

# Regression Methods

*ECON 5753 — University of Arkansas*

Prof. Kyle Butts

February 2025

## Theoretical Results from Last Time

1. Derive the OLS estimator  $\hat{\beta}_{\text{OLS}} = (\mathbf{W}'\mathbf{W})^{-1}\mathbf{W}'y$

2. Derive the sample distribution of  $\hat{\beta}_{\text{OLS}}$

$$\hat{\beta}_{\text{OLS}} \sim \mathcal{N}(\beta_0, (\mathbf{W}'\mathbf{W})^{-1}\mathbf{W}'\Sigma\mathbf{W}(\mathbf{W}'\mathbf{W})^{-1})$$

3. Forecasting  $w'\hat{\beta}_{\text{OLS}}$  and  $w'\hat{\beta}_{\text{OLS}} \sim \mathcal{N}(w'\beta_0, w'V_{\hat{\beta}_{\text{OLS}}}w)$

4. Marginal (predictive) effects,  $\frac{\partial}{\partial x_\ell} \hat{f}(X) = \sum_{k=1}^K \frac{\partial}{\partial x_\ell} g_k(X) \hat{\beta}_{\text{OLS},k}$

## This time

The rest of this topic will cover various practical uses of regression:

- Interpreting regression coefficients
- Different explanatory variables to include in  $W_i$  and how to interpret them
  - Polynomials, Indicators, Discrete Variables, Bins, Splines
- Interactions
- Partially linear model
- log transformed variables

Basically, this set of slides is the *applied* half of cross-sectional regression.

# Multivariate Regression – “All Else Equal”

## Regressors

Polynomials

Indicators

Multi-valued discrete variables

Interactions

Bins

Splines

log **transformations**

0/1 outcomes

## Marginal Effects with Multiple Variables

Say we have two variables in our linear model  $\hat{y} = \beta_0 + \beta_1 X_1 + \beta_2 + X_2$ .

Our predictions for unit  $i$  and unit  $j$  differ by

$$\hat{y}_j - \hat{y}_i = \beta_1(X_{1,j} - X_{1,i}) + \beta_2(X_{2,j} - X_{2,i})$$

## Marginal Effects with Multiple Variables

$$\hat{y}_j - \hat{y}_i = \beta_1(X_{1,j} - X_{1,i}) + \beta_2(X_{2,j} - X_{2,i})$$

Let's think about a simple version. Take two individuals with the same  $X_2$ , but  $X_1$  that differs by 1 unit (say  $X_{1,j} - X_{1,i} = 1$ )

## Marginal Effects with Multiple Variables

$$\hat{y}_j - \hat{y}_i = \beta_1(X_{1,j} - X_{1,i}) + \beta_2(X_{2,j} - X_{2,i})$$

Let's think about a simple version. Take two individuals with the same  $X_2$ , but  $X_1$  that differs by 1 unit (say  $X_{1,j} - X_{1,i} = 1$ )

Then, our estimated change is

$$\Delta\hat{y} = \beta_1$$

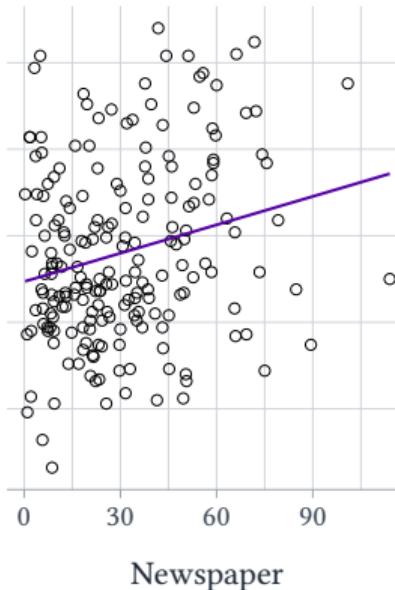
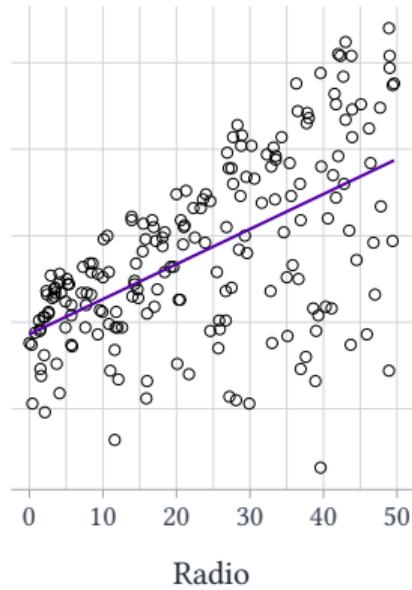
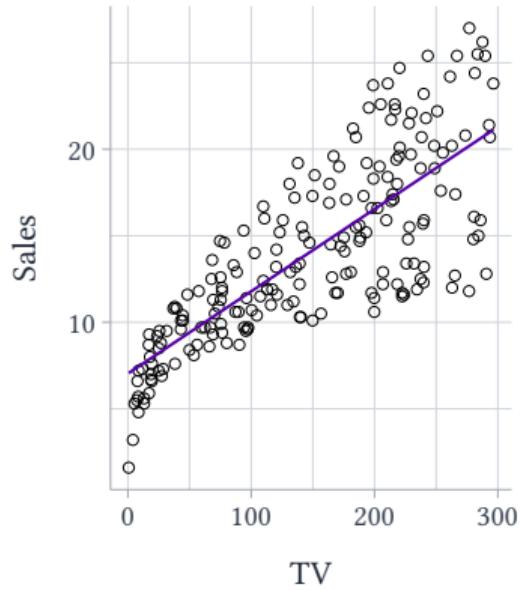
We refer to the  $\beta_k$  as “marginal effects”, i.e. the predicted change in  $y$  from a 1 unit increase in  $X$ , holding *fixed* all the other variables

## Advertising Example – Why Multiple Variable Regression?

Let's give an example. Say you're a business and you want to use advertising to boost sales. You have a bunch of different markets (e.g. cities) and you have data on how you've spent your advertising budget (TV, Radio, and Newspaper) and the sales in that market

Compare the results of regressing sales on all three variables to bivariate regressions

# Advertising Example – Bivariate Regressions



Dependent Variable:		Sales		
Constant		7.033*** (0.4578)	9.312*** (0.3882)	12.35*** (0.6338)
TV		0.0475*** (0.0027)		0.0458*** (0.0019)
Radio			0.2025*** (0.0217)	0.1885*** (0.0108)
Newspaper				0.0547*** (0.0186) -0.0010 (0.0064)
R <sup>2</sup>		0.61188	0.33203	0.05212
Signif. Codes: ***: 0.01, **: 0.05, *: 0.1				

## Interpreting regression results

Notice on the last table, in the bivariate regression, higher newspaper ads spending is associated with higher sales

*But*, after controlling for TV and Radio ads, there is no relationship between Newspaper ads and sales

- It seems like the bivariate relationship is being driven by newspaper ads being correlated with TV and Radio sales.
- Holding them fixed removed the relationship between Newspaper ads and sales

## “All Else Equal”

The latin word you'll sometimes see is *ceteris parabus* meaning “other things being equal”. This refers to variables *included* in your model

- In this case, we are “holding fixed” TV and Radio advertising
- Other things not included in the model still are moving as you compare areas with different Newspaper spending

For example, maybe cities with more newspaper spending have an older population than those with less

- OLS is assigning ‘credit’ to Newspaper for the effect of different age distributions in a city on sales

# Multivariate Regression – “All Else Equal”

## Regressors

Polynomials

Indicators

Multi-valued discrete variables

Interactions

Bins

Splines

log **transformations**

0/1 outcomes

# Multivariate Regression – “All Else Equal”

## Regressors

Polynomials

Indicators

Multi-valued discrete variables

Interactions

Bins

Splines

log **transformations**

0/1 outcomes

## Quadratic terms

Say you have the following model with wages as a quadratic function of age

$$w_i = \beta_0 + \beta_1 \text{age}_i + \beta_2 \text{age}_i^2 + u_i$$

Before we were discussing the idea of changing one variable, while holding the rest "equal"

→ How does one change age without changing age<sup>2</sup>?

This was why I focused on separating the variables you condition on  $X_i$  and the variables you include in your explanatory variables  $W_i = (g_1(X_i), \dots, g_K(X_i))'$

# Marginal Effects

$$\frac{\partial}{\partial x_\ell} \hat{f}(X) = \sum_{k=1}^K \frac{\partial}{\partial x_\ell} g_k(X) \hat{\beta}_{\text{OLS},k}$$

This is holding fixed all the other variables at the original covariate values:  $x_{1,i}, \dots, x_{K,i}$  and only changing  $x_\ell$

- If a given  $g_k$  does not depend on  $x_\ell$ , then  $\frac{\partial}{\partial x_\ell} g_k(X) = 0$
- multiple  $W_k$  can change from changing a particular  $x_\ell$

## Marginal effects with quadratic terms

$$\hat{w}_i = \hat{\beta}_0 + \hat{\beta}_1 \text{age}_i + \hat{\beta}_2 \text{age}_i^2$$

From calculus, we know that the partial derivative of  $\hat{w}_i$  with respect to  $\text{age}_i$  is given by

$$\frac{\partial}{\partial \text{age}_i} \hat{w}_i = \hat{\beta}_1 + 2\hat{\beta}_2 \text{age}_i$$

