

Research Design -- Topic 9

Fundamentals of Bivariate Regression and Correlation

© 2010 R. C. Gardner, Ph.D.

Bivariate regression (b) - - defining formulae

Bivariate correlation (r) - - defining formulae

Test of significance for regression

An example showing the distinction between b and r

Interpretations of correlation

Three limited truths

Factors that influence the magnitude of r

Special cases of the Pearson correlation

Tests of significance

Correlations with simple aggregates

1

Bivariate Regression and Correlation

Bivariate regression refers to an equation that relates a dependent variable to an independent variable, or a criterion to a predictor. The fundamental equation in raw score form is:

$$Y' = a + b_{yx}X$$

with a and b determined such that $\Sigma(Y - Y')^2 = a$ minimum.

and
$$a = \bar{Y} - b_{yx}\bar{X}$$

$$b_{yx} = \frac{\sum(X - \bar{X})(Y - \bar{Y})}{\sum(X - \bar{X})^2}$$

The formula in standard score form is:

$$Z'_Y = r_{XY}Z'_X$$

where r is as defined on the next slide

2

Bivariate correlation refers to covariation between two variables, X and Y . The most common measure is the Pearson product-moment correlation coefficient defined as:

$$r_{XY} = \frac{\sum(X - \bar{X})(Y - \bar{Y})}{nS_{b_x}S_{b_y}} = \frac{\sum(X - \bar{X})(Y - \bar{Y})}{(n-1)S_{u_x}S_{u_y}}$$

$$= \frac{\sum Z_X Z_Y}{n} = \frac{\sum Z_X Z_Y}{n-1}$$

using biased (S_b) and unbiased (S_u) estimates of the standard deviations respectively.

Or alternatives:

$$\frac{\sum(X - \bar{X})(Y - \bar{Y})}{\sqrt{\sum(X - \bar{X})^2 \sum(Y - \bar{Y})^2}} = \frac{N \sum XY - \sum X \sum Y}{\sqrt{(N \sum X^2 - (\sum X)^2)(N \sum Y^2 - (\sum Y)^2)}}$$

3

Given $Y = Y' + (Y - Y')$, we can compute:

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

$$SS_{TOTAL} = SS_{REGRESSION} + SS_{RESIDUAL}$$

And with some algebra, we can construct the following summary table

Source	df	Sums of Squares	
Regression	1	$r^2 SS_{TOTAL}$	$F = \frac{r^2 SS_{TOTAL}}{\frac{SS_{TOTAL}(1-r^2)}{n-2}}$
Residual	$n - 2$	$SS_{TOTAL}(1 - r^2)$	
Total	$n - 1$		$= \frac{r^2}{(1-r^2)/(n-2)}$

4

Consider the sample data set:

	X	Y	Z _x	Z _y
	3	3	-1.50	-1.50
	4	5	-.75	-.50
	4	5	-.75	-.50
	4	3	-.75	-1.50
	5	7	0	.50
	5	6	0	0
	5	7	0	.50
	6	9	.75	1.50
	7	7	1.50	.50
	7	8	1.50	1.00
Mean	5.0	6.0	.00	.00
S _u	1.33	2.00	1.00	1.00

5

Computing Regression Coefficients and Correlation

$$b_{yx} = \frac{\sum(X - \bar{X})(Y - \bar{Y})}{\sum(X - \bar{X})^2} = \frac{20}{16} = 1.25$$

