

Week 10

Questions of Description

Paul Laskowski and Alex Hughes

January 12, 2023

UC Berkeley, School of Information

Questions of Description

QUESTIONS OF DESCRIPTION



What is the shape of the relationship between a country's economic output and Internet access?

QUESTIONS OF DESCRIPTION

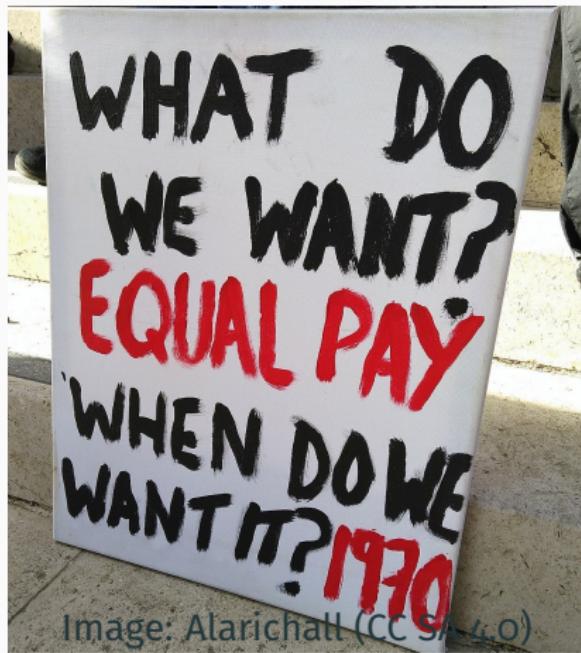


Image: Alarichall (CC SA 4.0)

How big is the pay gap in the United States?

QUESTIONS OF DESCRIPTION



How does the pay gap depend on the age of the worker?

Description

PLAN FOR THE WEEK

Preamble

- Three modes of model building

Three sections about descriptive modeling

1. Capturing nonlinear relationships
2. Measurement with controls
3. Modeling conditional effects

PLAN FOR THE WEEK (CONT.)

At the end of this week, you will be able to:

- Understand how three major modes of model building lead to very different models
- Balance design goals when creating a model for description
- Plan a set of model specifications for a regression table

What Is Linear Model Building?

WHAT IS LINEAR?

Model	Linear?
$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \epsilon$	
$Y = \beta_0 + \beta_1 X + \beta_2 X^2 + \epsilon$	
$\ln Y = \beta_0 + \beta_1 \ln X + \epsilon$	
$Y = \beta_0 X^{\beta_1} + \epsilon$	

WHAT IS LINEAR MODEL BUILDING?

How do you select from all the possible linear models?

- Which variables to include, which to exclude
- Whether and how to transform each variable
- Whether to create new variables
- Whether to multiply variables together

THE MODEL-BUILDING PROCESS

Propose Model



Evaluate Model

- Modeling goals
- Observations from data
- Background knowledge
- Background theory
- Diagnostic plots
- Measures of fit
- Tests of model assumptions

Modes of Model Building

AN IMPORTANT QUESTION TO KEEP IN MIND:

What are your goals?

THREE MODES OF MODEL BUILDING

1. Predictive modeling: What value can we expect for data we haven't seen?
2. Descriptive modeling: How can we make sense of the patterns in data?
3. Explanatory modeling: How can we measure effects inside a causal theory?

Modes of Model Building

Predictive Modeling

PREDICTIVE MODELING

Prediction: making guesses for unknown values

- For new people, new pictures of cats, or other units
- For future time periods

A key focus of machine learning and time series analysis

KEY GOALS OF PREDICTIVE MODELING

1. Accurately predict values
2. Be interpretable by humans (usually less important)

Implications for linear regression

- Hundreds of variables, or more
- Variable selection by algorithm
 - False discovery: Out of many variables, some will look important by chance.
- Coefficients usually don't have meaning

Modes of Model Building

Descriptive Modeling

DESCRIPTIVE MODELING

Description: summarizing or representing data in a compact, human-understandable way

- Gain understanding by interpreting the model's internal structure
- Popular in statistics, but “not commonly used for theory building and testing in other disciplines”
—Shmueli, Galit. *To Explain or to Predict?*

KEY GOALS OF DESCRIPTIVE MODELING

1. Measure complex concepts
2. Capture features of the world
3. Simplify phenomena to make them understandable
4. Highlight associations
5. Generate hypotheses

Implications for linear regression

- Parsimonious models
- Use of logarithms and other transforms with clear interpretation
- Long build process relying on EDA

Modes of Model Building

Explanatory Modeling

EXPLANATORY MODELING

Explanation: using data to test or estimate parameters in a causal theory

- Explanation lets us reason about actions
- Common mode in economics, political science, epidemiology, psychology, environmental science...
- Causal assumptions are required to make causal claims

KEY GOALS OF EXPLANATORY MODELING

1. Measure a causal effect.
2. Evaluate a theory.
3. Predict the consequences of potential actions.

Implications for linear regression

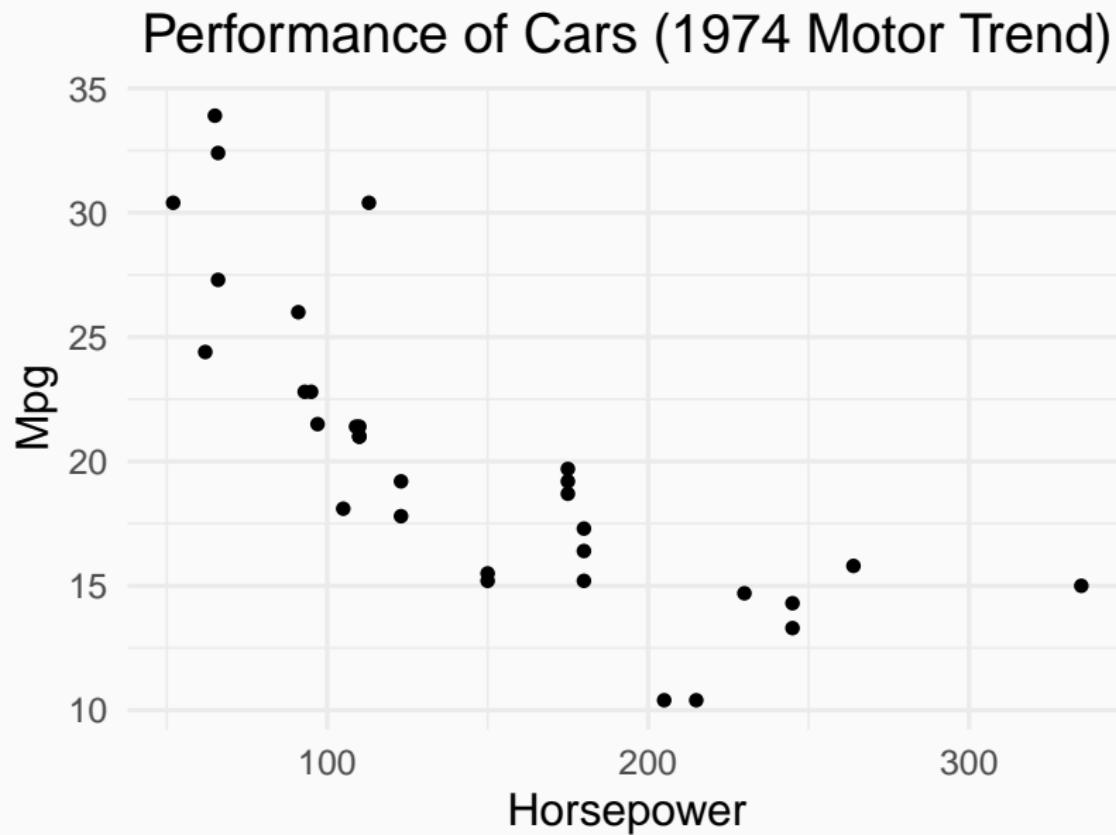
- Variable selection is guided by causal theory
- Models oriented to perform a specific measurement
- Key challenge is the operationalization gap between theoretical constructs and variables

THREE MODES OF MODEL BUILDING

1. Predictive modeling: What value can we expect for data we haven't seen?
2. Descriptive modeling: How can we make sense of the patterns in data?
3. Explanatory modeling: How can we measure effects inside a causal theory?

Introduction to Descriptive Modeling

How WOULD YOU DESCRIBE THESE DATA?



INTRODUCTION TO DESCRIPTIVE MODELING

Description: summarizing or representing data in a compact, human-understandable way

Linear regression is a tool that can help us answer:

- What patterns exist in a dataset?
- What is the shape of a specific relationship?
- What is the size of a feature or effect?

WHAT DOES IT TAKE TO BUILD A DESCRIPTIVE MODEL?

