

Welcome!

Nameplates please. And technology encouraged today!

All TF materials are available at github.com/nolankav/api-202.

If you want to follow along, download the dataset here:

In R: df <- read.csv ("http://tinyurl.com/api-202-tf-1")

In Excel: http://tinyurl.com/api-202-tf-2

EXCEL

What's the deal with regression?

API 202: TF Session 1

Nolan M. Kavanagh
January 26, 2024



Goals for today

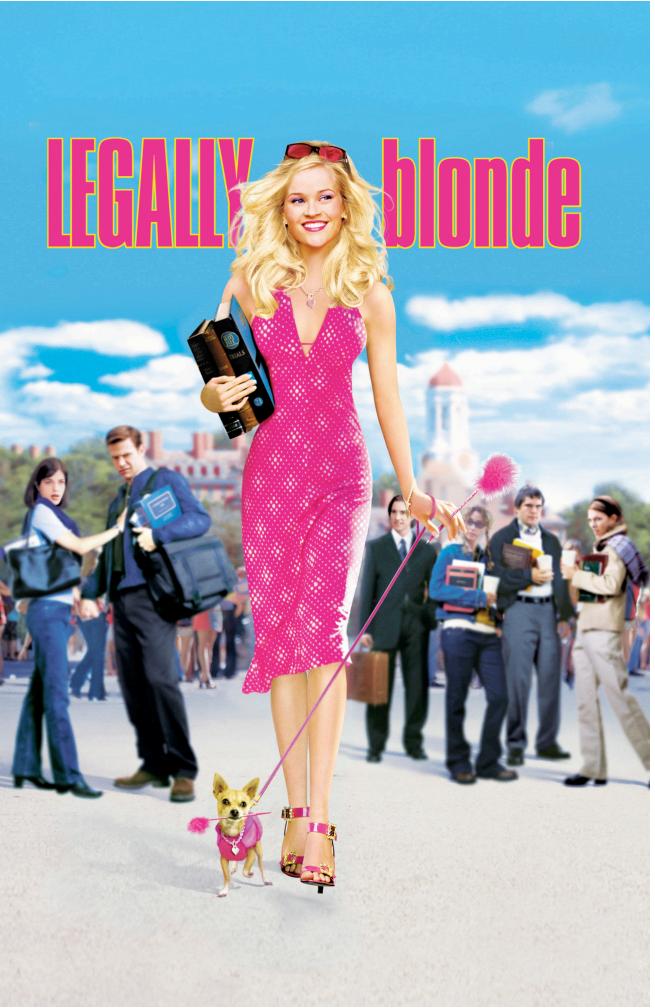
- 1. Get to know one another a little better.**
- 2. Review regression notation, including the PRF vs. SRF.**
- 3. Learn how to graph bivariate relationships.**
- 4. Learn how to run regressions.**
- 5. Review how to interpret regressions.**

We'll treat this session like a workshop with interactive examples.



MD/PhD student in health policy.

**“I don’t need backups.
I’m going to Harvard.”**



Hi, I'm Nolan.



GO BLUE!

My research is on the politics of health.

American Political Science Review (2021) 115, 3, 1104–1109
doi:10.1017/psr.2020.00065 © The Author(s). 2021. Published by Cambridge University Press on behalf of the American Political Science Association. This is an Open Access article, distributed under the terms of the Creative Commons Attribution licence (<http://creativecommons.org/licenses/by/4.0/>), which permits unrestricted re-use, distribution, and reproduction in any medium, provided the original work is properly cited.

Letter Does Health Vulnerability Predict Voting for Right-Wing Populist Parties in Europe?

NOLAN M. KAVANAGH *University of Pennsylvania*
ANIL MENON *University of Michigan*
JUSTIN E. HEINZE *University of Michigan*

Why do voters in developed democracies support right-wing populist parties? Existing research focuses on economic and cultural vulnerability, but little attention has been paid to the role of health vulnerability. We analyze data from the European Social Survey (2002–2020). Our findings suggest that voters with worse self-reported health were significantly more likely to vote for right-wing populist parties. The relationship persists even after accounting for measures of cultural and economic vulnerability, as well as voter satisfaction with both their personal lives and their country’s health system. The influence of health on support for right-wing populist parties appears to be greater than that of income and self-reported economic insecurity, while less than that of gender and attitudes about immigration. Our findings suggest that policies affecting public health could shape not only health outcomes but also the political landscape.

INTRODUCTION

Right-wing populist parties are surging in popularity across the Western world (Norris and Inglehart 2010). Why do voters in developed democracies support such parties? A great body of research has identified economic insecurity and cultural backlash as potential drivers of recent populist successes (Algan et al. 2007; Hochschild 2016; Inglehart and Norris 2016; Kavanagh and Menon 2018; Rodrik 2018; Smith and Hanley 2018). According to these explanations, once-dominant socioeconomic groups perceive an erosion of their economic opportunities or a threat to their privileged status in society. These threats cause voters’ perceived vulnerability in motivating them to support parties that promise to restore their socioeconomic standing through anti-multiculturalism, antiglobalisation, and anti-immigration (Inglehart and Norris 2016).

We argue that voters’ perceived health may similarly contribute to populist support via a similar mechanism. The development of illness and disability often produces frustration with one’s physical and

emotional limitations, and it prompts people to compare themselves with their peers (Banks, Gibson, and Baum 2013; Martz and Liao 2007). This experience may increase an individual’s sense of personal vulnerability regardless of their socioeconomic background. They may then blame their misfortune on the political system, and seek to change or dismantle existing political and economic structures (Lacouture 2005; Nussbaum 2018). If true, individuals who suffer poorer health and more disability would be desirous of changing the political status quo. This desire to change the political system may draw them toward parties that campaign for a fundamental restructuring of a “biased and broken” system.

As such, health-related vulnerability may contribute to the rise of right-wing populism. Indeed, some research has associated declining population health with right-wing populist voting. U.S. counties that experienced the greatest rise in mortality over recent decades, especially among whites, were most likely to vote for President Donald Trump in the 2016 election (Kahan et al. 2016; Bilal, Knapp, and Cooper 2018; Bor 2017). Similar associations have been shown for rates of chronic opioid use (Goodwin et al. 2018) and other markers of poor health (Wasylyshyn, et al., 2017; Wasylyshyn et al. 2017). In the U.K., areas that experienced greater rises in “deaths of despair” due to suicide or drug overdose in the previous decade were more likely to vote for Brexit (Koltai et al. 2019). However, the relationship between poor health and right-wing populist voting remains to be tested, and the causal link will require appropriate controls for economic and cultural vulnerability.

Understanding how poor health influences right-wing populist support could have important implications for

Nolan M. Kavanagh Medical student, Perelman School of Medicine, University of Pennsylvania; Lecturer, Department of Periodontics and Oral Medicine, School of Dentistry, University of Michigan; nolan.m.kavanagh@perleman.drexel.edu
Justin E. Heinze Assistant Professor, Department of Political Science, University of Michigan; jheinze@umich.edu
Received: June 26, 2020; revised: February 27, 2021; accepted: March 23, 2021. First published online: April 26, 2021.
<https://doi.org/10.1017/psr.2020.00065> Published online by Cambridge University Press

My go-to karaoke song is “Since U Been Gone.”