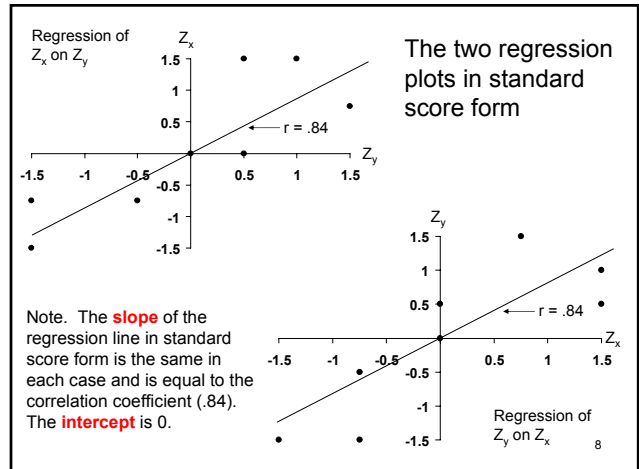
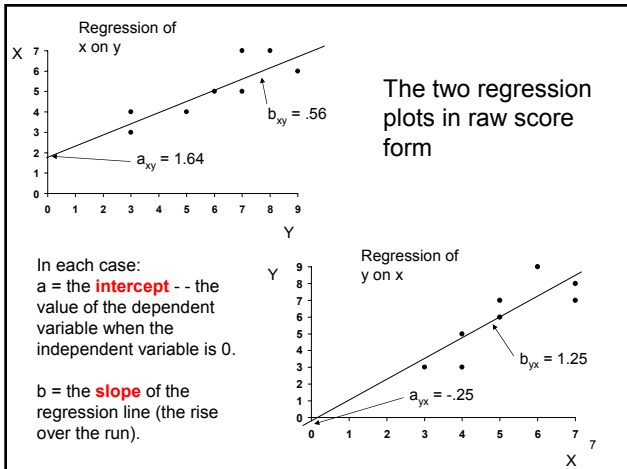
$$a_{yx} = \bar{Y} - b_{yx}\bar{X} = 6.0 - (1.25)(5.0) = -.25$$

$$b_{xy} = \frac{\sum(X - \bar{X})(Y - \bar{Y})}{\sum(Y - \bar{Y})^2} = \frac{20}{36} = .56$$

$$a_{xy} = \bar{X} - b_{xy}\bar{Y} = 5.0 - (.56)(6.0) = 1.64$$

$$r_{XY} = \frac{\sum(X - \bar{X})(Y - \bar{Y})}{(n-1)S_{u_x}S_{u_y}} = \frac{20}{9(1.33)(2.00)} = .84$$

6



Different Interpretations of Correlation

1. Correlation is a measure of the linear relation between y and y' :

$$r_{yy'} = \frac{\sum(y - \bar{y})(y' - \bar{y}')}{\sqrt{\sum(y - \bar{y})^2 \sum(y' - \bar{y}')^2}}$$

where: $y' = a + bx$

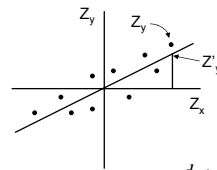
and $a = \bar{y} - b\bar{x}$

$$\therefore y' = \bar{y} + b(x - \bar{x})$$

$$= \frac{\sum(y - \bar{y})b(x - \bar{x})}{\sqrt{\sum(y - \bar{y})^2 b^2 \sum(x - \bar{x})^2}} = \frac{\sum(y - \bar{y})(x - \bar{x})}{\sqrt{\sum(y - \bar{y})^2 \sum(x - \bar{x})^2}} = r_{xy}$$

9

2. Correlation is a measure of the slope of the regression line in standard score form:



$$r = \text{slope} = \frac{Z_y'}{Z_x} \therefore Z_y' = rZ_x$$

Best fit line $\sum(Z_y - rZ_x)^2 = \text{minimum}$

$$\sum(Z_y - rZ_x)^2 = \sum Z_y^2 + r^2 \sum Z_x^2 - 2r \sum Z_x Z_y$$

$$\frac{d}{dr} (\sum Z_y^2 + r^2 \sum Z_x^2 - 2r \sum Z_x Z_y) = 0$$

$$0 + 2r \sum Z_x^2 - 2 \sum Z_x Z_y = 0$$

Note:

$$\frac{\sum Z_x^2}{n} = S_{Z_x}^2 = 1$$

$$r = \frac{\sum Z_x Z_y}{\sum Z_x^2} = \frac{\sum Z_x Z_y}{n} = r_{xy}$$

10

3. Correlation is a measure of the accuracy of predicting y given x :

Given $y = y' + (y - y')$

$$S_y^2 = S_{y'}^2 + S_{y-y'}^2 \quad \text{where } y' \text{ and } (y - y') \text{ are independent}$$

$$\therefore S_{y'}^2 = S_y^2 - S_{y-y'}^2$$

Defining $r_{xy}^2 = \frac{S_{y'}^2}{S_y^2} = \frac{S_y^2 - S_{y-y'}^2}{S_y^2}$

$$r_{xy}^2 = 1 - \frac{S_{y-y'}^2}{S_y^2}$$

$$\therefore r_{xy} = \pm \sqrt{1 - \frac{S_{y-y'}^2}{S_y^2}}$$

11

Three Limited Truths*

1. The Pearson product-moment correlation varies from -1 to +1. True, only under very specific circumstances.

Proof: Given $S_{Z_x}^2 = \frac{\sum Z_x^2}{N} = 1$

and $r = \frac{\sum Z_x Z_y}{N}$

r can equal +1, only if $Z_x = Z_y$ and -1, only if $Z_x = -Z_y$

Thus, for this to be true, the standardized distributions of x and y must be:

- Identical
- Symmetrical (not necessarily normal)

* Adapted from Gardner (2000).

12

2. Given a large enough sample size, the correlation will always be significant. True, only because of artifacts.

Proof: Given $X = T_X + E_{XR} + E_{XM}$ $Y = T_Y + E_{YR} + E_{YM}$

(i.e., the measures of X and Y consist of true scores (T_X & T_Y), random error (E_{XR} and E_{YR}) and measurement error (E_{XM} & E_{YM})).

Given: $\rho_{T_X T_Y} = 0$,

it is possible that $\rho_{XY} \neq 0$.

because the correlations

$\rho_{T_X E_{YM}}$, $\rho_{T_Y E_{XM}}$ and $\rho_{E_{XM} E_{YM}}$ are not 0

Thus, even with two variables that are truly independent, the correlation between measures of those variables may not be 0, and given a large enough sample size it may be significant.

3. Correlation does not mean causation. This is not a limitation of the statistic, but rather the nature of the underlying design.

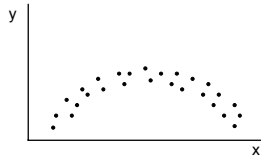
Consider an experiment on the effects of the amount of alcohol consumed in the afternoon and number of hours slept that night. This study could be run in controlled conditions with careful attention to detail, etc.

The correlation between the two could be considered an index of the linear effects of alcohol on hours slept (and an indication of causality) if the amount consumed was randomly determined and administered by the experimenter.

The correlation between the two would simply be an index of the covariation between the two if the amount consumed was not determined randomly. The regression equation would describe the nature of the linear relationship.

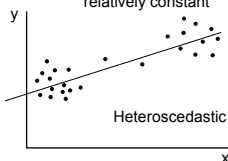
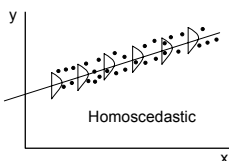
Factors That Affect the Pearson Product Moment Correlation Coefficient

1. Non-linear relationships.



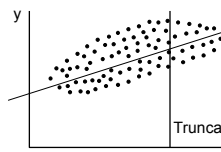
r_{xy} will approach 0, even though there is a non-linear relationship between the variables.

