

DISSERTATION PART 2

Model Selection with Multicollinear Explanatory Variables

Ella Foster-Molina

Abstract

Multicollinearity can obscure the effect of variables in regression analyses because it inflates standard errors and reduces statistical significance. Yet the effect of multicollinearity does not only depend on the degree of collinearity present; it also depends on how each explanatory variable affects the dependent variable. I use the collinearity between education and income to demonstrate three cases. In the first, the independent effect of each variable is statistically significant despite high collinearity. This is the best case scenario for a scholar interested in understanding the impact of education and income, or any other collinear variables. In the second case, exactly one of the collinear variables is statistically significant. In the third, a joint significance test rejects the null hypothesis that both coefficients are zero. Yet the independent impact of each is obscured because the model cannot reject the null hypothesis that either are differentiable from zero. I discuss two methods that respecify the model in a way that clarifies the joint effect of education and income. Both methods require reinterpreting the impact of omitted variables in the models. I then present a decision rule for determining how to handle multicollinearity in each case, as well as a method to visualize the impact of multicollinearity and omitted variable bias. I apply this method to a series of examples demonstrating the countervailing effects of education and income in range of political outcomes.

Contents

1	INTRODUCTION	4
2	EXISTING APPROACHES to MULTICOLLINEARITY	6
3	MULTICOLLINEARITY and OBSCURED RESULTS	10
4	REINTERPRETING OMITTED VARIABLE BIAS	12
5	MODEL SELECTION with COLLINEARITY	15
6	COLLINEARITY: EDUCATION and INCOME	21
7	EXAMPLE: LEGISLATIVE SUCCESS	23
7.1	Model Selection	32
8	ADDITIONAL EXAMPLES of EDUCATION and INCOME in POLITICS	41
8.1	Party Elected to Congress	42
9	CONCLUSION	58
	References	61

1 INTRODUCTION

The nature of research into class relations and politics often involves explanatory variables that are highly collinear with each other. Some examples include:

- partisanship \sim income + education
- educational outcomes \sim race + poverty
- turnout \sim unemployment + education

The goal is to interpret both the statistical and substantive effect of each collinear variable. Yet multicollinearity can interfere with this goal unless careful attention is paid to model selection.

Multicollinearity can obscure the effect of variables in regression analyses because it inflates standard errors and reduces statistical significance. Yet the effect of multicollinearity does not only depend on the degree of collinearity present; it also depends on how each explanatory variable affects the dependent variable. One solution to multicollinearity is to interpret the combined effect of the variables, but this is not always necessary. If it is necessary, it must be done with care. It is easy to end up misinterpreting the meaning of the combined impact. It is commonly misinterpreted as the independent impact of a single variable instead of the combined impact of all omitted multicollinear variables. This misinterpretation is commonly called omitted variable bias. Given a goal of interpreting statistical significance and substantive meaning, there is no one-size-fits-all model in the face of multicollinear variables.

I introduce a set of guidelines for model selection when there is joint statistical significance between two collinear variables. Case 1 is the best case scenario, where the large standard errors on the collinear variables is not so high that the model fails

to reject the null hypothesis for the independent impact of both collinear variables. In this case, the full model that includes both collinear variables provides the most information about the impact of collinear variables on the model. Case 2 is only slightly worse, as one variable is statistically significant but not the other. In this case the full model or a model that omits the statistically insignificant variable produces the most information. Of course, omitting the statistically significant variable would confound the remaining variable that was statistically insignificant in the full model. The last case is the trickiest, as the variables are jointly significant but neither are individually significant. I present two models that allow the joint effect to be interpreted for both statistical significance and substantive meaning when they have a theoretically relevant joint impact on the dependent variable.

Choosing between multiple models can be time intensive. In order to facilitate the choice, I develop a method to visualize the impact of multicollinearity and the joint effect of two variables. It also highlights the substantive impact of each collinear variable.

The example used throughout is the effect of education and income on political outcomes, both at the aggregate and individual level. These two variables help clarify more than the mathematical issue with multicollinearity. They also highlight the different, and occasionally opposing, effects of education and income in politics. That is, the well educated do not always behave in the same way as those who have high incomes. For example, among strong partisans, the highly educated tend to be Democratic while those with high incomes tend to be Republican. This is reflected in aggregate political outcomes: districts with many highly educated constituents tend to elect Democrats, while districts with many high income constituents tend to elect Republicans. I hypothesize some reasons why these differences should appear in

politics despite the idea that both income and education are benchmarks for similar types of cultural achievements. I also highlight the fact that any analysis that omits education to focus on income necessarily ends up focusing on the combined impact of education and income. Whether this substantively effects the interpretation of the results depends on whether education has a meaningful effect on the model even though it was not explicitly included. I present multiple examples where it is substantively important to correctly interpret the coefficients in the model as a combined effect.

The goal throughout out this paper is to find methods that allow the researcher to uncover:

- The statistical significance and standard errors for of a theoretically relevant variable or group of explanatory variables in a model.
- The substantive impact of those variables. For example, it is important to know whether a district where the median income is \$100,000 per year is twice as likely to elect a Republican as a district where the median income is \$50,000 per year, or 1% more likely.

2 EXISTING APPROACHES to MULTICOLLINEARITY

The statistical and inferential dangers of multicollinearity have been well examined (Belsley, 1991; Montgomery, Peck, & Vining, 2012). Multicollinearity, also called collinearity, occurs when one variable can be linearly predicted by another set of variables with a substantial degree of accuracy. When all multicollinear variables are included in a model, the standard errors of the estimated coefficients for the collinear

variables increase as the collinearity increases. This can result in situations where the joint effect of the collinear variables is statistically significant, but the individual coefficients are not. Short of collecting more data, there is no statistical solution that would help reveal the independent impact of each variable (Arceneaux & Huber, 2007).

When some multicollinear variables are omitted from the model, confounding can occur. That is, the coefficients on the retained variables no longer reflect the independent impact of each collinear variable. Instead they represent the combined impact of the retained variable with the omitted variables (Obrien, 2007). This combined effect can be substantially different from the independent effect, and can be misinterpreted if the researcher does not carefully interpret the meaning of the coefficient on the retained variable (Dormann et al., 2013).

A variety of remedies have been suggested. They can be divided into two camps: those that improve predictive modeling in the face of multicollinearity, and those that improve hypothesis testing. I will touch on the solutions proposed by those concerned about predictive modeling, as that is where the bulk of the remedies for multicollinearity exist. Yet the focus of this paper is on hypothesis testing, and the solutions used by the predictive modeling community do not adequately address the needs of those who are focused on hypothesis testing.

A brief note on the difference between these two priorities will help clarify why certain solutions do not apply to those engaged in hypothesis testing. Predictive modelers are trying to extrapolate their models to new data. This requires precision for the overall model instead of precision about importance of individual variables within a model. These differences are significant. Predictive modelers focus on mean squared errors and other measures of goodness of fit. This helps them predict new

observations more accurately. Hypothesis testers are trying to determine whether specific variables affect the dependent variable in a meaningful way. For example, the researcher may want to know if high income respondents in a survey are more likely to support tax cuts. In order to determine this, they will look for statistical significance and a meaningfully large change in support for tax cuts for a given change in income. They may also consider the 95% confidence interval for the size of the effect to determine how important the effect is at the upper and lower edge of confidence. Solutions that do not help improve statistical significance or interpreting the magnitude of the effect of a given variable will not help a hypothesis tester.

The predictive modeling camp has focused on improving the goodness of fit of the model in the face of multicollinearity. They focus on techniques such as ridge regression and LASSO (Dormann et al., 2013). I choose not to use these techniques for multiple reasons. For one, these techniques induce biased estimates with lower variances. Because of this bias, computing standard errors and p-values is not recommended (Park & Casella, 2008). As an extra complication the coefficients have, by design, reduced magnitudes (Dormann et al., 2013). Finally, I focus on cases where the multicollinear variables are jointly significant because the null hypothesis that $\beta_1 = \beta_2 = 0$ is rejected. In this case it is possible to observe theoretically meaningful, statistically significant, and interpretable magnitudes for either the joint impact of the multicollinear variables, or even the independent impact of these variables. None of this is possible for ridge regression and related techniques because that is not their aim.

Within the hypothesis testing camp, the remedies for multicollinearity focus on dropping certain multicollinear variables or combining them into a single variable that represents some version of the joint impact of all multicollinear variables

(Dormann et al., 2013; Graham, 2003; Obrien, 2007). As I will show, even omitting a variable has the effect of creating a coefficient that reflects its joint impact with the retained collinear variables. Proxying clusters of multicollinear variables with a specific variable is a well defined method for choosing which variables to omit (Dormann et al., 2013). As with all methods that omit certain multicollinear variables, the remaining coefficients implicitly include the joint impact of all omitted variables.

Other techniques explicitly combine the multicollinear variables in ways that reflect their joint impact using techniques like principle component analysis, residual regressions, or structural equation modeling. Principle component analysis uses a set of linear transformations to create a set of latent variables (Graham, 2003). Residual and sequential regressions create a sequence of new variables that combine the original variables in theoretically meaningful ways. Specifically, the first new variable in the sequence will reflect the independent impact of the first original variable as well as the joint effect of all other multicollinear variables. The second new variable will reflect the independent impact of a second original variable, as well as the joint impact of all original variables that have not already been absorbed by the first new variable. The procedure continues for as many new variables as there are multicollinear variables. Structural equation modeling uses theoretical grounds to develop a combination of the collinear variables (Dormann et al., 2013; Graham, 2003).

I will focus on two methods for interpreting the combined effect of collinear variables: omitting a variable and explicitly combining the collinear variables using principal component analysis. For clarity, my example model has only two collinear variables with a theoretically meaningful joint effect on the model.

3 MULTICOLLINEARITY and OBSCURED RESULTS

Multicollinearity is well known to distort the effects of variables when two or more collinear variables are included in a model.¹ Less well known are the limitations of that distortion. Multicollinearity can obscure the true effect of a variable by artificially inflating the standard errors. It cannot create false statistically significant results. The standard methods for dealing with inflated standard errors is to combine the variables into one, or omit at least one variable. This model choice must be made based on theoretical considerations, and the resulting coefficients must be interpreted with care.

One misconception is that the $\hat{\beta}$ estimator, ie the regression coefficients, can be biased by collinearity. This is false. The proof that $\hat{\beta}$ is unbiased does not require the dependent variables be uncorrelated. It will require that the error terms have a mean of 0, are identically distributed, and are uncorrelated, but it will not put any requirements on the collinearity of the matrix of independent variables, X . That is, the expected value of $\hat{\beta}$ is identical to the true value of β regardless of collinearity.

The confusion about collinearity's effect on regression estimates comes from a separate yet potentially confusing proof: the least squares distance between the true values of β and the $\hat{\beta}$ values is inflated when collinearity is present. It is an artifact of calculating distance, which involves squaring the difference between the true value beta and the estimated beta hat coefficients. This can be interpreted as indicating the unsquared $\hat{\beta}$ is too large, but it is not. Multicollinearity does not

¹The manual for the statistical package NCSS says: "Multicollinearity can create inaccurate estimates of the regression coefficients, inflate the standard errors of the regression coefficients, deflate the partial t-tests for the regression coefficients, give false, non-significant, p-values, and degrade the predictability fo the model (and that's just for starters)." (NCSS, n.d.)

cause $\hat{\beta}$ estimates to be biased.