## Marginal effects with quadratic terms

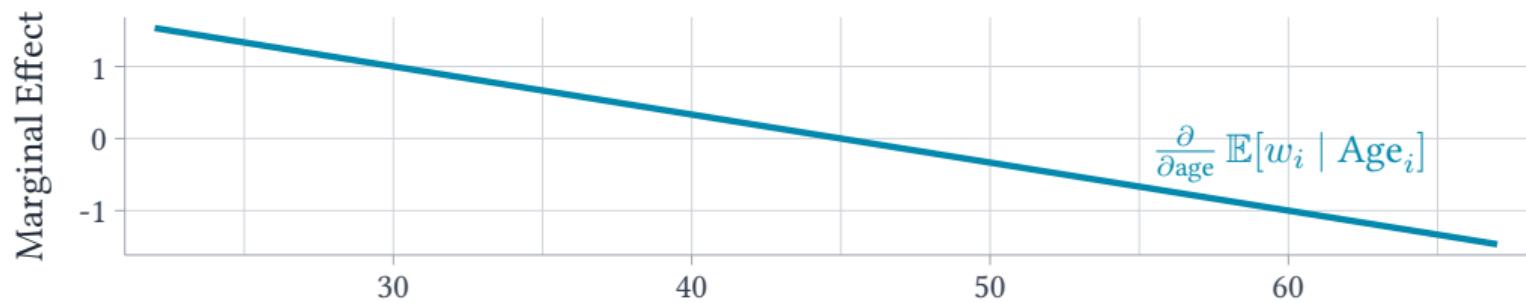
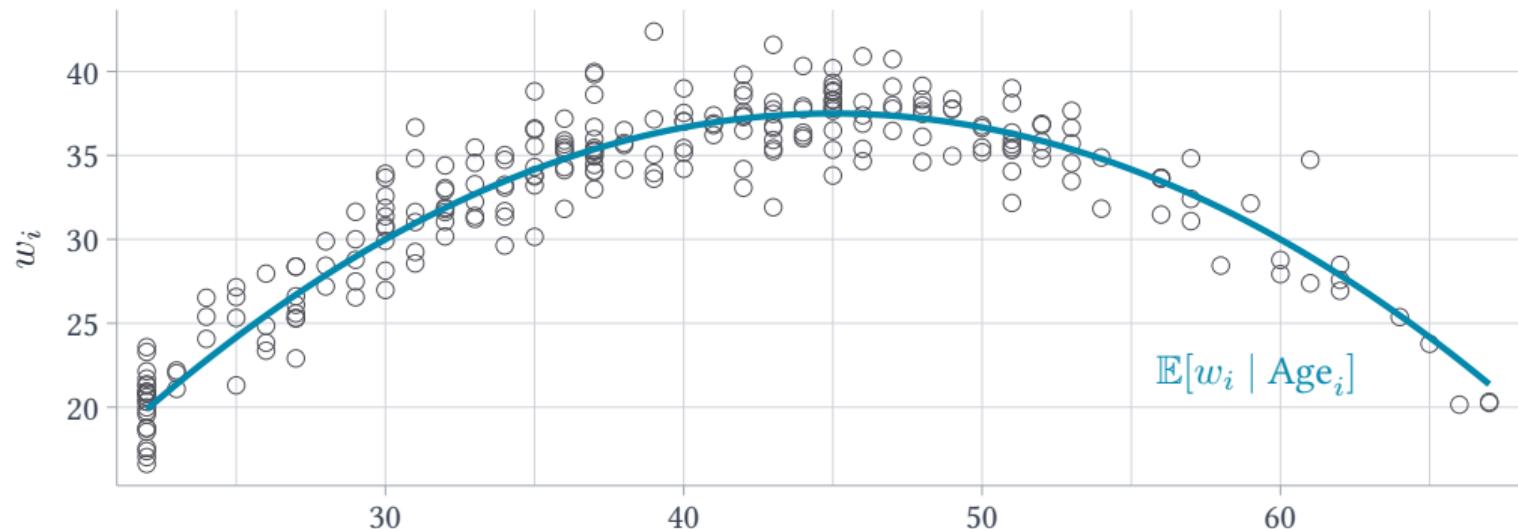
$$\hat{w}_i = \hat{\beta}_0 + \hat{\beta}_1 \text{age}_i + \hat{\beta}_2 \text{age}_i^2$$

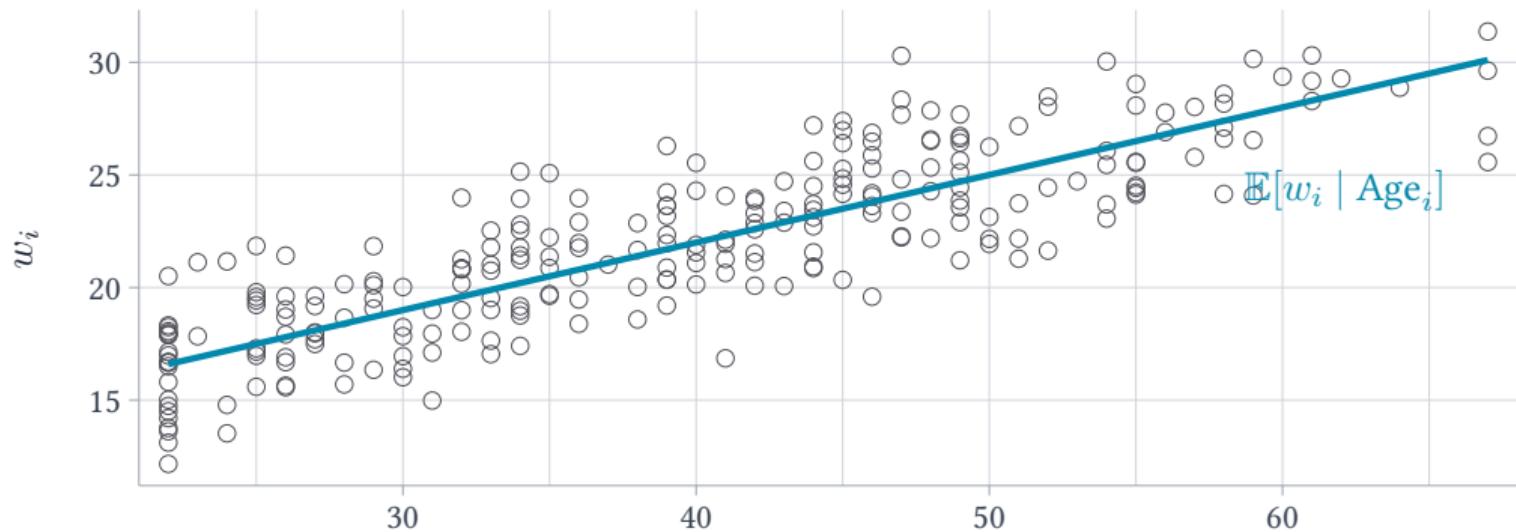
From calculus, we know that the partial derivative of  $\hat{w}_i$  with respect to  $\text{age}_i$  is given by

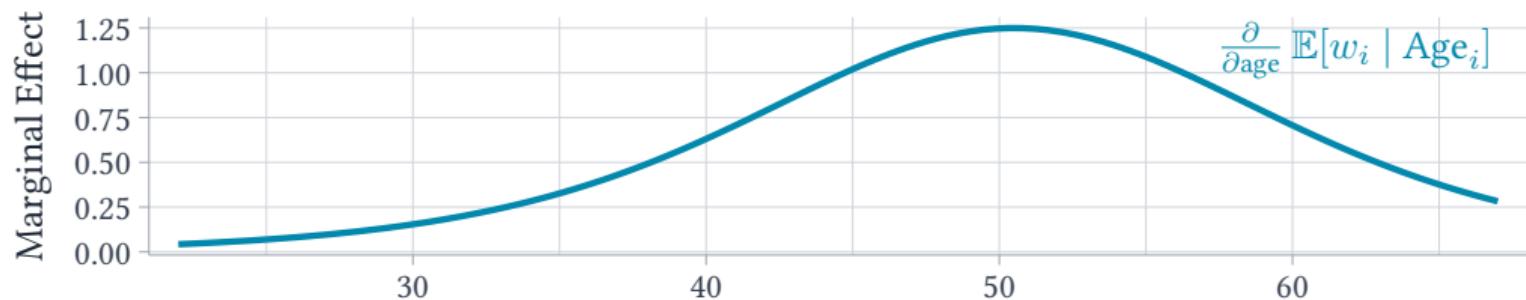
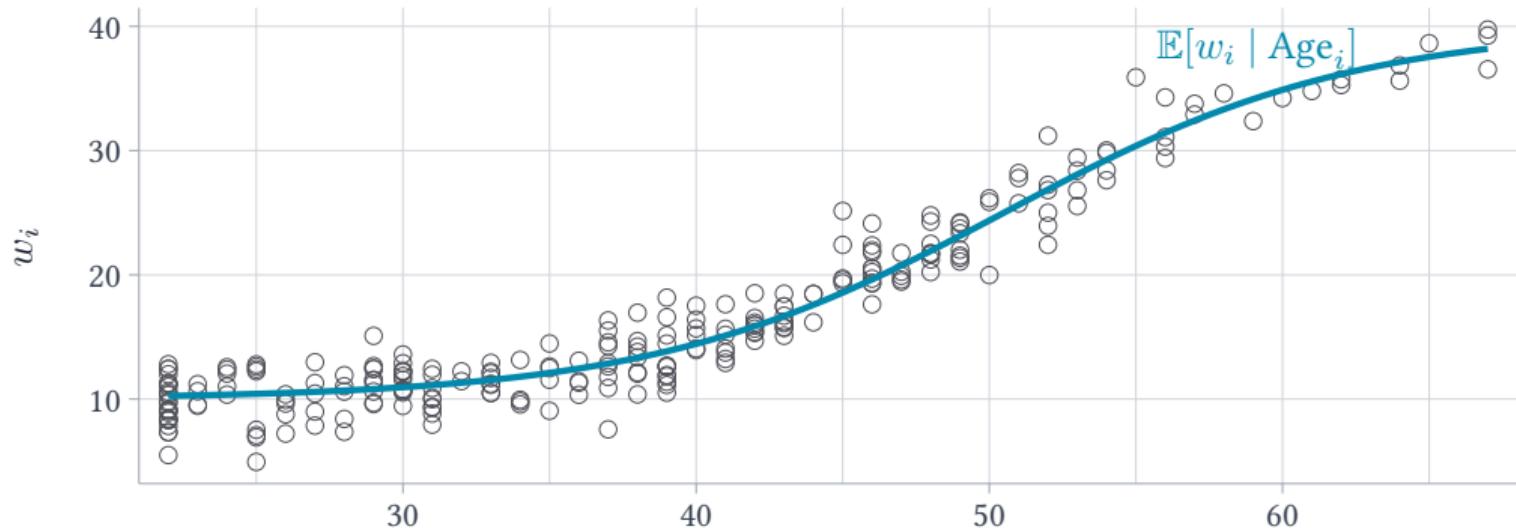
$$\frac{\partial}{\partial \text{age}_i} \hat{w}_i = \hat{\beta}_1 + 2\hat{\beta}_2 \text{age}_i$$

