Non-Linear Relationships Section 5.4

Cathy Poliak, Ph.D. cathy@math.uh.edu Office in Fleming 11c

Department of Mathematics University of Houston

Lecture 14 - 2311

- Fill in all of the proper bubbles.
- Make sure your ID number is correct.
- Make sure the filled in circles are very dark.
- This is popper number 10.

Calcuating LSLR

• Given the statistics: \bar{x} , \bar{y} , s_x , s_y , and r. Use the following to calulate $\hat{y} = a + bx$.

$$b = r rac{s_y}{s_x}$$

 $a = ar{y} - bar{x}$

• The residual is the difference between the actual y-value and the predicted y-value

$$resid = y - \hat{y}$$

 The coefficient of determination is the sqared value of the correlation coefficent; R². This is the percent (fraction) of the variation of y that can be explained by the LSLR.

Example

The following data was collected comparing score on a measure of test anxiety and exam score:

| | | | | | | | | \frown | |
|-------------------------|----|----|----|----|----|----|------|----------|----|
| Measure of test anxiety | 23 | 14 | 14 | 0 | 7 | 20 | 20/ | 15 | 21 |
| Exam score | 43 | 59 | 48 | 77 | 50 | 52 | 46 \ | 51 | 51 |

We will use R to:

- Construct a scatterplot.
- Find the LSRL and fit it to the scatterplot.
- Find r and r².
- Does there appear to be a linear relationship between the two variables? Based on wha you found, would you characterized the relationship as positive or negative? Strong or weak?
- Draw the residual plot.
- What does the residual plot reveal?
- Answer the question: Is it reasonable to conclude that test anxiety caused poor exam performance?

> anxiety=c(23,14,14,0,7,20,20,15,21) > exam=c(43,59,48,77,50,52,46,51,51) > plot(anxiety,exam) ح عدمد الإربان > lm(exam~anxiety) ح لجلار

Call: Im(formula = exam ~ anxiety)

Coefficients: (Intercept) anxiety 68.838 -1.064 & y = 68838 - 1.064 x X=15 y = 68.838 - 1.064 (15) = 52.878 Residual for X=15 : y = 51 residual = 51 - 52.878 = -1.818

Residual Plot

> plot(anxiety,resid(Im(exam~anxiety))) > abline(0,0)



Caution about correlation and regression

- An outlier is an observation that lies outside the overall pattern of the other observations. Points that are outliers in the y direction of sacatterplot have large regression residuals.
- Influential Observations: a data point whose removal causes the regression equation (and the line) to change considerably. Points that are extreme in the *x* direction of a scatterplot are often influential for the least-squares regression line.
- Data points do not appear to be scatter about a straight line: First draw a scatter diagram if the points do not appear to be scattered about a straight line, do not determine a regression line.

Caution about correlation and regression

- Extrapolation: Using the regression equation to make predictions for values of the predictor variable outside the range of the observed values of the predictor variable. (Grossly incorrect predictions can result from extrapolation.)
- Lurking Variable: A variable that is not among the explanatory or response variables in a study and yet may influence the interpretation of relationships among those variables.
- Association does not imply causation.

Popper 10 Questions $b = \frac{5}{5\pi}$ $a = \sqrt{5} - b \overline{x}$ $b = 6.9325(\frac{280843}{25.8790})$ a = 71.09133The following are the statistics to calculate a least squares line (LSLR) to predict the final grade (score) by the average quiz score. Score: mean = 71.09433, sd = 28.0842; Quiz: mean = 71.17697, sd = 25.87901; Correlation: r = 0.93257

1. Determine the LSLR to predict the final **score** based on the **quiz** score.

a.
$$\hat{y} = -0.939326 + 1.012x$$
b. $\hat{y} = 10.0826 + 0.8593x$ c. $\hat{y} = 70.913$ d. $\hat{y} = 1 + 0.85x$

2. A person has a **quiz** average for the semester at 86, has a final score of 93. Determine the residual for this student. x = 86 y = 9.3 $x = y - x_3 = 9.5$ y = 9.5 b. 86.1 c. -6.9 d. 0 x = 8.5

