CORRELATION & REGRESSION

CHAPTER 10

CORRELATION

10-2

BIVARIATE DATA

- DATA ON EACH OF TWO VARIABLES, WHERE EACH VALUE OF ONE OF THE VARIABLES IS PAIRED WITH A VALUE OF THE OTHER VARIABLE.
- APPENDIX B DATA SET 14 (P754)

| Car weight (lbs) | Highway MPG |
|------------------|-------------|
| 2560 | 34 |
| 2895 | 33 |
| 3320 | 28 |
| 3465 | 28 |
| 3835 | 26 |
| 4180 | 24 |

SCATTERPLOT

- USE THE STATISTICAL SOFTWARE
 ON YOUR CALCULATOR
- LIST EDITOR
- STAT PLOT
- ZOOM STAT







CORRELATION

- **CORRELATION** EXISTS BETWEEN TWO VARIABLE WHEN THE VALUES OF ONE VARIABLE ARE ASSOCIATED WITH THE VALUES OF THE OTHER VARIABLE.
- LINEAR CORRELATION EXISTS BETWEEN TWO VARIABLES WHEN THERE IS A CORRELATION AND THE PLOTTED POINTS OF PAIRED DATA RESULT IN A PATTERN THAT CAN BE APPROXIMATED BY A STRAIGHT LINE.

SCATTERPLOTS LINEAR CORRELATION COEFFICIENT, r, MEASURES THE STRENGTH OF THE LINEAR RELATIONSHIP



 $-1 \le r \le 1$

MEASURING CORRELATION COEFFICIENT

• 3 METHODS

• SUMS OF THE SQUARES & SQUARES OF THE SUMS

• Z-SCORES

• TECHNOLOGY

CORRELATION COEFFICIENT FORMULA 10-1 (P499)

 $n(\sum xy) - (\sum x)(\sum y)$ $\sqrt{n(\sum x^2) - (\sum x)^2} \sqrt{n(\sum y^2) - (\sum y)^2}$

$$r = \frac{n(\sum xy) - (\sum x)(\sum y)}{\sqrt{n(\sum x^2) - (\sum x)^2}\sqrt{n(\sum y^2) - (\sum y)^2}}$$

6(572325) - (20245)(173)

 $= \frac{1}{\sqrt{6(70066275) - (20245)^2}} \sqrt{6(5065) - (173)^2}$

| | x (car weight) | y (MPG) | x^2 | y^2 | ху |
|------|----------------|---------|----------|------|--------|
| | 2560 | 34 | 6553600 | 1156 | 87040 |
| | 2895 | 33 | 8381025 | 1089 | 95535 |
| | 3320 | 28 | 11022400 | 784 | 92960 |
| | 3465 | 28 | 12006225 | 784 | 97020 |
| | 3825 | 26 | 14630625 | 676 | 99450 |
| | 4180 | 24 | 17472400 | 576 | 100320 |
| sums | 20245 | 173 | 70066275 | 5065 | 572325 |

r = -0.982

CORRELATION COEFFICIENT FORMULA 10-2 (P499)

 $r = \frac{\sum(z_x z_y)}{n-1}$

 z_x denotes the z score for an individual sample value x z_y denotes the z score for an individual sample value y

 $r = \frac{\sum (z_x z_y)}{n-1}$ $r = \frac{-4.909}{5}$

| | x (car weight) | zx | y (MPG) | zy | zx*zy |
|--------|----------------|------------|-----------|--------------|------------|
| | 2560 | -1.3737361 | 34 | 1.318027211 | -1.8106215 |
| | 2895 | -0.8084936 | 33 | 1.06292517 | -0.8593682 |
| | 3320 | -0.0913949 | 28 | -0.212585034 | 0.0194292 |
| | 3465 | 0.1532623 | 28 | -0.212585034 | -0.0325813 |
| | 3825 | 0.760687 | 26 | -0.722789116 | -0.5498163 |
| | 4180 | 1.3596753 | 24 | -1.232993197 | -1.6764704 |
| mean | 3374.16667 | | 28.833333 | | -4.9094285 |
| st dev | 592.666 | | 3.92 | | |
| | | | | | |

r = -0.9818

USING TECHNOLOGY

- TI-84
- LIST EDITOR
- Stat
 - CALC
 - LINREGRESSION

TI-84 Plus





CORRELATION FOR LARGE DATA SETS

- APPENDIX B DATA SET 14 (P754)
- What if we use all of the data in data set 14
- *n* = 21

TI-84 Plus Texas Instruments



CORRELATION FOR LARGE DATA SETS



SMALL V. LARGE DATA SETS



n = 21r = -0.7927

SMALL V. LARGE DATA SETS

- Are we suggesting that the larger data set has a weaker Linear relationship?
- Use hypothesis testing to test the claim of a linear correlation between two variables