The change in predicted wage of a worker as they age is given by  $\hat{\beta}_1 + 2\hat{\beta}_2 \text{age}_i$  which depends on their age

- In words, how much predicted wages change as a worker gets a year older changes over a worker's lifetime







## Testing quadratic term

$$\hat{w}_i = \hat{\beta}_0 + \hat{\beta}_1 \text{age}_i + \hat{\beta}_2 \text{age}_i^2$$

In general, a linear term is much easier to interpret than a quadratic term: ‘a one unit change is associated with a  $\hat{\beta}_1$  unit change in  $w$ ’

→ We want to test whether the quadratic term is necessary or not

## Testing quadratic term

$$\hat{w}_i = \hat{\beta}_0 + \hat{\beta}_1 \text{age}_i + \hat{\beta}_2 \text{age}_i^2$$

You can test whether there is a significant quadratic relationship by testing the null that  $H_0 : \beta_2 = 0$ . Since we know  $\hat{\beta}_2$  has a normal distribution, can use a  $t$ -test:

$$\hat{t} \equiv \frac{\hat{\beta}_2 - 0}{\text{SE}(\hat{\beta}_2)}$$

If the  $p$ -value associated with this is larger than the level of significance, then we can not reject the simple linear model.

## Higher-order polynomials

This logic can be extended to higher-order polynomials to better estimate “wiggly” relationships between a given  $X$  and  $y$

- A key property of polynomials is that they are **smooth** functions
- Any smooth CEF can be approximated well by a polynomial of high-enough order (Taylor expansion)

## Higher-order polynomials

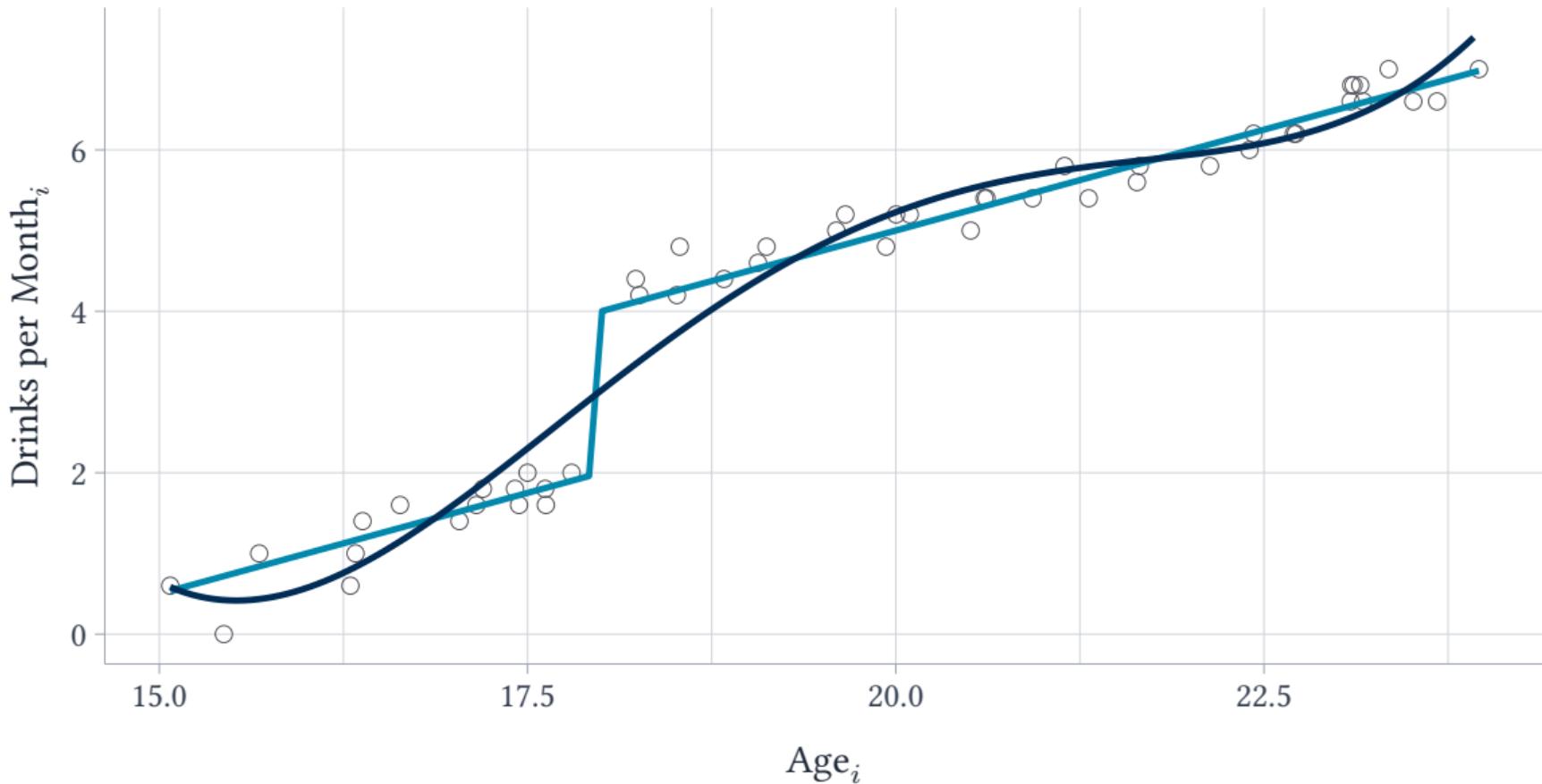
This logic can be extended to higher-order polynomials to better estimate “wiggly” relationships between a given  $X$  and  $y$

- A key property of polynomials is that they are **smooth** functions
- Any smooth CEF can be approximated well by a polynomial of high-enough order (Taylor expansion)

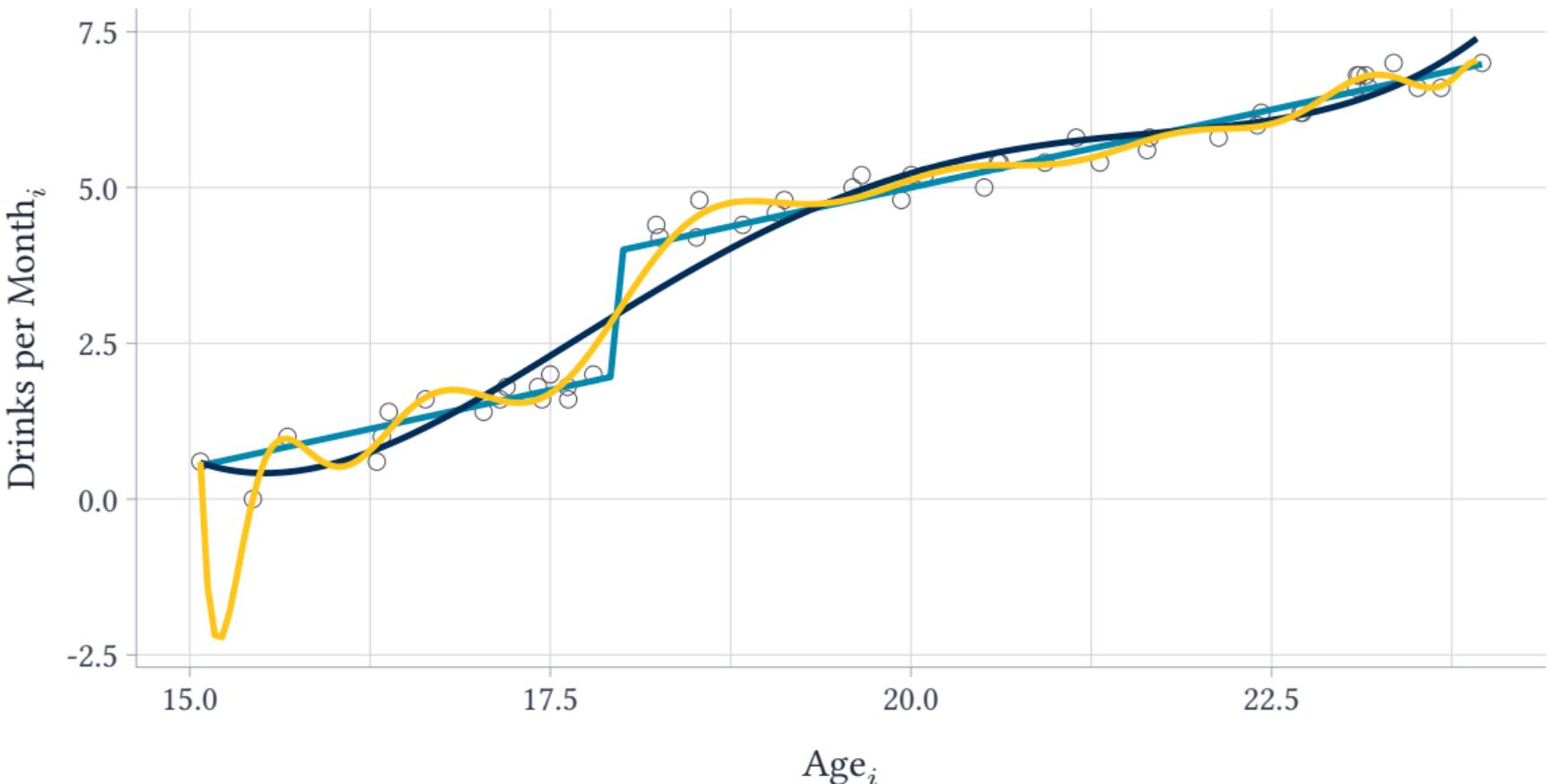
The smoothness can be a feature if you think outcomes evolve smoothly with a given  $X$