The primary influence of multicollinearity is on the confidence intervals. This affects the precision of the estimates, not the accuracy.

In effect, severe multicollinearity inflates the variances of the regression coefficients, and this increases the probability that one or more regression coefficients will have the wrong sign (Montgomery et al., 2012, p. 121).

Of course, the prediction about a coefficient having the wrong sign is only valid when the variances are inflated. That is, the coefficient may be both not statistically significant and have the wrong sign. Collinearity alone will not produce a statistically significant coefficient with the wrong sign. If the results are statistically significant, then it is at least as correct as any effect seen in the absence of multicollinearity.

As Belsley succinctly put it:

Thus, if an investigator is only interested in whether a given coefficient is significantly positive and is able, even in the presence of collinearity, to accept that hypothesis on the basis of the relevant t-test, then collinearity has caused no problem. Of course, the resulting forecasts or estimates may have wider confidence intervals than would be needed to satisfy a more ambitious researcher, but for the limited purpose of the test of significance initially proposed, collinearity has caused no practical harm. These cases serve to exemplify the pleasantly pragmatic philosophy that collinearity doesn't hurt so long as it doesn't bite (Belsley, 1991, p. 73).

The best case scenario for a researcher is the one in which multicollinearity does not bite. Although the statistical significance of the estimated coefficients may be inflated, the accuracy of $\hat{\beta}$ is not impacted and the null hypothesis is rejected. The researcher can estimate the substantive effect of the variables under the normal guidelines of hypothesis testing.

4 REINTERPRETING OMITTED VARIABLE BIAS

Another effect of multicollinearity is that omitting a collinear variable can confound the results. As I will show, this confounding can range from mildly altering the size of the coefficients on the remaining multicollinear variables to changing the sign and statistical significance. There are two ways to deal with this confounding. The first is to not drop any collinear variable, which is particularly valuable if multicollinearity does not bite and the coefficients are statistically significant. If multicollinearity bites, the second solution may be appropriate if it is theoretically grounded (Obrien, 2007). If the joint effect of the collinear variables are meaningful, the researcher can omit one or more of the multicollinear variables. Critically, they must then reinterpret the coefficients on the remaining variables as the joint effect of each retained multicollinear variable with the omitted variable. When these coefficients are not interpreted correctly, omitted bias is said to have occurred.

To show how this coefficients can be reinterpreted, I examine the simplest case where there are only two explanatory variables in the model. Greene (2003) discusses the more general case. Let the true relationship be represented by Equation 1, where the β terms are the coefficients of the explanatory variables X_1 and X_2 . Y is the dependent variable, c is the constant and μ is the error term.

$$Y = \beta_1 X_1 + \beta_2 X_2 + c + \mu \tag{1}$$

If X_1 and X_2 are collinear, then X_1 can be linearly predicted by X_2 with a substantial

degree of accuracy. This is be represented by Equation 2.

$$X_2 = \beta_{1a}X_1 + d + \epsilon \quad (2)$$

The effect of dropping X_2 from Equation 1 can be demonstrated by substituting Equation 2 into Equation 1. This substitution is presented in Equation 3.

$$Y = (\beta_1 + \beta_{1a}\beta_2)X_1 + (\beta_2d + c) + (\beta_2\epsilon + \mu) \quad (3)$$

Let the new coefficient on X_1 be

$$\beta_{OV} = \beta_1 + \beta_{1a}\beta_2 \quad (4)$$

where OV stands for the omitted variable model. The approach used here highlights that β_{OV} is not the direct effect of the independent impact of X_1 on Y . Instead, it is the sum of the direct effect, β_1 plus the indirect effect of X_2 , β_2 , channeled through the correlation between X_1 and X_2 , β_{1a} . Note that the independent effect of the omitted variable X_2 is captured in both the new coefficient on X_1 and in the error term $\beta_2\epsilon + \mu$.

This also makes it clear how confounding can happen. If the magnitude of $\beta_{1a}\beta_2$ is large relative to β_1 , then the influence of β_2 can move the value of β_{OV} far from the value of β_1 . If the sign of β_2 or β_{1a} opposes the sign of β_1 , then the sign of β_{OV} can be different from the sign of β_1 . This is a substantial problem if β_{OV} is assumed to be β_1 , as it will lead to the incorrect inference that the independent effect of X_1 on Y is positive when the true value of β_1 is negative, or vice versa. This misinterpretation is also referred to as omitted variable bias (Clarke,

2005). Yet if considered properly as distinct from β_1 and if the joint impact of the multicollinear variables is theoretically meaningful, β_{OV} is interpretable for both statistical significance and substantive meaning.

This coefficient is interpreted as the sum of the direct effect of X_1 on Y and the indirect effect of X_2 on Y through the relationship of X_1 and X_2 . In less technical terms, it is the joint effect of both X_1 and X_2 on Y . In a more complex case, it would be the joint effect of all dropped variables on the remaining variables in the model.

To restate, when a variable, call it X_2 , is omitted from the model but is collinear with a retained variable X_1 , the coefficient of the retained variable is no longer reflective of its independent effect on Y . If the assumption is that the coefficient reflects this independent effect, then it is fair to say that the coefficient is biased away from the independent effect. Yet the coefficient is not biased away from the joint effect of X_1 and X_2 on Y ; it is an unbiased estimate of Equation 4, $\beta_1 + \beta_{1a}\beta_2$ (Clarke, 2005). Thus, a solution to multicollinearity biting is to drop a variable and reinterpret the coefficient on the remaining variable as the joint effect of the dropped and retained variable on Y .

The decision of which variable to keep can be made based on which variable is more theoretically relevant. Yet it is critical to interpret the coefficient on the retained variable as the combined effect of both variables. Specifically, it the sum of the direct effect of the retained variable and the indirect effect of the omitted variable.

Another solution to finding the joint effect of collinear variables is to use principal component analysis. This will find one or more common latent variables for the collinear variables. This is useful under two conditions:

1. The effect of the collinear variables is thought to be due to common underlying causes, called common latent variables (Graham, 2003). For example, a desire to achieve success in society is an underlying cause of both high income and high educational attainment, and can cause a greater desire to influence politics. Thus, desire for cultural success is a latent cause of political engagement.
2. The latent variables are interpretable. The standard is to use the first set of principle components that account for between 80 and 90% of the variance in the data (Vajargah, n.d.; Dormann et al., 2013). However, if the primary goal is to understand the joint impact of the collinear variables,² only those components that have a relevant interpretation should be used. The simplest case is one in which the first component explains over 80% of the variation and has substantive meaning, so no other components need to be included.³

Two other methods for creating a combination of multicollinear variables were discussed in Section 2: residual and sequential regression, and structural equation modeling (Graham, 2003). At their core, they are all variations on estimating coefficients that reflects the joint impact of all or some of the multicollinear variables.

5 MODEL SELECTION with COLLINEARITY

Multicollinearity can make it difficult to interpret the magnitude of effect of multicollinear variables that are jointly statistically significant yet individually statistically insignificant. I present a guideline for selecting a model that provides information about the significance and substantive importance of the jointly significant

²If predictive power is also desired, then more components will help improve the R^2 value.

³By design, the components extracted from principle component analysis are orthogonal to each other; that is, they are uncorrelated with each other.

multicollinear variables, while avoiding misinterpreting coefficients as the independent impact of a given variable when it is in fact the joint impact of that variable combined with omitted collinear variables. For simplicity, I focus on the case where only two variables in the model have substantial collinearity.

Let X_1 and X_2 be two highly correlated variables. Let X_3 be a set of controls that do not have substantial multicollinearity with each other or X_1 and X_2 . Assume that the null hypothesis $\beta_1 = \beta_2 = 0$ is rejected. I present guidelines for model selection in three cases, each drawn from one of the three possible combinations of statistically significant coefficients for X_1 and X_2 in the following model:

$$Y \sim \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 \tag{5}$$

This guideline to model selection allows the magnitude of effect of both X_1 and X_2 , either individually or combined, to be interpreted under traditional rules of statistical significance.

Case 1

Scenario The independent effects of both variables are statistically significant despite high collinearity. This is the best case scenario for a scholar interested in understanding the impact of education and income, or any other collinear variables.

Model Selection The model should include both highly collinear variables.

$$Y \sim \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3$$

Omitting one would require reinterpreting the meaning of the remaining

variables, as the coefficients will now reflect the combined impact of both variables. Such an interpretation may be theoretically relevant in some cases, but should not be the default method.

Case 2

Scenario One of the two coefficients, β_1 , is statistically significant and the other, β_2 , is not.⁴

Model Selection One of the following two models should be used:

$$Y \sim \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 \quad (6)$$

$$Y \sim \beta_1 X_1 + \beta_3 X_3 \quad (7)$$

The interpretation of the coefficients in these two models remains important. The full model in Equation 6 includes both collinear variables. This provides the benefit of being able to interpret the independent effect of the statistically significant variable. The reduced model in Equation 7 drops the insignificant variable X_2 . Because statistical insignificance does not imply that β_2 is small, the new β_1 has to be interpreted to include the combined effect even in the case where X_2 is insignificant. If β_2 is small, then the joint effect reflected in the reduced model β_1 will be close to the value of β_1 in the full model.

Case 3

Scenario The independent impacts of both X_1 and X_2 are obscured because the model cannot reject the null hypothesis that either β_1 or β_2 is differ-

⁴A similar method applies if β_2 is statistically significant but not β_1 .

entiable from zero.

Model Selection In this scenario, it is statistically impossible to differentiate the independent effect of X_1 from X_2 . Yet it is still possible to derive meaningful interpretations, so long as the goal is to understand the joint impact of both variables. Three standard ways to do this are as follows:

- Use a model that omits X_1 :

$$Y \sim \beta_2 X_2 + \beta_3 X_3 \quad (8)$$

The key is that β_2 represents the combined impact of X_1 and X_2 , and should be interpreted as such.

- Use a model that omits X_2 :

$$Y \sim \beta_1 X_1 + \beta_3 X_3 \quad (9)$$

The interpretation is similar to dropping X_1 .

- Use principal component analysis, or another similar technique, to create a new variable X_{12} that combines the impact of X_1 and X_2 .

$$Y \sim \beta_{12} X_{12} + \beta_3 X_3$$

This will typically be a linear combination of X_1 and X_2 . This offers the benefit of being explicit that β_{12} is a combined effect, which can be difficult to remember when simply omitting either X_1 or X_2 . Principal component analysis offers the additional benefit that $\hat{\beta}_{12}$ may be more precise than either $\hat{\beta}_1$ in Equation 9 or $\hat{\beta}_2$ in Equation 8. This is because the weights in the linear combination of X_1 and X_2

that create X_{12} are reflective of the transformation that accounts for as much of the variability between X_1 and X_2 as possible.

