

Linear Regression &
Correlation

$f(x) = y = m x + b$ (Algebra)

$\hat{y} = b_0 + b_1 x$ (Statistics)

$r = \left\{ \begin{array}{l} \text{strength \& direction} \\ \text{of line} \end{array} \right\}$ (chapter 3)

$\left\{ \begin{array}{l} \text{dependent variable} \\ \text{predicted variable} \end{array} \right\} = y = \hat{y} = f(x)$

$\left\{ \begin{array}{l} \text{independent variable} \\ \text{predictor variable} \end{array} \right\} = x$

$\left\{ \text{y-intercept} \right\} = b = b_0 = \text{INTERCEPT}$
Excel function

$\left\{ \text{slope} \right\} = \left\{ \begin{array}{l} \text{how much y moves} \\ \text{for every 1 unit} \\ \text{of x} \end{array} \right\} = m = b_1 = \text{SLOPE}$
Excel function

$r = \left\{ \begin{array}{l} \text{strength \& direction} \\ \text{of linear} \\ \text{equation} \end{array} \right\} = r = \text{PEARSON}$
Excel function

$r^2 = \left\{ \begin{array}{l} \text{coefficient of} \\ \text{determination} \\ \text{"goodness of fit"} \\ \text{of equation"} \end{array} \right\} = r^2 = \text{RSQ}$
Excel function

chapter 12 : Simple Linear Regression

P.2

What we already know about "Relationship between 2 quantitative variables" from Math & chapter 3:

① Independent variable (X)

predictor variable.

② Dependent variable (y or $f(x)$ or $E(y)$ or \hat{y})

Variable that is predicted or estimated

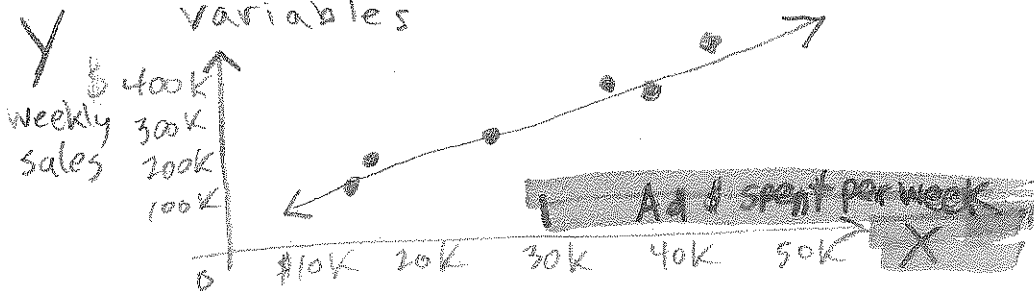
Example :

collected data from past :

X Ad \$ spent per week	Y weekly sales
14,000	97,000
27,000	185,000
40,000	260,000
17,000	143,000
34,000	270,000
43,000	398,000

③ Scatter Diagram/chart/Plot

Graphical technique to show relationship between 2 quantitative variables



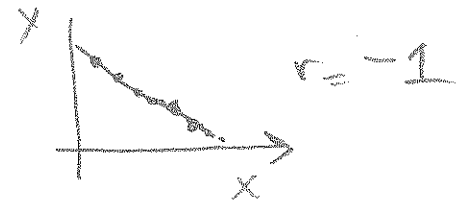
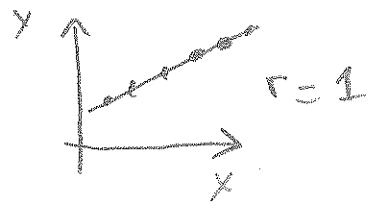
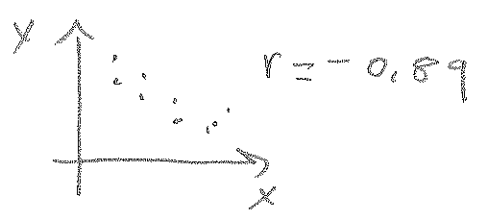
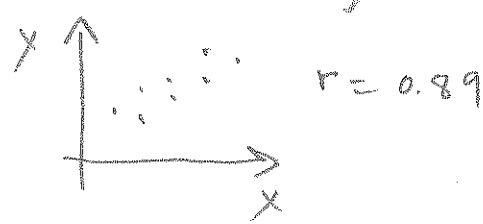
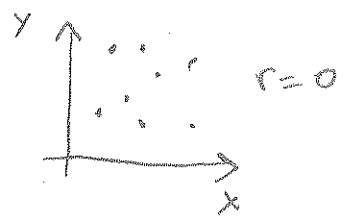
4

Coefficient of correlation (r) (Interval or Ratio Level Data)

P3

- ① Measure of the strength & direction of the linear relationship (between -1 and 1).
- ② 0 = No correlation.
- ③ 0.5 or -0.5 = Moderate correlation.
- ④ Near -1 or 1 = strong correlation.