- But, it can be a problem if you think there's a discontinuity present (e.g. turing 21 has a large jump in rates of DUI)

## $\mathbb{E}[\text{Drinks per Month}_i \mid \text{Age}_i]$ , 4th-order Polynomial fit



$\mathbb{E}[\text{Drinks per Month}_i \mid \text{Age}_i]$ , 4th-order Polynomial fit, 15th-order Polynomial fit



## Polynomials and forecasting

One thing to note about polynomials is that they will always shoot off to  $\pm\infty$  as you move  $X \rightarrow \infty$  and  $-\infty$

When you forecast using your model *outside* of the values of  $X$  you train on, this is called **extrapolation**

- Even for relatively small values outside the samples' **domain** of  $X$  can have very strange forecasts

# Polynomials and forecasting

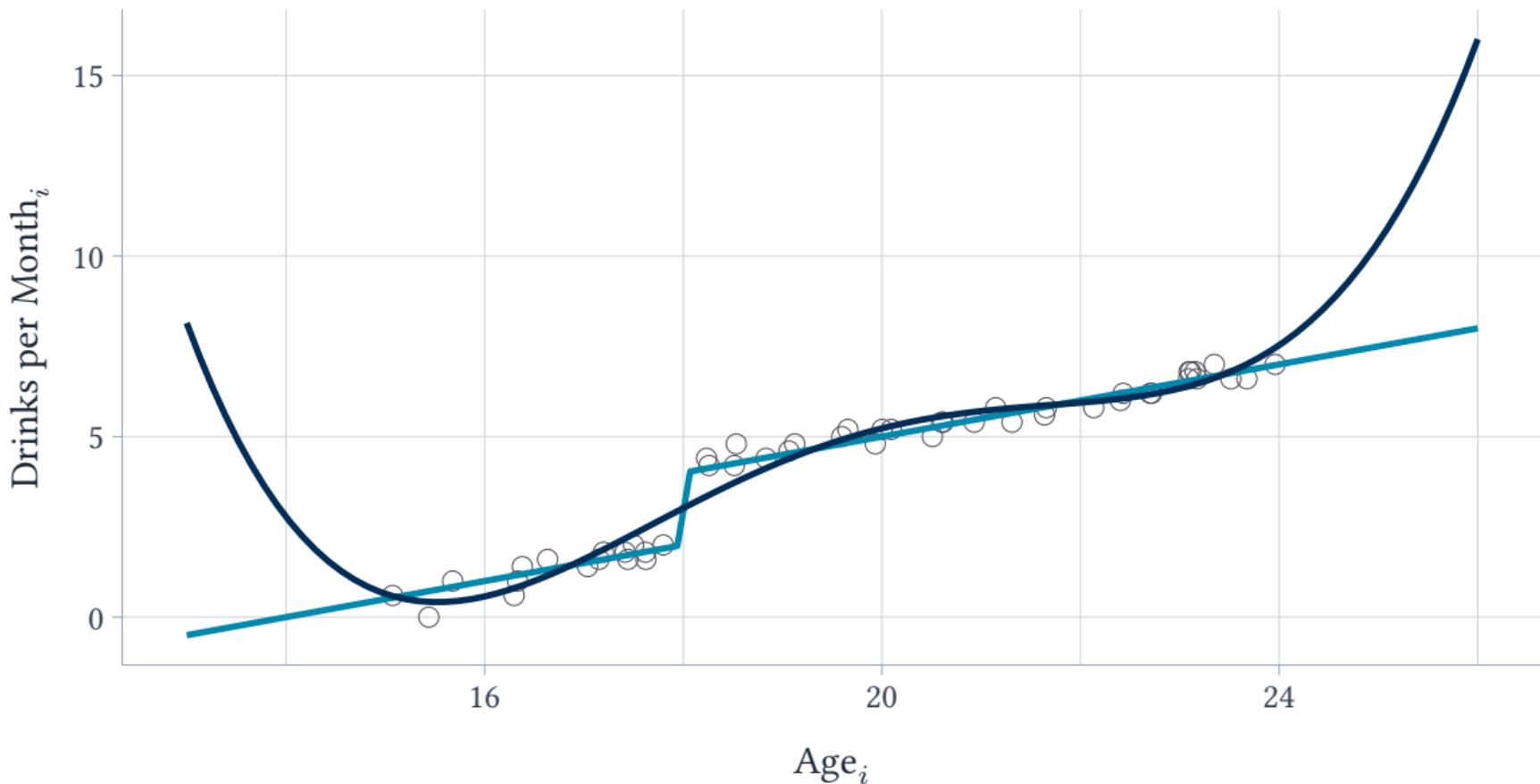
One thing to note about polynomials is that they will always shoot off to  $\pm\infty$  as you move  $X \rightarrow \infty$  and  $-\infty$

When you forecast using your model *outside* of the values of  $X$  you train on, this is called **extrapolation**

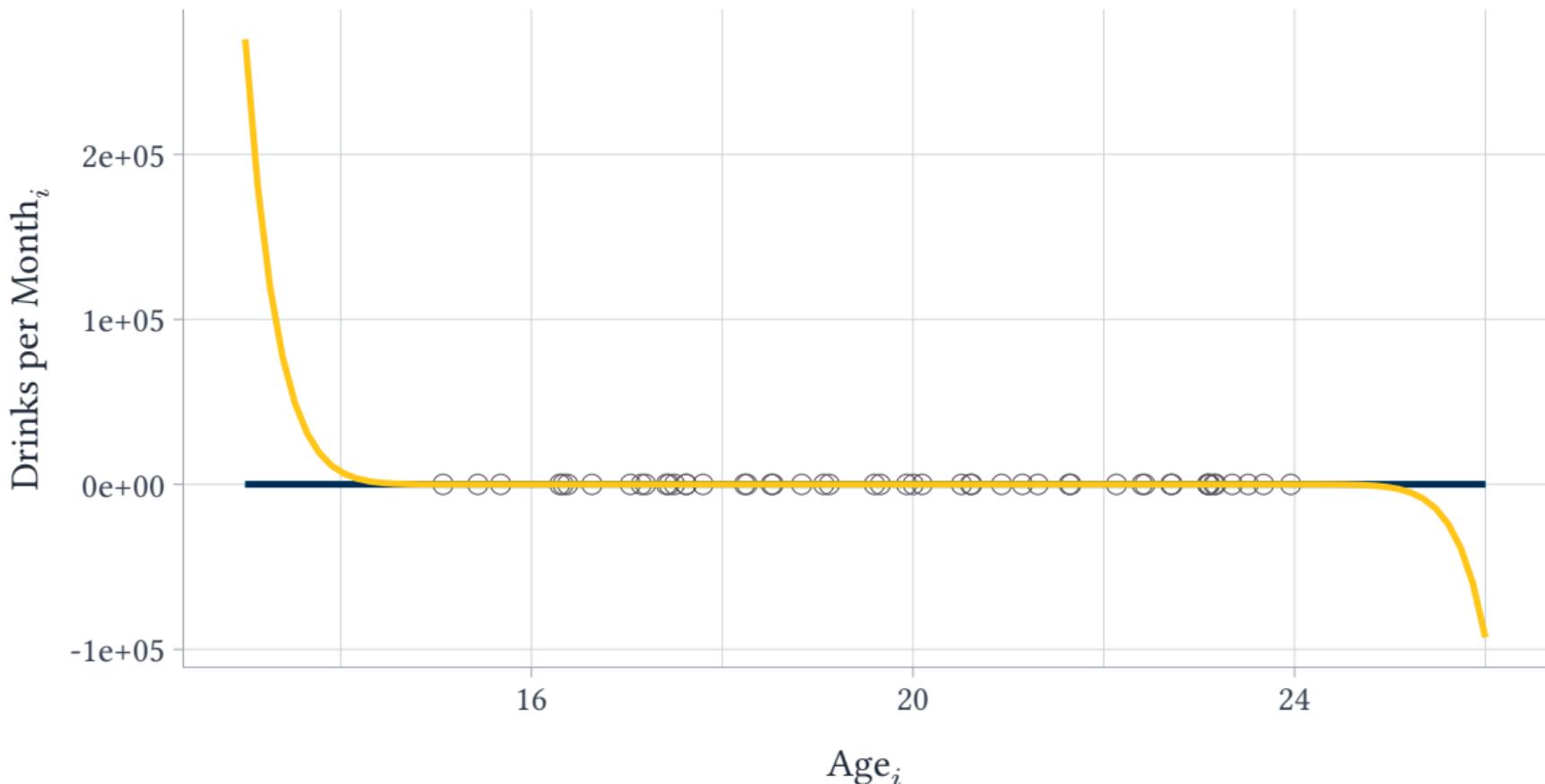
- Even for relatively small values outside the samples' **domain** of  $X$  can have very strange forecasts

