

Variance & Covariance

I. Measures of Central Tendency

Mode: Most frequent observation

Median: Middle of rank ordered X_i

Mean: Arithmetic = $\bar{X} = \frac{1}{n} \sum_{i=1}^n (X_i)/N$

Geometric = $\sqrt[n]{\prod_{i=1}^n (X_i)}$

Harmonic = $\frac{N}{\sum_{i=1}^n (1/X_i)}$

II. Measures of Dispersion

Range: maximum - minimum

Interquartile range 75% - 25%

average absolute deviation from median

deviation score = $x = X - \bar{X}$

mean deviation = $\frac{1}{n} \sum_{i=1}^n (x_i)/N =$

$$\frac{1}{n} \sum_{i=1}^n (X - \bar{X})/N = \frac{1}{n} \sum_{i=1}^n (X)/N - \bar{X} = 0$$

standard deviation = $\sigma_x = \text{root mean square deviation}$

variance = mean square deviation = σ^2

Standard Deviation = σ_x = root mean square deviation

$$\sigma_x = \sqrt{\frac{1}{N} \sum_{i=1}^n (x_i^2)}$$

Variance = mean square deviation = $\sigma^2 = \frac{1}{N} \sum_{i=1}^n (x_i^2)$

unbiased estimate of variance from a sample =

$$\frac{1}{N-1} \sum_{i=1}^n (x_i^2)$$

$$V_x = \frac{\frac{1}{N} \sum_{i=1}^n (x_i^2) - \left(\frac{1}{N} \sum_{i=1}^n x_i \right)^2}{N-1}$$

Sensitivity to transformations:

$$M(X+C) = M(X) + C$$

$$V(X+C) = V(X)$$

$$V(XC) = C^2 V(X)$$

Standard Score = deviation score / standard deviation

$$z = \frac{x - M_x}{S_x} = \frac{(X-M)}{S_x} = (\text{a unit free index of dispersion})$$

$$M_z = 0 \quad V_z = 1 \quad S_z = 1$$

Coefficient of variation = S_x/M_x (\Rightarrow ratio measurement)

Variance of Composites

$$V(X+Y) = V(x+y) = \frac{\sum_{i=1}^n ((x_i+y_i)^2)/(N-1) =}{(N-1)} \frac{\sum_{i=1}^n ((x_i)^2) + \sum_{i=1}^n ((y_i)^2) + 2 \sum_{i=1}^n ((x_i * y_i))}{(N-1)}$$

$$V(X+Y) = V_x + V_y + 2 Cov_{xy}$$

Covariance of x and y =

$$\frac{\sum_{i=1}^n ((X_i * Y_i)) - \sum_{i=1}^n ((X_i)) * \sum_{i=1}^n ((Y_i))}{(N-1)}$$

V(x+y) a visual representation

	x	y
x	V _x	C _{xy}
y	C _{xy}	V _y

Variance of Composites: an example

Standard deviation of GRE Verbal = 100

Standard deviation of GRE Quant = 100

Variance of Verbal = 100 * 100 = 10,000

Variance of Quant = 100 * 100 = 10,000

Covariance of GRE Q and V = 6,000

Variance of GRE (V + Q) = V_v + V_q + 2 C_{vq}

	Verbal	Quantitative
Verbal	10,000	6,000
Quantitative	6,000	10,000

V(V+Q) = 32,000 ==> SD(V+Q) = 179

Generalization of Variance of Composites to N variables:

$$\begin{aligned}
V(x_1 + x_2 + \dots + x_n) = & \\
Vx_1 + Vx_2 + \dots + Vx_n & + 2(C_{x_1x_2} + C_{x_1x_3} + \dots + C_{x_ix_j} + \dots \\
) & \\
(n \text{ terms}) & \qquad (n * (n-1) \text{ terms})
\end{aligned}$$

Variance of N variables: (figural representation)

	x_1	x_2	x_3		$\dots x_i$		$\dots x_j$		$\dots x_n$
x_1	V_1	C_{12}	C_{13}		$\dots C_{1i}$				$\dots C_{1n}$
x_2	C_{21}	V_2	C_{23}		$\dots C_{2i}$				$\dots C_{2n}$
x_3	C_{31}	C_{32}	V_3		$\dots C_{3i}$				$\dots C_{3n}$
\dots					\dots				\dots
x_i					$\dots V_i$		$\dots C_{ij}$		$\dots C_{in}$
\dots					\dots				\dots
x_j	C_{j1}	C_{j2}	C_{j3}	\dots	$\dots C_{ji}$	\dots	$\dots V_j$		$\dots C_{jn}$
\dots									
x_n	C_{n1}	C_{n2}	C_{n3}		$\dots C_{ni}$		$\dots C_{nj}$		$\dots V_n$

A total of n variance terms on the diagonal and $n * (n-1) = n^2 - n$ covariance terms off the diagonal.