Skill, art, and a lot of iteration

- No causal theory to constrain model choice
- Often many ways to operationalize concepts
- Need to balance many competing goals
- Takes time to build intuition in many dimensions

Logarithms

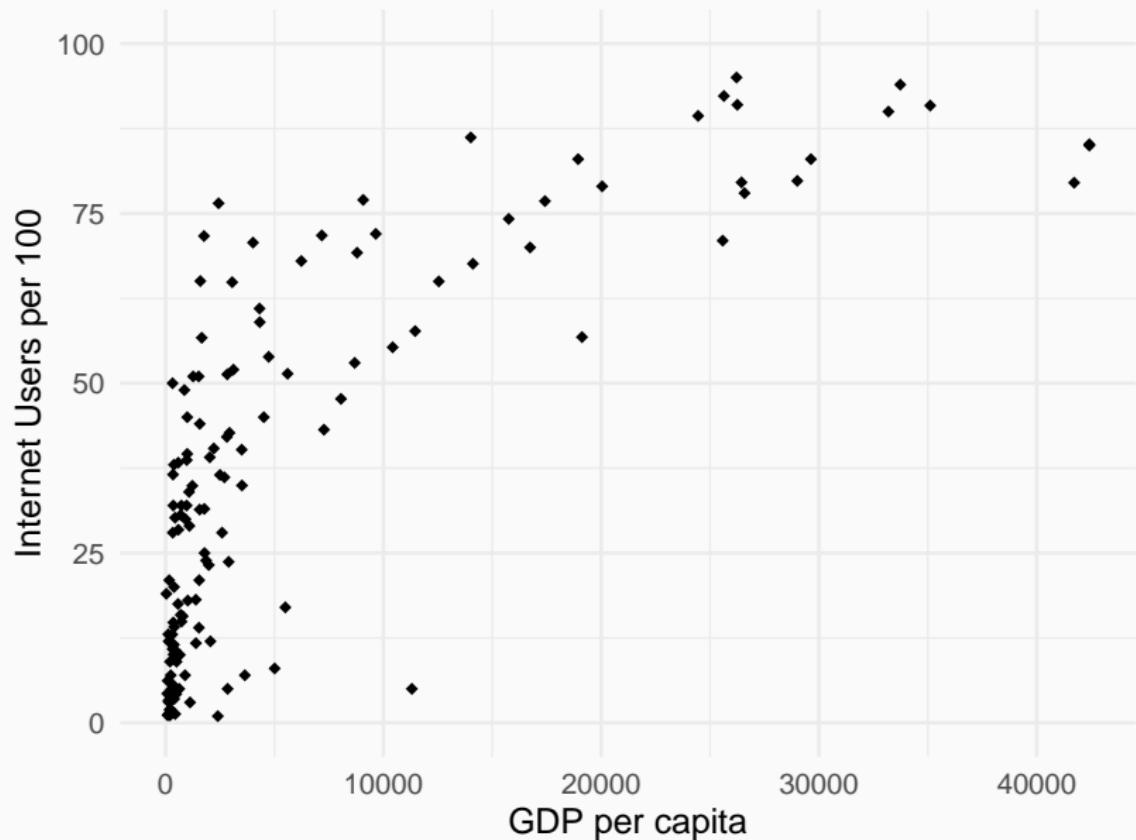
A QUESTION OF DEVELOPMENT

How does the economic productivity of a country relate to Internet access?

Data taken from the World Bank

- GDP per capita
- Internet users per 100

Is OLS REGRESSION APPROPRIATE?



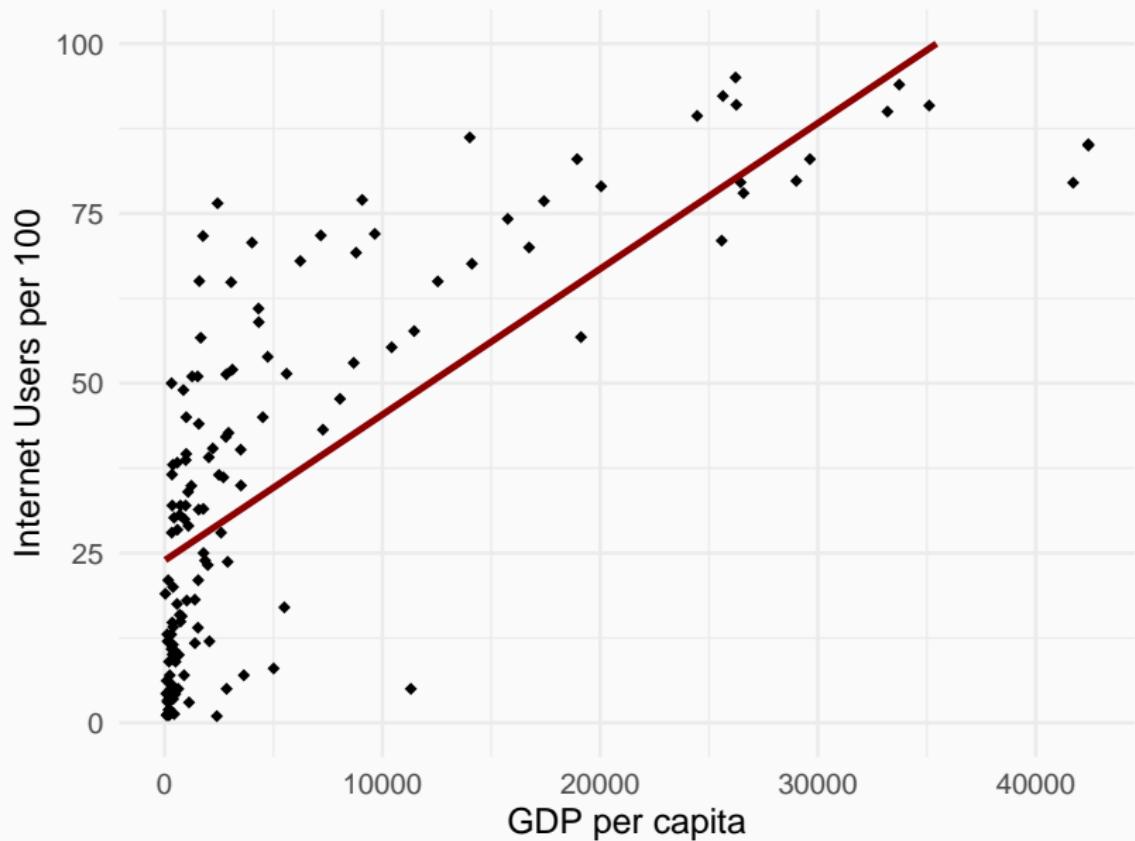
INTERPRETING THE LEVEL-LEVEL MODEL

$$\widehat{\text{users}} = 2.4 + .0021 \text{ GDP}$$

Interpretation 1: A country with an extra \$1 of GDP per capita is predicted to have another .0021 Internet users per 100.

Interpretation 2: A country with an extra \$1,000 of GDP per capita is predicted to have another 2.1 Internet users per 100.

DOES OLS REGRESSION MEET OUR GOALS?