The higher the order of the polynomial, the worse this gets

## $\mathbb{E}[\text{Drinks per Month}_i \mid \text{Age}_i]$ , 4th-order Polynomial fit



$\mathbb{E}[\text{Drinks per Month}_i \mid \text{Age}_i]$ , 4th-order Polynomial fit, 15th-order Polynomial fit



# Multivariate Regression – “All Else Equal”

## Regressors

Polynomials

Indicators

Multi-valued discrete variables

Interactions

Bins

Splines

log **transformations**

0/1 outcomes

## Intercept-only regression estimates sample mean

```
feols(y ~ 1, data = df)
```

Running this regression is useful because it will estimate the mean of  $y$  and give us the standard error estimate  $\frac{\hat{\sigma}}{\sqrt{n}}$

→ This makes inference easier: hypothesis testing and confidence intervals

## Comparing means using indicator variables

An *indicator variable* is a variable that can only equal 0 and 1, splitting units into groups

- $X$  "indicates" when a unit is of type 0 or type 1
- Often, it is written as  $\mathbb{1}[\cdot]$  where ‘ $\cdot$ ’ is a true/false condition

E.g. include

- being born male ( $= 1$ ) or female ( $= 0$ )
- being White ( $= 1$ ) or not ( $= 0$ )
- having a high-school degree ( $= 1$ ) or not ( $= 0$ )
- $\mathbb{1}[\text{Height}_i \geq 6]$ , being at least 6 foot tall

## Indicator variable

Let's work through some properties of an indicator variable,  $D_i$ . First, the sample mean of an indicator variable is the proportion of units with a 1:

$$\frac{1}{n} \sum_{i=1}^n D_i = \frac{\# \text{ of } 1s}{n} = \% \text{ of sample with } D_i = 1$$

## Indicator variable

Let's work through some properties of an indicator variable,  $D_i$ . First, the sample mean of an indicator variable is the proportion of units with a 1:

$$\frac{1}{n} \sum_{i=1}^n D_i = \frac{\# \text{ of } 1s}{n} = \% \text{ of sample with } D_i = 1$$

Define  $\pi$  as the (population) fraction of units with  $D_i = 1$ . Using the formula for variance, we can derive:

$$\text{Var}(D_i) = \pi(1 - \pi)$$

## Covariance with an indicator variable

What is  $\text{Cov}(D_i, Y_i)$ ? Rembmer

$$\text{Cov}(D_i, Y_i) = \mathbb{E}[D_i Y_i] - \mathbb{E}[D_i] \mathbb{E}[Y_i]$$

## Covariance with an indicator variable

What is  $\text{Cov}(D_i, Y_i)$ ? Rembmer

$$\text{Cov}(D_i, Y_i) = \mathbb{E}[D_i Y_i] - \mathbb{E}[D_i] \mathbb{E}[Y_i]$$

Again, skipping the math:

$$\text{Cov}(D_i, Y_i) = \pi(1 - \pi) (\mathbb{E}[Y_i | D_i = 1] - \mathbb{E}[Y_i | D_i = 0])$$

## Covariance with an indicator variable

Math details (if you're curious):

$$\begin{aligned}\text{Cov}(D_i, Y_i) &= \mathbb{E}[D_i Y_i] - \mathbb{E}[D_i] \mathbb{E}[Y_i] \\&= \mathbb{E}[D_i Y_i] - \pi (\pi \mathbb{E}[Y_i | D_i = 1] + (1 - \pi) \mathbb{E}[Y_i | D_i = 0]) \\&= \pi \mathbb{E}[Y_i | D_i = 1] - \pi \pi \mathbb{E}[Y_i | D_i = 1] - \pi (1 - \pi) \mathbb{E}[Y_i | D_i = 0] \\&= \pi(1 - \pi) (\mathbb{E}[Y_i | D_i = 1] - \mathbb{E}[Y_i | D_i = 0])\end{aligned}$$

## Regression with an indicator variable

Say you have a regression of  $Y_i = \beta_0 + \beta_1 * D_i + u_i$ . What does  $\hat{\beta}_0$  and  $\hat{\beta}_1$  equal?

$$\hat{\beta}_1 = \frac{\text{Cov}(D_i, Y_i)}{\text{Var}(D_i)}$$

## Regression with an indicator variable

Say you have a regression of  $Y_i = \beta_0 + \beta_1 * D_i + u_i$ . What does  $\hat{\beta}_0$  and  $\hat{\beta}_1$  equal?

$$\begin{aligned}\hat{\beta}_1 &= \frac{\text{Cov}(D_i, Y_i)}{\text{Var}(D_i)} = \frac{\pi(1 - \pi) (\mathbb{E}[Y_i | D_i = 1] - \mathbb{E}[Y_i | D_i = 0])}{\pi(1 - \pi)} \\ &= \mathbb{E}[Y_i | D_i = 1] - \mathbb{E}[Y_i | D_i = 0]\end{aligned}$$

The coefficient  $\hat{\beta}_1$  tells me the difference in sample means between the group with  $D_i = 1$  and the group with  $D_i = 0$

## Regression with an indicator variable

Say you have a regression of  $Y_i = \beta_0 + \beta_1 * D_i + u_i$ . From the last slide, we have:

$$\hat{\beta}_1 = \mathbb{E}[Y_i | D_i = 1] - \mathbb{E}[Y_i | D_i = 0]$$

Solving our other first-order condition for  $\hat{\beta}_0$ , we have:

$$\begin{aligned}\beta_0 &= \mathbb{E}[Y_i] - \hat{\beta}_1 \mathbb{E}[D_i] \\ &= \mathbb{E}[Y_i] - \hat{\beta}_1 \pi \\ &= \underbrace{\pi \mathbb{E}[Y_i | D_i = 1] + (1 - \pi) \mathbb{E}[Y_i | D_i = 0]}_{-\pi \mathbb{E}[Y_i | D_i = 1] - \pi \mathbb{E}[Y_i | D_i = 0]} \\ &= \mathbb{E}[Y_i | D_i = 0]\end{aligned}$$

## Marginal effects

Our forecast is  $\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 D_i$ . Since  $D_i$  can only be 0 or 1, our marginal effect just compares these values directly:

$$D_i = 0 : \quad \hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 * 0 = \hat{\beta}_0$$

$$D_i = 1 : \quad \hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 * 1 = \hat{\beta}_0 + \hat{\beta}_1$$

$\hat{\beta}_0$  is our predicted value for  $Y_i$  for the group with  $D_i = 0$  and  $\hat{\beta}_0 + \hat{\beta}_1$  is our predicted value for  $Y_i$  for the group with  $D_i = 1$

→ This makes  $\hat{\beta}_1$  is the *difference in the means* between those with  $D_i = 1$  compared to  $D_i = 0$

## Interpreting coefficient on an indicator

Our forecast is  $\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 D_i$ . We can interpret  $\hat{\beta}_1$  as follows:

“On average, someone with  $D_i = 1$  has a  $\hat{\beta}_1$  larger/smaller value of  $\hat{Y}$  compared to the units with  $D_i = 0$  *holding all else equal*”

→ Where you add the last part if you include additional covariates

## Example

Let's revisit our example with the `mtcars` dataset. There is an indicator variable, `am` for being an automatic ( $= 1$ ) or manual ( $= 0$ ). Regress the miles per gallon a car gets, `mpg`, on `am`.

→ In `fixest`, we can use `i(am)` to make it print out more nicely

```
with(mtcars, mean(mpg[am == 0]))  
#> [1] 17.14737  
with(mtcars, mean(mpg[am == 1]) - mean(mpg[am == 0]))  
#> [1] 7.24494  
  
feols(mpg ~ i(am), data = mtcars)  
#> OLS estimation, Dep. Var.: mpg  
#> Observations: 32  
#> Standard-errors: IID  
#> Estimate Std. Error t value Pr(>|t|)  
#> (Intercept) 17.14737 1.12460 15.24749 1.1340e-15 ***  
#> am::1 7.24494 1.76442 4.10613 2.8502e-04 ***
```

# Multivariate Regression – “All Else Equal”

## Regressors

Polynomials

Indicators

Multi-valued discrete variables

Interactions

Bins

Splines

log **transformations**

0/1 outcomes

## Multi-valued discrete variables

This intuition will extend directly to settings where we have a discrete variable that obtains  $K$  distinct values:

- E.g. race/ethnicity, 10-year bins of age, number of cylinders in engine

## Multi-valued discrete variables

This intuition will extend directly to settings where we have a discrete variable that obtains  $K$  distinct values:

- E.g. race/ethnicity, 10-year bins of age, number of cylinders in engine

We can construct a *set of* indicator variables for each value that  $X$  can obtain. For  $k = 1, \dots, K$

$$X_{ik} \equiv \mathbb{1}[X_i = x_k]$$

There are  $K$  such variables:  $X_{i1}, \dots, X_{iK}$

## Multi-valued variable regression

$$y_i = \sum_{k=1}^K X_{ik}\beta_k + u_i$$

From the same intuition as before, we have  $\hat{\beta}_k$  is the sample average of  $y_i$  for individuals with  $X_i = x_k$

- We are in a very special case since these variables are mutually exclusive (only one of them is non-zero per unit), so this is easy to show with matrix algebra

## Example

Let's revisit our example with the `mtcars` dataset. Let's see if `mpg` differs based on the number of cylinders a car has, `cyl`.