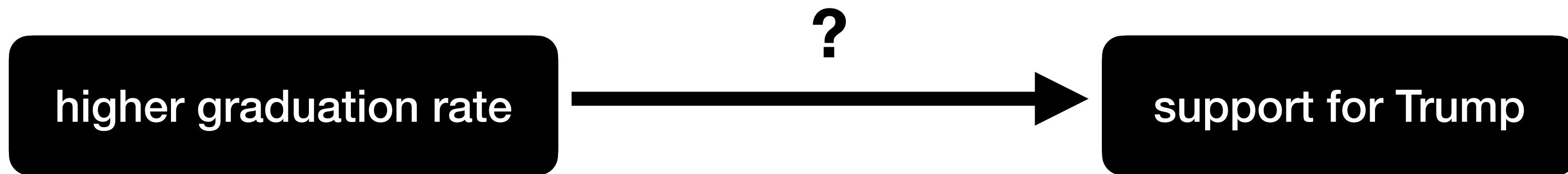
Overview of our sample data

Dataset of U.S. county-level characteristics in 2020

state	State of county	<i>Administrative</i>
county_fips	County FIPS identifier	<i>Administrative</i>
pc_under_18	Percent of county under age 18	<i>American Community Survey (2016–2020)</i>
pc_over_65	Percent of county over age 65	<i>American Community Survey (2016–2020)</i>
pc_male	Percent of county that is male	<i>American Community Survey (2016–2020)</i>
pc_black	Percent of county that is Black	<i>American Community Survey (2016–2020)</i>
pc_latin	Percent of county that is Hispanic/Latino	<i>American Community Survey (2016–2020)</i>
pc_hs_grad	Percent of county that graduated high school	<i>American Community Survey (2016–2020)</i>
unemploy_rate	County unemployment rate (%)	<i>American Community Survey (2016–2020)</i>
median_income	County median income (\$)	<i>American Community Survey (2016–2020)</i>
pc_uninsured	Percent of county without health insurance	<i>American Community Survey (2016–2020)</i>
pc_trump	Percent of county votes for Trump in 2020	<i>MIT Election Lab</i>

Tell me a story.

Let's say we're interested in the relationship between high school graduation and support for Trump.

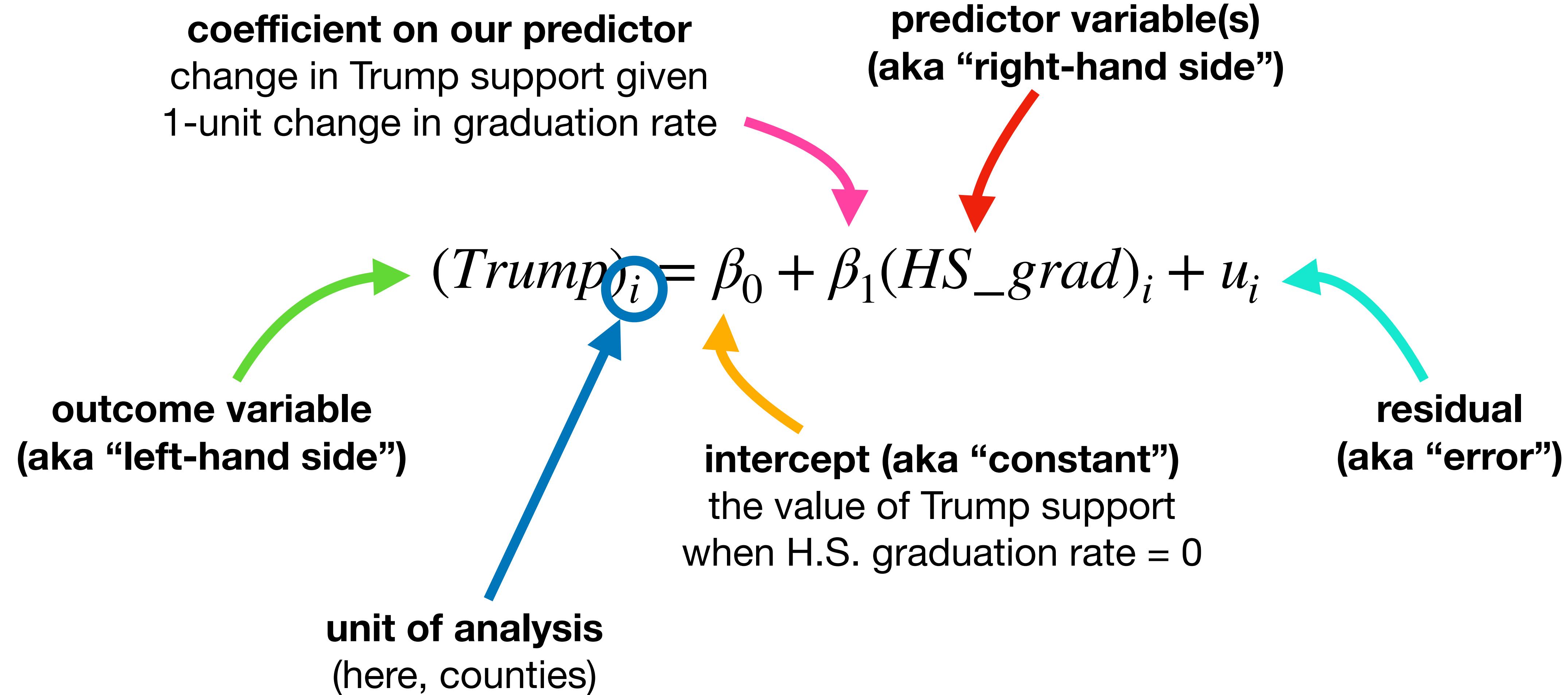


What might be the mechanism?

More education = liberal values = prefer multiculturalism?

More education = more income = prefer lower taxes?