HYPOTHESIS TESTING FOR LINEAR CORRELATION

- r The sample correlation coefficient
- ρ (RHO) The population correlation coefficient
- $H_0: \rho = 0$ (there is no linear correlation)
- H_A : $\rho \neq 0$ (THERE IS A LINEAR CORRELATION)

THE TEST STATISTIC



CRITICAL VALUE CAN BE FOUND IN TABLE A-3
n-2 degrees of freedom

HYPOTHESIS TEST

 $H_0: \rho = 0$ $H_A: \rho \neq 0$ r = -0.7927n = 21df = 19 $\alpha = 0.05$ $t_{\alpha/2} = 2.093$

$$=\frac{r}{\sqrt{\frac{1-r^2}{n-2}}}$$

t

t =

r

$$= \frac{-0.7927}{\sqrt{\frac{1 - (-0.7927)^2}{21 - 2}}}$$

t = -5.66

ALTERNATIVE METHOD

- INSTEAD OF USING THE TEST STATISTIC, COMPARE THE SAMPLE CORRELATION COEFFICIENT TO THE CRITICAL VALUES OF THE PEARSON CORRELATION COEFFICIENT *r*
- TABLE A-6 ON PAGE 732
- r = -0.7927

To test $H_0: \rho = 0$ against $H_A: \rho \neq 0$, reject H_0 if the absolute value of r is greater than the critical value in the table.

• *n* = 21

. . .

12-52

. . 31.

161

 $H_0: \rho = 0$ $H_A: \rho \neq 0$

r = -0.7927n = 21

> 0.7927 > 0.444 $reject H_0$

| in a se | Coefficie | ent r |
|---------|----------------|----------------|
| n | $\alpha = .05$ | $\alpha = .0$ |
| 4 | .950 | .990 |
| 5 | .878 | .959 |
| 6 | .811 | .917 |
| . 7 | .754 | .875 |
| 8 | .707 | .834 |
| .9 | .666 | .798 |
| 10 | .632 | .765 |
| 11 | .602 | ° .735 |
| 12 | .576 | .708 |
| 13 | .553 | .684 |
| 14 | .532 | .661 |
| 15 | .514 | .641 |
| 16 | .497 | .823 |
| 17 | .482 | .606 |
| 18 | .468 | .590 |
| 19 | .456 | .575 |
| 20 | .444 | .561 |
| 25 | .396 | .505 |
| 30 | .361 | .463 |
| 35 | .335 | .430 |
| 40 | .312 | .402 |
| 45 | .294 | .378 |
| 50 | .279 | .361 |
| 60 | .254 | .330 |
| 70 | .236 | .305 |
| 80 | 000 | C Ob Bro Louis |

5

a mandatak.

1947 22223

Ne

984 1710

HOMEOWRK

• P513 #13-16, 24, 29

ADDITIONAL THINGS ABOUT CORRELATION

PROPERTIES OF r

- The value of r is always between -1 and 1 inclusive
- The value of r is not affected by the choice of x or y
- r only measures the strength of a linear relationship
- r is very sensitive to outliers

COMMON ERRORS

- CORRELATION DOES NOT MEAN CAUSATION
- ERRORS ARISE WHEN DATA IS BASED ON AVERAGES OR RATES
- *r* is only a test for linear Correlation. Just because paired data are not related linearly does not mean that they aren't related in some Other way.

r^2 THE PROPORTION OF VARIATION

r = -0.7927
r² = 0.6284
This means that 62.84% of the variation in the MPG can be explained by the linear relationship

Linear Regression

Regression line

The straight line that "best" fits a set of bivariate data
\$\hat{y} = b_0 + b_1 x\$
\$x\$ - predictor variable, independent variable
\$\hat{y}\$ - response variable, or dependent variable
\$b_0\$ - y-intercept
\$b_1\$ - slope

Linear Regression

Ouse the linear regression function in the calculator to find the equation for the line of best fit



Requirement check

• Data is assumed to be simple random sample

• A scatterplot suggests a linear pattern

OThere is a linear correlation OThere are no outliers



scatterplots





Linear regression

• Use three significant digits • $\hat{y} = 50.4 - 0.006x$



Storing the regression equation



Making predictions

 $\hat{y} = 50.4 - 0.006x$

- Find the expected Highway Fuel efficiency for a car weighing 5000 lbs.
- Find the fuel efficiency of the Hummer H2 weighing 6,400 lbs
- Find the fuel efficiency of the Smart Car weighing 1,500 lbs

• How much should a car weight to get 60 MPG??

What if there is no linear correlation!?

Other the value of \hat{y} is assumed to be \bar{y} for any predictor value of x

OWhy?

OBecause \bar{y} is the expected value of \hat{y}

Homeowrk

OP529 #13-16, 24, 29