Which of the three cases applies depends on the magnitude and direction of each variable in the model, which is impossible to know prior to using the model. Thus, when high multicollinearity is present, multiple models may be examined and compared against each other.

In the best case scenario, Case 1 applies. In Case 1, no additional models are needed as the full model provides enough information to interpret the independent effects of both X_1 and X_2 under the traditional rules of statistical significance.

In Case 2, it is appropriate to compare the estimates $\hat{\beta}_1$, $\hat{\beta}_2$, and $\hat{\beta}_{12}$ from equations (1) and (2) below. In Case 3, it is appropriate to compare the estimates from all four of the following models:

$$Y \sim \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 \quad (10)$$

$$Y \sim \beta_1 X_1 + \beta_3 X_3 \quad (11)$$

$$Y \sim \beta_2 X_2 + \beta_3 X_3 \quad (12)$$

$$Y \sim \beta_{12} X_{12} + \beta_3 X_3 \quad (13)$$

The model in Equation 10 must be run simply to determine that Case 1 does not apply, so the model selection guidelines from Case 2 or 3 must be used. The coefficients $\hat{\beta}_1$, $\hat{\beta}_2$, and $\hat{\beta}_{12}$ estimated from Equations 11, 12, and 13 all reflect the combined impact of both X_1 and X_2 . Under conditions of high collinearity, the magnitude of each coefficient⁵ will be comparable if X_1 and X_2 are measured on

⁵They could have opposing signs depending on whether X_1 and X_2 are positively or negatively correlated with each other

the same scale.⁶ The benefit from using the coefficient that represents the effect of the combination of both X_1 and X_2 is that the magnitude of the coefficient can be interpreted, which is not possible by simply using a joint significance test.

Since comparing coefficients for multiple variables across multiple models is not simple, it is very useful to have tools that help visualize and understand the results. In the examples presented below, I develop a visualization that allows the reader to quickly examine both the magnitude of the effect of the coefficient as well as the statistical significance. The statistical significance is relevant to determine how multicollinearity impacted the various models, and therefore which models to are appropriate. Showing the magnitude of the effect allows the reader to examine the impact of omitted variable bias, as well as the substantive effect of each explanatory variable.

In order to allow the magnitudes to be easily compared between different variables, I focus on the expected change in Y caused by a change in X instead of the coefficient β . While it is possible to plot the coefficients themselves on the same graph, the coefficients for one variable are likely to be drastically bigger or smaller than coefficients for the other variable. Therefore I compare the expected Y value for a high value of X against the expected Y value for a low value of X . I choose to use the third quartile of X as a proxy for a high value, and the first quartile as a proxy for a low value. Using this method, the units being compared across different variables is the same: a change in the Y value. An added benefit is that this method highlights whether the β value is substantively meaningful.

The following sections apply these guidelines and implement these graphics using the collinearity between education and income.

⁶In the discussion of visualizing the coefficients, I demonstrate a method to compare the coefficients of variables that are not measured in the same units.

6 COLLINEARITY: EDUCATION and INCOME

I show evidence from four political outcomes that the independent impact of education . Analyzing each on its own has been the standard, and has produced many useful results. Yet including one or both in a regression analysis without sufficient understanding of the ways their collinearity can impact the interpretation of the model can easily lead to misinterpreted effects. Specifically, the following misinterpretations are common:

- If education is omitted from the model but has an independent impact, then the coefficient for income must be interpreted as the joint effect of income and education.
- If income and education are retained in the model but are not statistically significant, they may be misinterpreted to have no effect on the dependent variable. Yet to truly determine this, their joint significance should be examined. If they are jointly significant, it can be substantively meaningful to interpret that joint effect.

To show this I examine data on legislative success, partisan affiliation, and replicate results on country turnout rates (Burden & Wichowsky, 2014) and ideology (McCarty, Poole, & Rosenthal, 2006). I reveal insights by examining education and income with a better theoretical grounding on the conflicting requirements of multicollinearity and omitted variable bias.

Unlike income, race, and gender, education does not have a strong line of scholarly study. It is so closely tied to income, which has been rightly viewed as the predominant driver of political outcomes in the modern era, that it seems to have been dismissed as inconsequential. Not many studies include it in their models, and

when they do it is rarely analyzed in any detail. (McCarty et al., 2006) show that Republicans have lost many educated voters while they have gained high income voters over the past few decades. This could imply that they have ceased to be responsive to the educated as well, and this is precisely what I find in the third part of this dissertation. Viewed from the perspective of national legislative outcomes, as the responsiveness of Republicans to high income districts has increased, their responsiveness to the educated has decreased. That is, Republicans are more successful in legislation for richer districts, and less successful in educated districts.

There are theoretical grounds to believe that education and income will have disparate effects on policymaking. While Republicans tend to win the votes of the rich, Democrats are known as the party of Ivy League intellectuals (Ansolabehere, Rodden, & Snyder, 2006; Gilens, 2012). Republicans are thought to be more responsive to economic interests of the rich, while Democrats are more responsive to the intellectual elite. Therefore Democrats should respond to the highly educated more than Republicans.

Yet it is also plausible that education and income have the same effect in other areas of politics. High education and high income are both indicators of achievement and engagement in society. Individuals engaged in society can be expected to be more engaged in politics and more likely to demand policy congruence from their representatives. There are many questions to be answered about when the effect of income and education should have similar impacts on politics, and when they should be different.

In practice, the lack of theoretical grounding for the independent impact of education relative to income on political outcomes means that the potential independent impact of education needs to be carefully explored whenever income is the

primary explanatory variable. If both variables are included in the model and are statistically significant impact, then there are not interpretation issues due to the collinearity between these two variables. If both are included and neither are individually statistically significant, it is possible to draw the incorrect conclusion that they have no meaningful impact on the model. If education is omitted, it is possible to incorrectly assume the coefficient on income reflects the independent impact of income instead of the joint impact of income and education. In the worst case, this can lead to conclusions that income has a positive (negative) effect on the dependent variable, when in fact it has a negative (positive) and statistically significant impact. I show an example of this in voter turnout in a later section.

The datasets I use reveal the confounding influence of education on income in a variety of political contexts. When they do not confound each other, it is relatively common for them to be jointly statistically significant yet not individually significant. The effects of the collinearity between education and income occasionally clash with the need to avoid omitted variable bias. I demonstrate one useful method to untangle the disparate effects of income and education, and apply it to replications of McCarty et al. (2006) and Burden and Wichowsky (2014). I show that educational and economic characteristics of a district have different effects on legislative outcomes, partisan preferences, and representational preferences.

7 EXAMPLE: LEGISLATIVE SUCCESS

The primary example I use theorizes a connection between district demographics, constituent preferences over policy, and legislative outcomes. The theoretical connection between district demographics and legislative activity lies in both representational and policy preferences. People prefer policy from their members of

Congress when they are richer, as seen in Part 1. Republicans provide policy that is more congruent with the preferences of their constituents when those constituents have high incomes (Rhodes & Schaffner, 2017; Grossmann & Williams, 2018; Lax, Phillips, & Zelizer, 2018). Republicans are the party ideologically aligned with the interests of the rich, while Democrats are the party ideologically aligned with the interests of the educated and working class. Thus, Republican members of Congress who represent high income districts should find it easier to create successful legislation, as the legislation supported by their party and their ideology is congruent with the preferences of their constituents. Part 3 of this dissertation provides additional evidence for the mechanism behind this link.

The purpose of this section is to discuss and apply methods for detecting collinearity, implement the model selection criteria laid out in this part of the dissertation, and present the visualization of the statistical significance and substantive effect of income and education in each model. I use two examples from the data to highlight the model selection criteria. In the first, drawn from data on Democratic legislative behavior from the 1980s, I show a situation where Case 1 applies. The graphics will clearly show why using the combined impact obscures the independent effect of both income and education. The second example is drawn from Republican legislative behavior in the 2000s. Here, Case 3 applies. Neither income nor education are independently significant, but they are jointly statistically significant. In this case, it is possible to draw meaningful inferences about the joint effect of income and education.

Additional examples provided in the next section will examine situations where Case 2 applies, as well as the different ways income and education can impact political outcomes.

Data

The primary independent variables are district income, district education, and a combined measure of district socioeconomic status. District income is measured as the percent of a legislator’s district that earned over \$75,000 per household per year in 2016 dollars. This captures the percent of the district that is high income. The cutoff for this is not always perfectly \$75,000, as the value changes according to inflation and the income brackets used by the census. It is always in the range of \$65,000-\$75,000 in inflation adjusted 2016 dollars, and always falls above the median income of the nation at the time. The results are consistent across a variety of income measures. The secondary independent variable is district education, which is measured as the percent of the district with at least a bachelor’s degree. All parts of the analyses separate Republicans from Democrats because the ideologies of each party create different kinds of responses across district economic and educational levels. In order to explicitly capture the joint impact of income and education, I find the principle component behind education and income and call it socioeconomic status (SES). The socioeconomic variable is a linear combination of education and income. This method assumes a latent variable behind district education and district income that causes representatives to be more successful legislatively.

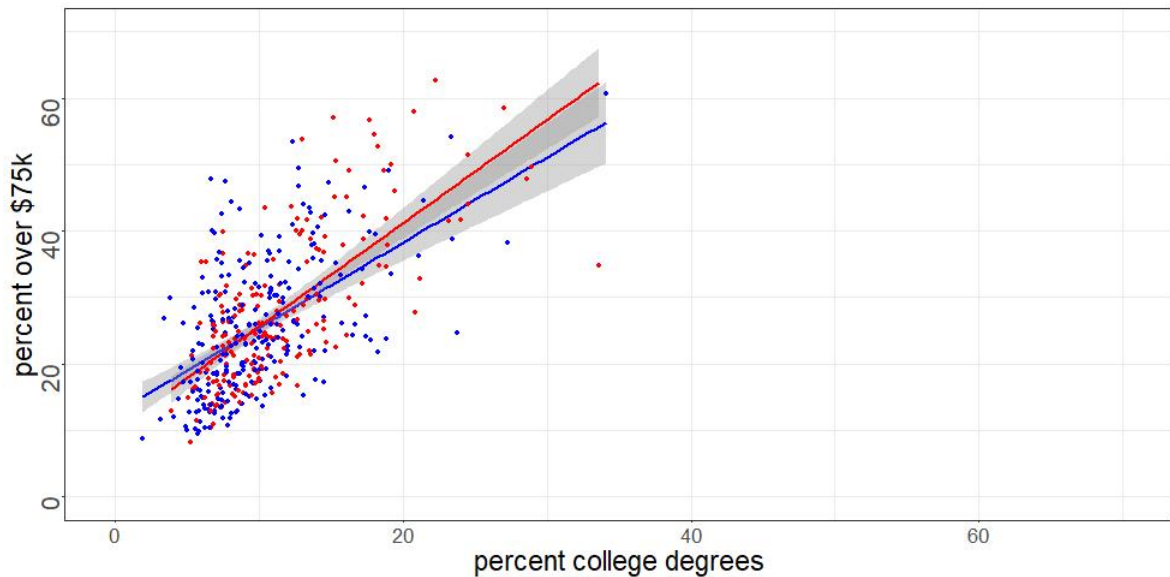
The primary dependent variable is the number of bills a representative sponsors that are approved of in a House vote.