Population regression function



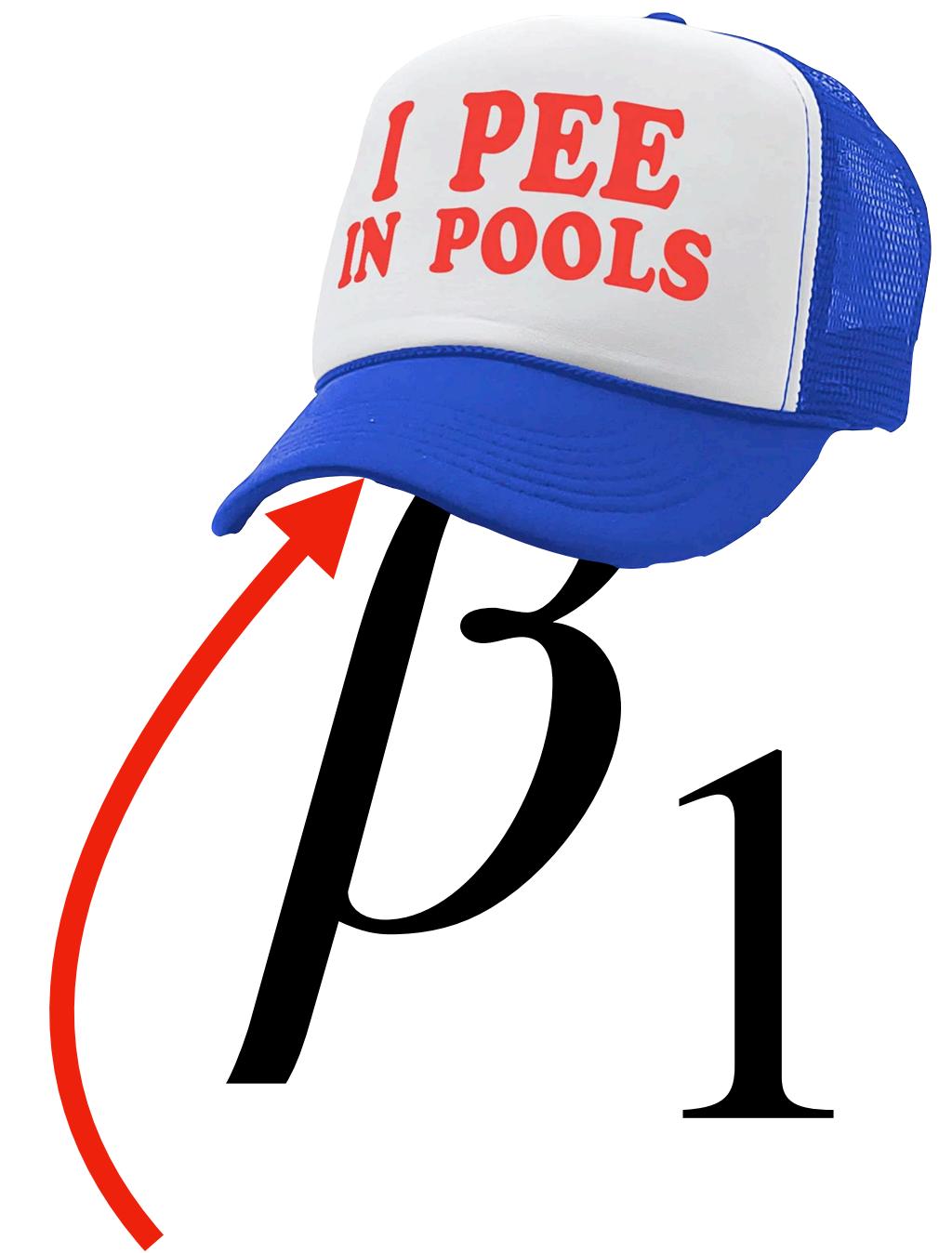
Population regression function

$$(Trump)_i = \beta_0 + \beta_1(HS_grad)_i + u_i$$

Sample regression function

$$(Trump)_i = \hat{\beta}_0 + \hat{\beta}_1(HS_grad)_i + \hat{u}_i$$

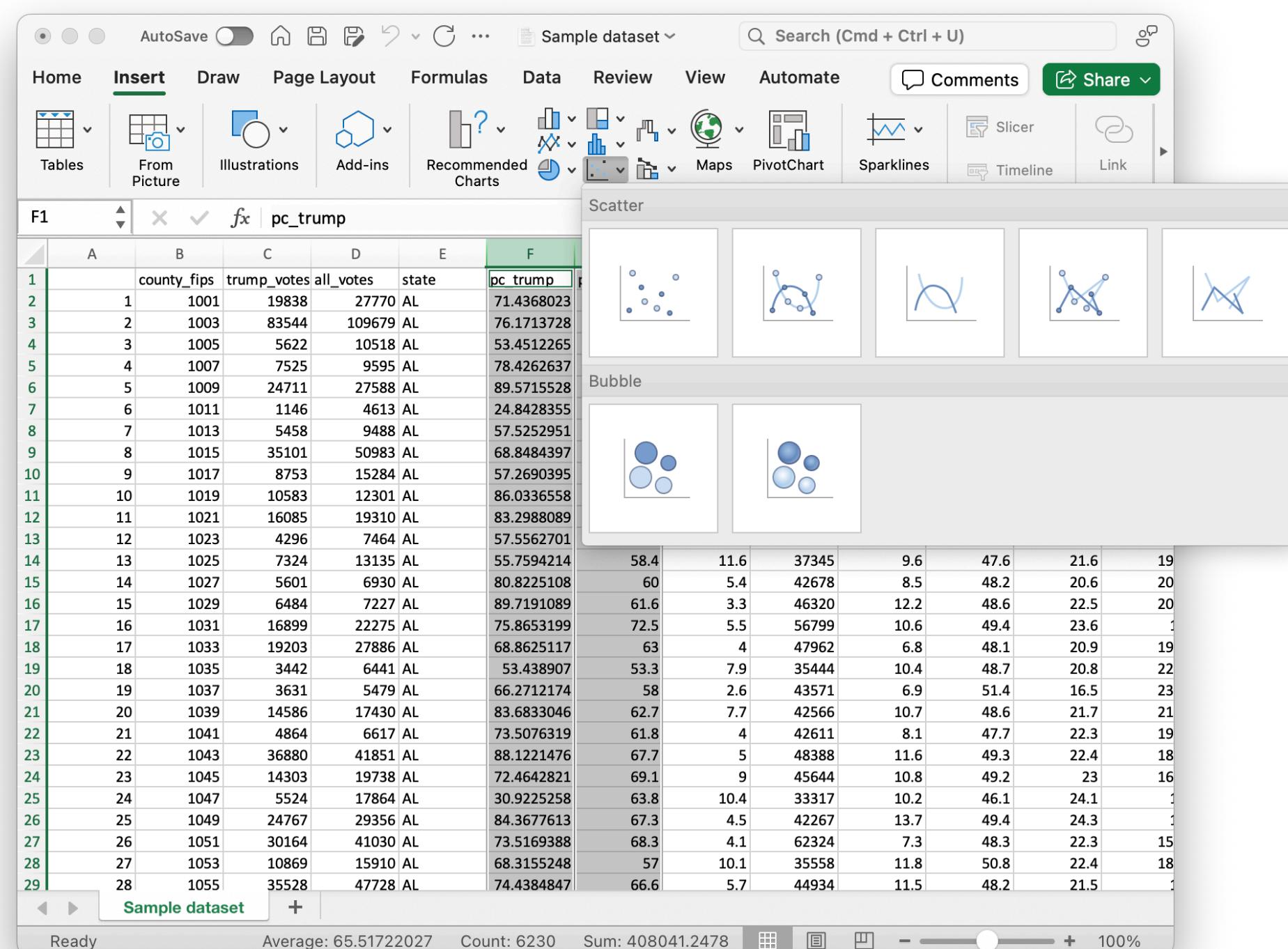
We add “hats” to signify estimated values in our sample.



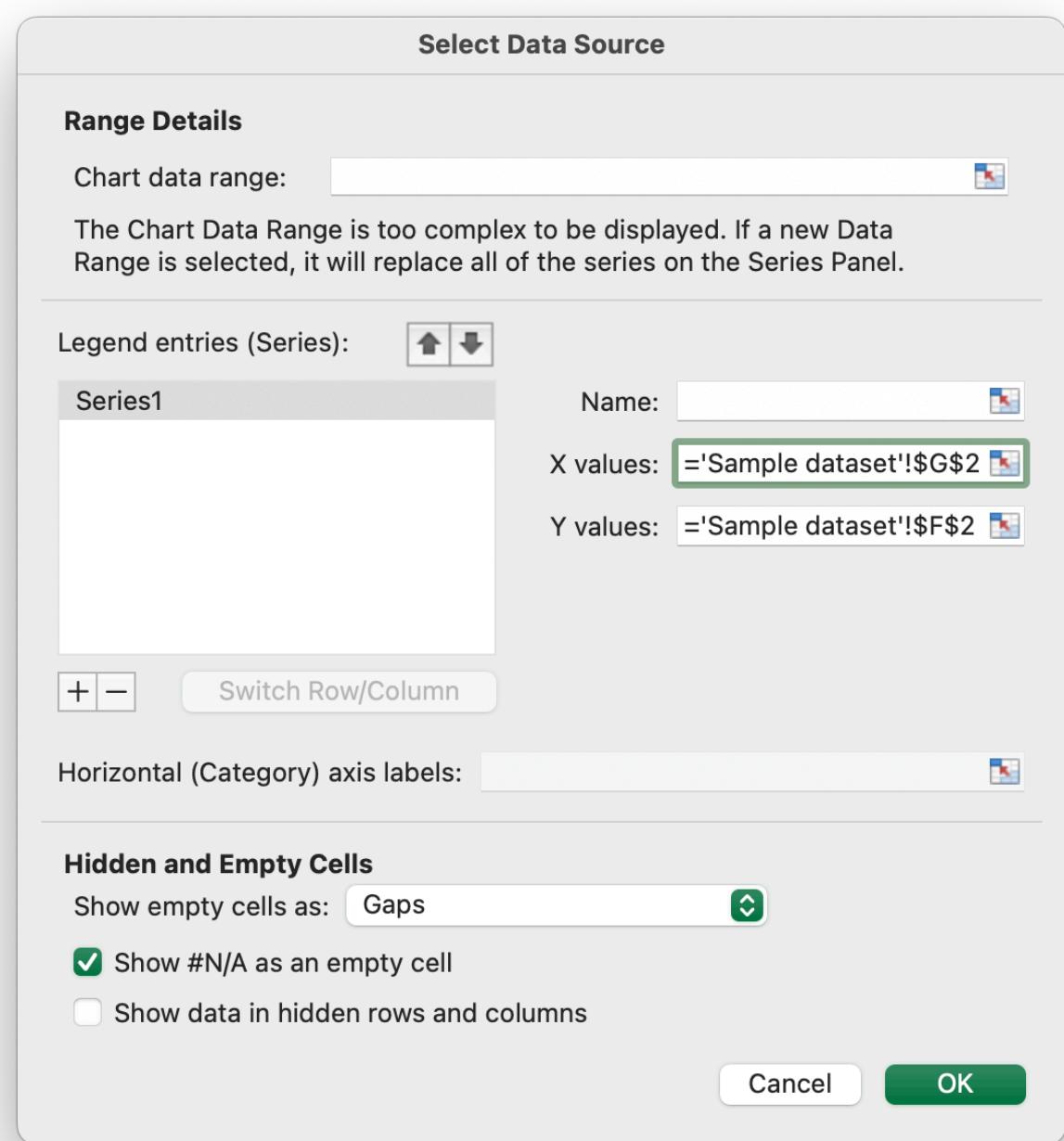
Only a specific sample would ever wear this hat.