- In `fixest`, we can use `i(cyl)` to make indicators for each value of a variable
- Otherwise, we could for 4, 6, and 8 create the indicator variables with e.g.

```
mtcars$cyl4 = (mtcars$cyl == 4)
```

Interpret these coefficients:

```
feols(mpg ~ 0 + i(cyl), data = mtcars)
#> OLS estimation, Dep. Var.: mpg
#> Observations: 32
#> Standard-errors: IID
#>           Estimate Std. Error t value Pr(>|t|)
#> cyl::4   26.6636   0.971801 27.4373 < 2.2e-16 ***
#> cyl::6   19.7429   1.218217 16.2064 4.4933e-16 ***
#> cyl::8   15.1000   0.861409 17.5294 < 2.2e-16 ***
#> ---
#> Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
#> RMSE: 3.0683  Adj. R2: 0.704784
```

## Intercept and Multicollinearity

In reality, we will want to include many predictor variables (beyond our single multi-valued discrete variable). In this case, we want to include an *intercept*

$$y_i = \beta_0 + \sum_{k=1}^K X_{ik} \beta_k + u_i$$

# Multicollinearity

Our matrix  $\mathbf{W}$  look like this. Note the 3 cyl indicator variables sum to the intercept

```
#> (Intercept) cyl::4 cyl::6 cyl::8  
#> 1 0 1 0  
#> 1 0 1 0  
#> 1 1 0 0  
#> 1 0 1 0  
#> 1 0 0 1  
#> 1 0 1 0  
#> 1 0 0 1  
#> 1 1 0 0
```

# Multicollinearity

$$\hat{y}_i = \hat{\beta}_0 + \sum_{k=1}^K X_{ik} \hat{\beta}_k$$

It turns out that we face a non-uniqueness problem because of the **multicollinearity** we identified

→ We can add 10 to  $\hat{\beta}_0$  and subtract 10 from  $\hat{\beta}_4$ ,  $\hat{\beta}_6$ , and  $\hat{\beta}_8$  and get the exact same  $\hat{y}$

Therefore, we will typically need to drop one of the  $X_{ik}$  variables (or R will do it for you)

```
feols(mpg ~ 1 + i(cyl), data = mtcars)
#> OLS estimation, Dep. Var.: mpg
#> Observations: 32
#> Standard-errors: IID
#>           Estimate Std. Error t value Pr(>|t|)
#> (Intercept) 26.66364  0.971801 27.43735 < 2.2e-16 ***
#> cyl::6      -6.92078  1.558348 -4.44110 1.1947e-04 ***
#> cyl::8      -11.56364 1.298623 -8.90453 8.5682e-10 ***
#> ---
#> Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
#> RMSE: 3.0683 Adj. R2: 0.714009
```

## Interpreting Multicollinearity

In the previous example, we dropped  $\mathbb{1}[X_i = 4]$ . This is the **omitted group**. What happened to our coefficient estimates?

Just like in the indicator variable case,  $\hat{\beta}_0$  estimated the mean of mpg for cars with  $X_i = 4$

## Interpreting Multicollinearity

In the previous example, we dropped  $\mathbb{1}[X_i = 4]$ . This is the **omitted group**. What happened to our coefficient estimates?

Just like in the indicator variable case,  $\hat{\beta}_0$  estimated the mean of mpg for cars with  $X_i = 4$

The coefficients on the other  $\hat{\beta}_k$  now represent the *difference* in means between the group for  $X_i = 6$  and the ‘omitted group’  $X_i = 4$ .

→ The mean for  $X_i = 6$  is  $19.742 = 26.663 - 6.921$

## Specifying ref option

```
feols(mpg ~ i(cyl, ref = 6), data = mtcars)
#> OLS estimation, Dep. Var.: mpg
#> Observations: 32
#> Standard-errors: IID
#>             Estimate Std. Error t value Pr(>|t|)
#> (Intercept) 19.74286    1.21822 16.20636 4.4933e-16 ***
#> cyl::4       6.92078    1.55835  4.44110 1.1947e-04 ***
#> cyl::8      -4.64286    1.49200 -3.11182 4.1522e-03 **
#> ---
#> Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
#> RMSE: 3.0683 Adj. R2: 0.714009
```

## Significance with indicator variables

```
#>           Estimate Std. Error t value Pr(>|t|)  
#> (Intercept) 26.66364   0.971801 27.43735 < 2.2e-16 ***  
#> cyl::6       -6.92078   1.558348 -4.44110 1.1947e-04 ***  
#> cyl::8       -11.56364  1.298623 -8.90453 8.5682e-10 ***  
#> ---  
#> Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Including an intercept also helps with certain statistical inference. The estimates test that the average of  $y$  for the omitted group is *the same* for the other groups

→ Rejecting this ( $p\text{-value} < \alpha$ ) rejects the null that the two means are the same

## Interpreting coefficient on indicators for a discrete variable

Our forecast is  $\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_6 X_{i6} + \hat{\beta}_8 X_{i8}$ . We can interpret  $\hat{\beta}_6$  as follows:

“On average, a car with 6 cylinders has a  $\hat{\beta}_6$  larger/smaller value of mpg compared to cars with 4 cylinders *holding all else equal*”

→ Where you add the last part if you include additional covariates

# Multivariate Regression – “All Else Equal”

## Regressors

Polynomials

Indicators

Multi-valued discrete variables

Interactions

Bins

Splines

log **transformations**

0/1 outcomes

# Why interactions

## *Wages Example*

Consider a model where we want to understand how wages are influenced by both being a female and being a college-educated worker. We can write the model as:

$$w_i = \beta_0 + \beta_1 \text{female}_i + \beta_2 \text{college}_i + \beta_3 (\text{female}_i \times \text{college}_i) + u_i$$

Here,  $\beta_3$  captures the interaction effect between female and college-education status on wages

- This means that the effect of females on wages may differ depending on whether the worker has a college-degree, and vice versa.

# Interactions

## *Wages Example*

Consider the difference in predicted wages for non-college educated male vs. non-college educated workers:

$$w_{i,NC,F} - w_{i,NC,M} = (\beta_0 + \beta_1) - \beta_0 = \beta_1$$

Compare this to the difference in predicted wages for college educated male vs. college educated workers:

$$w_{i,C,F} - w_{i,C,M} = (\beta_0 + \beta_1 + \beta_2 + \beta_3) - (\beta_0 + \beta_2) = \beta_1 + \beta_3$$

# Interactions

## *Wages Example*

Wage-gap for college educated workers is  $\beta_1 + \beta_3$  and the wage-gap for non-college educated workers is  $\beta_1$

→  $\beta_3$  represents the difference in wage gaps of college-educated vs. non-college-educated workers.

# Interactions

## *Wages Example*

Wage-gap for college educated workers is  $\beta_1 + \beta_3$  and the wage-gap for non-college educated workers is  $\beta_1$

→  $\beta_3$  represents the difference in wage gaps of college-educated vs. non-college-educated workers.

More generally, can interpret interactions how one variable changes the marginal effect of another variable

→ This is similar to when you have a quadratic function of  $X$ , the marginal effect depends where you are along the distribution of  $X$ .

# Interactions

## *Partial Derivative*

$$\hat{w}_i = \hat{\beta}_0 + \hat{\beta}_1 \text{female}_i + \hat{\beta}_2 \text{college}_i + \hat{\beta}_3 (\text{female}_i \times \text{college}_i)$$

We can derive this result using partial derivatives:

$$\frac{\partial \hat{w}_i}{\partial \text{female}_i} = \hat{\beta}_1 + \hat{\beta}_3 \text{college}_i$$

# Interactions

## *Partial Derivative*

$$\hat{w}_i = \hat{\beta}_0 + \hat{\beta}_1 \text{female}_i + \hat{\beta}_2 \text{college}_i + \hat{\beta}_3 (\text{female}_i \times \text{college}_i)$$

We can derive this result using partial derivatives:

$$\frac{\partial \hat{w}_i}{\partial \text{female}_i} = \hat{\beta}_1 + \hat{\beta}_3 \text{college}_i$$

In this case, it's a little weird to think of "small" changes in the female variable. Instead, we will think of it as a 1 unit change (from 0 to 1)

## Continuous interacted with a discrete variable

Let  $X_i$  be a continuous variable and  $D_i$  be a dummy variable and consider the regression

$$y_i = \beta_0 + D_i\beta_1 + X_i\beta_2 + X_iD_i\beta_3 + u_i$$

## Continuous interacted with a discrete variable

Let  $X_i$  be a continuous variable and  $D_i$  be a dummy variable and consider the regression

$$y_i = \beta_0 + D_i\beta_1 + X_i\beta_2 + X_iD_i\beta_3 + u_i$$

In this case, the marginal effect of  $X_i$  is given by  $\frac{\partial \hat{y}_i}{\partial X_i} = \hat{\beta}_2 + D_i\hat{\beta}_3$

- The marginal effect for group  $D_i = 0$  is given by  $\hat{\beta}_2$
- The marginal effect for group  $D_i = 1$  is given by  $\hat{\beta}_2 + \hat{\beta}_3$

## Continuous interacted with a discrete variable

Let  $X_i$  be a continuous variable and  $D_i$  be a dummy variable and consider the regression

$$y_i = \beta_0 + D_i\beta_1 + X_i\beta_2 + X_iD_i\beta_3 + u_i$$

In this case, the marginal effect of  $X_i$  is given by  $\frac{\partial \hat{y}_i}{\partial X_i} = \hat{\beta}_2 + D_i\hat{\beta}_3$

- The marginal effect for group  $D_i = 0$  is given by  $\hat{\beta}_2$
- The marginal effect for group  $D_i = 1$  is given by  $\hat{\beta}_2 + \hat{\beta}_3$
- Therefore,  $\hat{\beta}_3$  is the difference in marginal effects between  $D_i = 1$  relative to  $D_i = 0$

## Continuous interacted with a discrete variable

$$y_i = \beta_0 + D_i\beta_1 + X_i\beta_2 + X_iD_i\beta_3 + u_i$$

Exercise:

- In words, how would you interpret a  $t$ -test for the null that  $\hat{\beta}_2 = 0$ ?
- In words, how would you interpret a  $t$ -test for the null that  $\hat{\beta}_3 = 0$ ?

## mtcars example