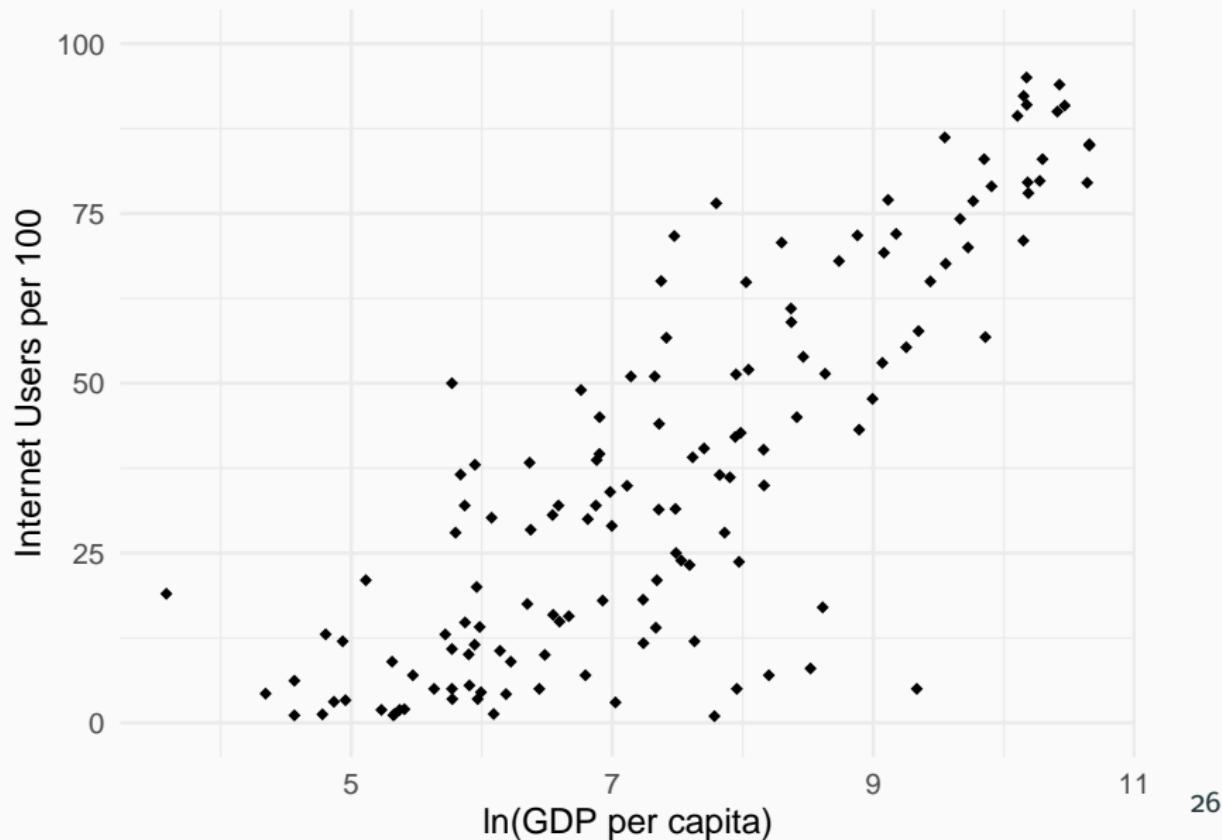
VARIABLE TRANSFORMATION

Variable transformation

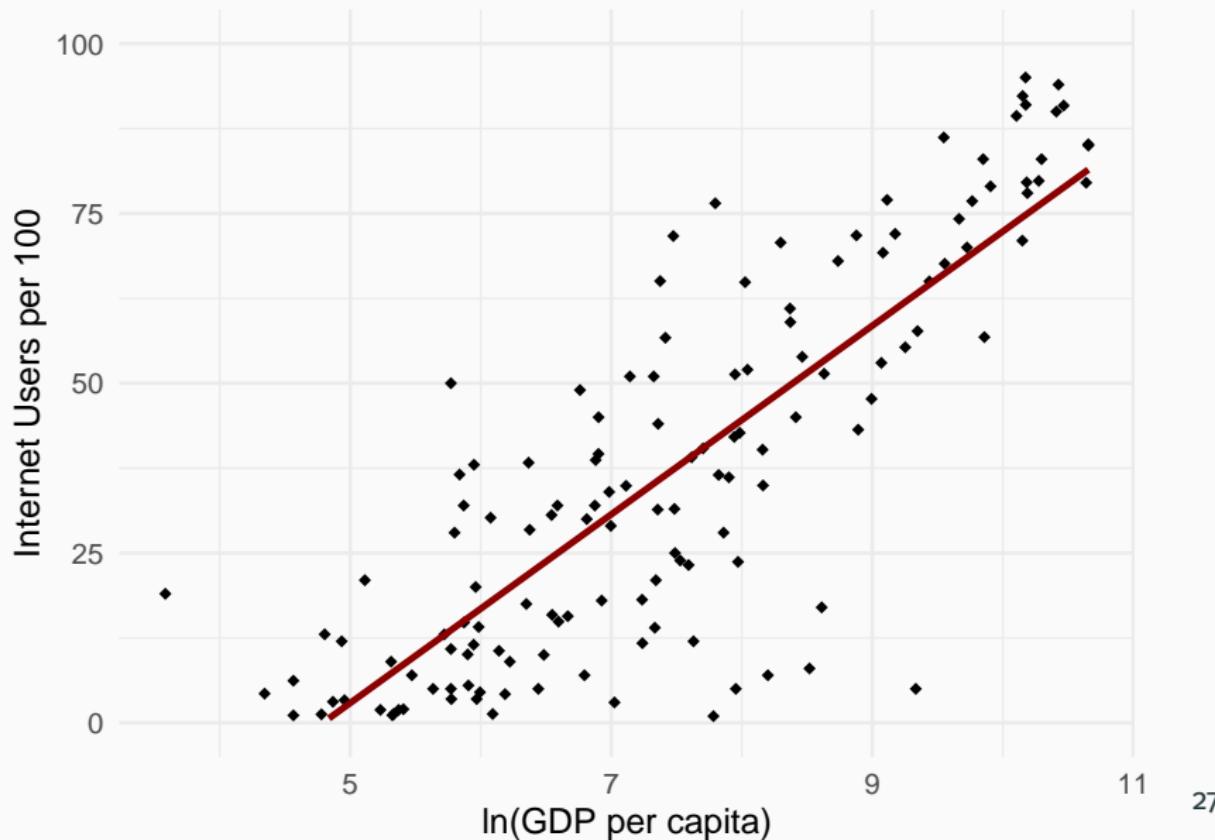
Replace X with $f(X)$ for some $f : \mathbb{R} \rightarrow \mathbb{R}$

- $\ln(X)$
- $\log_{10}(X)$
- Indicator functions
- X^a
- Polynomials(X)