Let's graph our data.

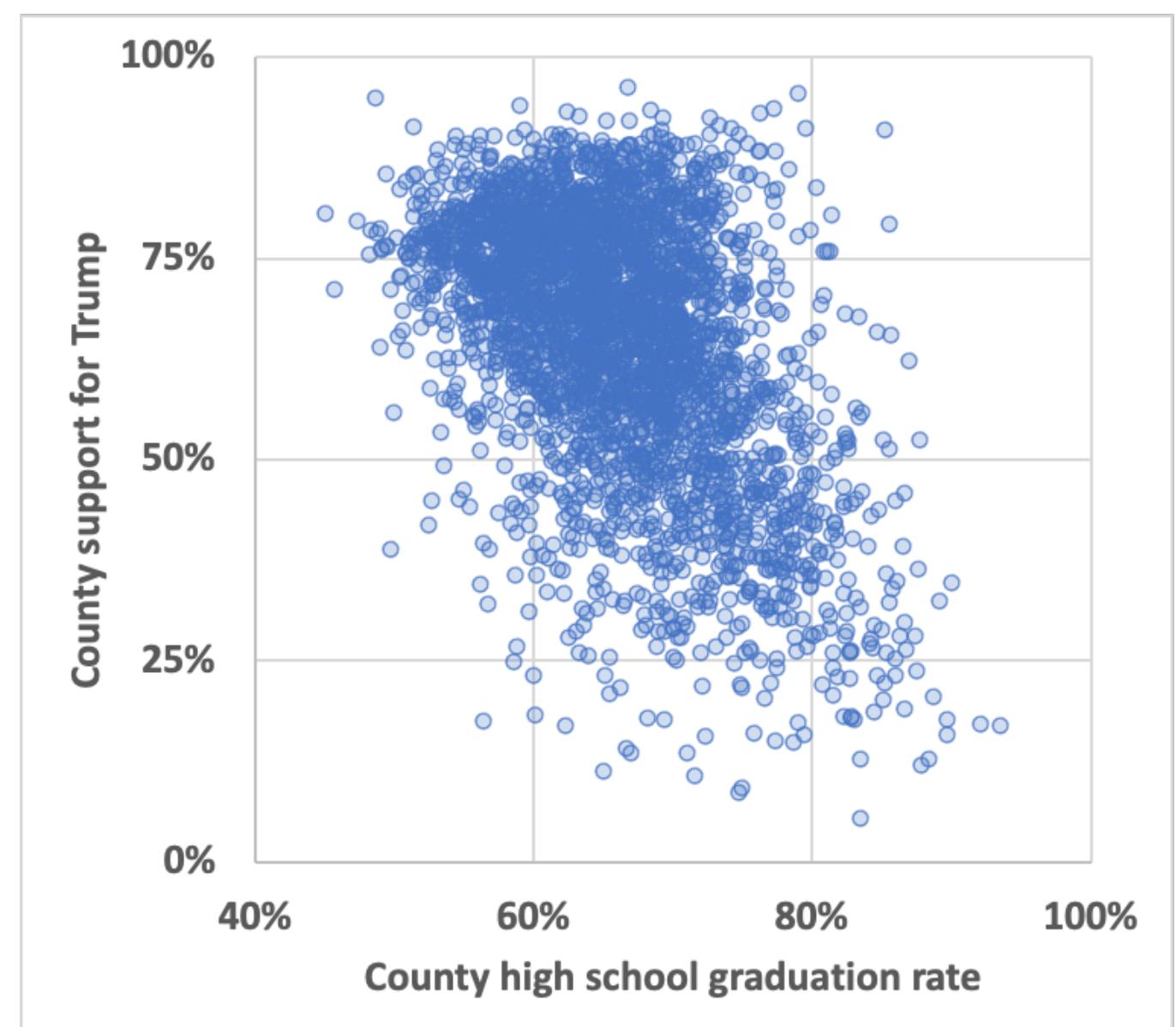
1. Insert a scatterplot.



2. Modify X and Y data ranges, as necessary.



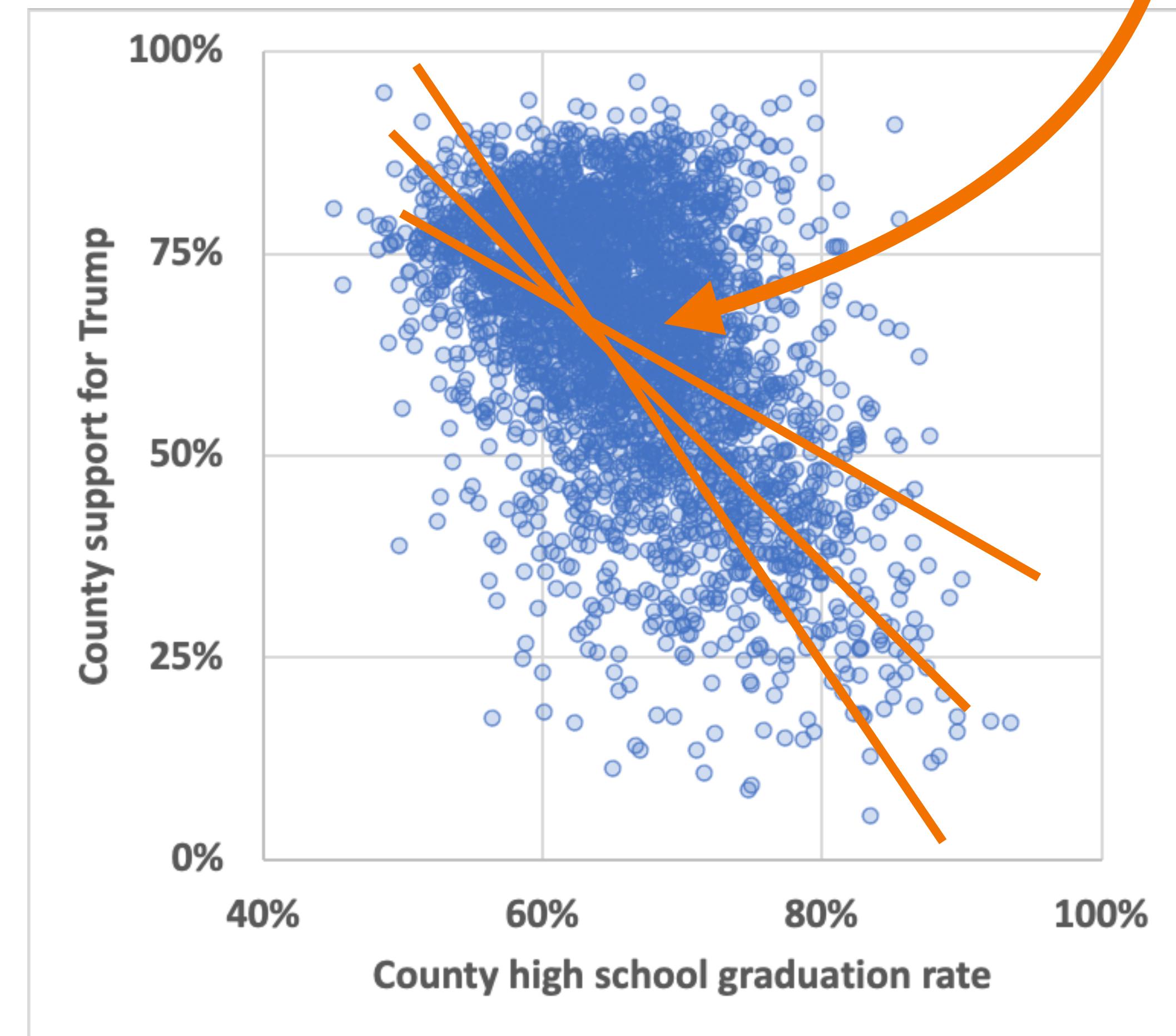
3. Chart Design > Add Chart Element, like axis labels.