Coefficient of correlation = r



5

Coefficient of determination (r^2)

A measure of the goodness of fit of the estimated regression equation. The proportion of the total variation in the dependent variable, y , that is explained, or accounted for, by the variation in the independent variable, x . It does NOT say anything about causation, that is, it does not say that x causes y .

$$\left\{ \begin{array}{l} \text{coefficient} \\ \text{of} \\ \text{correlation} \end{array} \right\} = r = \frac{\sum (X - \bar{X}) * (Y - \bar{Y})}{(n-1) * S_x * S_y}$$

How to calculate coefficient of correlation
Example 1

n = 6
n-1 = 5

$\bar{X} = 29,150$

$\bar{Y} = 225,500$

$S_x = 11,938.80$

$S_y = 107,596.93$

only 2 independent variables

X = particular value

\bar{X} = mean of Xs

Y = particular value

\bar{Y} = mean of Ys

n = count of observed pairs

S_x = standard deviation of Xs

S_y = standard deviations of Ys

X	Y	$X - \bar{X}$	$Y - \bar{Y}$	$(X - \bar{X}) * (Y - \bar{Y})$
14,000	97,000	$14,000 - 29,150 = -15,150$	$97,000 - 225,500 = -128,500$	$+ 1,946,775,000$
17,000	143,000	$17,000 - 29,150 = -12,150$	$143,000 - 225,500 = -82,500$	$+ 1,002,725,000$
27,000	185,000	$27,000 - 29,150 = -2,150$	$185,000 - 225,500 = -40,500$	$+ 87,075,000$
34,000	270,000	$34,000 - 29,150 = 4,850$	$270,000 - 225,500 = 44,500$	$+ 215,825,000$
39,900	260,000	$39,900 - 29,150 = 10,750$	$260,000 - 225,500 = 34,500$	$+ 370,875,000$
43,000	398,600	$43,000 - 29,150 = 13,850$	$398,600 - 225,500 = 173,100$	$+ 2,389,125,000$
		$\Sigma = 0$	$\Sigma = 0$	$\Sigma = 6,012,050,000$

$$\left\{ \begin{array}{l} \text{Strength} \\ \& \\ \text{direction} \end{array} \right\} = r = \frac{6,012,050,000}{5 * 11,938.8 * 107,596.93} = 0.936034653$$

What does .936 mean?
 strength = Strong
 direction = direct, as X increase, Y increase.

In Excel: \rightarrow PEARSON or CORREL function!
 \rightarrow =PEARSON(yrange, Xrange)

$$\left\{ \begin{array}{l} \text{Coefficient} \\ \text{of} \\ \text{Correlation} \end{array} \right\} = r = .936034653$$

P. 5

Ad dollars spent & sales earned Example 1:

.936034653 indicates that the relationship between Ad dollars spent & sales earned is very strong. In addition the relationship is direct, which means that as Ad dollars increase sales increase. However "strong" is not numerically precise. But coefficient of determination is numerically precise.

$$\left\{ \begin{array}{l} \text{Coefficient} \\ \text{of} \\ \text{Determination} \end{array} \right\} = r^2 = \left\{ \begin{array}{l} \text{Influence X} \\ \text{has on Y,} \\ \text{(Not causation)} \end{array} \right\} = .936034^2 = .876$$

Goodness of Fit

Ad dollars spent & sales earned Example 1:

We can say that 87.6% of the variation in y ($f(x)$) can be explained by the variation in x (not causation). Correlation does not mean

causation!! (from textbook: As population of

donkeys decreases, number of doctoral degrees increases.)

This is called "spurious correlations." conclusion: correlation is great for building models we can use to make predictions, but correlation does not mean causation.

⑥ Estimated Simple Linear Regression Equation (p. 6)

An equation that expresses the linear relationship between 2 variables

$$y = mx + b \quad (\text{Algebra})$$

$$\hat{y} = b_0 + b_1 X \quad (\text{Statistics})$$

$y = \hat{y}$ = predicted variable = dependent variable

$X = X$ = predictor variable = independent variable

$M = b_1$ = slope = How much y moves for 1 unit of X

$b = b_0$ = y -intercept

x_i = particular X

y_i = particular y

$\bar{x} = \bar{X}$ = sample mean

$\bar{Y} = \bar{Y}$ = sample mean

$$b_1 = \frac{\sum (x_i - \bar{x}) * (y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

$$b_0 = \bar{y} - b_1 * \bar{x}$$

$$\bar{X} = 29150 \quad n = 6$$

$$n - 1 = 5$$

(p.7)

$$\bar{Y} = 225,500$$

$$S_x = 11,938.80$$

$$S_y = 107,596.93$$

From bottom:

$$\hat{y} = b_0 + b_1 X$$

$$\hat{y} = -20,406 + 8.44X$$

can use to estimate!!

X	Y	$(X - \bar{X}) * (Y - \bar{Y})$	$(X - \bar{X})^2$
14,000	97,000	$(14000 - 29150) * (97000 - 225500)$ = 1,946,775,000	$(14000 - 29150)^2$ = 229,522,500
27,000	185,000	$(27000 - 29150) * (185000 - 225500)$ = 87,075,000	$(27000 - 29150)^2$ = 4,622,500
39,900	260,000	$(39900 - 29150) * (260000 - 225500)$ = 370,875,000	$(39900 - 29150)^2$ = 115,562,500
17,000	143,000	$(17000 - 29150) * (143000 - 225500)$ = 1,002,375,000	$(17000 - 29150)^2$ = 147,622,500
34,000	270,000	$(34000 - 29150) * (270000 - 225500)$ = 215,825,000	$(34000 - 29150)^2$ = 23,522,500
43,000	398,000	$(43000 - 29150) * (398000 - 225500)$ = 2,389,125,000	$(43000 - 29150)^2$ = 191,922,500
Total :		6,012,050,000	712,675,000

$$\text{slope} = m = b_1 = \frac{\sum (X - \bar{X}) * (Y - \bar{Y})}{\sum (X - \bar{X})^2} = \frac{6,012,050,000}{712,675,000} = 8.44$$

$$\text{Intercept} = b = b_0 = \bar{Y} - b_1 * \bar{X} = 225,500 - 8.44 * 29150 = -20,406.28$$

$$\hat{y} = -20,406.28 + 8.44 X$$

P. 8

Use to make predictions

If we spend \$55,000 on weekly ads,
we can expect the sales Revenue
to be :

$$f(x) \hat{y} = -20,406.28 + 8.44 * 55,000$$

$$f(x) = \hat{y} = 443,567.83 = f(55,000)$$