The variance of the composite of n variables = the sum of the n variances and the $n * (n-1)$ covariances.

Correlation and Regression

The problem of predicting y from x :

Linear prediction $y = bx + c$ $Y = b(X - M_x) + M_y$
 $+ M_y$

error in prediction = predicted y - observed y

problem is to minimize the squared error of prediction

minimize the error variance = $[\sum (y_p - y_o)^2] / (N-1)$

$$V_e = V(bx - y) = \frac{\sum (bx - y)^2}{(N-1)} =$$

$$\frac{\sum (b^2x^2 - 2bxy + y^2)}{(N-1)} =$$

$$b^2 \frac{\sum x^2}{(N-1)} - 2b \frac{\sum xy}{(N-1)} + \frac{\sum y^2}{(N-1)} \implies$$

$$V_e = b^2 V_x - 2b C_{xy} + V_y$$

V_e is minimized when the first derivative (w.r.t. b) = 0 ==
>

$$\text{when } 2bV_x - 2C_{xy} = 0 \implies$$

$$b_{y.x} = C_{xy} / V_x$$

Similarly, the best $b_{x.y}$ is C_{xy} / V_y

The Pearson Product Moment Correlation Coefficient (PPMC C) is the geometric mean of these two slopes:

$$r_{xy} = \frac{C_{xy}}{\sqrt{V_x V_y}} = \frac{C_{xy}}{S_x S_y}$$

$r_{xy} =$

Error!

$$r_{xy} = \frac{C_{xy}}{\sqrt{V_x V_y}} = \frac{C_{xy}}{S_x S_y} = r_{xy} = \frac{\sum xy}{\sqrt{\sum x^2 \sum y^2}}$$

1) If x and y are continuous variables, then $r = \text{Pearson } r$

2) if x and y are rank orders, then $r = \text{Spearman } r$

3) if x is continuous and y is dichotomous $r = \text{point biserial}$

4) if x and y are both dichotomous, then $r = \text{phi} = \sqrt{\frac{\text{chi}^2 \text{ square}}{N}}$

5) **Tetrachoric correlation** is an estimate of continuous (Pearson) based upon dichotomous data. This assumes bivariate normality.

6) **Biserial correlation** estimates continuous based upon one dichotomous and one continuous. It also assumes normality.

Calculating formulae:

$$r_{xy} = \frac{\sum xy}{\sqrt{\sum x^2 \sum y^2}}$$

$$\text{Covariance } xy = (\sum XY - \sum X \sum Y / N) / (N-1)$$

$$\text{Variance } X = (\sum X^2 - (\sum X)^2 / N) / (N-1)$$

$$\text{Variance } Y = (\sum Y^2 - (\sum Y)^2 / N) / (N-1)$$

$$\text{Correlation} = \frac{\text{Covariance}}{\sqrt{(\text{Variance } X)(\text{Variance } Y)}} =$$

$$\frac{(\sum XY - \sum X \sum Y / N) / (N-1)}{\sqrt{[(\sum X^2 - (\sum X)^2 / N) / (N-1)][(\sum Y^2 - (\sum Y)^2 / N) / (N-1)]}}$$

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

Correlation

1) **Slope of regression** ($b_{xy} = C_{xy}/V_x$) reflects units of x and y but the correlation $\{r = C_{xy}/(S_x S_y)\}$ is unit free.

2) **Geometrically**, $r = \cos(\text{angle between test vectors})$

3) **Correlation as prediction:**

Let y_p = predicted deviation score of y = predicted Y - M
 $y_p = b_{xy}x$ and $b_{xy} = C_{xy}/V_x = rS_y/S_x \implies y_p/S_y = r(x/S_x) \implies$
 predicted z score of y (z_{y_p}) = r_{xy} * observed z score of x (z_x)
 predicted z score of x (z_{x_p}) = r_{xy} * observed z score of y (z_y)

4) Amount of **error variance** (residual or unexplained variance) in y given x and r

$$V_e = \bullet e^2/N = \bullet (y - bx)^2/N = \bullet \{y - (r \cdot S_y \cdot x/S_x)\}^2$$

$$V_y + V_y \cdot r^2 - 2(r \cdot S_y \cdot C_{xy})/S_x$$

(but $S_y \cdot C_{xy}/S_x = V_y \cdot r$)

$$V_y + V_y \cdot r^2 - 2(r^2 \cdot V_y) = V_y(1 - r^2) \implies$$

$$V_e = V_y(1 - r^2) \iff V_{y_p} = V_y(r^2)$$

Residual Variance = Original Variance * (1 - r²)