To modify an existing element, double-click on it and a formatting window will pop up on the right.

Let's graph our data.

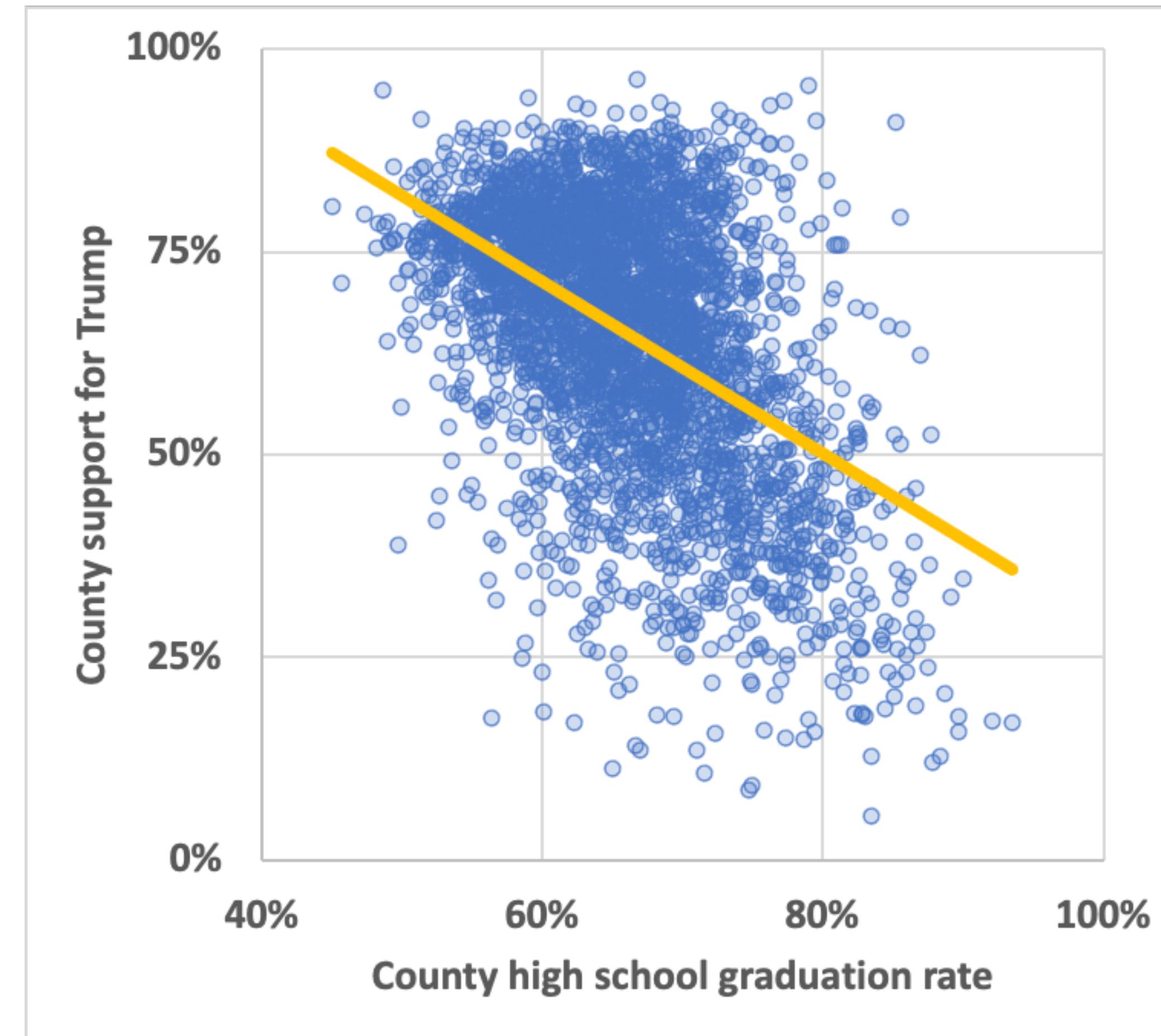
Which line fits best?



Let's graph our data.

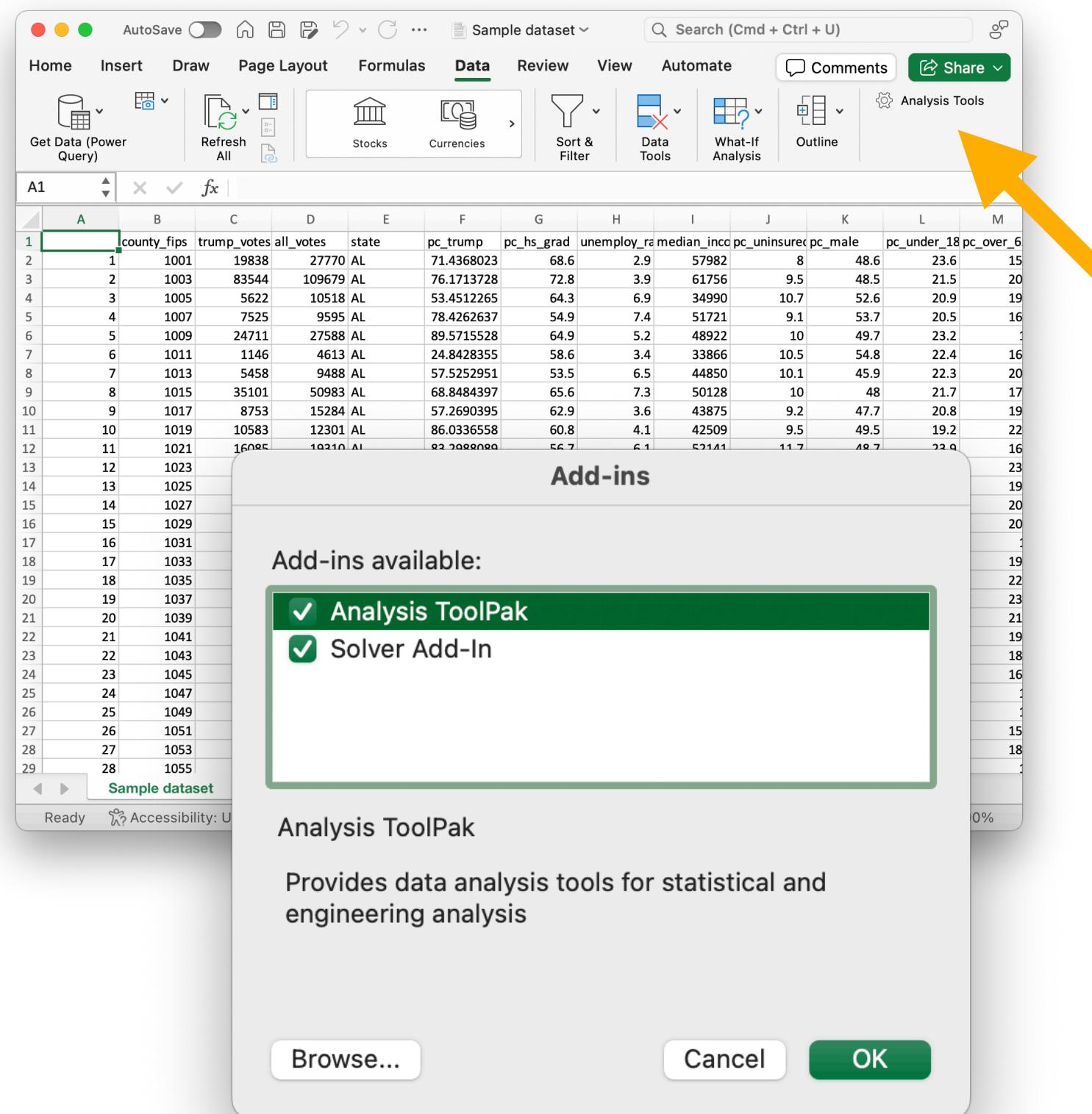
We can add a linear model under Add Chart Element > Trendline > Linear.