Detecting Collinearity

As can be seen in Figure 1 and Figure 2, income and education are strongly correlated. The graph plots the district income and education for each representative in the 97th and 113th Congresses. Other Congresses have very similar correla-

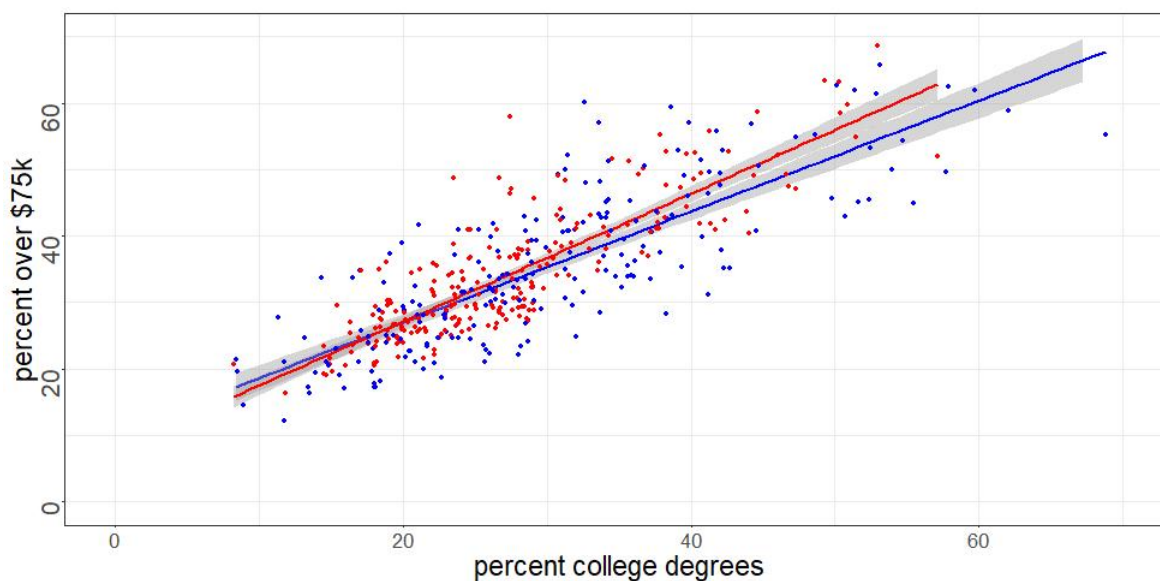
tions. Red dots represent districts represented by a Republican, and blue represents Democrats. The regression line is also depicted, in blue for Democrats and red for Republicans. Clearly, there is a correlation between district education and income for both Democrats and Republicans. This collinearity affects the ways omitted variable bias appears. As discussed in the prior section, it fortunately cannot artificially create statistically significant results in a model that includes all multicollinear variables.

Figure 1: Collinear Relationship between Education and Income, 1981-82



These two Congresses were chosen because they show the underlying correlation between district education and income for two different cases for model selection. For Democrats between 1981 and 1986, Case 1 holds. For Republicans between 2011 and 2014, Case 3 holds. One pattern of note in these two graphs is that the range for spread of college degrees has substantially increased since 1981. The percent of a district with college degrees ranges between 15 and 65% today. In 1981 and 1982, is ranged between 3 and 35%. Another pattern of note is the

Figure 2: Collinear Relationship between Education and Income, 2013-2014



fact that the tie between district income and educational attainment is higher today than it was in the past. These trends are consistent across all Congresses in the dataset.

The first example examines Democratic legislative behavior between 1981 and 1986. I will show how to determine that Case 1 applies to this example using diagnostics for multicollinearity. That is, the diagnostics will show that there is not much multicollinearity for this data. However, the key diagnostic is not whether the collinearity diagnostic tests indicate a problem, but whether the regression models fail to reject the null hypothesis because of multicollinearity.

Three diagnostic tools are commonly used to assess multicollinearity: the condition index, variance inflation factor, and a perturbation of the Y values to see if the coefficients substantially change. Table 1 and Table 2 depict these three different diagnostics for collinearity. They indicate that Case 1 applies. That is, the coefficients on income and education are both statistically significant in the full model

Table 1. Democratic 1981-1986 Collinearity Diagnostics: Condition Index and Variance Inflation Factor

Condition Index	Variance Decomposition Proportions						Weights
	intercept	income	education	year	conservativeness	seniority	percent black
1.000	0.000	0.001	0.002	0.000	0.003	0.005	0.000
3.390	0.000	0.010	0.007	0.000	0.001	0.007	0.668
4.483	0.000	0.012	0.027	0.000	0.000	0.850	0.002
6.701	0.001	0.000	0.144	0.001	0.347	0.028	0.002
7.066	0.000	0.016	0.200	0.000	0.645	0.102	0.188
10.272	0.000	0.830	0.255	0.001	0.000	0.005	0.121
43.091	0.026	0.002	0.034	0.062	0.003	0.000	0.013
123.566	0.973	0.129	0.332	0.937	0.000	0.002	0.001

Democratic 1981-1986 Variance Inflation Factor					
income	education	year	conservativeness	seniority	percent black
1.320129	1.510961	1.288464	1.087829	1.021415	1.099920

Table 2. Democratic 1981-1986 Collinearity Diagnostics: Perturbations

Impact of Perturbations on Coefficients				
	mean	std deviation	minimum coefficient	maximum coefficient
Intercept	1.843	0.074	1.674	2.035
income	-0.018	0.001	-0.016	-0.019
education	0.016	0.001	0.017	0.014
year	-0.023	0.001	-0.025	-0.022
conservativeness	-0.299	0.006	-0.312	-0.284
seniority	0.103	0.000	0.102	0.103
percent black	-0.003	0.000	-0.004	-0.003

that includes both, so both should be retained in the final model selected.

The first diagnostic is the condition index, which is a measure of the degree to which the principle components of a variance-covariance matrix are unequal (Montgomery et al., 2012). Different sources provide slightly different cutoff values: Montgomery et al. (2012) sets the cutoff for high collinearity at 100, while the R package Perturb sets it at 30. Collinearity affects those variables that have a high variance decomposition proportion, set at 0.50 by Belsley (1991), or 50% of the variance inflation. These values are bolded in Table 1. While the condition index does rise above the cutoff values, it does not affect both education and income at the same time. Therefore, education and income do not have sufficient collinearity for Democrats between 1981 and 1986 to obscure the results. Instead, the high condition index is influenced by the intercept and the year. Collinearity between the intercept and year has no meaningful effect on this analysis.

The variance inflation factor is calculated from the correlation matrix (Montgomery et al., 2012). Again, the threshold for high collinearity varies, but it is common to set the cutoff at values higher than five or ten. None of the variance inflation factors

exceed this cutoff. This supports the conclusion that the collinearity in the model does not meaningfully obscure the results.

The most useful diagnostic comes from the Perturb package in R. This package induces small changes to the variables to see if they unduly influence $\hat{\beta}$. It is considered the best test to see how much collinearity affects regression results. The results are presented in Table 2. It presents the lower bound of the 95% confidence interval for each coefficient, as well as the upper bound. If the coefficient does not change signs, it is a strong indication that collinearity does not affect the direction of the coefficients. For this group of observations, Democrats between 1981 and 1986, collinearity does not appear to influence the direction of the effects of income or education, nor meaningfully change the magnitude of the coefficient.

Overall, these diagnostics indicate that for Democrats between 1981 and 1986, Case 1 will apply. The model that best reflects the underlying effect of income and education is likely to be one that includes both. Yet these are merely indications, not conclusive.

Turning to an example of Case 3, I examine the same group of variables for Republicans between 2011 and 2014. As I will show, the diagnostics indicate that the collinearity between income and education bites in this case. Here, district income and education are both associated with increased legislative success. However, this time the effect of income (education) controlling for education (income) are artificially obscured by the inherent collinearity between the two variables.

The impact of collinearity can be seen in the diagnostics for Republicans between 2011 and 2014, shown in Table 3 and Table 4. While the variance inflation factor does not signal any major issues for the coefficients on income and education, the condition index and the perturbations do. For the condition index of 24.271,

Table 3. Republican 2011-14 Collinearity Diagnostics: Condition Index and Variance Inflation Factor

Condition Index	Variance Decomposition Proportions						Weights
	intercept	income	education	year	conservativeness	seniority	percent black
1.000	0.000	0.000	0.000	0.000	0.001	0.005	0.000
3.829	0.000	0.000	0.000	0.000	0.001	0.392	0.000
4.477	0.000	0.002	0.003	0.000	0.007	0.533	0.001
8.405	0.000	0.055	0.088	0.000	0.079	0.014	0.009
12.983	0.000	0.000	0.008	0.000	0.863	0.045	0.043
24.271	0.000	0.936	0.898	0.000	0.017	0.011	0.006
30.667	0.001	0.004	0.000	0.001	0.031	0.000	0.933
848.367	0.999	0.002	0.001	0.999	0.001	0.000	0.003
							0.008

Variance Inflation Factor					
income	education	year	conservativeness	seniority	percent black
1.958373	1.939159	1.003454	1.023099	1.034675	1.010547

Table 4. Republican 2011-14 Collinearity Diagnostics: Perturbations

	Impact of Perturbations on Coefficients			
	mean	std deviation	minimum coefficient	maximum coefficient
Intercept	-18.229	0.188	-18.632	-17.780
income	0.002	0.002	-0.005	0.008
education	0.013	0.003	0.006	0.021
year	0.166	0.002	0.162	0.169
conservativeness	-0.502	0.010	-0.526	-0.480
seniority	0.050	0.000	0.050	0.051
percent black	-0.025	0.000	-0.025	-0.024

both education and income have very high variance decomposition values, well over the 0.50 or 50% threshold. This is reflected in the values of the coefficients when small perturbations are introduced, as described by Table 4. Here, the sign of the coefficient on income in the model that has a control for education flips between -0.005 and 0.008. Thus, collinearity can be expected to effect the sign of the coefficient for income in addition to inflating the standard errors.

These results indicate that the coefficient on education, and possibly that for income, will not be statistically significant. Yet, as with the previous diagnostics, they are not conclusive. The only conclusive test for whether multicollinearity will obscure the independent impact of variables is if the multicollinear variables are jointly significant but individually insignificant.