```
OLS estimation, Dep. Var.: mpg
```

	Estimate	Std. Error	t value	Pr(> t )	
(Intercept)	26.624848	1.346754	19.769644	< 2.2e-16	***
am::1	5.217653	2.324898	2.244250	3.2904e-02	*
hp	-0.059137	0.008957	-6.602265	3.6781e-07	***
am::1:hp	0.000403	0.013362	0.030152	9.7616e-01	
---					
Signif. codes:	0	'***'	0.001	'**'	0.01
		'*'	0.05	'. '	0.1
				' '	1

Exercise: Does the estimate relationship between being a car's horsepower and miles per gallon depend on whether it is an automatic? How do you know?

## Interaction terms always should have ‘main effects’

When including an interaction term, it is important to (almost) always include the **main effects**

The **main effects** are the variables by themselves. E.g. if you interact gender with race, you want to include race and gender as separate terms as well

- The main effects are what let us interpret the interaction term as the ‘difference’ in marginal effects

## Continuous-Continuous interactions

Now consider two continuous variables being interacted:

$$y_i = \beta_0 + X_{1,i}\beta_1 + X_{2,i}\beta_2 + X_{1,i}X_{2,i}\beta_3 + u_i$$

This is common when you think there are complementarities between variables

- E.g.  $y$  is crop-yield,  $X_1$  is the amount of fertilizer applied, and  $X_2$  is the amount of water given. Does it help to do more of both (complements)?
- $y$  is a measure of job performance,  $X_1$  is a measure of on-the-job experience, and  $X_2$  is a measure of training

## Continuous-Continuous interactions

$$\frac{\partial \hat{y}_i}{\partial X_{1,i}} = \hat{\beta}_1 + X_{2,i}\hat{\beta}_3 \quad \text{and} \quad \frac{\partial \hat{y}_i}{\partial X_{2,i}} = \hat{\beta}_2 + X_{1,i}\hat{\beta}_3$$

Can interpret it in two ways:

1. The marginal effect of  $X_1$  grows/shrinks with the value of  $X_2$  (depending on the sign of  $\hat{\beta}_3$ )