...but we still want to run the regression.

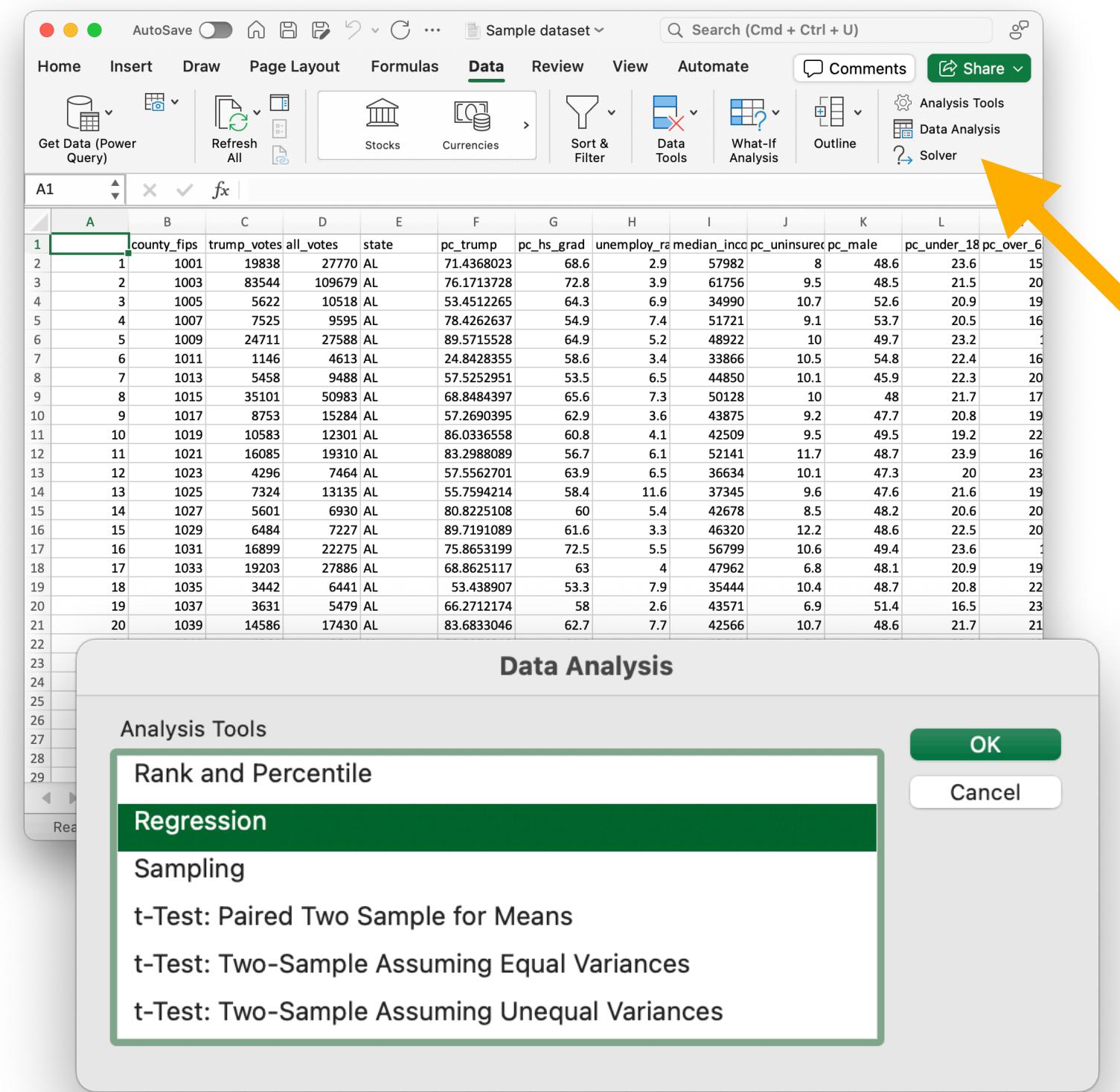


OK fine, let's run a regression.

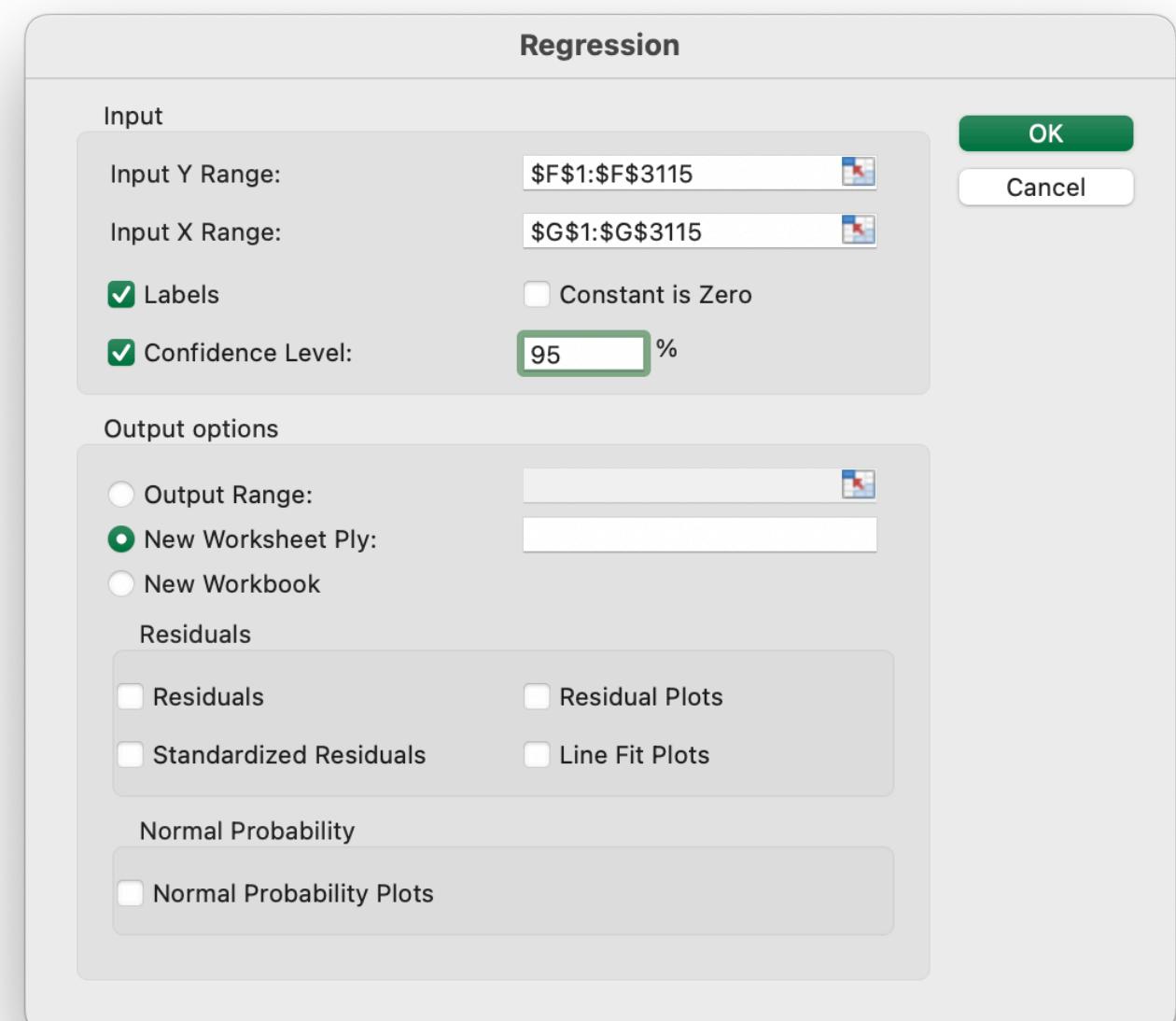
1. Data > Analysis Tools.
Add the Analysis ToolPak.



