

Lecture 27 – Centered Vectors and the Sample Correlation R

Centering

Define the new vectors $\dot{X} = X - \bar{X} = \begin{bmatrix} X_1 - \bar{X} \\ X_2 - \bar{X} \\ \vdots \\ X_n - \bar{X} \end{bmatrix}$ and $\dot{Y} = Y - \bar{Y} = \begin{bmatrix} Y_1 - \bar{Y} \\ Y_2 - \bar{Y} \\ \vdots \\ Y_n - \bar{Y} \end{bmatrix}$. The new vectors \dot{X} and \dot{Y}

are the **Centered** versions of the original vectors X and Y . Centering the data has many advantages and is often done prior to running a regression. I will focus on just a couple of the applications of centering.

Centering and the Sample Correlation

One of the coolest applications of centering involves the sample correlation r (or R). The sample correlation R is an estimator of the correlation between the variables X and Y , $\rho_{X,Y} = \frac{Cov(X,Y)}{\sigma_X \sigma_Y}$, where

$Cov(X,Y) = E(X - \mu_X)(Y - \mu_Y)$ is the covariance of X and Y . The sample correlation R is computed as

$$R = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\sum_{i=1}^n (X_i - \bar{X})^2} \sqrt{\sum_{i=1}^n (Y_i - \bar{Y})^2}}. \text{ Now, rewriting this using centered vectors, } R = \frac{\sum_{i=1}^n \dot{X}_i \dot{Y}_i}{\sqrt{\sum_{i=1}^n \dot{X}_i^2} \sqrt{\sum_{i=1}^n \dot{Y}_i^2}}.$$

The last equation can be rewritten $R = \frac{\dot{X}^T \dot{Y}}{\sqrt{\dot{X}^T \dot{X}} \sqrt{\dot{Y}^T \dot{Y}}} = \frac{\dot{X}^T \dot{Y}}{\|\dot{X}\| \|\dot{Y}\|}$. The final expression may not look

familiar unless you've had Multivariate Calculus, so let me introduce the relevant theorem.

Theorem: The dot product of vectors u and v obeys the equality $u \cdot v = \|u\| \|v\| \cos \theta$, where θ is the angle between the vectors. The corollary is $\cos \theta = \frac{u \cdot v}{\|u\| \|v\|}$.

Returning to the transformed formula for the sample correlation, $R = \frac{\dot{X}^T \dot{Y}}{\|\dot{X}\| \|\dot{Y}\|}$, we rewrite it as $R = \cos \theta$,

where θ is the angle between the centered vectors \dot{X} and \dot{Y} in \mathbb{R}^n . This is cool enough already, but the real coolness is that we've just proven one of the main properties of the sample correlation, that $-1 \leq R \leq 1$. Normally, the proof would require a fair amount of ordinary algebra, and doesn't provide the geometrical insight acquired from the linear algebraic proof above.

Special Cases: $R = 0$ and $R = \pm 1$

Two special cases of the sample correlation are when $R = 0$ and when $R = \pm 1$. Although they are almost never encountered in practice, they help to visualize the geometry of correlation in simple regression.

Special Case: $R = 0$

From the transformed formula $R = \frac{\dot{\mathbf{X}}^T \dot{\mathbf{Y}}}{\|\dot{\mathbf{X}}\| \|\dot{\mathbf{Y}}\|}$, $R = 0$ if and only if $\dot{\mathbf{X}}^T \dot{\mathbf{Y}} = 0$ if and only if $\dot{\mathbf{Y}} \perp \dot{\mathbf{X}}$.