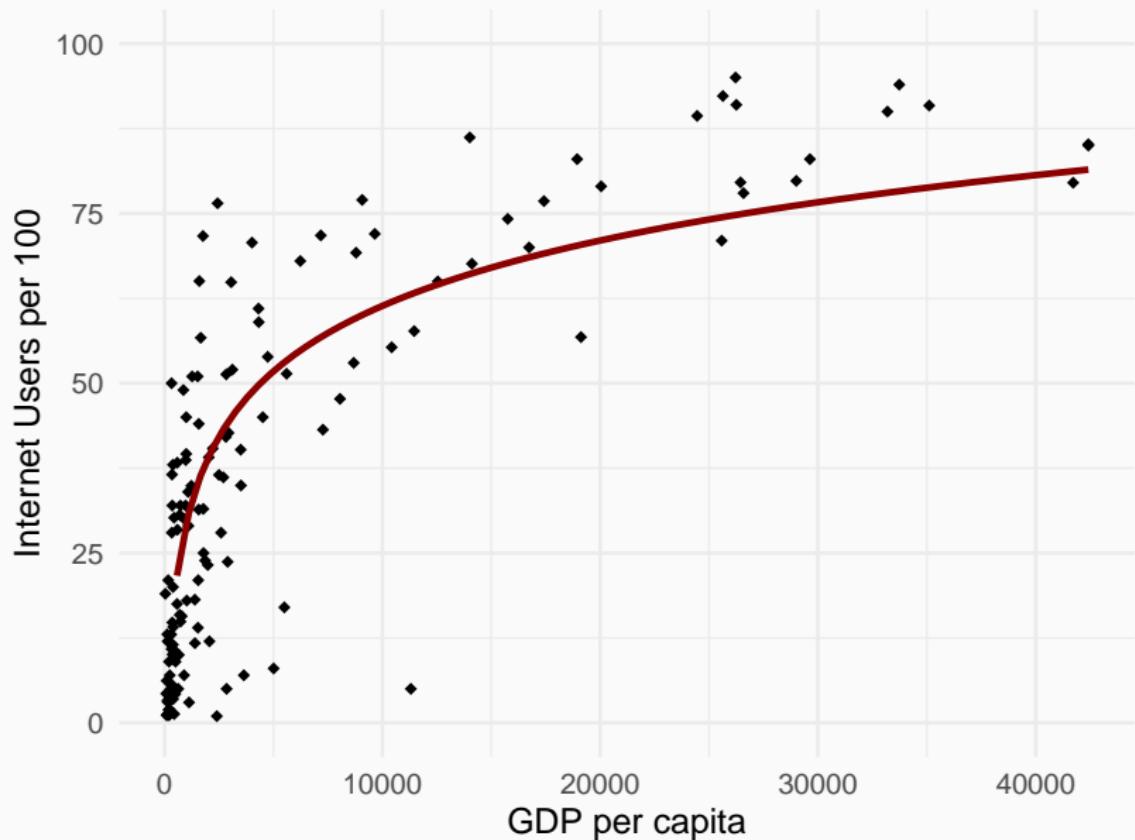
TAKING THE LOG OF GDP



TAKING THE LOG OF GDP



MODEL: $\widehat{\text{users}} = \beta_0 + \beta_1 \ln GDP$ $R^2 = .68$



Interpreting and Applying Logarithms

INTERPRETING LOGARITHMS

Base 10 log

$$\widehat{\text{users}} = -66.5 + 32.0 \log_{10} GDP$$

Interpretation: adding a 0 to GDP associated with 32 more Internet users per 100

Base e log

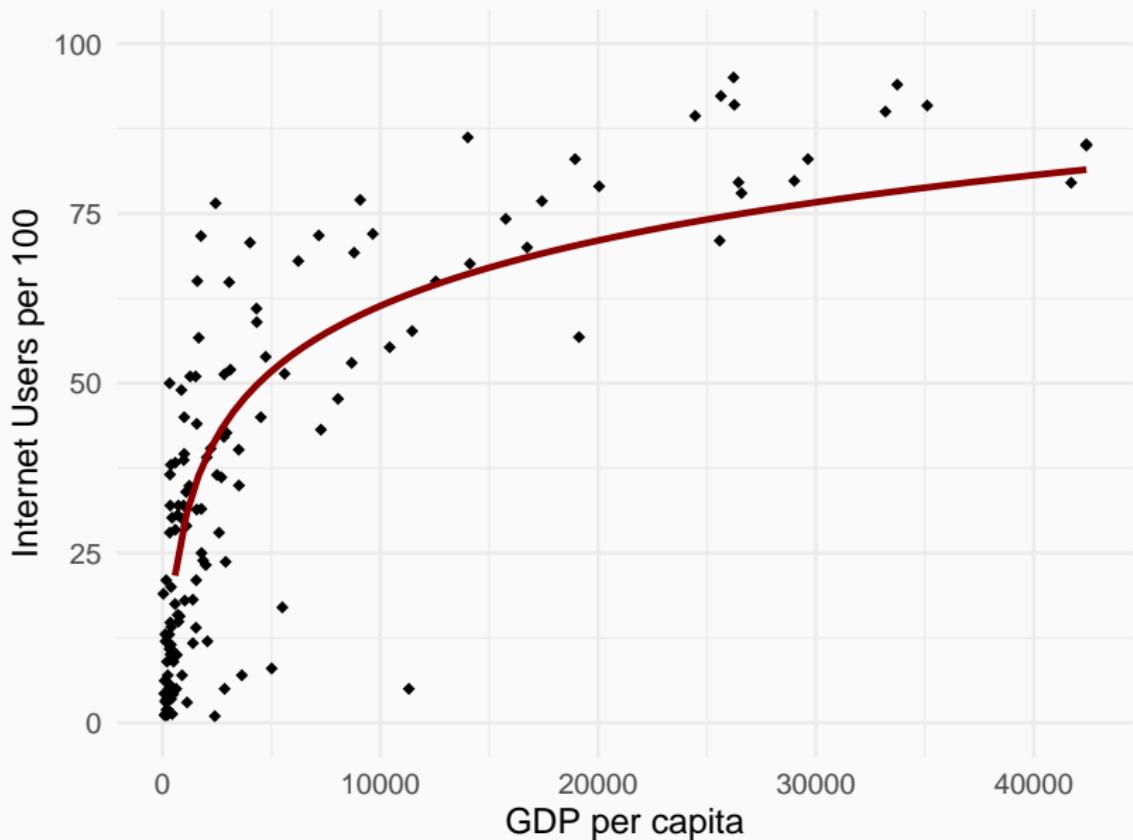
$$\widehat{\text{users}} = -66.5 + 13.9 \ln GDP$$

Interpretation: multiplying GDP by a small $1 + \alpha$ associated with 13.9α more Internet users per 100

INTERPRETING LOGARITHMS (CONT.)

$$\widehat{\frac{\partial \text{users}}{\partial GDP}} = \frac{\partial}{\partial GDP}[-66.5 + 13.9 \ln GDP]$$

MODEL: $\widehat{\text{users}} = \beta_0 + \beta_1 \ln GDP$ $R^2 = .68$



OTHER USES OF LOGARITHMS

Log-linear form: $\ln Y = \beta_0 + \beta_1 X$

- Interpretation: $\frac{\Delta Y}{Y} \approx \beta_1 \Delta X$
- Example: add 0.1 to $X \rightarrow$ add $0.1\beta_1 \cdot Y$ to Y

Log-log form: $\ln Y = \beta_0 + \beta_1 \ln X$

- Constant elasticity model
- Interpretation: $\frac{\Delta Y}{Y} \approx \beta_1 \frac{\Delta X}{X}$
- Example: increase X by 1% \rightarrow increase Y by $\beta_1\%$

WHEN TO APPLY LOGS

What makes a variable a good candidate for a log transform?

- Always positive
 - Don't add a constant to make variable positive
- Spans multiple orders of magnitude
- Clustered near zero with high outliers
- Percent changes are meaningful

POWER TRANSFORMATIONS

Power transformation: Replace X with X^α for some $\alpha \in \mathbb{R}$.

- \sqrt{X} behaves similarly to $\ln X$.
- X^2, X^3, \dots spreads values far from zero further apart.

Interpretation is more difficult than with logs.

WHEN TO APPLY TRANSFORMATIONS

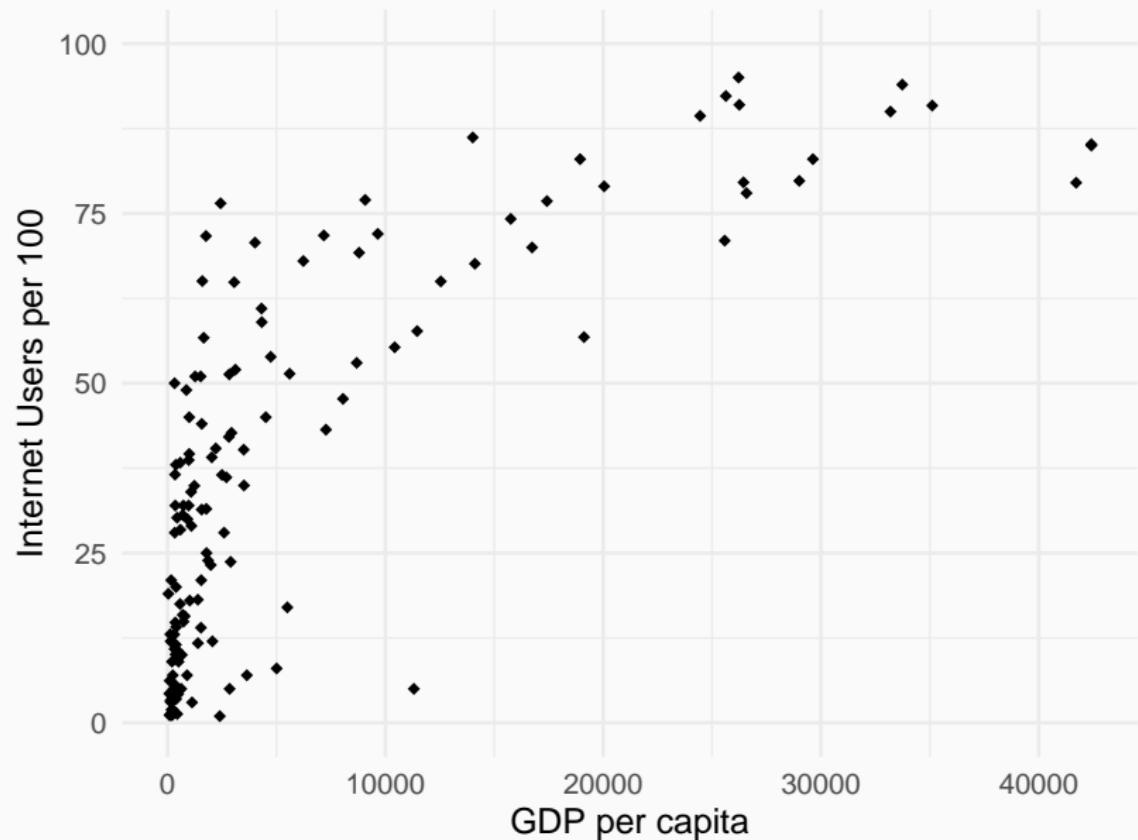
Question: Should you always transform your variables to make them normal?

- Normal variables are **never** a requirement for OLS regression.
 - Even in the classical linear model, *errors* are normal, not variables.

⇒ Don't focus too much on normality. Capturing relationships is more important.

Polynomials

HOW CAN WE IMPROVE PREDICTIVE ACCURACY?



