

8 📜 🗏

Q

Home ► All Journals ► Mathematics, Statistics & Data Science ► Journal of Statistics Education ► List of Issues ► Volume 10, Issue 1 ► More on Venn Diagrams for Regression

Journal of Statistics Education > Volume 10, 2002 - Issue 1

Free access

3,43510ViewsCrossRef citations to dateAltmetric

Listen

**Original Articles** 

# More on Venn Diagrams for Regression

E. Kennedy Peter 🔽

| Published online: 01 Dec 2017

#### **L** Cite this article **I** https://doi.org/10.1080/10691898.2002.11910547



were not included in Ip's presentation, and to demonstrate some ways of using these applications to improve student understanding of ordinary least squares (OLS) regression. Because these alternative applications are not new to the literature, the main contribution of this paper consists of suggestions for how this approach can be used effectively in teaching.

The first use of Venn diagrams in regression analysis appears to be in the textbook by <u>Cohen and Cohen (1975)</u>. A major difficulty with its use occurs in the presence of suppressor variables, a problem discussed at length by <u>Ip (2001)</u>. No one denies that Venn diagrams can mislead, just as no one denies that ignoring friction in expositions of physical phenomena misleads, or using Euclidian geometry misleads because the surface of the earth is curved. Such drawbacks have to be weighed against the pedagogical benefits of the "misleading" expository device. As recognized by Ip, in the case of applying the Venn diagram to regression analysis, reasonable instructors could disagree on the pedagogical value of the Venn diagram because of the suppressor variable problem.

Ip's article is confined to the use of Venn diagrams for analyzing the coefficient of determination R<sup>2</sup>, partial correlation, and sums of squares. In these cases, exposition is compromised in the presence of suppressor variables. But there are other concepts in thought by many to be of considerably more importance than R<sup>2</sup>, regressic which ar X ng bias and tion of Venn variance these diagram diagram 2. Ar ariance in <u>Kenn</u> + ε. Here the cont the depe S represer d error ε. The usu id that X is fixed in e measured dovi liagramo

labeled y represents "variation" in y, and the circle labeled X represents "variation" in X, where, for pedagogical purposes, "variation" is not explicitly defined but is left as an intuitive concept. The overlap between the y and X circles, the purple area in Figure 1, is interpreted as "variation" that y and X have in common - in this area y and X "move together." This co-movement is used by the OLS formula to estimate  $\beta_x$ , the slope coefficient of X.

Figure 1. Venn diagram for regression.



The purple area is interpreted as information; the black area interpretation is quite different - its magnitude reflects the magnitude of a parameter estimate.

To analyze multiple regression, Kennedy adopted the three intersecting circles diagram of <u>Cohen and Cohen (1975)</u>, which they named the "Ballantine" because of its resemblance to the logo of a brand of beer; Kennedy reinterpreted the areas as described above, and to emphasize this new interpretation changed its spelling to "Ballentine." Such a diagram is shown in Figure 2, in which a new circle marked W is added (with an associated slope  $\beta_w$ ), representing variation in another explanatory variable. In the presentations of <u>Cohen and Cohen (1975)</u> and <u>Ip (2001)</u>, the overlap between the y circle and the X and W circles represents the variation in y explained by variation in X and in W. The ratio of this area (the blue plus red plus green area in Figure 2) to the y circle is interpreted as the R<sup>2</sup> from regressing y on X and W. Trouble happens in the presence of suppressor variables because the red area may have to be negative.

Figure 2. Ballentine Venn diagram.



experience are some particularly effective ways of using this diagram to teach the properties of OLS estimates in the CLR model.

### 3. Teaching some Properties of OLS

I begin by carefully explaining the interpretations of the areas in Figure 1 as described above. Following this I put up Figure 2 on the overhead and ask the class what will happen using OLS when there is more than one explanatory variable, drawing to their attention that it is not obvious what role is played by the red area. I note that if we were to regress y on X alone, OLS would use the information in the blue plus red areas to create its estimate of  $\beta_x$ , and if we were to regress y on W alone OLS would use the information in the green plus red areas to create its estimate of  $\beta_{w}$ . I present three options for the OLS estimator when y is regressed on X and W together.