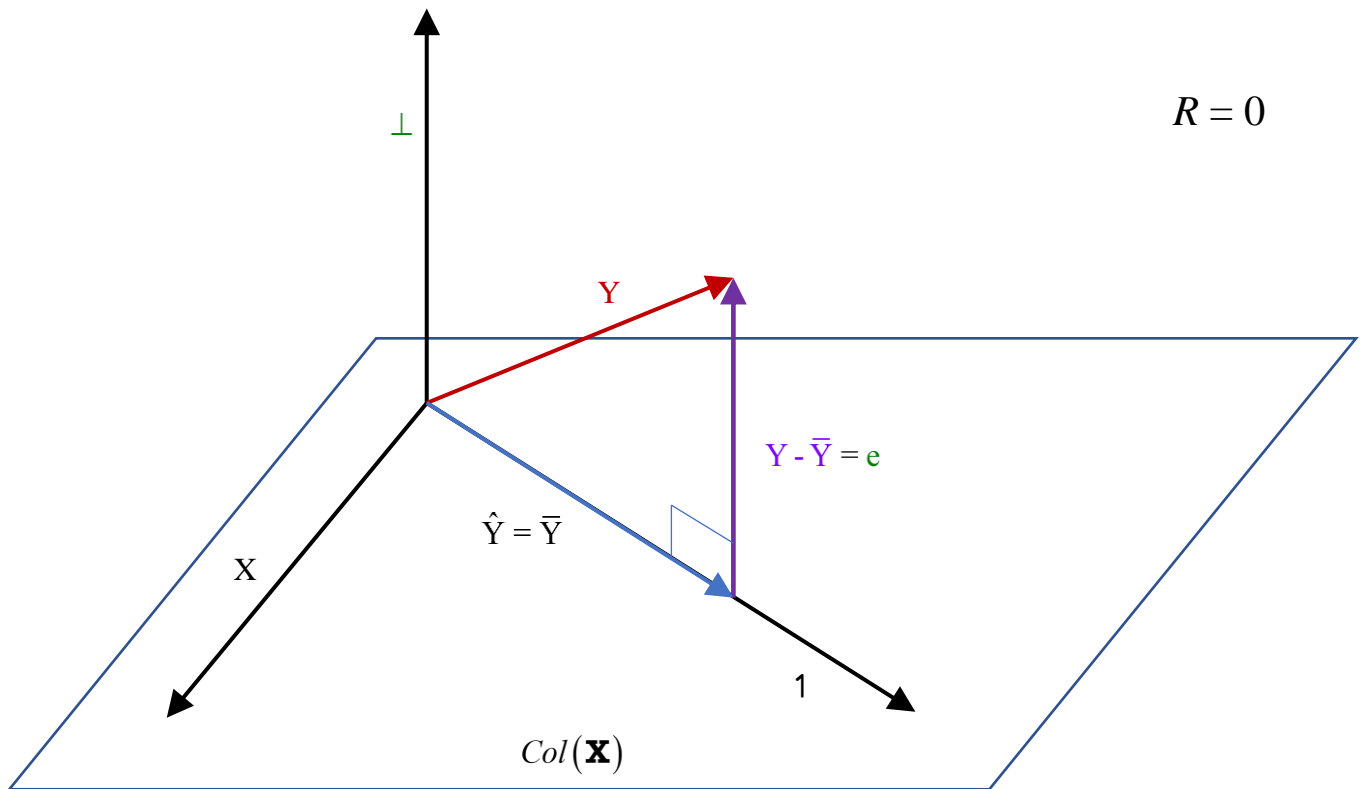
The following are always true in regression:

- $\dot{\mathbf{Y}} \perp \mathbf{1}$, where $\mathbf{1}$ is the n -vector of ones, because $\dot{\mathbf{Y}}^T \mathbf{1} = (\mathbf{Y} - \bar{Y})^T \mathbf{1} = \sum (Y_i - \bar{Y}) = 0$
- Similarly, $\dot{\mathbf{X}} \perp \mathbf{1}$ because $\sum (X_i - \bar{X}) = 0$

In addition, when the sample correlation is zero, $\dot{\mathbf{Y}} \perp \dot{\mathbf{X}}$.