POLYNOMIAL REGRESSION

Motivation: Polynomials can approximate any continuous function (e.g. Weierstrass theorem)

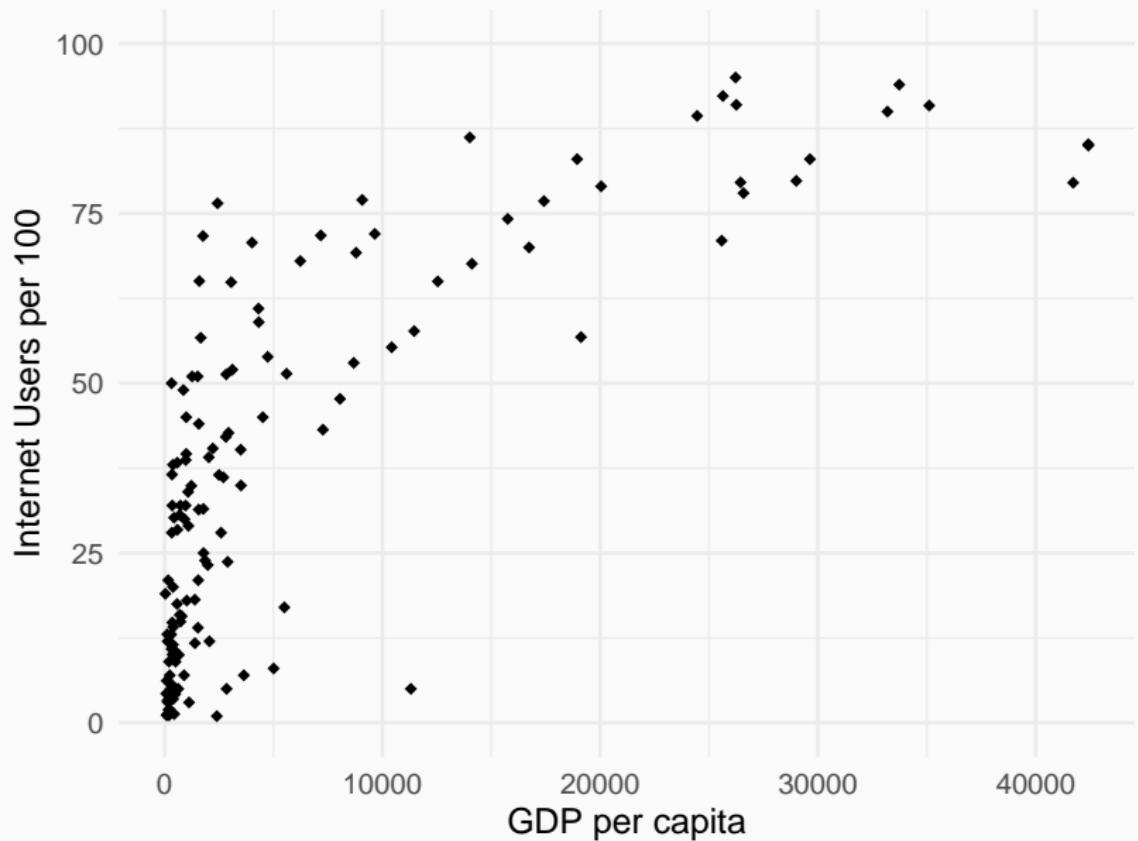
Quadratic: $\widehat{\text{Users}} = \beta_0 + \beta_1 \text{GDP} + \beta_2 \text{GDP}^2$

- OLS finds the parabola that minimizes MSE.

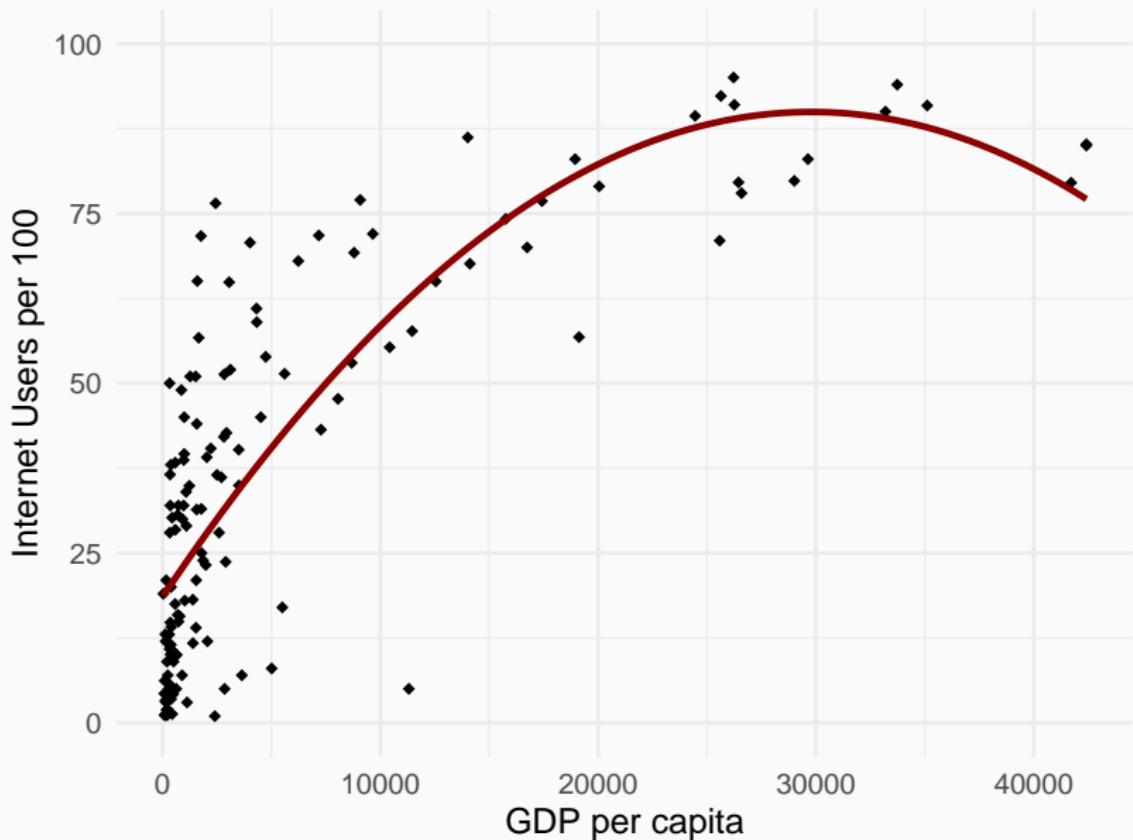
Cubic: $\widehat{\text{Users}} = \beta_0 + \beta_1 \text{GDP} + \beta_2 \text{GDP}^2 + \beta_3 \text{GDP}^3$

- OLS finds the cubic function that minimizes MSE.

A PARABOLIC PATTERN?



MODEL: $\widehat{\text{users}} = \beta_0 + \beta_1 \text{GDP} + \beta_2 \text{GDP}^2$ $R^2 = .66$



INTERPRETING THE QUADRATIC SPECIFICATION

$$\widehat{\text{Users}} = 18.5 + .0048 \text{ GDP} - 8.0 \cdot 10^{-8} \text{ GDP}^2$$

$$\frac{\partial \widehat{\text{Users}}}{\partial \text{GDP}} = .0048 - 1.6 \cdot 10^{-7} \text{ GDP}$$

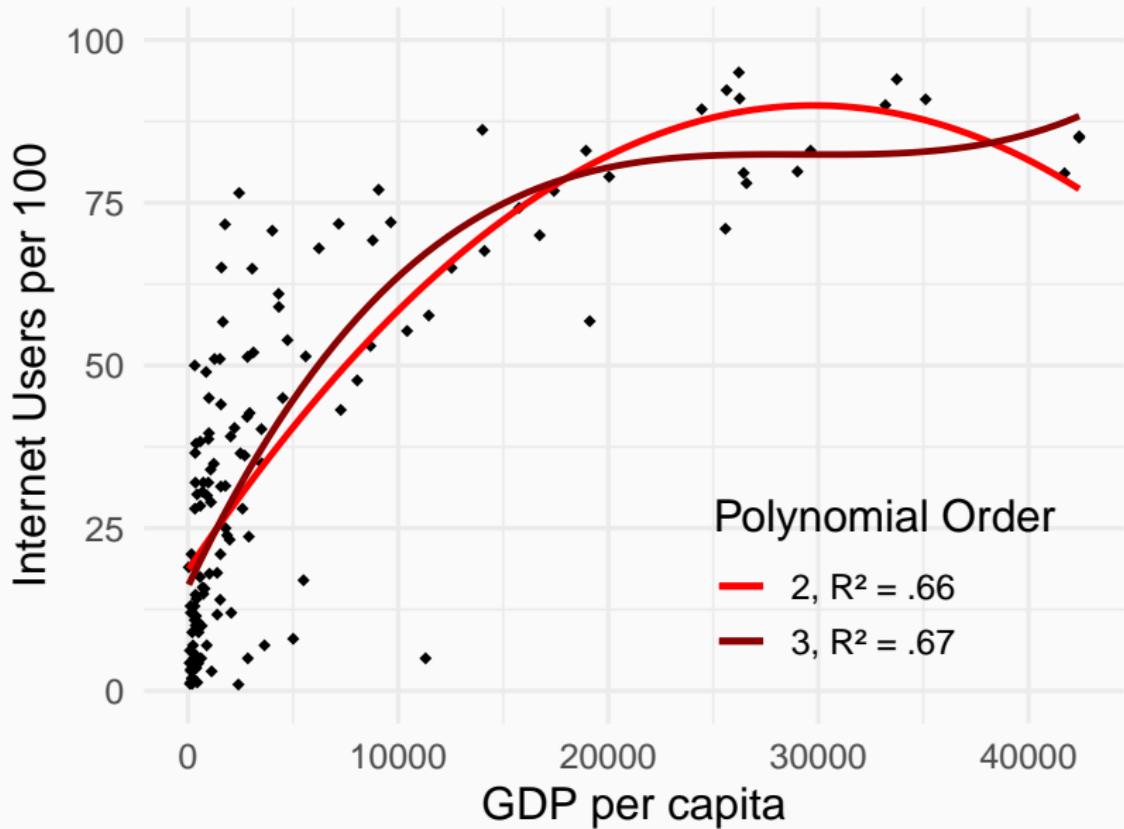
Example 1: If $\text{GDP} = 10,000$, then $\frac{\partial \widehat{\text{Users}}}{\partial \text{GDP}} = .0032$.