7.1 Model Selection

The goal is to understand the statistical significance and substantive impact of income and education. When possible, their independent effects are examined. When this is not possible because multicollinearity prevents rejecting the null hypothesis

that their independent effects are zero, the statistical significance and substantive impact of their joint impact can be examined. The four models that reveal these effects are:

1. Model 1 includes both income and education. Select this model in Case 1 or Case 2, where the independent effects of one or both of them are statistically significant in Model 1.

$$\text{legislative success} \sim \beta_1 * \text{income} + \beta_2 * \text{education} + \beta_3 * \text{controls}$$

2. Model 2 includes income while omitting education. This model can be used in Case 2 or Case 3, but β_1 must be interpreted as the joint effect of income and education.

$$\text{legislative success} \sim \beta_1 * \text{income} + \beta_3 * \text{controls}$$

3. Model 3 includes education while omitting income. This model can be used in Case 2 or Case 3, but β_2 must be interpreted as the joint effect of income and education.

$$\text{legislative success} \sim \beta_2 * \text{education} + \beta_3 * \text{controls}$$

4. Model 4 replaces both income and education with a socioeconomic status variable created from income and education through principle component analysis. This model can be used in Case 2 or Case 3, but β_{12} must be interpreted as the joint effect of income and education.

$$\text{legislative success} \sim \beta_{12} * \text{SES} + \beta_3 * \text{controls}$$

The diagnostics above indicated that the association between district demographics and legislation for Democrats between 1981 and 1986 would land in Case 1, where the full model is appropriate. They also indicated that a similar association for Republicans would land in Case 3, where the coefficients on income and education are jointly, but not individually, statistically significant. Model 1 is used to check this along with joint significance tests.

The results for Democrats presented in Table 5 show that the coefficients for income and education are statistically significant in Model 1. Thus, Case 1 applies and no further models are required.

Table 5. Collinearity Examples, Democrats 1981-86

	Model 1	Model 2	Model 3	Model 4
income	-0.015* (0.006)	-0.006 (0.005)		
education	0.025** (0.009)		0.014 (0.008)	
SES				0.001 (0.004)
conservativeness	-0.730* (0.305)	-0.820** (0.307)	-0.601* (0.299)	-0.714* (0.306)
seniority	0.112*** (0.012)	0.110*** (0.012)	0.114*** (0.012)	0.112*** (0.012)
percent black	-0.068* (0.031)	-0.061* (0.031)	-0.057 (0.031)	-0.056 (0.031)
Black Caucus	-0.375 (0.216)	-0.420 (0.217)	-0.330 (0.216)	-0.379 (0.218)
comm. chair	0.591*** (0.163)	0.575*** (0.165)	0.600*** (0.165)	0.586*** (0.165)
powerful comm.	-0.673*** (0.112)	-0.687*** (0.112)	-0.678*** (0.112)	-0.686*** (0.112)
Observations	768	768	768	768
Akaike Inf. Crit.	2,550.349	2,555.514	2,553.602	2,556.678

Note:

*p<0.05; **p<0.01; ***p<0.001

Table 6. Collinearity Examples, Republicans 2011-14

	Model 1	Model 2	Model 3	Model 4
income	0.004 (0.009)	0.016*** (0.005)		
education	0.016 (0.010)		0.020*** (0.005)	
SES				0.011*** (0.003)
conservativeness	-0.603* (0.280)	-0.534 (0.278)	-0.615* (0.279)	-0.579* (0.278)
seniority	0.026* (0.012)	0.024* (0.012)	0.027* (0.011)	0.025* (0.011)
percent black	-0.212*** (0.052)	-0.207*** (0.053)	-0.213*** (0.052)	-0.210*** (0.052)
comm. chair	0.901*** (0.144)	0.911*** (0.145)	0.896*** (0.144)	0.907*** (0.144)
powerful comm.	-0.231* (0.116)	-0.204 (0.115)	-0.237* (0.115)	-0.220 (0.114)
Observations	480	480	480	480
Akaike Inf. Crit.	1,614.676	1,615.063	1,612.844	1,613.087

Note:

*p<0.05; **p<0.01; ***p<0.001

Democrats who represented highly educated districts produced more policy, and they produces less when they represented high income districts. This is congruent with the modern policy platforms of the Democratic party. It is also congruent with modern voting behavior: the highly educated tend to vote Democratic, while the rich tend to vote Republican (Gelman, 2009; McCarty et al., 2006).

The results presented in Table 6 show that multicollinearity bites. As seen in Model 1, neither coefficient for education or income is statistically significant. Indeed, the standard error on income is substantially larger than the estimated coefficient. Yet they are jointly significant, as seen in both an F-test and the coefficients for the joint impacts visible in Models 2-4. This means that Case 3 applies, and it is appropriate to use any of Models 2-4 to interpret the joint effect of income and education on legislative success.

The substantive interpretation of the results in Table ?? is that the highly educated and high income districts tend to have Republican representatives who are more legislatively successful between 2011 and 2014. It entails that Republicans have an easier time producing legislation when their constituents are some combination of highly educated and high income. The independent impacts of income and education are impossible to uncover with the given data. It is possible that with more data, both the independent impacts of income and education would gain statistical significance. Yet it is also possible that a latent variable captured by income and education is causing both high income districts, high education districts, and legislative productivity. If that is the case, then they may not have a true independent impact, merely a joint impact.

The regression results presented provide sufficient information to select a model that uncovers the joint and/or independent effects of income and education. Yet it

requires examining information from multiple models and coefficients. A graphical solution to highlight the relevant information would help streamline the model selection process. This is what I present in the next subsection: a graphic that quickly reveals that statistical significance and size of each coefficient for income and education in the four models. This graphic not only helps with model selection in the face of multicollinearity, but it also highlights the substantive impact of both district education and income on legislative success.

Visualizing Effects of Collinearity

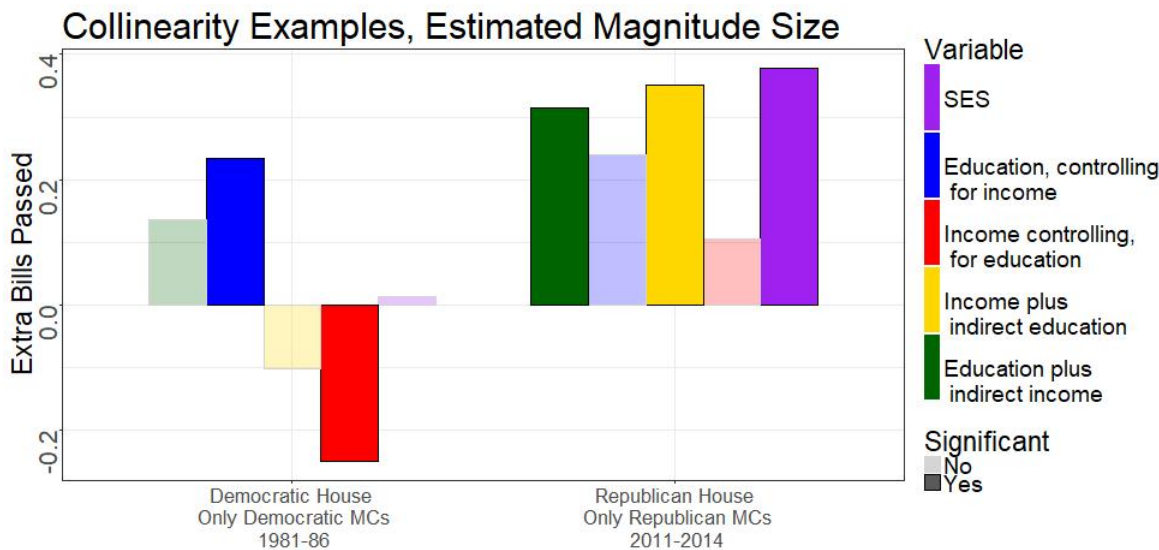
Figure 3 graphically summarizes these two examples of the effects of collinearity. The y-axis shows the estimated number of extra bills passed by an individual member of Congress based on the socioeconomic characteristics of their district. This magnitude is estimated by comparing the number of successful bills sponsored by a legislator who represents the top quartile of the independent variable versus the number of successful bills sponsored by a legislator who is statistically identical except for representing a district at the bottom quartile. For example, the purple bar on the far right of Figure 3 shows that Republicans in the Republican controlled House between 2011 and 2014 who represent a district at the top quartile of the socioeconomic distribution sponsor 0.36 more successful bills than an equivalent representative from a district at the bottom quartile of the socioeconomic distribution. The effect is statistically significant, so the bar is bolded and outlined in black.

Each bar represents the magnitude for one of five coefficients derived from Models 1-4: the combined impact of district income and education channeled through education (green), education controlling for income (blue), the combined impact of

district income and education channeled through income (yellow), income controlling for education (red), and a combined socioeconomic measure (purple). Note that the estimates based on the education variables are blue and green. The estimates based on the income variables are yellow and red. The socioeconomic variable is purple, which is the color wheel result from combining of blue and red.

Another benefit of this graphic is that the magnitudes of effects calculated from the coefficients for education and income in each model are comparable. In a typical regression table, coefficients often cannot be compared because they are measured in different units. For the graphic, I calculated the expected effect of each explanatory variable on the dependent variable. The result is that all magnitudes are measured in the same units as the dependent variable: number of successful bills sponsored. This allows the effect of district educational attainment to be directly compared against the effect of district income.⁷

Figure 3: Effects of Collinearity



⁷As it so happens, both income and education are measured as a percentage of the district, so the units are the same. Yet it is unusual for variables to be measured in the same units.

The colors are faded when they are not statistically significant at the $\alpha = 0.05$ level, and bolded with a black outline when they are statistically significant. For example, for Democrats in the Democratic Houses between 1981 and 1986, the coefficient for education when income has been controlled for, in blue, is statistically significant. So is the coefficient for income when education is controlled for, in red. The other three bars are not statistically significant. Specifically, the coefficients for education without controlling for income (green), income without controlling for education (yellow), and overall socioeconomic status (purple) are not statistically significant. The opposite is true for Republicans in Republican Houses of 2011-2014.

Each group of five bars was chosen as a case study in the two main ways collinearity can effect the results. The left hand group of five bars shows the magnitude of the effect for each of the five coefficients for 1981-86, focusing on Democrats. The House was controlled by Democrats at this time. The right hand group of five bars shows the magnitudes of each coefficient for 2011-14, focusing on Republicans in the Republican controlled House.

One take away from Figure 3 is that for Democrats between 1981 and 1986, Case 1 applies. Multicollinearity does not bite and the independent impacts of district income (red) and education (blue) are statistically significant. As I discuss in the third part of this dissertation, the magnitudes described are also substantively meaningful because each Democratic member of Congress in this time period passed an average of only 1.72 bills per Congressional session. Thus, the blue and red bars represent the coefficients of Model 1, which provides the most information regarding the impact of education and income on legislative success.

The faded bars on the left all represent the combined impact of income and education. They also provide an object lesson for why the coefficients in Models 2

and 3 must be interpreted as the joint effect of income and education. The faded yellow bar on the left is negative, but much closer to zero than the bold red bar. This indicates that when the joint impact of income and education is channeled through income, the combined effect appears to be negative but is statistically insignificant. Yet the joint effect of income and education, when channeled through education, is positive and not statistically significant. Clearly, the magnitude of the effect the combined measure represented by the green bar is lower than the independent effect of education because the coefficient for the combined measure incorporates the countervailing effect of the income.