- Continue to use blue plus red to estimate  $\beta_x$  and green plus red to estimate  $\beta_w$ .
- Throw away the red area and just use blue to estimate  $\beta_x$  and green to estimate  $\beta_w$
- Divide the red information into two parts in some way, and use blue plus part of red to estimate  $\beta_x$  and green plus the other part of red to estimate  $\beta_w$ .



to W, so to be on the safe side we should throw away this information.

• If only the blue area information is used to estimate  $\beta_x$  and the green area information is used to estimate  $\beta_w$ , unbiased estimates are produced, because the blue area corresponds to variation in y uniquely attributable to X and the green area corresponds to variation in y uniquely attributable to W.

An instructor may wish to elaborate on point b by noting that the variation in y in the red area is actually due to joint movements in X and W because the red area corresponds to X and W "moving together" as well as to y and X moving together and y and W moving together. Suppose that in the red area when X changes by one W changes by two, so that a joint movement of one by X and two by W gives rise to a movement in y of  $\beta_x + 2\beta_w$ . If  $\beta_x = 5$  and  $\beta_w = 7$ , this would be a movement of 19. If we were to match this 19 movement in y with a unit movement in X we would get a  $\beta_x$ estimate badly off the true value of  $\beta_x = 5$ . When this is combined with the unbiased estimate coming from the blue area information, a biased estimate results. Similarly, if the 19 movement in y were matched with a two movement in W we would get an estimate of 9.5 for  $\beta_w$ , badly off its true value of 7.

Instructors presenting an algebraic version of this material can demonstrate this result by working through the usual derivation of the OLS estimate of  $\beta_x$  as  $(X^*X^*)^{-1}X^*y^*$ 



Ask the students how a greater degree of multicollinearity would manifest itself on the Venn diagram. They should be able to guess that it is captured by increasing the overlap between the X and W circles, as shown by moving from Figure 3a to Figure 3b. Next ask them if the higher collinearity causes bias. I go around the class and ask everyone to commit to yes, no, or don't know. (Those answering don't know are not asked later to explain their rationale.) Once this has been done, have someone who voted with the majority offer an explanation, and then have someone who voted with the minority offer a counter-explanation. Work with this until everyone sees that because the OLS formula continues to use the information in the blue and green areas to estimate the slopes of X and W, these estimates remain unbiased - the blue and green areas continue to correspond to variation in y uniquely attributable to X and W, respectively.

Figures 3a and 3b. Ballentine Venn diagrams displaying modest and considerable collinearity.



Return to Figure 2, reproduced here as Figure 4, specifying that both X and W affect y, and ask how the properties of the OLS estimate of  $\beta_x$  would be affected if W is omitted from the regression, perhaps because a researcher isn't aware that it belongs, or has no way of measuring it. There are three questions of interest here.

Figure 4. Ballentine Venn diagram.



this voting, the discussion should continue until everyone sees that because the blue plus red area information is used (instead of just the blue area information), more information is being used, so that variance should be smaller. What if the omitted explanatory variable W is orthogonal to X? In this case the OLS estimator continues to use just the blue information, so variance is unaffected.

At this stage the instructor might want to note that by omitting a relevant explanatory variable it should be clear that bias is created, a bad thing, but variance is reduced, a good thing, and comment that the mean square error criterion becomes of interest here because it is a way of trading off bias against variance. A good example to use here is the common procedure of dropping an explanatory variable if it is highly collinear with other explanatory variables. Ask the students how the results developed above could be used to defend this action. They should be able to deduce that omitting a highly-collinear variable can markedly reduce variance, and so may (but may not!) reduce the mean square error.

Third, what can we say about our estimate of  $\sigma^2$  (the variance of the error term)? I ask everyone to commit to unbiased, biased upward, biased downward, or don't know. After this voting, the discussion should continue until everyone sees that the OLS procedure uses the magnitude of the yellow plus green area to estimate the magnitude of the yellow area, so the estimate will be biased upward. The instructor can follow up by asking if X gonal to X. In summ oefficient ance of the estimate error ter timation remains nated. 3.3 De trend is Supp removed detrended X? This but I find it he students better to obtain s and that the OLS slop 's of ogged real quarterl Article contents

CANSIM (<u>www.statcan.ca/english/CANSIM</u>) data base; data values have been logged to reflect the functional form specification usually adopted in this context.