(Extra \$1,000 in GDP associated with 3.2 extra users)

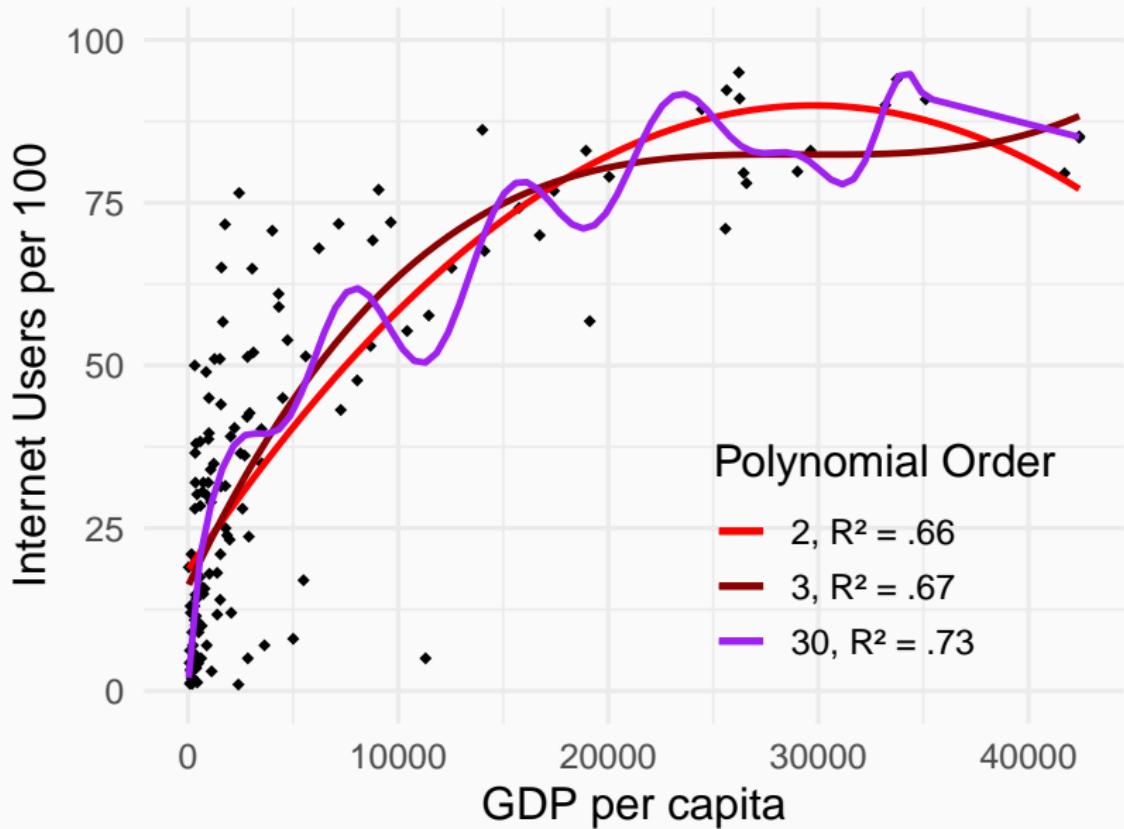
Example 2: If $\text{GDP} = 20,000$, then $\frac{\partial \widehat{\text{Users}}}{\partial \text{GDP}} = .0016$.

(Extra \$1,000 in GDP associated with 1.6 extra users)

HIGHER-ORDER POLYNOMIALS



HIGHER-ORDER POLYNOMIALS



Measurement With Controls

MEASUREMENT WITH CONTROLS

Note: This is a placeholder slide for an introduction that we will provide to the section. We're just placing it here for organization.

Interpreting Indicator Variables

EXAMPLE: MEASURING THE WAGE GAP

78 cents on the dollar: The facts about the gender wage gap

US women made economic strides in 2018, but pay gap persists

The gender pay gap is even worse if you're a woman with a college degree



EXAMPLE: MEASURING THE WAGE GAP (CONT.)

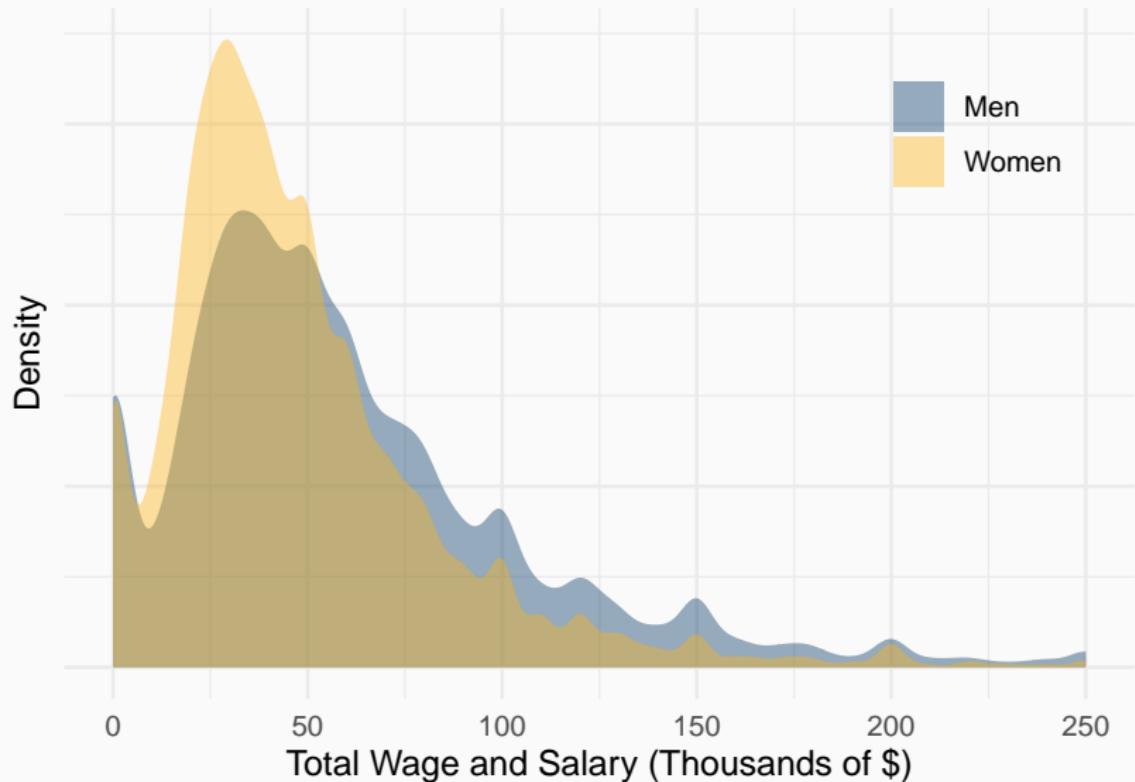
Data: Current Population Survey, 2019 Annual Social and Economic Supplement (ASEC)

Key variables:

- Status: full-/part-time work status
 - 1: not in labor force
 - 2: full-time hours (35 or more), usually full-time
 - ⋮
- Pay: total wage and salary earnings
- Sex
 - 1: male
 - 2: female

OVERALL COMPARISON

Women Earn Less than Men



DIRECT APPLICATION OF T-TEST

Mean for men: \$68,700

Mean for women: \$52,100

Mean difference: \$16,600

(76 cents to the dollar)

$$t = 27.938, df = 61,733, p < 2.2e - 16$$

THE SAME ANALYSIS IN A REGRESSION FRAMEWORK, PART I

$$\widehat{\text{Pay}} = \beta_0 + \beta_1 \text{Female}$$

For men, $\widehat{\text{Pay}} = \beta_0 + \beta_1(0) = \beta_0.$

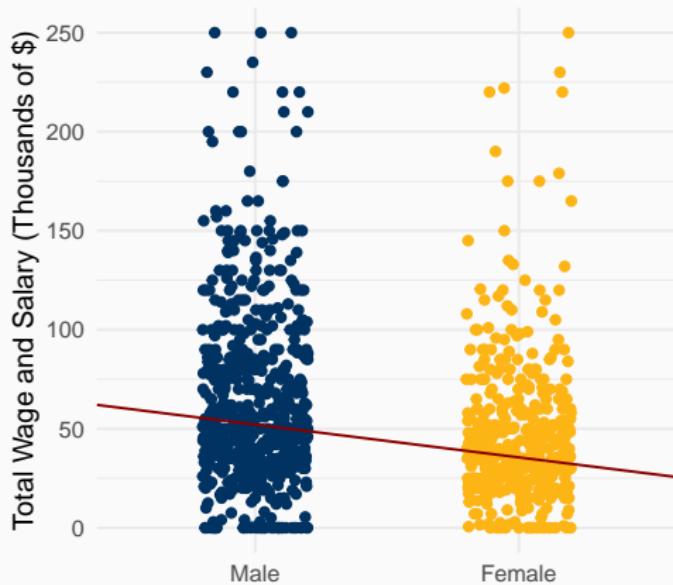
For women, $\widehat{\text{Pay}} = \beta_0 + \beta_1(1) = \beta_0 + \beta_1.$

$\implies \beta_1$ represents the difference between groups.