Another take away from Figure 3 is that for Republicans between 2011 and 2014, Case 3 applies. Multicollinearity bites, and although income and education are jointly significant, it is impossible to distinguish their independent effects (faded red and blue). Yet the combined effect represented by the green, yellow, and purple bars is still substantively meaningful and statistically significant. The direction of the effect and statistical significance for all three are identical, and the magnitudes of effects are similar. Indeed, the magnitudes of the effects are larger than those evident for the independent impacts of income and education for Democrats on the left hand side.⁸

This graphic summarize a large amount of information to help determine a useful model specification in the face of collinearity: statistical significance, direc-

⁸A word of caution is relevant for comparing coefficients across groups. If the distribution of the number of bills passed is substantially different for Republicans between 2011 and 2014 versus Democrats in 1981-1986, then the differences between the groups will be not be due to the differences caused by income and education, but due to the differences due to the group. Fortunately, both time periods focus on the majority party, and members of the majority party have similar rates of legislative success across time and party. To use a separate example, it would not be meaningful to compare how weight changes the heart rate of a mouse as compared to how it changes the heart rate of a giraffe. Yet comparisons may still be possible for a different measure, such as the percentage change in heart rate.

tion of effects, and effect sizes for five coefficients derived from four models for each of two groups. Again, the goal is to capture the statistical significance and substantive importance of education and income. This part of the dissertation shows that this is possible so long as the joint impact of income and education is statistically significant. It demonstrates a method to uncover the substantive effect by appropriately interpreting coefficients based on whether they capture the independent or joint effects of income and education.

8 ADDITIONAL EXAMPLES of EDUCATION and INCOME in POLITICS

The remainder of this part of the dissertation will examine three other political outcomes in which the collinearity between education and income influences the interpretation of the coefficients in Models 1-4. These are the party elected to Congress based on district education and income levels, how conservative those members of Congress are, and voter turnout. There are some theoretical expectations through these examples, but the primary purpose is to show the application of the model selection procedure and to demonstrate that failing to accurately account for the confounding influence of education can lead to incorrect inferences. Specifically, I show that

- Failing to include education in a model predicting whether Republicans are elected can lead to the incorrect inference that district income does not appear to change who is elected. The correct inference is that district income increases the probability of electing a Republican, but the combined impact of income and education has no statistically significant effect on which party is elected.

- Failing to include education in a model predicting the ideology of members of Congress can lead to the incorrect inference that district income is associated with more polarized members of Congress. The correct inference for is that highly educated districts are associated with more polarized members of Congress in the time period examined, but the independent effect of income is not statistically significant.
- Failing to include education in a model predicting voter turnout can lead to the incorrect inference that counties with high incomes or low unemployment rates have lower levels of turnout. The correct inferences is that counties with high unemployment rates, and correspondingly low incomes, have higher levels of turnout.

8.1 Party Elected to Congress

There has been some debate over the influence of income on which party is elected to Congress. This section will show that if the effect of education is not controlled for, then it will appear as if income has no statistically significant impact on the party elected. That is, the countervailing independent effect of education is substantial enough shrink the combined effect of income and education down toward zero.

This example uses the same dataset examined in the prior section, but this time the focus is on which party is elected to Congress based on district demographics. In order to determine which model best captures the effect of district education and income, I present the results of Models 1-4 in Table 7 and Figure 4.

Model 1 shows that income (red bar) and education (blue bar) have statistically significant independent impacts on which party is elected to Congress. Thus, Case 1 applies and Model 1 should be used. The magnitude of the coefficient is on the

Table 7. Probability Republican Elected 2007-2010

	Model 1	Model 2	Model 3	Model 4
district income	0.051*** (0.012)	0.012 (0.006)		
district education	-0.053*** (0.013)		-0.006 (0.007)	
district SES				0.002 (0.004)
percent black	-0.290*** (0.071)	-0.315*** (0.070)	-0.334*** (0.070)	-0.327*** (0.070)
Observations	896	896	896	896
Akaike Inf. Crit.	1,156.620	1,171.687	1,174.016	1,174.582

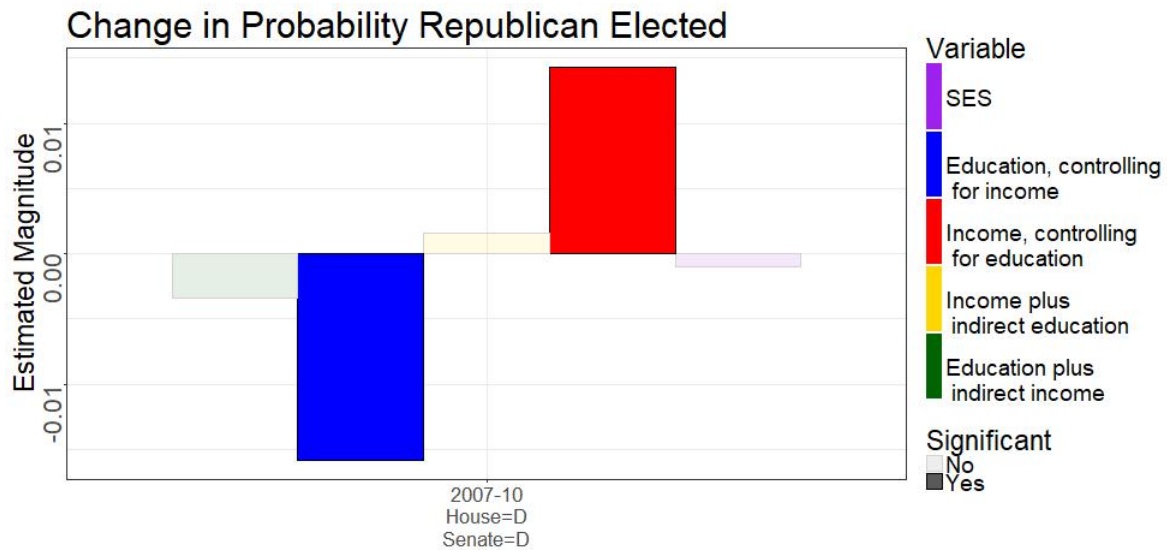
Note: Logistic regression

*p<0.05; **p<0.01; ***p<0.001

small side. A district at the top quartile of educational attainment (blue bar) is 1.6% more likely to elect a Democrat to Congress than an equivalent district at the bottom quartile of educational attainment. Meanwhile, a district at the top quartile of the district income (red bar) spectrum is 1.4% more likely to elect a Republican. This means that we could expect to see around just over one fewer Republicans elected from districts at the bottom end of the educational attainment spectrum, and a bit over one more elected from districts at the top end of the district income spectrum.

Yet although Model 1 best reflects the impacts of education and income, Models 2-4 still provide some interesting insights. They highlight the fact that the combined impact of income and education has no statistically significant effect on the party of elected to Congress. That is, the three measures of the joint impact of income and education (faded green, faded yellow, and faded purple) are not statistically significant. This shows that if education is excluded from the model, the true positive effect of income on the probability a Republican is elected will not be

Figure 4: Effect of District Demographics on Party Elected to Congress 2007-2010



visible.

One additional result from Table 7 is that, unsurprisingly, districts with large numbers of black constituents are unlikely to elect a Republican. Finally, AIC indicates that the Model 1 is a more informative model.

Ideology over Time

Another example of the effects of the collinearity between education and income comes from replicating a study on the influence of district income and education on a representative's conservativeness. I examine how members of Congress vote on all bills by examining their first dimension DW-Nominate scores. DW-Nominate scores are a measure developed in Poole and Rosenthal (1997), and are commonly used to capture how conservative or liberal a member of Congress is.

This section replicates and expands on results from McCarty et al. (2006). As they show, district income and education have clearly distinguishable effects on

the ideology of members of Congress. They also show that high income districts have become increasingly conservative over the past 60 years. By using the methods presented in this section of the dissertation, I show that this effect is due to Republicans taking control of Congress and because Republicans now lean more conservative when they represent highly educated districts.

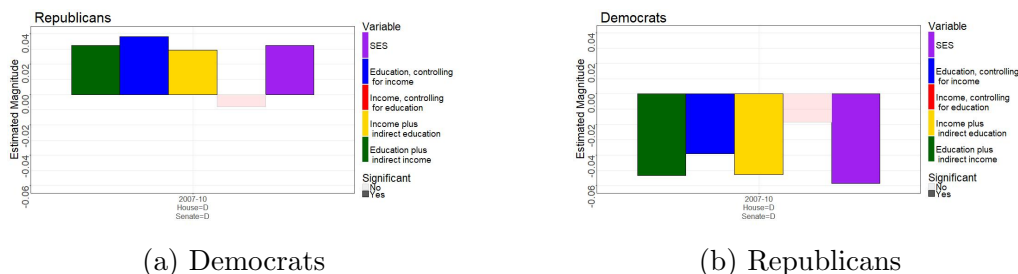
The dependent variable is the ideology of legislators, as measured by first dimension DW-Nominate scores. The more conservative a legislator, the closer his or her ideological score is to 1. The more liberal legislators have ideological scores close to -1. The economic and educational independent variables all range between 0 and 1, so their coefficients can be compared relatively directly.⁹

Before I examine the changing impact of district education and income over time, I will focus on a simpler case from 2007-2010. These results are presented in Figure 5. It is clear from the graphic that Case 2 applies: in Model 1, represented by the blue for the independent impact of education and faded red for the independent impact of income, only one is statistically significant. Case 2 indicates that either the full model, represented by the red and blue bars in the graphic, or the model that omits the variable for income, represented by the green bar, should be selected. If the model that omits income is selected, then the coefficient on the education variable (green) must be interpreted as the joint effect of income and education.

The graphics show that Republicans are more conservative (positive values) when they represent both highly educated and jointly highly educated and high income districts. Democrats are more liberal under the same conditions. In other words, highly educated districts have more ideologically polarized representatives,

⁹The combined factor for socioeconomic status, SES, was normalized to range between 0 and 1. Education and income are both percentages that inherently are bounded between 0 and 1. Education reflects the percent of the district that has a college degree, and income reflects the percent of the district that earns over ~\$75,000 in 2009 inflation adjusted dollars.

Figure 5: Impact of Socioeconomic Characteristics on Conservativeness 2007-2010



as do districts that are both highly educated and high income. Nothing can be said about the independent impact of district income, as it is statistically insignificant in both models.

Note that it would be easy to choose the model that omits education (yellow), especially since education has not widely been theorized to be a more important factor in political outcomes than income. The graphic shows what happens if Model 2 is selected (yellow). The yellow bar shows that the combined impact of income and education is polarizing. Yet it would be easy to misinterpret this variable as the independent impact of income and conclude that high income districts have polarized representatives.

I now turn to examining the changing impact of district demographics across time. McCarty et al. (2006) show that connection between district income and legislator conservativeness has increased over the past 40 years. They argue that this is partially due to an increase in magnitude of the coefficient, but also because districts themselves are facing larger inequality.