The vectors $\mathbf{1}$ and $\dot{\mathbf{X}}$ are in $Col(\mathbf{X})$, and are independent because they are orthogonal. Furthermore, both are orthogonal to $\dot{\mathbf{Y}}$ when $R = 0$, so $\dot{\mathbf{Y}}$ is orthogonal to the column space of the design matrix, $Col(\mathbf{X})$, when $R = 0$. This immediately implies all of the following. (See sketch below)

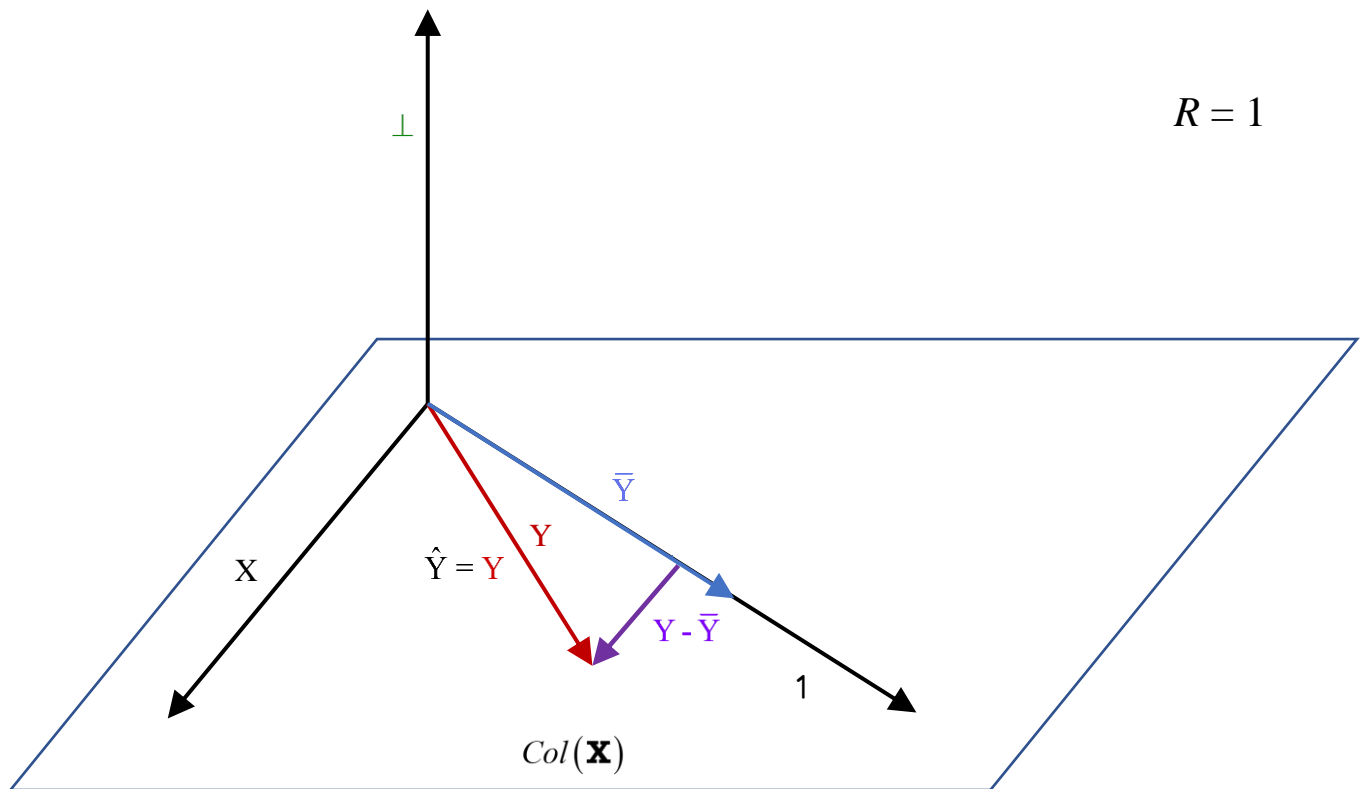
- $\mathbf{Y} - \bar{Y} = \mathbf{e}$, i.e., the residuals are just the deviations, $e_i = Y_i - \bar{Y}$.
- $\hat{\mathbf{Y}} = \bar{Y}$, i.e., the fitted values are just the average value of the response for the sample because $\hat{\beta}_1 = 0$ and $\hat{\beta}_0 = \bar{Y}$ when $R = 0$. ($\hat{Y}_i = \bar{Y}$ for all X_i , and the fitted regression line is horizontal.)
- The vector of observations \mathbf{Y} lies directly above the vector of ones, $\mathbf{1}$, because the vector \mathbf{X} plays no role in fitting values when $\hat{\beta}_1 = 0$.