THE SAME ANALYSIS IN A REGRESSION FRAMEWORK, PART

II

Regression View of Wage Gap
Random Subset of 1,000 Data Points



$$\Delta Y = \frac{\text{slope}}{\Delta X}$$

THE SAME ANALYSIS IN A REGRESSION FRAMEWORK, PART III

Dependent Variable: Pay	
Female	-16,561*** (618)
Intercept	68,667*** (408)
Observations	62,110

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

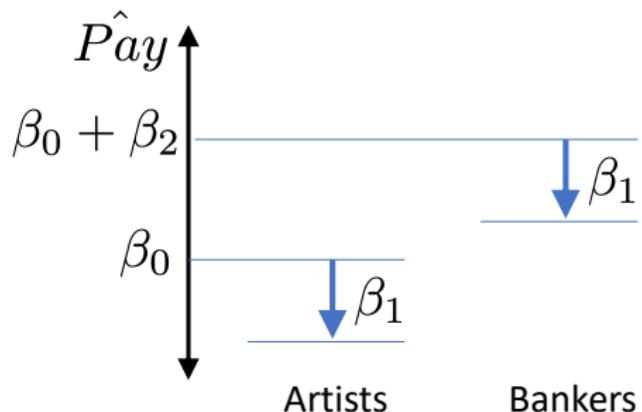
Indicator Variables as Controls

CONTROLLING FOR OCCUPATION, PART I

How can we compare men and women in the same occupation?

CONTROLLING FOR OCCUPATION, PART II

$$\widehat{Pay} = \beta_0 + \beta_1 Female + \beta_2 Banker + \beta_3 Engineer + \dots$$



CONTROLLING FOR OCCUPATION, PART III

Dependent Variable: Pay		
	(1)	(2)
Female	-16,561*** (618)	-15,035*** (679)
Intercept		183,292*** (3,440)
occup20		84,255*** (3,714)
occup40		76,881*** (13,132)
occup50		108,112*** (3,753)
:	:	:

CONTROLLING FOR OCCUPATION, PART IV

Dependent Variable: Pay		
	(1)	(2)
Female	-16,561*** (618)	-15,035*** (679)
Intercept	68,667*** (408)	7,110** (2,321)
Occupation FE	No	Yes
Observations	62,110	62,110

Note:

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

Metric Variables as Controls

AGE AS A COVARIATE

Ways to add age to a specification

Linear effect: $\widehat{Pay} = \beta_0 + \beta_1 Female + \beta_2 Age$

Polynomial effect: $\widehat{Pay} = \beta_0 + \beta_1 Female + \beta_2 Age + \beta_3 Age^2$

UNDERSTANDING THE LINEAR AGE EFFECT

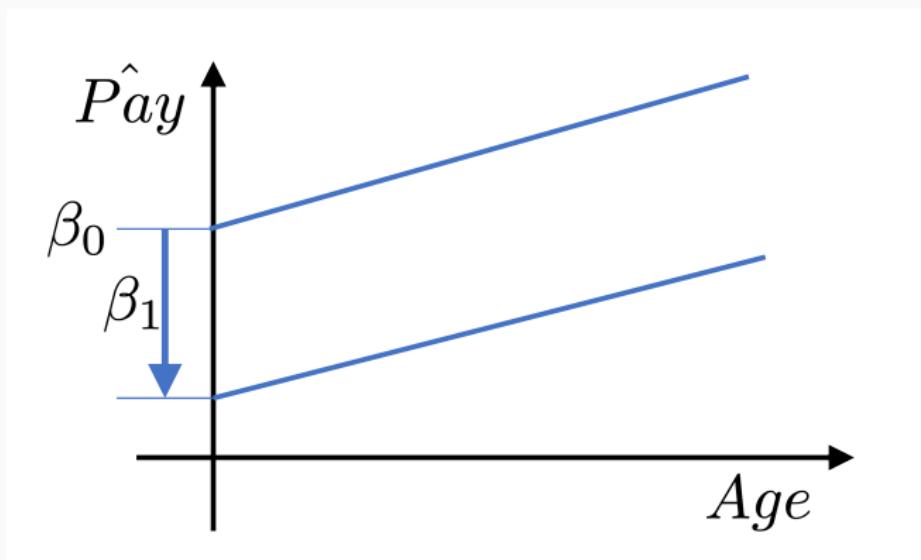
For men: $\widehat{Pay} = \beta_0 + \beta_1(0) + \beta_2 Age = \beta_0 + \beta_2 Age$

For women: $\widehat{Pay} = \beta_0 + \beta_1(1) + \beta_2 Age = \beta_0 + \beta_1 + \beta_2 Age$

UNDERSTANDING THE LINEAR AGE EFFECT

For men: $\widehat{Pay} = \beta_0 + \beta_1(0) + \beta_2 Age = \beta_0 + \beta_2 Age$

For women: $\widehat{Pay} = \beta_0 + \beta_1(1) + \beta_2 Age = \beta_0 + \beta_1 + \beta_2 Age$



CONTROLLING FOR AGE

Dependent Variable: Pay			
	(1)	(2)	(3)
Female	-16,561*** (618)	-15,035*** (679)	-15,063*** (677)
Age			469*** (22)
Intercept	68,667*** (408)	7,110** (2,321)	-10,140*** (2,447)
Occupation FE	No	Yes	Yes
Observations	62,110	62,110	62,110

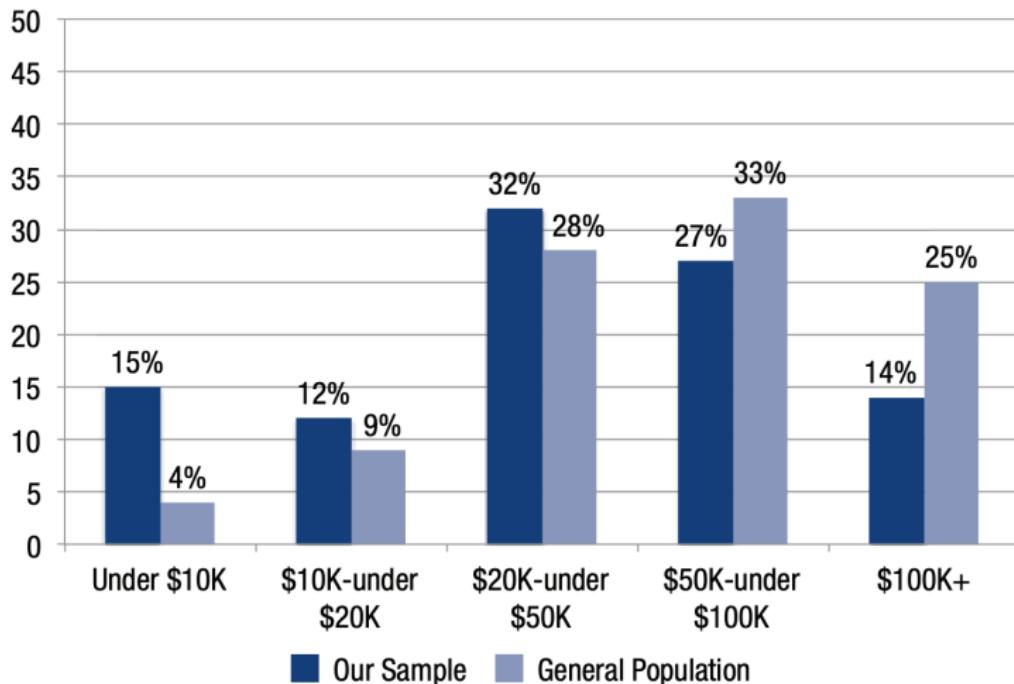
Note:

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

Measuring With More Categories

NATIONAL TRANSGENDER DISCRIMINATION SURVEY

Household Incomes of Respondents³



EXAMPLE OF A MORE INCLUSIVE QUESTION

How would you describe yourself?

- Female
- Male
- Transgender female
- Transgender male
- Non-binary/nonconforming
- Prefer not to answer

ANALYZING MORE GENDER CATEGORIES

$$\widehat{\text{Pay}} = \beta_0 + \beta_1 \text{Female} + \beta_2 \text{Transgender_Male} + \beta_3 \text{Transgender_Female} + \beta_4 \text{Nonbinary}$$

⇒ Report each β to describe the pay gap for a different gender identity.