Tables 8 and 9 show this is also due to the fact that Republicans are now more conservative when they represent districts with high socioeconomic status, not because all legislators have become more conservative when representing districts with high socioeconomic status. Democrats are still more liberal when they represent

high socioeconomic status districts, just as they were in the 1970s and 1980s.¹⁰

Table 8. Republican Ideology by District Demographics 1972-2014

	Model 1	Model 2	Model 3	Model 4
income	0.014* (0.006)	-0.015*** (0.004)		
income*year	-0.0001* (0.0001)	0.0001*** (0.00004)		
education	-0.050*** (0.008)		-0.036*** (0.005)	
education*year	0.0005*** (0.0001)		0.0003*** (0.0001)	
SES				-0.015*** (0.003)
SES*year				0.0001*** (0.00003)
year	0.024*** (0.001)	0.025*** (0.001)	0.023*** (0.001)	0.024*** (0.001)
majority	0.015* (0.007)	0.022** (0.007)	0.020** (0.007)	0.022** (0.007)
seniority	-0.008*** (0.001)	-0.008*** (0.001)	-0.008*** (0.001)	-0.008*** (0.001)
percent black	-0.010*** (0.002)	-0.010*** (0.002)	-0.010*** (0.002)	-0.010*** (0.002)
comm. chair	-0.001 (0.014)	-0.003 (0.014)	-0.002 (0.014)	-0.003 (0.014)
powerful comm.	-0.023*** (0.006)	-0.021*** (0.006)	-0.022*** (0.006)	-0.022*** (0.006)
Observations	4,134	4,134	4,134	4,134
Akaike Inf. Crit.	-3,205.318	-3,168.685	-3,202.619	-3,185.130

Note: OLS, *p<0.05; **p<0.01; ***p<0.001

Figure 6 highlights the changing effects of income, socioeconomic status, and education. As in the prior graphics, each bar in the graphics represent one measure of socioeconomic status from a regression that looked at that group of legislators. This time, the height of the bar reflects the size of the coefficient instead of the magnitude of the effect. For example, the green bar for Democrats between 1973 and 1980 shows the coefficient for education on Democrats ideology in that time

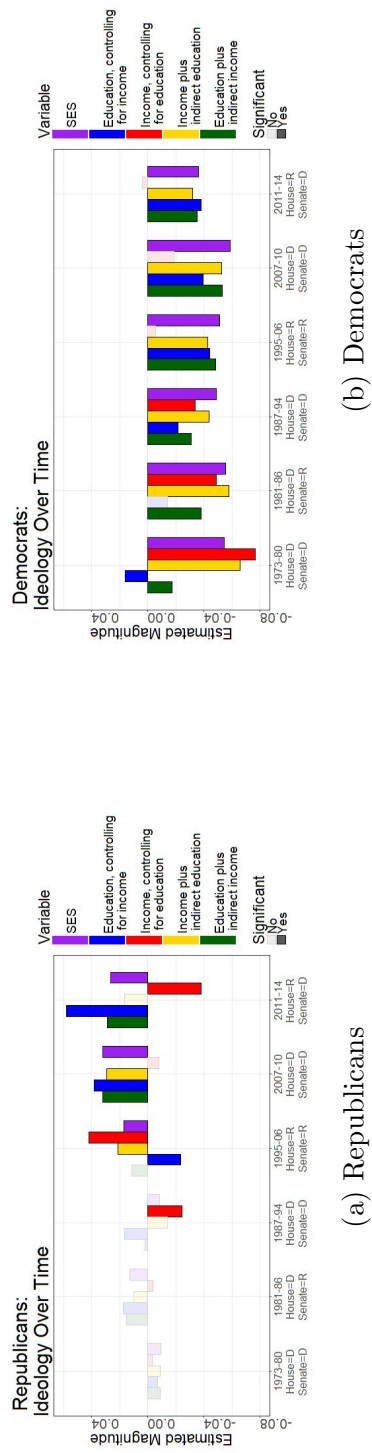
¹⁰The effect has declined somewhat for Democrats, but it is still strongly negative.

Table 9. Democratic Ideology by District Demographics 1972-2014

	Model 1	Model 2	Model 3	Model 4
income	-0.037*** (0.004)	-0.022*** (0.003)		
income*year	0.0003*** (0.00004)	0.0002*** (0.00003)		
education	0.015** (0.005)		-0.010* (0.004)	
education*year	-0.0002** (0.0001)		0.0001 (0.00004)	
SES				-0.016*** (0.002)
SES*year				0.0001*** (0.00002)
year	-0.007*** (0.001)	-0.007*** (0.001)	-0.001 (0.001)	-0.004*** (0.001)
majority	-0.031*** (0.006)	-0.043*** (0.006)	-0.022*** (0.006)	-0.033*** (0.005)
seniority	-0.005*** (0.0005)	-0.005*** (0.0005)	-0.005*** (0.001)	-0.005*** (0.0005)
percent black	0.017*** (0.002)	0.017*** (0.002)	0.022*** (0.002)	0.019*** (0.002)
Black Caucus	-1.653*** (0.113)	-1.611*** (0.112)	-1.447*** (0.113)	-1.530*** (0.111)
comm. chair	-0.024** (0.009)	-0.024** (0.009)	-0.023* (0.009)	-0.024** (0.009)
powerful comm.	-0.032*** (0.005)	-0.034*** (0.005)	-0.032*** (0.005)	-0.032*** (0.005)
Black Caucus*year	0.014*** (0.001)	0.013*** (0.001)	0.012*** (0.001)	0.012*** (0.001)
Observations	5,060	5,060	5,060	5,060
Akaike Inf. Crit.	-5,717.505	-5,661.439	-5,527.723	-5,684.847

Note: OLS, *p<0.05; **p<0.01; ***p<0.001

Figure 6: Impact of Socioeconomic Characteristics on Conservativeness



period, controlling for the all non-economic or educational variable in Table 9.¹¹ So between 1973 and 1980, Democrats who represented highly educated districts were more liberal.

The impact of district demographics on Democratic ideology is highly consistent across time. Democrats who represent districts with many constituents who are of high socioeconomic status are consistently more liberal, whether socioeconomic status is measured by income, education, or both. The one exception is the effect of education when controlling for income in 1973-1980. For that time period and that measure, Democrats were more conservative when they represented districts with high education levels relative to their income level. The effects are almost always statistically significant. The increasing impact of income on conservative ideology is not due to Democratic behavior.

Once again reflecting the trends in legislative success, the impact of district demographics on Republican ideology is less consistent over time. In fact, prior to the 1994 Republican take over of the House, district socioeconomics had almost no impact on Republican ideology. The fact that high income districts overall create more conservative legislators is entirely driven by the changes over time for Republicans, as well as their newfound control over the House of Representatives. Note, however, that while income (yellow bar), education (green bar), and the combined socioeconomic status variable (purple bar) are always positive and usually statistically significant, the effect sometimes becomes negative when education or income are controlled for. Between 1995 and 2006, the independent effect of education (blue bar) for Republican ideology was negative. That is, for two districts with similar numbers of high income inhabitants, the district with more educated inhabitants

¹¹Majority party is also not controlled for, as this time period was always controlled by Democrats.

would be expected to have a legislator who was more liberal. Similarly, between 1987 and 1994, and between 2011 and 2014, the independent effect of income (red bar) was associated with more liberal Republican legislators. This discrepancy is currently unexplained, and merits future investigation.

Overall, the increasing association between district socioeconomics and legislator ideology is driven by Republicans, not Democrats. We can see this in Figure 6. This graphic highlights the differences between Democrats and Republicans. This is a very similar pattern to the one revealed in the analysis of legislative success, and is consistent with the theory presented by (Barker & Carman, 2012) and the voting patterns described by (Gelman, 2009). Republicans have changed their ideological grounding, particularly since 1990, and that appears to be reflected in the demographic ties to how they vote and how they create legislation.

Note that yet again the disparate effects of income and education for Republican ideology are highly relevant. Generally, Republicans are more conservative when they represent socioeconomically elite districts. Yet the opposite effect can be obscured if they are not examined both together and separately. For example, look at Democrats in a Republican House, as seen in Figure ???. It is impossible to tell whether the effect of income when controlling for education is obscured by collinearity or if the apparent effect of income is entirely due to the collinearity with education. Yet we can see that overall socioeconomic status, as well the stand-alone impact of income, are both tied to more liberal Democrats. As another example, focus on the impact of income for Republicans in a Democratic House. Here, the independent effect of income once education is controlled for shows a correlation with more liberal members of Congress. Yet if income is looked at alone, without controlling for education, it could be misinterpreted to indicate that the district

income is correlated with more conservative Republicans.

A few other results from these tables stand out. The first is that Republicans with large numbers of black constituents tend to be more liberal, while Democrats with large numbers of black constituents tend to be more conservative. This is an unexpected finding. Future work could examine whether this effect is due to Southern Democrats prior to the Republican revolution of 1994. However, a black Democratic member of Congress is, as expected, more liberal than a white Democrat.¹²

In general, Republicans have become more conservative over time. The trend for Democrats is less clear, as it depends highly on which demographic variables are controlled for. Intriguingly, when facing a Republican majority Democrats become more liberal while Republicans become more conservative.

Voter Turnout

The confounding effect of education is also apparent in a replication of the results presented in Burden and Wichowsky (2014). They claim that a common understanding used to be that counties with high levels of unemployment would have lower turnout rates. They make the argument that, once adequate controls have been included in the model, counties with higher turnout rates in fact have lower levels of unemployment. I will demonstrate that a key confounding variable is the educational attainment of the county.

This is another example where Case 1 applies. That is, when the full model is used education and unemployment both have a statistically significant and countervailing effects. I will show that the same holds for an analysis where income is a primary explanatory variable instead of unemployment, and unemployment is

¹²This result may be strongly influenced by whether the number of black constituents in the district is controlled for.

omitted from the model.

The following two tables demonstrate the effect of collinearity between income and education in voter turnout rates. The analysis in the original paper refutes conventional wisdom that county unemployment rates have a negative relationship to voter turnout, and provides a theoretical grounding for the opposite effect. I replicated these results and focused on the impact of education, unemployment, and income.

Table 10 reveals one potential origin of conventional wisdom. Namely, when education is not controlled for as in Model 4, it appears that there is a statistically significant negative correlation between unemployment rates and voter turnout. However, controlling for education reveals that this effect was driven by education rates, and once education is controlled for, as in Models 1-3, the effect of unemployment reverses while remaining statistically significant. That is, counties with high unemployment rates do have lower turnout for voting, but this effect is driven by the fact that counties with low high school graduation rates have both high unemployment rates and low turnout rates. Low educational attainment drives low voter turnout, while high rates of unemployment help increase turnout once educational attainment has been controlled for.

Of note, the effect of the collinearity between education and unemployment is not explicitly addressed in Burden and Wichowsky (2014). They acknowledge that their findings oppose tradition wisdom, but do not explain that traditional wisdom was confounded by omitting the influence of education. No other variable reversed the effect of county unemployment when omitted from the model; for this analysis, the most important collinear variable is education.