Special Case: $R = \pm 1$

From $R = \cos \theta$, $R = \pm 1$ implies that $\dot{Y} \parallel \dot{X}$. Since $\bar{Y} = \bar{y}1$ always lies in $Col(\mathbf{X})$, if $\dot{Y} = (Y - \bar{Y})$ also lies in $Col(\mathbf{X})$, then so must Y . This makes sense when we reflect that all points in the scatterplot line up along the regression line when $R = \pm 1$, so the residuals are all 0. Finally, when $R = \pm 1$, the fitted values are just the observations themselves, i.e., $\hat{Y} = Y$. Contrast this with the situation when $R = 0$ and $\hat{Y} = \bar{Y}$.

Because $\dot{Y} \parallel \dot{X}$ and Y lies in $Col(\mathbf{X})$ when $R = \pm 1$, it's tempting to assume that Y is parallel to X in this case, but it need not be. However, if the regression line passes through the origin and $R = \pm 1$, then $Y = \hat{\beta}_1 X$ and the vector of observations Y is parallel the vector X .



Centering and Regression

Centering can also be used to obtain the estimated regression coefficients. There are several ways to accomplish this depending upon your goals, but we'll look at one simple approach.

Consider the following simple regression model: $\dot{Y} = \hat{\beta}_1 \dot{X} + \varepsilon$. So, what happened to the intercept term $\hat{\beta}_0$? Well, one effect of centering is that the fitted regression line will always pass through the origin, so we don't need an intercept term in the model. (But have no fear, because $\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X}$ as before.)

The Design Matrix in Simple Regression Using Centering

Centering results in a different design matrix $\dot{\mathbf{X}}$ because the first column of ones is absent (because there is no intercept term in the model) and the remaining columns involve a centered vector for each predictor variable. In simple regression there is only one predictor, so the design matrix is particularly simple, it is

the simply the vector $\dot{\mathbf{X}} = \begin{bmatrix} \dot{X}_1 \\ \dot{X}_2 \\ \vdots \\ \dot{X}_n \end{bmatrix}$. The equations leading to the least-squares estimators of the

regression coefficients remains the same minus the missing $\hat{\beta}_0$, $\hat{\beta} = (\dot{\mathbf{X}}^T \dot{\mathbf{X}})^{-1} \dot{\mathbf{X}}^T \dot{Y}$.

In simple regression only $\hat{\beta}_1$ needs to be estimated this way, to the equation becomes $\hat{\beta}_1 = (\dot{\mathbf{X}}^T \dot{\mathbf{X}})^{-1} \dot{\mathbf{X}}^T \dot{Y}$,

but $\dot{\mathbf{X}}^T \dot{\mathbf{X}}$ is just a number, so this becomes $\hat{\beta}_1 = \frac{\dot{\mathbf{X}}^T \dot{Y}}{\dot{\mathbf{X}}^T \dot{\mathbf{X}}}$.

If this looks familiar to those who've had vector calculus, this states that the vector of fitted values \hat{Y} is the vector projection of \dot{Y} onto \dot{X} , i.e., $\hat{Y} = \dot{\mathbf{X}} \hat{\beta}_1 = \dot{\mathbf{X}} \frac{\dot{\mathbf{X}}^T \dot{Y}}{\dot{\mathbf{X}}^T \dot{\mathbf{X}}} = \left(\frac{\dot{\mathbf{X}}^T \dot{Y}}{\dot{\mathbf{X}}^T \dot{\mathbf{X}}} \right) \dot{\mathbf{X}} = \text{proj}_{\dot{\mathbf{X}}} \dot{Y}$. How cool is that?

Those of you who suffered through Math 32 with professor Hauser now understand why I put vector projections on every vector calculus exam.