Variance of predicted scores = original variance * r²

5) Correlation of x with predicted y = $C_x(bx)/(S_x \cdot S_{y_p})$

but $C_x(bx) = bV_x = V_x \cdot r \cdot S_y/S_x$ and $S_{y_p} = rS_y$ and therefore $r(x \text{ with predicted } y) = 1$

	x	y	y_p	residual
Variance	V_x	V_y	$V_y(r^2)$	$V_y(1 - r^2)$
Correlation with x	1	r_{xy}	1	0
Correlation with y	r_{xy}	1	r_{xy}	$\sqrt{(1 - r^2)}$

Multiple Correlation

$$r_{xy} = \frac{C_{xy}}{\sqrt{V_x V_y}}$$

The problem of predicting y from x₁, x₂:

Linear prediction $y = b_1 x_1 + b_2 x_2 + c$

Just as the optimal b weights in regression are $b_{y.x} = C_{xy} / V_x$, so are the optimal b weights in multiple regression, however, they are corrected for the effect of the other variables:

In the two variable case the b weights (betas) are:

$$b_1 = \frac{C_{yx_1.x_2}}{V_{x_1.x_2}} \qquad b_2 = \frac{C_{yx_2.x_1}}{V_{x_2.x_1}}$$

$$b_1 = \frac{r_{x_1 y} - r_{x_1 x_2} r_{x_2 y}}{1 - r_{x_1 x_2}^2} \qquad b_2 = \frac{r_{x_2 y} - r_{x_1 x_2} r_{x_1 y}}{1 - r_{x_1 x_2}^2}$$

The amount of variance accounted for by the model is the sum of the product of the betas and the zero order correlations:

$$R^2 = \sum \beta_j \cdot r_{x_j y}$$

Consider the following example:

extraversion with leadership $r = .56$ $r^2 = .31$
 dominance with leadership $r = .42$ $r^2 = .18$
 extraversion with dominance $r = .48$ $r^2 = .23$

$$\beta_1 = \frac{.56 - (.42)(.48)}{1 - .48^2} = .47 \qquad \beta_2 = \frac{.42 - (.56)(.48)}{1 - .48^2} = .20$$

$$R^2 = \beta_1 \cdot r_{x_1 y} + \beta_2 \cdot r_{x_2 y} = .47 \cdot .56 + .20 \cdot .42 = .35$$

Multiple Correlation as weighted linear composites

$$r_{xy} = \frac{C_{xy}}{\sqrt{V_x V_y}}$$

The problem is to find the Covariance of $(x_1 \ x_2)$ with y

	x_1	x_2	y
x_1	V_1	C_{12}	C_{1y}
x_2	C_{21}	V_2	C_{2y}
y	C_{y1}	C_{y2}	V_y

Covariance $((x_1 \ x_2) \ y) = C_{y1} + C_{y2}$

Variance $(x_1 \ x_2) = V_1 + C_{12} + C_{21} + V_2$

Variance $(y) = V_y$

Standardized solution:

	z_1	z_2	z_y
z_1	1	r_{12}	r_{1y}
z_2	r_{21}	1	r_{2y}
z_y	r_{y1}	r_{y2}	1

Covariance $((x_1 \ x_2) \ y) = r_{y1} + r_{y2}$

Variance $(x_1 \ x_2) = 1 + r_{12} + r_{21} + 1$

Variance $(y) = 1$

Multiple Correlation is optimally weighted composite:

$$r_{xy} = \frac{C_{xy}}{\sqrt{V_x V_y}}$$

	$b_1 z_1$	$b_2 z_2$	z_y
$b_1 z_1$	b_1^2	$b_1 b_2 r_{12}$	$b_1 r_{1y}$
$b_2 z_2$	$b_1 b_2 r_{21}$	b_2^2	$b_2 r_{2y}$
z_y	$b_1 r_{y1}$	$b_2 r_{y2}$	1

Covariance $((x_1 \ x_2) \ y) = b_1 r_{y1} + b_2 r_{y2}$

Variance $(x_1 \ x_2) = b_1^2 + b_1 b_2 r_{21} + b_1 b_2 r_{21} + b_2^2$

Variance $(y) = 1$

$$b_1 = \frac{r_{1y} - r_{12} r_{2y}}{1 - r_{12}^2} \qquad b_2 = \frac{r_{2y} - r_{12} r_{1y}}{1 - r_{12}^2}$$

$$C_{12,y} = \frac{(r_{1y} - r_{12} r_{2y}) r_{1y} + (r_{2y} - r_{12} r_{1y}) r_{2y}}{1 - r_{12}^2}$$