3. Determine the coefficient of determination.

a. 0.93257 b. 0.8967 c. 0.9215 d. 1.0852

Example

Can we predict the people per square mile based on year? The following is the population int he city of Houston from 1900 to 2010.

| Year | Population | | |
|------|------------|--|--|
| 1900 | 44,633 | | |
| 1910 | 78,800 | | |
| 1920 | 138,276 | | |
| 1930 | 292,352 | | |
| 1940 | 384,514 | | |
| 1950 | 596,163 | | |
| 1960 | 938,219 | | |
| 1970 | 1,233,505 | | |
| 1980 | 1,595,138 | | |
| 1990 | 1,631,766 | | |
| 2000 | 1,953,631 | | |
| 2010 | 2,100,263 | | |

Give the scatterplot, LSLR, R² and the residual plot of this data

Cathy Poliak, Ph.D. cathy@math.uh.edu Offic

Section 5.4

R code

```
year=c(1900,1910,1920,1930,1940,1950,1960,1970,1980,1990,2000,2010)
population=c(44633,78800,138276,292352,384514,596163,938219,1233505,
1595138,1631766,1953631,2100263)
plot(year, population)
pop.lm=lm(population~vear)
pop.lm
Call:
lm(formula = population ~ year)
Coefficients:
(Intercept)
                year
                                -39,649,04 + 20,749x
                             =
-39649140
               20749
cor(year, population)^2
[1] 0.9630567
plot(vear, resid(pop.lm))
```

Scatterplot



Cathy Poliak, Ph.D. cathy@math.uh.edu Offic

11/24

Residual Plot



- The coefficient of determination, R^2 is high.
- The scatter plot and residual plot shows a non-linear pattern.
- The least-square regression line might not be the "best fit" for this data.

- Many times a scatterplot reveals a curved pattern instead of a linear pattern.
- We can transform the data by changing the scale of the measurement that was used when the data was collected.
- In order to find a good model we may need to transform our x value or our y value or both.

Transform Year

In our example it appears like the relationship is $y = x^2$. If we transform the points (x, y) into (x, \sqrt{y}) we get the scatterplot. Plot(x, sqrt(y)



R code

```
tpop=sqrt(population)
plot(year,tpop)
tpop.lm=lm(tpop~year)
tpop.lm
```

```
Call:

lm(formula = tpop ~ year)

Coefficients:

(Intercept) year

-23267.19 12.34 => (5) = -23267.19 + (2.34 x
```

cor(year,tpop)^2
[1] 0.9854597
plot(year,resid(tpop.lm))
lines(year,rep(0,length(year)))

Residual Plot of Transformed Data



Calculating LSLR from Mosaic Date

Check mark the Mosaic Data in the programs list. The following is to create the scatterplot, LSLR, and the residual plot for predicting **temperature** in Celsius for water based on **time** in minutes that water was poured into a mug. Let X = time and Y = temp.

$$x = CoolingWater$timey = CoolingWater$tempplot(x,y)cor(x,y)^2[1] 0.778089 G R2 = 77.87water.lm=lm(y~x) Gwater.lmCall:lm(formula = y ~ x)Coefficients:(Intercept) x64.2766 -0.2164 $\frac{3}{3}$ = **G**4.2700 - **O**.2164
UNIVERSITY of HOUSION
DEVICE NOT HOUSION$$

Section 5.4

Scatterplot

Using this scatterplot, could we say that the LSLR is a good model to predict y?



NIVERSITY of HOUSTON

Residual Plot



Transformation

A good transformation for this is $y = \frac{1}{\sqrt{x}}$. We will transform the points (x, y) to $(x, \frac{1}{y^2})$. Plot (x, \sqrt{y}, \sqrt{y})



See your textbook for other examples of transformations.

UNIVERSITY of HOUSTON DEPARTMENT OF MATHEMATICS

Section 5.4

y2=1/y^2 twater.lm=lm(y2~x) plot(x,resid(twater.lm),ylim=c(-1,1)) cor(x,y2)^2 [1] 0.9986831 (= R²

Residual Plot



Other Possible Transfers



UNIVERSITY of HOUSTON DEPARTMENT OF MATHEMATICS

Cathy Poliak, Ph.D. cathy@math.uh.edu Offic

Section 5.4