From our example, the test of the regression coefficient is:

$$H_0: \quad b = 0$$
$$H_1: \quad b \neq 0$$

$$s_{y.x}^2 = \frac{\sum(Y - Y')}{N - k - 1} = \frac{ss_{res}}{N - k - 1} = \frac{107.7}{18} = 5.983$$

$$\sum x^2 = 40$$

$$t = \frac{b}{\sqrt{\frac{s_{y.x}^2}{\sum x^2}}} = \frac{0.75}{\sqrt{\frac{5.983}{40}}} = 1.94$$

■ Typing in "=tdist(1.94,18,2)" (2 for a two-tailed test), gives a p-value of 0.0683 - the same as before.

■ Recall $F = t^2$ from somewhere before...

# Testing the Regression Coefficient

- In our example, we tested against zero.
- In reality, any value can be tested against:

$$H_0: \quad b = B$$
$$H_1: \quad b \neq B$$

$$t = \frac{b - B}{\sqrt{\dfrac{s_{y.x}^2}{\sum x^2}}}$$

# Key Assumptions Behind Simple Linear Regression

- ■ Basic assumptions:
  - ◆ $X$ is fixed - measured without error.
  - ◆ The mean of $Y|X$ has a linear relationship with $X$

- ■ Error assumptions:
  - ◆ Mean of errors is zero.
  - ◆ All errors independent.
  - ◆ The variance of the errors does not change as a function of $X$ - homoscedasticity.
  - ◆ All errors are uncorrelated with $X$.
  - ◆ Statistically speaking, errors are IID: $e \sim N(0, \sigma_e^2)$.

- An easy way to check many error based assumptions a visualization:
  - ◆ A plot of the standardized residual versus the predicted value.
- The next class will focus on other diagnostic measures.

# Curvilinear Example

What if we were to fit a linear regression to these data?

# Curvilinear Example

Fitting the simple linear regression model gives:



Z = 7.75 + −0.55 * X
R−Square = 0.04

Not good?

# Curvilinear Example

Let's take a look at a plot of the standardized residuals versus the predicted values:

# When $X$ is Random

- Commonly, we have no control over our predictor variable.
- Statistically speaking, having $X$ being random adds:
  - ◆ The potential for measurement error in $X$ - reliability.
  - ◆ Limitations to the causal inferences that can be drawn from a regression analysis.
- All test statistics and measures are used equivalently.
- We can predict $Y$ from $X$ or $X$ from $Y$.
- Correlations are more common place: common variance rather than causal influence.

# Wrap Up Example

Example taken from Weisberg (1985, p. 230):

"Perhaps the most famous geyser is Old Faithful, in Yellowstone National Park, Wyoming. The intervals of eruptions of Old Faithful range from about 30 to 90 minutes. Water shoots to heights generally over 35 meters, with eruptions lasting from 1 to 5.5 minutes.

"Data on Old Faithful has been collected for many years by ranger/naturalists in the park, using a stopwatch. The duration measurements have been rounded to the nearest 0.1 minute or 6 seconds, while intervals reported are to the nearest minute. The National Park Service uses $x$ (values of the duration of an eruption) to predict $y$ (the interval to the next eruption)."

# Old Faithful Postcard

# Old Faithful Scatterplot

# Descriptive Statistics

| Variable | N | Mean | SD |
|---|---|---|---|
| $x$ - Duration | 107 | 3.4607 | 1.0363 |
| $y$ - Interval | 107 | 71.000 | 12.967 |

$$r_{xy} = 0.858$$
$$r_{xy}^2 = 0.736$$
$$\sum x^2 = 113.84$$
$$\sum y^2 = 17,822.00$$
$$\sum xy = 1,222.70$$

# Regression Line

$$b = \frac{\sum xy}{\sum x^2} = \frac{1,222.70}{113.84} = 10.74$$

$$a = \bar{Y} - b\bar{X} = 71.0 - (10.74) \times 3.4607 = 33.84$$

$$Y' = 33.84 + 10.74X$$

# Regression Line

Interval to Next Eruption = 33.83 + 10.74 * X
R−Square = 0.74

# Hypothesis Tests

$$ss_{reg} = b \sum xy = 10.74 \times 1,222.70 = 13,132.99$$

$$ss_{res} = \sum y^2 - ss_{reg} = 17,822.0 - 13,132.99 = 4,689.01$$

$$F = \frac{\frac{ss_{reg}}{k}}{\frac{ss_{res}}{N-k-1}} = \frac{13,132.99/1}{4,689.01/(107-1-1)} = 294.084$$

# Residual Plot

# Data

# Data

# Plots from Today

Scatterplots were made with the interactive graphs: Graphs...Interactive...Scatterplot.

# Simple Linear Regression

Analyze...Regression...Linear.

# Saving Residuals

Analyze...Regression...Linear → Save Button.

# Final Thought

- Simple linear regression can be simple and straight forward.
- Need to understand difference between statistical significance and practical relevance.



- The tests of significance and model estimates are only valid if all assumptions are met.
- Regression is fairly robust to violations of assumptions, but...
- Usually when assumptions are violated, something else is wrong with the model.
- Listen to your data, they are trying to tell you something.

# Next Time

- Regression diagnostics - Chapter 3.
- More SPSS.
- More practice with fitting models.
- QUESTION: Do you have any data you are willing to share with class?