Example

Using our evergreen example:

x	1	2	4	5
y	8	4	6	2

We create the centered vectors $\dot{\mathbf{X}} = \begin{bmatrix} -2 \\ -1 \\ 1 \\ 2 \end{bmatrix}$, $\dot{\mathbf{Y}} = \begin{bmatrix} 3 \\ -1 \\ 1 \\ -3 \end{bmatrix}$ (**Note:** A centered vectors is just the vector of

deviations.)

$$\text{Then, } \hat{\beta}_1 = \frac{\dot{\mathbf{X}}^T \dot{\mathbf{Y}}}{\dot{\mathbf{X}}^T \dot{\mathbf{X}}} = \frac{\begin{bmatrix} -2 & -1 & 1 & 2 \end{bmatrix} \begin{bmatrix} 3 \\ -1 \\ 1 \\ -3 \end{bmatrix}}{\begin{bmatrix} -2 & -1 & 1 & 2 \end{bmatrix} \begin{bmatrix} -2 \\ -1 \\ 1 \\ 2 \end{bmatrix}} = \frac{-10}{10} = -1 .$$

$$\text{Also, } \hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X} = 5 - (-1)(3) = 8 .$$

In Statgraphics, I've entered the centered data below.

XC	YC
-2	3
-1	-1
1	1
2	-3

Regressing YC on XC produces the following *Analysis* window,

Coefficients

	<i>Least Squares</i>	<i>Standard</i>	<i>T</i>	
<i>Parameter</i>	<i>Estimate</i>	<i>Error</i>	<i>Statistic</i>	<i>P-Value</i>
Intercept	0	1.11803	0	1.0000
Slope	-1.0	0.707107	-1.41421	0.2929

Analysis of Variance

<i>Source</i>	<i>Sum of Squares</i>	<i>Df</i>	<i>Mean Square</i>	<i>F-Ratio</i>	<i>P-Value</i>
Model	10.0	1	10.0	2.00	0.2929
Residual	10.0	2	5.0		
Total (Corr.)	20.0	3			

Correlation Coefficient = -0.707107

R-squared = 50.0 percent

R-squared (adjusted for d.f.) = 25.0 percent

Standard Error of Est. = 2.23607

Notice that the estimated intercept is zero, as expected. Alternatively, we could have unchecked the Include Constant box in the Analysis Options window.

Simple Regression Options

×

Type of Model

☒ Linear
☐ Square Root-Y
☐ Exponential
☐ Reciprocal-Y
☐ Squared-Y
☐ Square Root-X
☐ Double Square Root
☐ Log-Y Square Root-X
☐ Reciprocal-Y Square Root-X

☐ Squared-Y Square Root-X
☐ Logarithmic-X
☐ Square Root-Y Log-X
☐ Multiplicative
☐ Reciprocal-Y Log-X
☐ Squared-Y Log-X
☐ Reciprocal-X
☐ Square Root-Y Reciprocal-X
☐ S-Curve

☐ Double Reciprocal
☐ Squared-Y Reciprocal-X
☐ Squared-X
☐ Square Root-Y Squared-X
☐ Log-Y Squared-X
☐ Reciprocal-Y Squared-X
☐ Double Squared
☐ Logistic
☐ Log Probit

☐ Include constant

Alternative Fit

☒ None (least squares only)
☐ Minimize absolute deviations
☐ Use medians of 3 groups

OK

Cancel

Help

Below is the output from Statgraphics. Notice that the estimated slope and correlation are unchanged, but the Mean Square Error is different because we are working with a different model. In this case, the *MSE* from the previous output is the one we want, as can be seen by rerunning the regression one more time on the original, uncentered, data.

Coefficients

	<i>Least Squares</i>	<i>Standard</i>	<i>T</i>	
<i>Parameter</i>	<i>Estimate</i>	<i>Error</i>	<i>Statistic</i>	<i>P-Value</i>
Slope	-1.0	0.57735	-1.73205	0.1817

Analysis of Variance

<i>Source</i>	<i>Sum of Squares</i>	<i>Df</i>	<i>Mean Square</i>	<i>F-Ratio</i>	<i>P-Value</i>
Model	10.0	1	10.0	3.00	0.1817
Residual	10.0	3	3.33333		
Total	20.0	4			

Correlation Coefficient = -0.707107

R-squared = 50.0 percent

R-squared (adjusted for d.f.) = 50.0 percent

Standard Error of Est. = 1.82574