2. Heteroscedasticity.



An assumption underlying r_{xy} is homoscedasticity – viz., that the variation around the regression line is relatively constant

3. Restriction in Range.



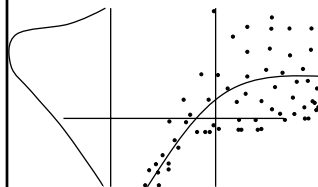
If selection on one variable has restricted its range, this could influence the magnitude of r_{xy} . To correct for this, calculate

$$r_{cor} = \frac{r_{xy} S_x / s_x}{\sqrt{1 - r_{xy}^2 + r_{xy}^2 S_x^2 / s_x^2}}$$

where S_x = stand. dev. uncurtailed
 s_x = stand. dev. curtailed

Entire Range \rightarrow r relatively high
Truncated Range \rightarrow r relatively low

4. Asymmetrical Distributions.



If the distributions aren't symmetrical, it can influence the value of r_{xy} .

concentration of cases

Phi Coefficient = Φ

Correlation between two dichotomous variables.

		x		
		0	1	
y	1	b	a	a + b
	0	d	c	c + d
		b + d	a + c	N = a + b + c + d

$$\sum x = \sum x^2 = a + c$$

$$\sum y = \sum y^2 = a + b$$

$$\sum xy = a$$

$$\therefore r = \frac{N \sum xy - \sum x \sum y}{\sqrt{[N \sum x^2 - (\sum x)^2][N \sum y^2 - (\sum y)^2]}} = \frac{ad - bc}{\sqrt{(a+b)(c+d)(b+d)(a+c)}} = \phi$$

and $\phi = \sqrt{\frac{x^2}{N}}$ for 2 x 2 tables

Effect Strength and Power

Cohen's (1988) definitions:

Small $\rho = .10$. "many relationships pursued in "soft" behavioral sciences are of this magnitude" (p. 79).

Medium $\rho = .30$. "this degree of relationship would be perceptible to the naked eye of a reasonably sensitive observer" (p. 80).

Large $\rho = .50$. "around the upper end of the range of (nonreliability) r's one encounters in those fields of behavioral science which use them extensively" (p.80).

Power can be calculated using Cohen (1988) or G*Power3. Note that G*Power3 has three routines that can be used for this purpose, one with t, one with F, and one with Exact tests.

Testing the significance of a single bivariate correlation coefficient

1. $H_0: \rho = 0$.

$$F = \frac{r^2}{(1-r^2)/(N-2)}$$

@ $df_1 = 1$; $df_2 = N - 2$.

or its equivalent:

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

@ $df = N - 2$.

2. $H_0: \rho = 0$. (For large N).

$$Z = r\sqrt{N-1}$$

Testing the significance of a single multiple correlation coefficient

3. $H_0: \rho = 0$.

$$F = \frac{R^2/p}{(1-R^2)/(N-p-1)}$$

@ $df_1 = p$; $df_2 = N - p - 1$

Testing the difference between two correlation coefficients

4. $H_0: \rho_1 = \rho_2$ for independent samples.

Fisher's Z $Z_{r_1} = \frac{1}{2} \log_e \frac{(1+r_1)}{(1-r_1)}$

$$Z_{r_2} = \frac{1}{2} \log_e \frac{(1+r_2)}{(1-r_2)}$$

and

$$Z = \frac{Z_{r_1} - Z_{r_2}}{\sqrt{\frac{1}{n_1 - 3} + \frac{1}{n_2 - 3}}}$$

25

Comparing two correlations from the same sample (with a common variable)

5. $H_0: \rho_{12} = \rho_{13}$ for correlated correlations.

1. Test proposed by Dunn and Clark (1969).

$$Z = \frac{(r_{12} - r_{13})\sqrt{N}}{\sqrt{(1-r_{12}^2)^2 + (1-r_{13}^2)^2 - 2r_{23}^3 - (2r_{23} - r_{12}r_{13})(1-r_{12}^2 - r_{13}^2 - r_{23}^2)}}$$