$$V_{12} = \frac{[(r_{1y} - r_{12} r_{2y})^2 + 2(r_{1y} - r_{12} r_{2y})(r_{2y} - r_{12} r_{1y}) + (r_{2y} - r_{12} r_{1y})^2]}{(1 - r_{12}^2)(1 - r_{12}^2)}$$

expand and collect terms ==>

$$V_{12} = Cov = R^2_{12,y} = \frac{r_{1y}^2 + r_{2y}^2 - 2 r_{1y} r_{12} r_{2y}}{1 - r_{12}^2}$$

if $r_{12} = 0$ then notice that $R^2_{12,y} = r_{1y}^2 + r_{2y}^2$

Unit Weights versus Multiple R

Multiple Correlation is optimally weighted composite:

$$r_{xy} = \frac{C_{xy}}{\sqrt{V_x V_y}}$$

But consider what happens if we equal (unit) weights rather than optimal weights

Standardized solution with unit weights

	z ₁	z ₂	z _y
z ₁	1	r ₁₂	r _{1y}
z ₂	r ₂₁	1	r _{2y}
z _y	r _{y1}	r _{y2}	1

Covariance ((x₁ x₂) y) = r_{y1} + r_{y2}

Variance (x₁ x₂) = 1 + r₁₂ + r₂₁ + 1

Variance (y) = 1

$$R = \frac{r_{y1} + r_{y2}}{\sqrt{1 + r_{12} + r_{21} + 1}} = \frac{r_{y1} + r_{y2}}{\sqrt{2 * (1 + r_{12})}}$$

Consider several examples:

rx1x2	rx1y	rx2y	beta 1	beta 2	R	R ²	Unit Wt	UW ²
0.0	0.5	0.5	0.50	0.50	0.71	0.50	0.71	0.50
0.3	0.5	0.5	0.38	0.38	0.62	0.38	0.62	0.38
0.5	0.5	0.5	0.33	0.33	0.58	0.33	0.58	0.33
0.7	0.5	0.5	0.29	0.29	0.54	0.29	0.54	0.29
0.3	0.5	0	0.55	-0.16	0.52	0.27	0.31	0.10
0.3	0.5	0.3	0.45	0.16	0.52	0.27	0.50	0.25

Partial Correlation

$$r_{xy} = \frac{C_{xy}}{\sqrt{V_x V_y}}$$

To find r_{xy} with w held constant (partial $r=r_{xy.w}$) or R_{xyw} (multiple R), we need to find the Covariance and Variances.

Conceptual solution:

- find residual x after predicting from w ($x.w$)
- find residual y after predicting from w ($y.w$)
- correlate these residual scores.

Variance of residual = (Variance of original) \cdot $(1-r^2)$
 Covariance of residuals =
 original covariance - covariance with control

$$z_{\text{predicted}} = r \cdot z_{\text{predictor}}$$

$$z_{\text{residual}} = z_{\text{original}} - r \cdot z_{\text{predictor}}$$

$$z_{x.w} = z_x - r_{xw} \cdot z_w \qquad z_{y.w} = z_y - r_{yw} \cdot z_w$$

$$\text{Covariance}(z_{x.w}, z_{y.w}) = \text{Cov}(z_x, z_y) - r_{xw} \cdot r_{yw} \text{ since}$$

$$\text{Covariance}(z_{x.w}, z_{y.w}) = \bullet (z_x - r_{xw} \cdot z_w) \cdot (z_y - r_{yw} \cdot z_w) / N =$$

$$\bullet (z_x - r_{xw} \cdot z_w) \cdot (z_y - r_{yw} \cdot z_w) / N =$$

$$\bullet (z_x \cdot z_y - r_{xw} \cdot z_w \cdot z_y - z_x \cdot r_{yw} \cdot z_w + r_{xw} \cdot z_w \cdot r_{yw} \cdot z_w) / N =$$

$$\{ \bullet (z_x \cdot z_y) - r_{xw} \cdot \bullet (z_w \cdot z_y) - r_{yw} \cdot \bullet (z_x \cdot z_w) + r_{xw} \cdot r_{yw} \cdot \bullet (z_w \cdot z_w) \} / N$$

$$\text{Cov}(z_x, z_y) - r_{xw} \cdot r_{yw} \text{ - } r_{yw} \cdot r_{xw} + r_{xw} \cdot r_{yw} \cdot \text{Var } z_w =$$

$$\text{Cov}(z_x, z_y) - r_{xw} \cdot r_{yw}$$

$$\text{Variance residual} = V_x \cdot (1 - r_{xw}^2)$$

$$\text{Partial } r_{xy.w} = \frac{\text{Cov}(z_x, z_y) - r_{xw} \cdot r_{yw}}{\sqrt{(1 - r_{xw}^2) \cdot (1 - r_{yw}^2)}}$$