2. Click Data Analysis.
Then, select Regression.



3. Modify X and Y ranges.
Add labels and 95% C.I.s.



Other tidbits

- Constant should not be 0.
- Can ask for residuals if you plan to use them for another analysis.
- For ranges, have to give first and last cell in the range: "\$F\$1:\$F\$3115".
Can't just name the column.

OK fine, let's run a regression.

If your X or Y columns have missing data (e.g. blank cells or “NAs”), you’ll get an error and the regression won’t run.

We have a few options for dealing with it:

- “Sort” the column(s) with missing data so that the NAs go to the bottom. Then, specify a range that excludes the missing rows.
- Use “Filter” and uncheck the NAs. Then, copy the remaining rows into a new sheet and run the regression on these rows.

The screenshot shows a Microsoft Excel window with a sample dataset. The ribbon is visible at the top, with the 'Data' tab selected. A yellow arrow points to the 'Data Tools' icon in the ribbon, which is part of the 'Data' tab group. A dropdown menu is open from this icon, showing options like 'Get Data (Power Query)', 'Refresh All', 'Stocks', 'Currencies', 'Sort & Filter', 'Data Tools' (which is the active option), 'Clear', 'Reapply', and 'Advanced'. The main area of the screen displays a table with various columns and rows of data, representing a sample dataset.

OK fine, let's interpret a regression.

SUMMARY OUTPUT								
Regression Statistics								
Multiple R	0.4832049							
R Square	0.23348698							
Adjusted R S	0.23324067							
Standard Err	14.1346772							
Observations	3114							
ANOVA								
	df	SS	MS	F	Significance F			
Regression	1	189388.892	189388.892	947.944069	6.07E-182			
Residual	3112	621743.678	199.7891					
Total	3113	811132.57						
	Coefficients	Standard Error	t Stat	P-value	Lower 95%	Upper 95%	Lower 95.0%	Upper 95.0%
Intercept	134.921115	2.28637329	59.0109742	0	130.438162	139.404068	130.438162	139.404068
pc_hs_grad	-1.0588226	0.03438998	-30.7887	6.07E-182	-1.126252	-0.9913933	-1.126252	-0.9913933

OK fine, let's interpret a regression.

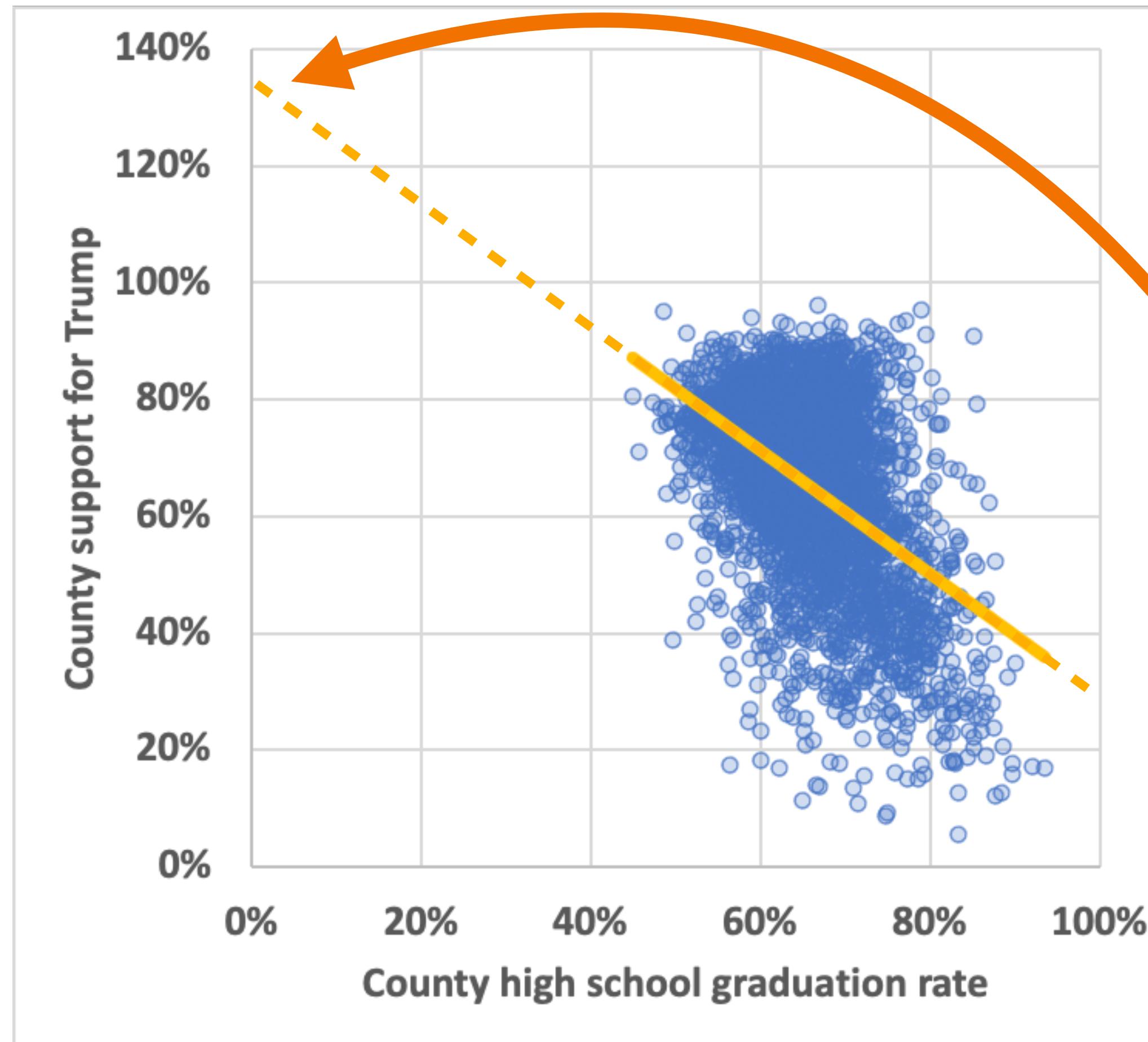
SUMMARY OUTPUT								
Regression Statistics								
Multiple R	0.4832049							
R Square	0.23348698							
Adjusted R S	0.23324067							
Standard Err	14.1346772							
Observations	3114							
ANOVA								
	df	SS	MS	F	Significance F			
Regression	1	189388.892	189388.892	947.944069	6.07E-182			
Residual	3112	621743.678	199.7891					
Total	3113	811132.57						
	Coefficients	Standard Error	t Stat	P-value	Lower 95%	Upper 95%	Lower 95.0%	Upper 95.0%
Intercept	134.921115	2.28637329	59.0109742	0	130.438162	139.404068	130.438162	139.404068
pc_hs_grad	-1.0588226	0.03438998	-30.7887	6.07E-182	-1.126252	-0.9913933	-1.126252	-0.9913933

We'll start with the intercept.

When 0% of a county has graduated high school, support for Trump is an estimated 134.9%.

(Is a graduation rate of 0% meaningful?)

OK fine, let's interpret a regression.



We'll start with the intercept.

When 0% of a county has graduated high school, support for Trump is an estimated 134.9%.

(Is a graduation rate of 0% meaningful?)