Additional analyses show that this effect is primarily created by counties with

Table 10. Voter Turnout by Unemployment, Education, Gubernatorial Election (1980-2008)

	<i>Voter Turnout</i>			
	Full	(2)	(3)	(4)
County unemployment	0.146*** (0.019)	0.187*** (0.018)	0.191*** (0.019)	-0.290*** (0.019)
High school graduation	1.506*** (0.107)	4.475*** (0.060)	4.510*** (0.060)	
Concurrent gubernatorial race	5.328*** (0.299)	1.747*** (0.141)		
State unemployment	0.510*** (0.035)			
Percent black	0.031 (0.017)			
Median income	-0.389 (0.215)			
Competitive presidential race	0.010*** (0.002)			
Concurrent senatorial race	0.634*** (0.060)			
Year fixed effects	Yes	No	No	No
County fixed effects	Yes	No	No	No
AIC	167704.3	203835.5	203986.1	209162.5
Observations	27,899	27,901	27,901	27,901

Note:

*p<0.05; **p<0.01; ***p<0.001

the lowest rates of educational attainment, as shown in Table 11. The interaction effect between unemployment and educational attainment is statistically significant.¹³ Counties with the lowest education levels produce higher turnout when they face high unemployment rates. The effect of unemployment disappears for counties in the top half of educational attainment.¹⁴

Yet their primary independent variable is unemployment, not income. To demonstrate that the effects are similar to the results in this dissertation, these results must be very similar when the analysis focuses on income instead of unemployment. In this case, because median income and county unemployment should produce similar effects on turnout rates, median income and unemployment rates should produce interchangeable statistical results. Table 12 shows that this effect holds when unemployment rates are omitted, instead focusing primarily on income. Namely, low income counties have overall lower turnout rates, but only when education is not accounted for. Low income counties of similar education levels have higher than expected turnout rates.

As the results from Table 10 show, when both income and unemployment are included in the model, only unemployment shows up as statistically significant. That is, unemployment does a better job of explaining turnout than does median

¹³Model with the interaction effect is not shown.

¹⁴Note that multiple variables change magnitude substantially based on educational attainment. County unemployment and median income both become smaller as educational attainment goes up. Larger black populations are associated with reduced turnout in low education counties, and with increased turnout in high education counties. Competitive presidential races are associated with lower voter turnout in low education counties, but are associated with higher voter turnout in high education counties. One possible theory that explains this rests on the educational attainment of Democrats versus Republicans. Highly educated counties will tend to be more Democratic. Democrats tend to promote the participation of those who are poorer and minorities, so highly educated Democratic districts should see an increase in voter turnout when they have more poor and minority members. Yet Republicans tend to fan the anger of unemployed white males, so counties with high unemployment, particularly white male unemployment, should see higher levels of turnout relative to other similar counties.

Table 11. Voter Turnout by Educational Attainment

	<i>Turnout based on county educational attainment quartiles</i>			
	Lowest quartile	2nd lowest quartile	2nd highest quartile	Highest quartile
County unemployment	0.173*** (0.032)	0.109** (0.042)	0.052 (0.041)	0.085 (0.045)
High school graduation	2.169*** (0.270)	1.322*** (0.384)	1.811*** (0.438)	0.983** (0.302)
Concurrent gubernatorial race	4.755*** (0.569)	5.197*** (0.553)	4.387*** (0.593)	3.516*** (0.794)
State unemployment	0.646*** (0.075)	0.708*** (0.076)	0.379*** (0.064)	0.331*** (0.067)
Percent black	-0.037 (0.033)	-0.181*** (0.042)	0.281*** (0.037)	0.420*** (0.040)
Median income	-2.283*** (0.688)	-2.610*** (0.611)	-1.501** (0.457)	0.421 (0.348)
Competitive presidential race	-0.010* (0.005)	-0.0002 (0.005)	0.018*** (0.004)	0.027*** (0.004)
Concurrent senatorial race	0.651*** (0.138)	0.672*** (0.125)	0.540*** (0.098)	0.414*** (0.105)
Year fixed effects	Yes	Yes	Yes	Yes
County fixed effects	Yes	Yes	Yes	Yes
AIC	43092.7	42571.7	39231	40248.8
Observations	6,979	6,975	6,978	6,979

Note:

*p<0.05; **p<0.01; ***p<0.001

Table 12. Voter Turnout by Income, Education, Gubernatorial Election (1980-2008)

	<i>Voter Turnout</i>			
	(1)	(2)	(3)	(4)
Median income	-0.599** (0.213)	-4.433*** (0.127)	-4.513*** (0.127)	1.894*** (0.113)
High school graduation	1.493*** (0.106)	5.723*** (0.069)	5.774*** (0.069)	
Concurrent gubernatorial race	5.335*** (0.300)	1.504*** (0.139)		
State unemployment	0.657*** (0.029)			
Percent black	0.037* (0.017)			
Competitive presidential race	0.010*** (0.002)			
Concurrent senatorial race	0.641*** (0.060)			
Year fixed effects	Yes	No	No	No
County fixed effects	Yes	No	No	No
AIC	167876.9	202901.6	203017	209276.8
Observations	27,919	27,921	27,921	27,921

Note: *p<0.05; **p<0.01; ***p<0.001

income. In more precise terms, unemployment captures almost all of the variance in voter turnout that would otherwise be attributed to median income. Thus, for Burden and Wichowsky, unemployment is a more important explanatory variable than is median income.

The results in Burden and Wichowsky point to yet another example of income related variables creating opposing political effects to education. I suggest that income and education frequently create opposing effects in politics. Any time that the impacts of education and income oppose each other, omitting one can dramatically obscure the independent effect of the other (Clarke, Kenkel, & Rueda, 2016).

9 CONCLUSION

Highly correlated explanatory variables can confound the interpretation of statistical models if not appropriately accounted for. On the one hand, omitting a highly correlated variable can substantially change the interpretation of the retained correlated variables. On the other hand, including all correlated and collinear variables will inflate the standard errors on the independent effects of each collinear variable. Multiple techniques have been proposed to account for these problems. I propose a solution that involves carefully considering all combinations of the collinear variables in models. The choice of the model depends on the joint significance of the variables and whether their independent effects can be examined under traditional rules of statistical significance.

I argue that this model selection criteria will help scholars avoid interpretation problems when faced with multicollinearity. A model that has jointly significant collinear variables with no independent significance can still be interpreted to reveal the effect of their joint impact. The solution is to omit a variable or use some other

technique that combines the impact of the joint variables. A model that omits a highly correlated variable must be interpreted to be the combined effect of the omitted and included variables whenever it is known that the omitted variable has a high correlation with the included variable. This does not account for all the unknown correlated variables excluded from the model, but I argue that when we know about the potential for a confounded model, we should be extra careful about model selection.

In the process of demonstrating this model selection criteria, I highlight the importance of accounting for education in models of political outcomes. Whether the model used examines the joint impact of education and income, or the independent impact of education and income, inferential errors are possible. For example, the analysis of the impact of district income and education on Republican legislative success shows that Case 3 applies. This case is prone to the inferential error that income and education have no statistically significant effect on the model. Yet the correct inference is that education and income have a joint effect that is statistically significant. The magnitude of that joint effect is easily calculated if the appropriate model is selected.

On the other hand, the analysis of the impact of county income on voter turnout shows the danger of analyzing the joint impact of income and education if desired outcome is to understand the independent impact of income. Specifically, without accounting for education, the effect of county income appears to be positive. This fits easily with theories of income and voter turnout, but is not accurate. The correct interpretation of that effect is that the joint effect of income and education is positive. Yet closer examination shows that the independent effect of income is negative. That is, counties with more high income individuals have statistically

significantly lower voter turnout rates. Burden and Wichowsky (2014) show that this is because counties with lower incomes have higher unemployment rates, and high unemployment rates are associated with higher voter turnout. I show that this effect is particularly strong for low education districts, which are more likely to be blue collar and rural.

The examples used throughout only involve two collinear variables. The next step would be to develop a model selection criteria with over two collinear variables. The number of combinations of multicollinear variables that would have to be examined would get large quickly, as the foundation for this model selection criteria was an examination of one model for each possible combination of the two collinear variables.

References

- Ansolabehere, S., Rodden, J., & Snyder, J. M. (2006). Purple America. *The Journal of Economic Perspectives*, 20(2), 97–118.
- Arceneaux, K., & Huber, G. A. (2007). What to do (and not do) with multicollinearity in state politics research. *State Politics & Policy Quarterly*, 7(1), 81–101.
- Barker, D. C., & Carman, C. J. (2012). *Representing red and blue: How the culture wars change the way citizens speak and politicians listen*. Oxford University Press.
- Belsley, D. A. (1991). *Conditioning diagnostics: Collinearity and weak data in regression* (No. 519.536 B452). Wiley.
- Burden, B. C., & Wichowsky, A. (2014). Economic discontent as a mobilizer: Unemployment and voter turnout. *The Journal of Politics*, 76(4), 887–898.
- Clarke, K. A. (2005). The phantom menace: Omitted variable bias in econometric research. *Conflict management and peace science*, 22(4), 341–352.
- Clarke, K. A., Kenkel, B., & Rueda, M. R. (2016). Omitted variables, countervailing effects, and the possibility of overadjustment. *Political Science Research and Methods*, 1–12.
- Dormann, C. F., Elith, J., Bacher, S., Buchmann, C., Carl, G., Carré, G., . . . others (2013). Collinearity: a review of methods to deal with it and a simulation study evaluating their performance. *Ecography*, 36(1), 27–46.
- Gelman, A. (2009). *Red state, blue state, rich state, poor state: Why Americans vote the way they do*. Princeton University Press.
- Gilens, M. (2012). Under the influence. *Boston Review*, July–August.
- Graham, M. H. (2003). Confronting multicollinearity in ecological multiple regres-

- sion. *Ecology*, 84(11), 2809–2815.
- Greene, W. H. (2003). *Econometric analysis*. Pearson Education India.
- Grossmann, M., & Williams, I. (2018). *Oligarchy or class war? Political parties and interest groups in unequal public influence on policy adoption*. (Working paper)
- Lax, J., Phillips, J., & Zelizer, A. (2018). *The party or the purse? Unequal representation in the US Senate*. (Working paper)
- McCarty, N., Poole, K. T., & Rosenthal, H. (2006). *Polarized america: The dance of ideology and unequal riches (vol. 5)*. MIT Press.
- Montgomery, D. C., Peck, E. A., & Vining, G. G. (2012). *Introduction to linear regression analysis* (Vol. 821). John Wiley & Sons.
- NCSS. (n.d.).
- Obrien, R. M. (2007). A caution regarding rules of thumb for variance inflation factors. *Quality & Quantity*, 41(5), 673–690.
- Park, T., & Casella, G. (2008). The bayesian lasso. *Journal of the American Statistical Association*, 103(482), 681–686.
- Poole, K. T., & Rosenthal, H. (1997). *Congress: A political-economic history of roll call voting*. Oxford University Press on Demand.
- Rhodes, J. H., & Schaffner, B. F. (2017). Testing models of unequal representation: Democratic populists and Republican oligarchs? *Quarterly Journal of Political Science*, 12(2), 185–204.
- Vajargah, K. F. (n.d.). Comparing ridge regression and principal components regression by monte carlo simulation based on mse.