2. Test proposed by Meng, Rosenthal & Rubin (1992).

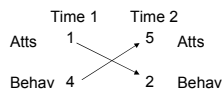
$$Z = (Z_{r_1} - Z_{r_2})\sqrt{\frac{N-3}{2(1-r_{23})h}}$$

where each: $Z_r = \frac{1}{2} \log_e \frac{(1+r)}{(1-r)}$

and $f = \frac{1-r_{23}}{2(1-(r_{12}^2 + r_{13}^2)/2)}$ $h = \frac{1-f(r_{12}^2 + r_{13}^2)/2}{1-(r_{12}^2 + r_{13}^2)/2}$ 26

Comparing two correlations from the same sample (with different variables)
(Cross Lagged Panel Analysis)

6. $H_0: \rho_{12} = \rho_{45}$ for correlated correlations.



$$Z = \frac{(r_{12} - r_{45})\sqrt{N}}{\sqrt{(1-r_{12}^2)^2 + (1-r_{45}^2)^2 - r_{12}r_{45}(r_{14}^2 + r_{15}^2 + r_{24}^2 + r_{25}^2) - 2(r_{14}r_{25} + r_{15}r_{24}) + C}}$$

where:

$$C = 2(r_{12}r_{14}r_{15} + r_{12}r_{24}r_{25} + r_{14}r_{24}r_{45} + r_{15}r_{25}r_{45})$$

27

Testing the Significance of an average correlation

7. $H_0: \rho_{av} = 0$

$$Z_{AV} = \frac{(n_1 - 3)Z_{r_1} + (n_2 - 3)Z_{r_2} + \dots + (n_k - 3)Z_{r_k}}{(n_1 - 3) + (n_2 - 3) + \dots + (n_k - 3)}$$

where:

$$Z_r = \frac{1}{2} \log_e \frac{1+r}{1-r}$$

then:

$$Z = Z_{AV} \sqrt{((n_1 - 3) + (n_2 - 3) + \dots + (n_k - 3))}$$

28

Testing the significance of a partial correlation

8. $H_0: \rho_{12.3} = 0$

$$t = \frac{r_{12.3}}{\sqrt{(1 - r_{12.3}^2)/(N - 3)}} \quad df = N - 3.$$

Testing the significance of a semipartial (part) correlation

9. $H_0: \rho_{1(2.3)} = 0$

$$F = \frac{(N - 3) r_{1(2.3)}^2}{1 - R_{1.23}^2} \quad df = 1, \quad N - 3$$

Note. These two statistics yield identical results, except that $F = t^2$ (both at $N-3$ df).

29

Comparing two bivariate regression coefficients

10. $H_0: b_1 = b_2$ in the population

$$t = \frac{b_1 - b_2}{S_{Db}}$$

where:

$$S_{Db} = C \sqrt{\frac{n_1 S_{y1}^2 (1 - r_1^2) + n_2 S_{y2}^2 (1 - r_2^2)}{(n_1 + n_2 - 4)}}$$

and:

$$C = \sqrt{\frac{1}{n_1 s_{x1}^2} + \frac{1}{n_2 s_{x2}^2}}$$

and $df = n_1 + n_2 - 4$.

30

Correlations Involving Aggregates

Raw Data

X_1	X_2	X_3	Y	T_x	T_z
6	10	32	100	48	-.51
8	9	25	97	42	-2.34
10	13	31	103	54	2.37
9	13	29	106	51	1.06
10	15	30	105	55	2.66
7	9	27	92	43	-2.15
6	10	29	85	45	-1.64
11	18	26	106	55	2.73
7	9	24	90	40	-3.28
9	12	30	93	51	1.10

For these data:

$$T_z = Z_{X_1} + Z_{X_2} + Z_{X_3}$$

$$r_{T_z, Y} = .769$$

$$T_x = X_1 + X_2 + X_3$$

$$r_{T_x, Y} = .754$$

31

Correlation Matrix

	Y	X_1	X_2	X_3
Y	1.0000	.7415	.7346	.2675
X_1	.7415	1.0000	.8688	.0259
X_2	.7346	.8688	1.0000	.1742
X_3	.2675	.0259	.1742	1.0000