**That looks
about right!**

**...but here, it's not meaningful. We don't
have any counties near 0% graduation rates,
and we can't ever have 135% support!**

OK fine, let's interpret a regression.

SUMMARY OUTPUT							
Regression Statistics							
Multiple R	0.4832049						
R Square	0.23348698						
Adjusted R S	0.23324067						
Standard Err	14.1346772						
Observations	3114						
ANOVA							
	df	ss	MS	F	Significance F		
Regression	1	189388.892	189388.892	947.944069	6.07E-182		
Residual	3112	621743.678	199.7891				
Total	3113	811132.57					
	Coefficients	standard Error	t Stat	P-value	Lower 95%	Upper 95%	Lower 95.0%
Intercept	134.921115	2.28637329	59.0109742	0	130.438162	139.404068	130.438162
pc_hs_grad	-1.0588226	0.03438998	-30.7887	6.07E-182	-1.126252	-0.9913933	-1.126252
							-0.9913933

We'll start with the intercept.

When 0% of a county has graduated high school, support for Trump is an estimated 134.9%.

(Is a graduation rate of 0% meaningful?)

The standard error is 2.3%. This gives us a 95% C.I. of $134.9 \pm 1.96 \cdot 2.3 = [130.4\% \text{ to } 139.4\%]$.

The t-statistic is 59.0. The p-value is <0.05.

Thus, we can conclude that the intercept is significantly different from 0.

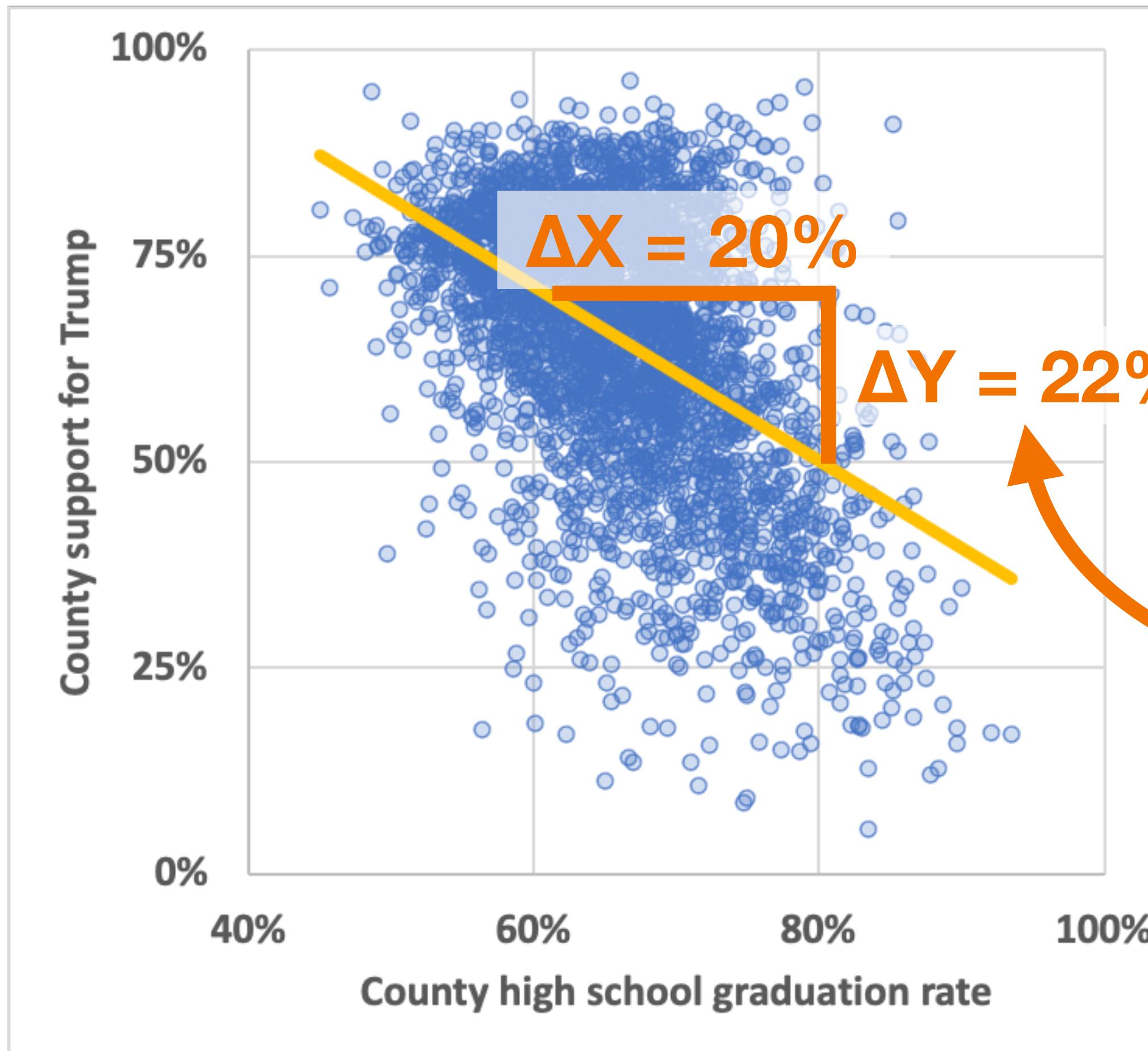
OK fine, let's interpret a regression.

SUMMARY OUTPUT								
Regression Statistics								
Multiple R	0.4832049							
R Square	0.23348698							
Adjusted R S	0.23324067							
Standard Err	14.1346772							
Observations	3114							
ANOVA								
	df	SS	MS	F	Significance F			
Regression	1	189388.892	189388.892	947.944069	6.07E-182			
Residual	3112	621743.678	199.7891					
Total	3113	811132.57						
	Coefficients	Standard Error	t Stat	P-value	Lower 95%	Upper 95%	Lower 95.0%	Upper 95.0%
Intercept	134.921115	2.28637329	59.0109742	0	130.438162	139.404068	130.438162	139.404068
pc_hs_grad	-1.0588226	0.03438998	-30.7887	6.07E-182	-1.126252	-0.9913933	-1.126252	-0.9913933

Now for the coefficient on pc_hs_grad.

For each 1 percentage point (pp) increase in a county's high school graduation rate, the estimated support for Trump decreases by 1.1 pp.

OK fine, let's interpret a regression.



Now for the coefficient on `pc_hs_grad`.

For each 1 percentage point (pp) increase in a county's high school graduation rate, the estimated support for Trump decreases by 1.1 pp.

**That looks
about right!**

OK fine, let's interpret a regression.

SUMMARY OUTPUT								
Regression Statistics								
Multiple R	0.4832049							
R Square	0.23348698							
Adjusted R S	0.23324067							
Standard Err	14.1346772							
Observations	3114							
ANOVA								
	df	SS	MS	F	Significance F			
Regression	1	189388.892	189388.892	947.944069	6.07E-182			
Residual	3112	621743.678	199.7891					
Total	3113	811132.57						
	Coefficients	Standard Error	t Stat	P-value	Lower 95%	Upper 95%	Lower 95.0%	Upper 95.0%
Intercept	134.921115	2.28637329	59.0109742	0	130.438162	139.404068	130.438162	139.404068
pc_hs_grad	-1.0588226	0.03438998	-30.7887	6.07E-182	-1.126252	-0.9913933	-1.126252	-0.9913933

Now for the coefficient on pc_hs_grad.

For each 1 percentage point (pp) increase in a county's high school graduation rate, the estimated support for Trump decreases by 1.1 pp.

The standard error is 0.03%. This gives us a 95% C.I. of $-1.06 \pm 1.96 \times 0.03 = [-1.13\%, -0.99\%]$.

The t-statistic is -30.8. The p-value is <0.05.

Thus, the coefficient is significantly different from 0. There's a negative association between graduation rates and Trump support.

That's all good. But is it causal?

We'll spend lots of time in API 202 asking this very question.

What problems with a causal interpretation come to mind?

What variables/influences are we missing?