REPORTING OVERALL EFFECT

Law of total variance:

$$V[Pay] = V[E[Pay|Gender]] + E[V[Pay|Gender]]$$

$$V[Pay] = \text{Gender-Explained Variance} + \text{Other Variance}$$

$$\eta^2 = \frac{\text{Gender-Explained Variance}}{\text{Total Variance}}$$

$$\text{Gender-Explained Standard Deviation} = \sqrt{V[E[Pay|Gender]]}$$

F-test: Null is that all genders have equal pay.

Planning Multiple Specifications

	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
Log working hours	-0.175*** (0.007)	-0.152*** (0.013)	-0.151*** (0.013)	-0.151*** (0.013)	-0.137*** (0.013)	-0.128*** (0.012)	-0.113*** (0.013)	-0.139*** (0.011)
Trend			0.001*** (0.000)	0.001*** (0.000)	0.001*** (0.000)	0.001*** (0.000)	0.001*** (0.000)	0.001*** (0.000)
<i>Day of the week dummies (reference: Monday)</i>								
Tuesday					0.004 (0.004)	0.007* (0.003)	-0.023 (0.024)	-0.452*** (0.045)
Wednesday					-0.000 (0.004)	0.001 (0.004)	0.029 (0.025)	-0.427*** (0.052)
Thursday					-0.001 (0.004)	0.001 (0.003)	0.066*** (0.025)	-0.609*** (0.063)
Friday					0.004 (0.004)	0.006* (0.004)	0.011 (0.026)	-0.424*** (0.050)
Saturday					0.109*** (0.008)	0.110*** (0.007)	-0.471*** (0.069)	-0.890*** (0.074)
Sunday					0.809*** (0.058)	0.417*** (0.063)	0.806*** (0.039)	0.608*** (0.040)
Age								-0.001*** (0.000)
Tenure								0.008*** (0.000)
Male								-0.003 (0.004)
Individual fixed effects	No	Yes	Yes	Yes	Yes	Yes	Yes	No
Team fixed effects	No	No	No	Yes	Yes	Yes	Yes	Yes
Hour-of-the-day dummies	No	No	No	No	No	Yes	Yes	Yes
Day fixed effects	No	No	No	No	No	No	Yes	Yes
R-squared	0.085	0.094	0.152	0.160	0.198	0.285	0.385	0.403
N	33,123	33,123	33,123	33,123	33,123	33,123	33,123	31,525
Individuals	332	332	332	332	332	332	332	332

WHY REPORT MULTIPLE SPECIFICATIONS?

Modeling decisions are tough.

- Is it worthwhile to control for education, even if standard errors increase?
- Is it worth switching from a log to a square root to get a much better fit?
- Should we include occupation, even if that absorbs some of the concept we want to measure?

Would your results be different with different choices?

⇒ Try different paths to see if your results are **robust**.

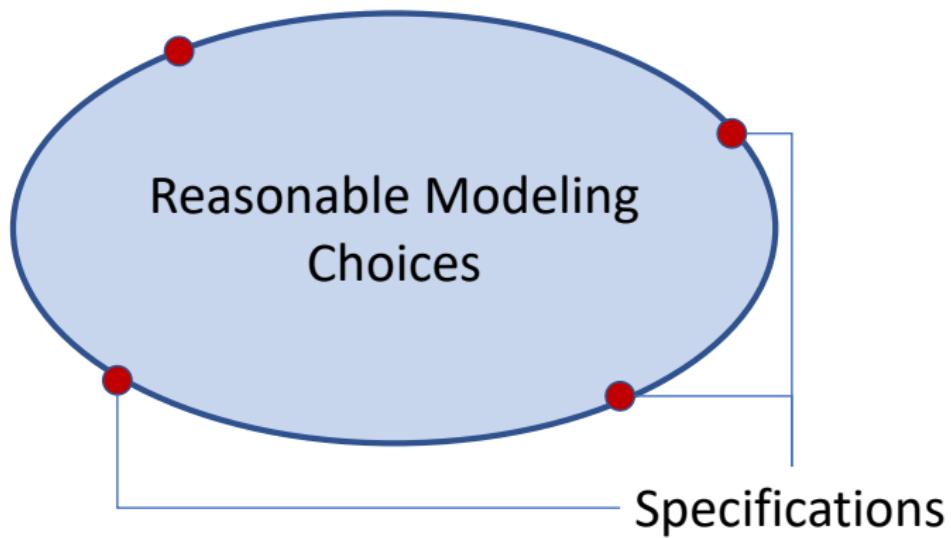
WHY REPORT MULTIPLE SPECIFICATIONS? (CONT.)

Error rate inflation: a tendency for effects to appear large (and p-values small), especially when a researcher uses the same data for exploration and testing

p-Hacking: a deliberate attempt to generate significant results by altering the model specification and other researcher degrees of freedom

⇒ Report multiple specifications to guard against error rate inflation.

How SHOULD YOU THINK ABOUT SPECIFICATIONS?



AN EXAMPLE SPECIFICATION TABLE

Dependent Variable: Pay			
	(1)	(2)	(3)
Female	-16,561*** (618)	-15,035*** (679)	-15,063*** (677)
Age			469*** (22)
Intercept	68,667*** (408)	7,110** (2,321)	-10,140*** (2,447)
Occupation FE	No	Yes	Yes
Observations	62,110	62,110	62,110

Note:

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

Modeling Conditional Effects

Conditional Effects

CONDITIONAL EFFECTS

Idea: The relationships we measure may be different for different people or units.

- A vaccine may reduce infection rates more in adults than in children.
- Network access may be more associated with loan availability in poor countries.
- The pay gap may be different for people of different ages.

CONDITIONAL EFFECTS OF AGE, PART I

Interaction term: a variable in a regression formed by multiplying two other variables together

$$\widehat{Pay} = \beta_0 + \beta_1 Female + \beta_2 Age + \beta_3 Female \cdot Age$$

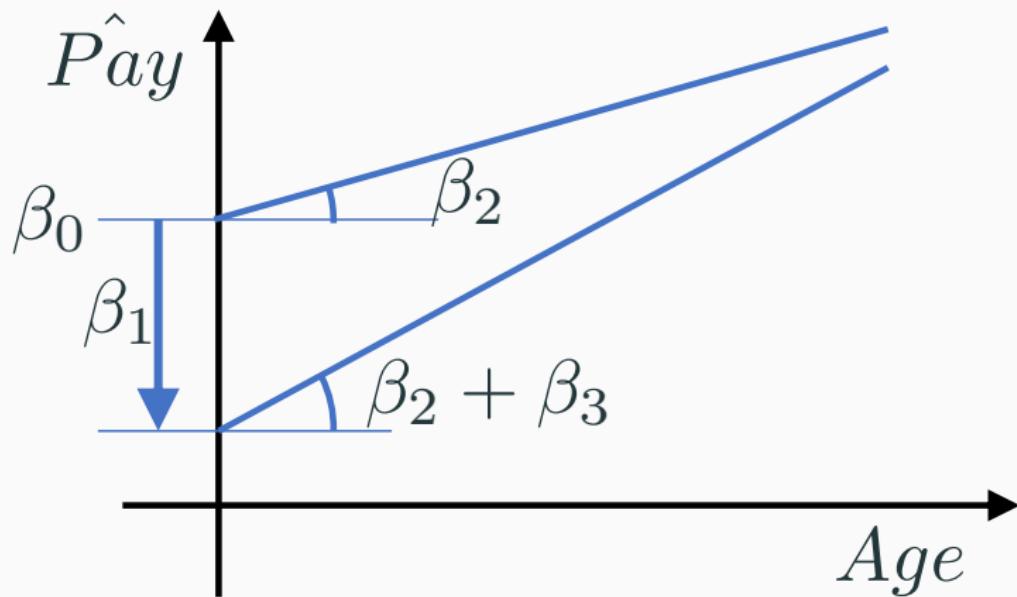
For men:

$$\widehat{Pay} =$$

For women:

$$\widehat{Pay} =$$

CONDITIONAL EFFECTS OF AGE, PART II



CONDITIONAL EFFECTS OF AGE, PART III

Dependent Variable: Pay	
Female	-6,230** (1,967)
Age	559*** (29)
Female:Age	-207*** (43)
Intercept	-13,971*** (2,575)
Occupation FE	No
Observations	62,110

Note:

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

Interaction Terms

INTERACTIONS FOR INDICATOR VARIABLES

Use old content: 12.9 Interaction Terms for Indicator Variables, Part 1

If there's extra time, I would redo this using the wage gap example.

GUIDELINES FOR POLYNOMIAL AND INTERACTION TERMS

Use old content: 12.12 Guidelines for Polynomial and Interaction Terms

Could redo some of the discussion of goals if there is time.