Aggregated Standard Scores:

$$r_{T_z, Y} = \frac{\sum_{j=1}^m r_{jy}}{\sqrt{\sum_{j=1}^m \sum_{k=1}^m r_{jk}}} = \frac{.7415 + .7346 + .2675}{\sqrt{1.000 + .8688 + \dots + .1742 + 1.000}}$$

$$= \frac{1.7436}{\sqrt{5.1378}} = \frac{1.7436}{2.2667} = .769$$

32

Covariance Matrix

	Y	X ₁	X ₂	X ₃
Y	55.503	9.778	16.473	5.321
X ₁	9.778	3.133	4.629	.122
X ₂	16.473	4.629	9.060	1.400
X ₃	5.321	.122	1.400	7.129

Aggregated Raw Scores:

$$r_{r_{xy}} = \frac{\sum_{j=1}^m \text{cov}_{jy}}{S_y \sqrt{\sum_{j=1}^m \sum_{k=1}^m \text{cov}_{jk}}} = \frac{9.778 + 16.473 + 5.321}{\sqrt{55.503 \sqrt{3.133 + 4.629 + \dots + 1.400 + 7.129}}} = \frac{31.572}{(7.450)\sqrt{31.624}} = \frac{31.572}{41.895} = .754 \quad 33$$

Correlations Involving Difference Scores

(1) Correlation of Initial Score with the Difference

$$r_{x(y-x)} = \frac{\sum(x - \bar{x})[(y - x) - (\bar{y} - \bar{x})]}{N S_x S_{y-x}} = \frac{M r_{xy} - 1}{\sqrt{1 + M^2 - 2M r_{xy}}}$$

$$\text{where: } M = \frac{S_y}{S_x}$$

34

(2) Correlation of one variable (A) with a Difference (y - x)

$$r_{A(y-x)} = \frac{\sum(A - \bar{A})[(y - x) - (\bar{y} - \bar{x})]}{N S_A S_{y-x}} = \frac{M r_{Ay} - r_{Ax}}{\sqrt{1 + M^2 - 2M r_{xy}}} \quad \text{where: } M = \frac{S_y}{S_x}$$

35

(3) Correlation between two Difference Scores

$$r_{(B-A)(y-x)} = \frac{\sum[(B - A) - (\bar{B} - \bar{A})][(y - x) - (\bar{y} - \bar{x})]}{N S_{B-A} S_{y-x}} = \frac{M(Lr_{By} - r_{Ay}) - (Lr_{Bx} - r_{Ax})}{\sqrt{[L^2 + 1 - 2Lr_{AB}][M^2 + 1 - 2M r_{xy}]}}$$

$$\text{where: } M = \frac{S_y}{S_x} ; L = \frac{S_B}{S_A}$$

36

References

- Cohen, J. (1988). *Statistical Power Analysis for the Behavioral Sciences (Second Edition)*. Hillsdale, NJ: Lawrence Erlbaum.
- Dunn, O.J. & Clark, V.A. (1969). Correlation coefficients measured on the same individuals. *Journal of the American Statistical Association*, *64*, 366-377.
- Gardner, R. C. (2000). Correlation, causation, motivation and second language acquisition. *Canadian Psychology*, *41*, 10-24.
- Gardner, R.C. & Erdle, S. (1984). Aggregating scores: To standardize or not to standardize? *Educational and Psychological Measurement*, *44*, 813-822.
- Gardner, R.C. & Neufeld, R.W.J. (1987). Use of the simple change score in correlational analyses. *Educational and Psychological Measurement*, *47*, 849-864.
- Meng, X-L, Rosenthal, R. & Rubin, D.B. (1992). Comparing correlated correlation coefficients. *Psychological Bulletin*, *111*, 172-175.
- Tabachnick, B. G. & Fidell, L.S. (2007). *Using Multivariate Statistics (Fifth Edition)*. Needham Heights, MA: Allyn & Bacon.

37