They are then told to perform the following three estimations, following which they are asked to use the Ballentine Venn diagram to explain their results:

- 1. Regress y on X and W to obtain  $\beta_x$ , and its estimated variance vb<sup>\*</sup>.
- 2. Regress X on W, save the residuals r, and regress y on r to get c<sup>\*</sup>, the estimate of the r coefficient, and its estimated variance vc<sup>\*</sup>.
- 3. Regress y on W, save the residuals s, and regress s on r to get d<sup>\*</sup>, the estimate of the r coefficient, and its estimated variance vd<sup>\*</sup>.

With their data my students obtain results reported in Table 1.



Students are surprised that these three estimates  $b^*$ ,  $c^*$ , and  $d^*$  are identical to six

decimal	×	this. First,
b <sup>*</sup> is the		e area in
Figure 5		he orange
plus blue		g the y
circle on		the same
informat		d third, s is
the part		ea. The
overl		yellow on
the o		e same
informat		to the other
two.		
Figure 5		
	Polatod rocoarch	



Trouble begins when they try to explain the estimated variance results. Because exactly the same information is being used to produce  $b^*$ ,  $c^*$  and  $d^*$ , they should all have exactly the same variance. But in Table 1 the three numbers are different. Students react to



the dependent variable not explained is the yellow area, so in this case  $\sigma^2$  is also estimated by the magnitude of the yellow area. But in the case of estimating c<sup>\*</sup> (by regressing the entire y circle on orange plus blue) the variation in y not explained is the yellow plus red plus green areas. As a result, in this case  $\sigma^2$  is overestimated. This overestimation causes vc<sup>\*</sup> to be larger than vb<sup>\*</sup> and vd<sup>\*</sup>.

Instructors may wish to supplement this explanation by noting that the formula for the variance of an OLS estimator involves both  $\sigma^2$  and variation in the explanatory variable data which is "independent" of variation in other explanatory variables. (In this example variance would be given by the formula  $\sigma^2(X^tM_wX)^{-1}$ .) In all three cases here, the "independent" variation in X is reflected by the blue plus orange areas, so the relative magnitudes of the estimated variances depend entirely on the estimates of  $\sigma^2$ .

How does all this relate to regressing on detrended data? If W is a time trend, then s is detrended y and r is detrended X, so that regressing s on r produces estimates identical to those of regressing on raw data. One concludes that it doesn't matter if one regresses on raw data including a time trend, or if one removes the linear trend from data and regresses on detrended data. Similarly, if W is a set of quarterly dummies, it doesn't matter if one regresses on raw data plus these dummies, or if one regresses on data that have been linearly deseasonalized. More generally, this reflects the well-known result that clone estimates are identical using raw data or appropriately

residuali	×
4. Con	
Ine Ball	kposits the
applicati	ing an
irrel	ble
estima	n-nested
hypothe	attention to
this use	f using this
diagram	
There do	ample, the
Ballentir	ng y on X
Article contents	Related research

from regressing X on W, obtain v by taking the residuals from regressing W on X, and then regress y on r and v. My experience with its use in the classroom has been overwhelmingly positive, however; confined to standard analyses, the advantages of this Venn diagram interpretation as a pedagogical device are too powerful to ignore.

## References

1. Cohen, J., and Cohen, P. (1975), Applied Multiple Regression/Correlation Analysis for the Behavioral Sciences, Hillside, NJ: Lawrence Erlbaum Associates.

#### Google Scholar

 Ip, E. H. S. (2001), "Visualizing Multiple Regression," Journal of Statistics Education [Online], 9 (1). (ww2.amstat.org/publications/jse/v9n1/ip.html)

Google Scholar



Related research (

People also read	Recommended articles	Cited by 1	
Information for	Open access		
Authors	Overview		
R&D professionals	Open journals	Open journals	
Editors	Open Select	Open Select	
Librarians	Dove Medical Pre	Dove Medical Press	
Societies	F1000Research	F1000Research	
Opportunities	Help and inform	nation	
Reprints and e-prints	Help and contact	Help and contact	
Advertising solutions	Newsroom	Newsroom	
Accelerated publication	All journals	All journals	
Corporate access solutions	Books		

Keep up Register t by email Sigr Copyright Accessibi X

or & Francis Group