# Multivariate Regression – “All Else Equal”

## Regressors

Polynomials

Indicators

Multi-valued discrete variables

Interactions

Bins

Splines

log **transformations**

0/1 outcomes

## Key explanatory variables

In a forecast model, we might have a bunch of different variables we want to use as predictors

Many times, we really want to see how the forecasted outcome varies with a *particular* variable of interest

- Want to “control” for a bunch of variables and then model the relationship between  $y$  and key variable,  $X_1$ , flexibly

## Partially-linear model

The **Partially linear model** mixes high model flexibility in a key variable we care about and linear model for the rest of the covariates:

$$y_i = \mu(X_{1i}) + W_i' \beta + u_i$$

- $\mu(X_{1i})$  is a highly flexible function
- $W_i$  is a set of *linear* control variables

This allows you to prevent the curse of dimensionality by linearly controlling for most of the variables. Allows a flexible model for the key variable of interest,  $X_i$ , that is good for graphing.

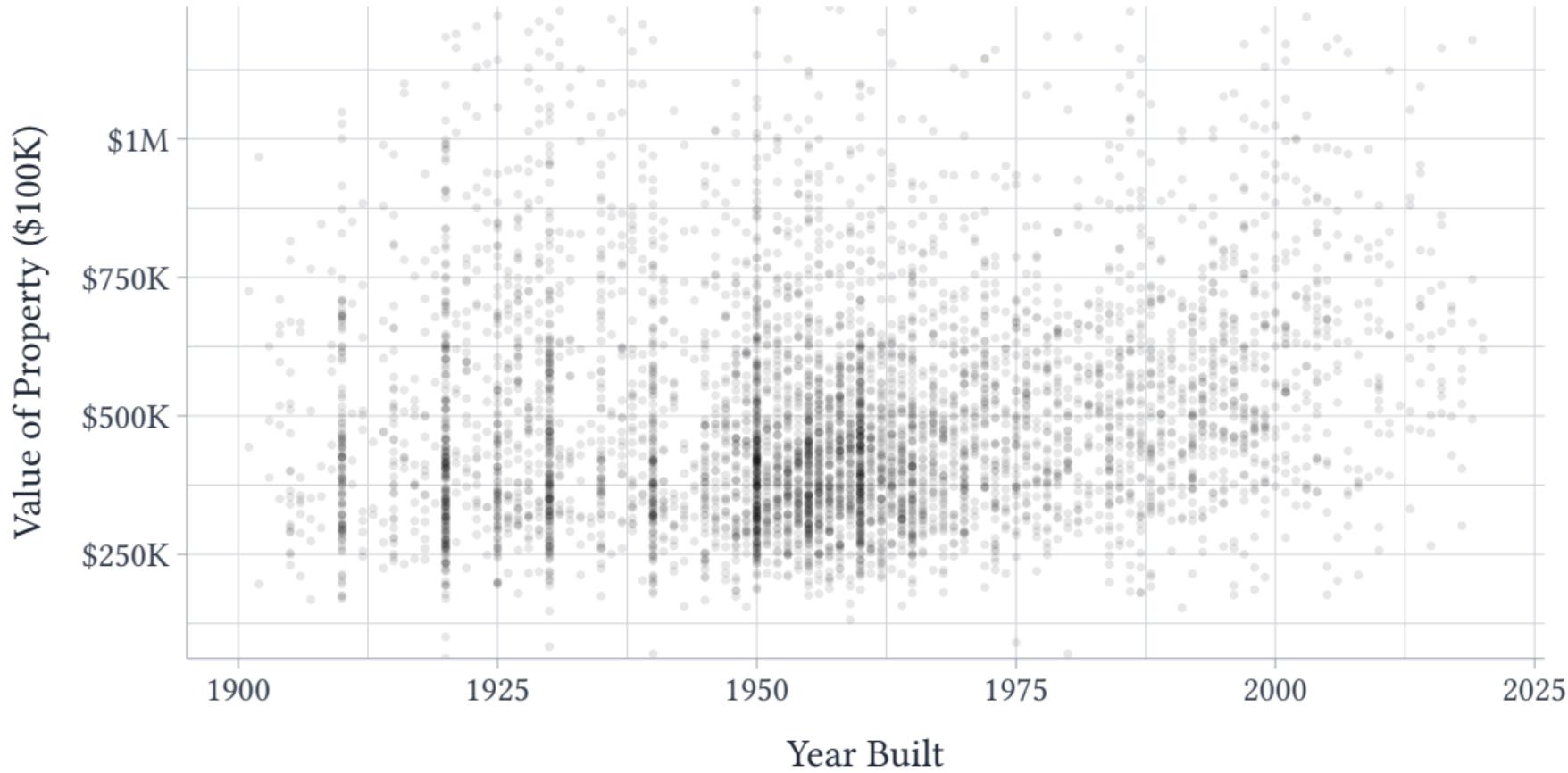
## Chopping $X_1$ into bins

The first approach, is to take  $X_{1i}$  and make it into a discrete set of bins:

- E.g. Person's age into 20 – 24, 25 – 29, 30 – 34, ...

Then, we can treat this as a multi-valued discrete variable

- Indicators for each bin and one omitted category
- We estimate sample means for each bin (relative to omitted group's mean)



Value of Property (\$100K)

\$1M

\$750K

\$500K

\$250K

1900

1925

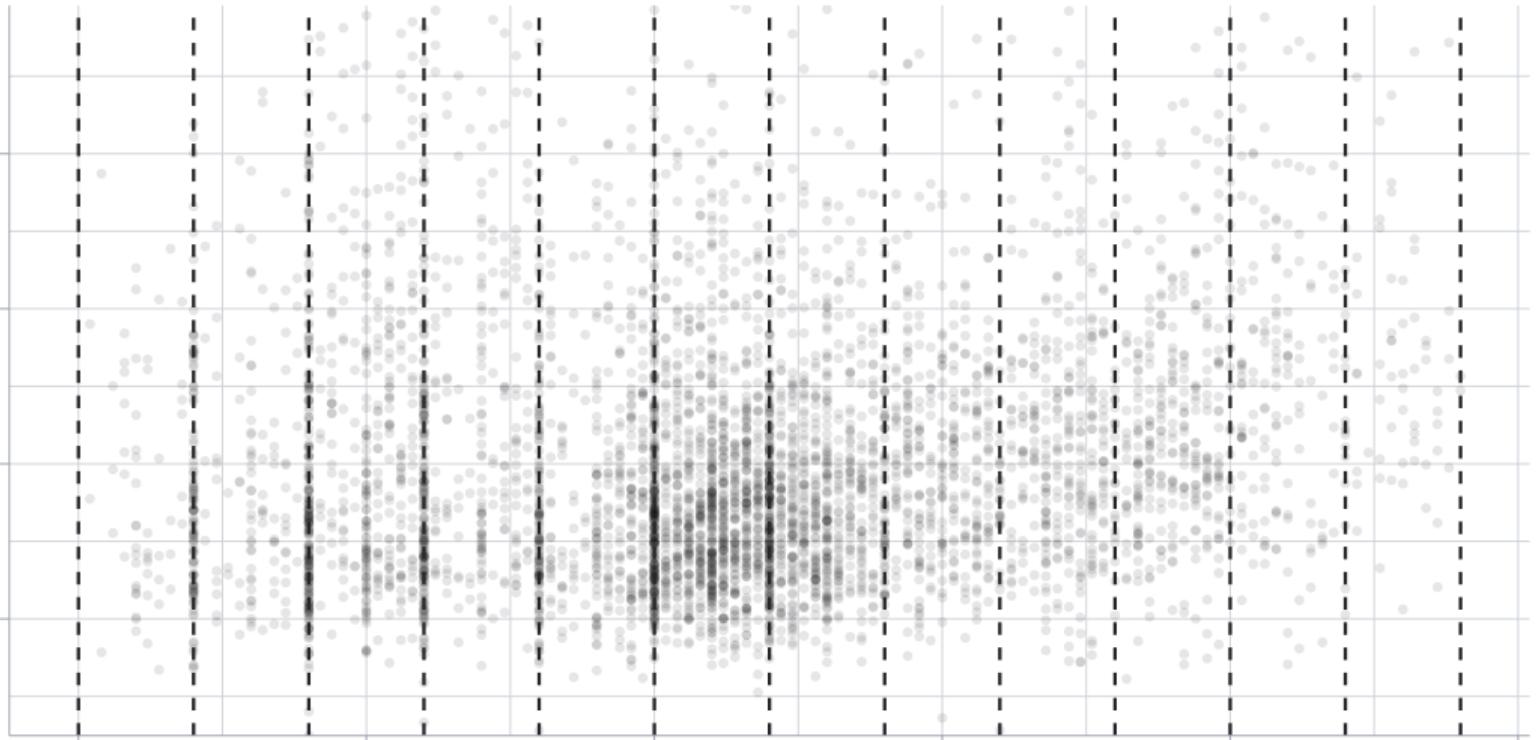
1950

1975

2000

2025

Year Built



Value of Property (\$100K)

\$1M

\$750K

\$500K

\$250K

1900

1925

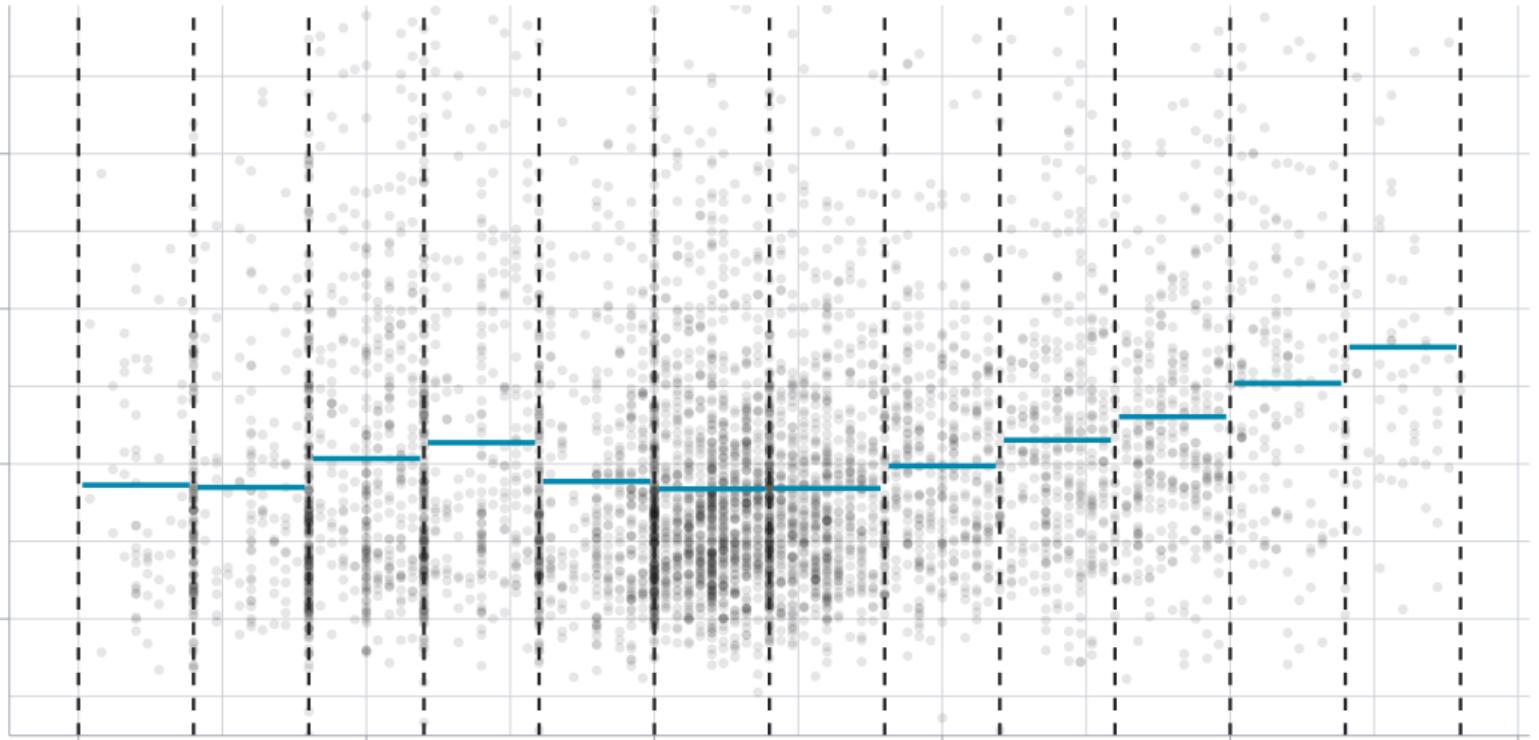
1950

1975

2000

2025

Year Built



Value of Property (\$100K)

\$1M

\$750K

\$500K

\$250K

1900

1925

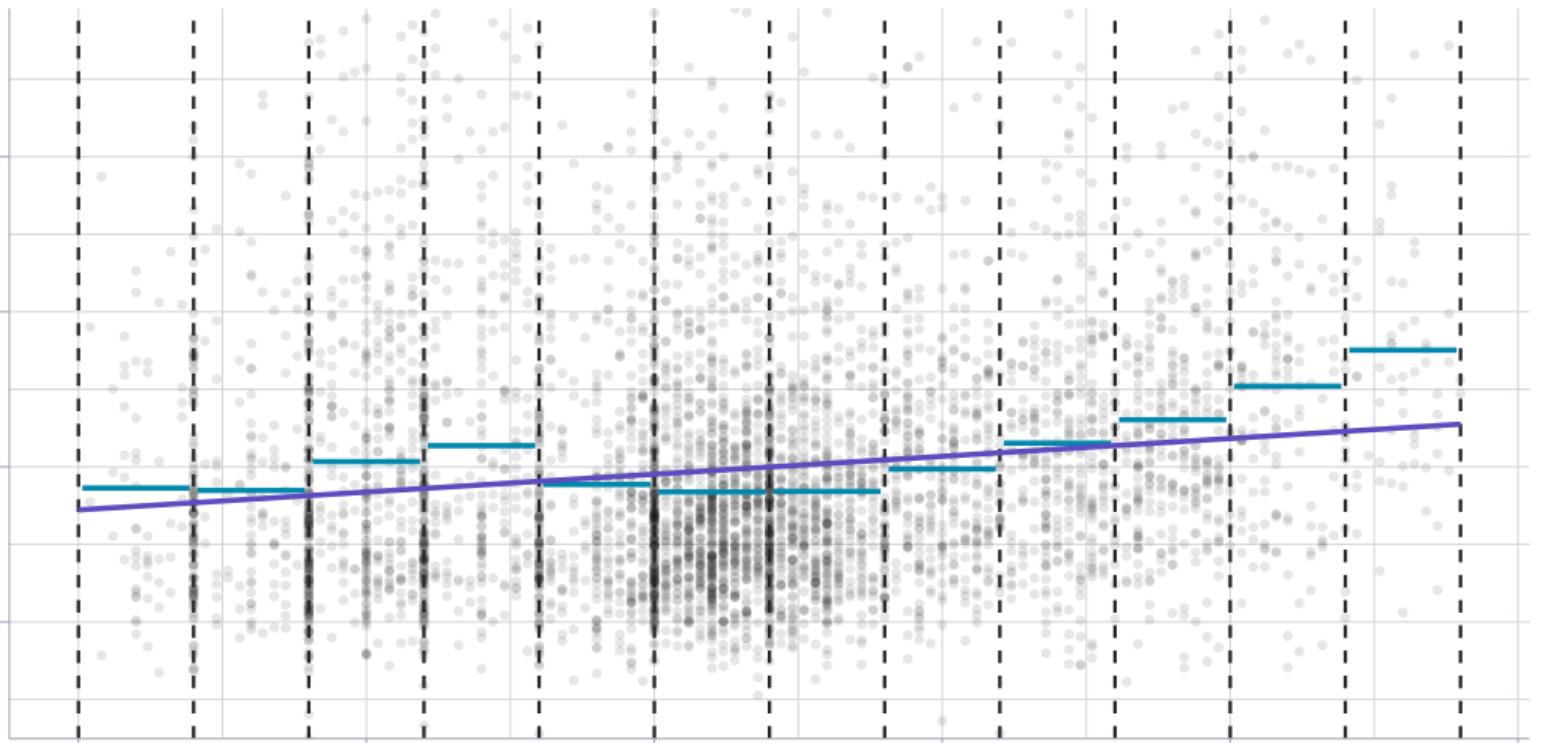
1950

1975

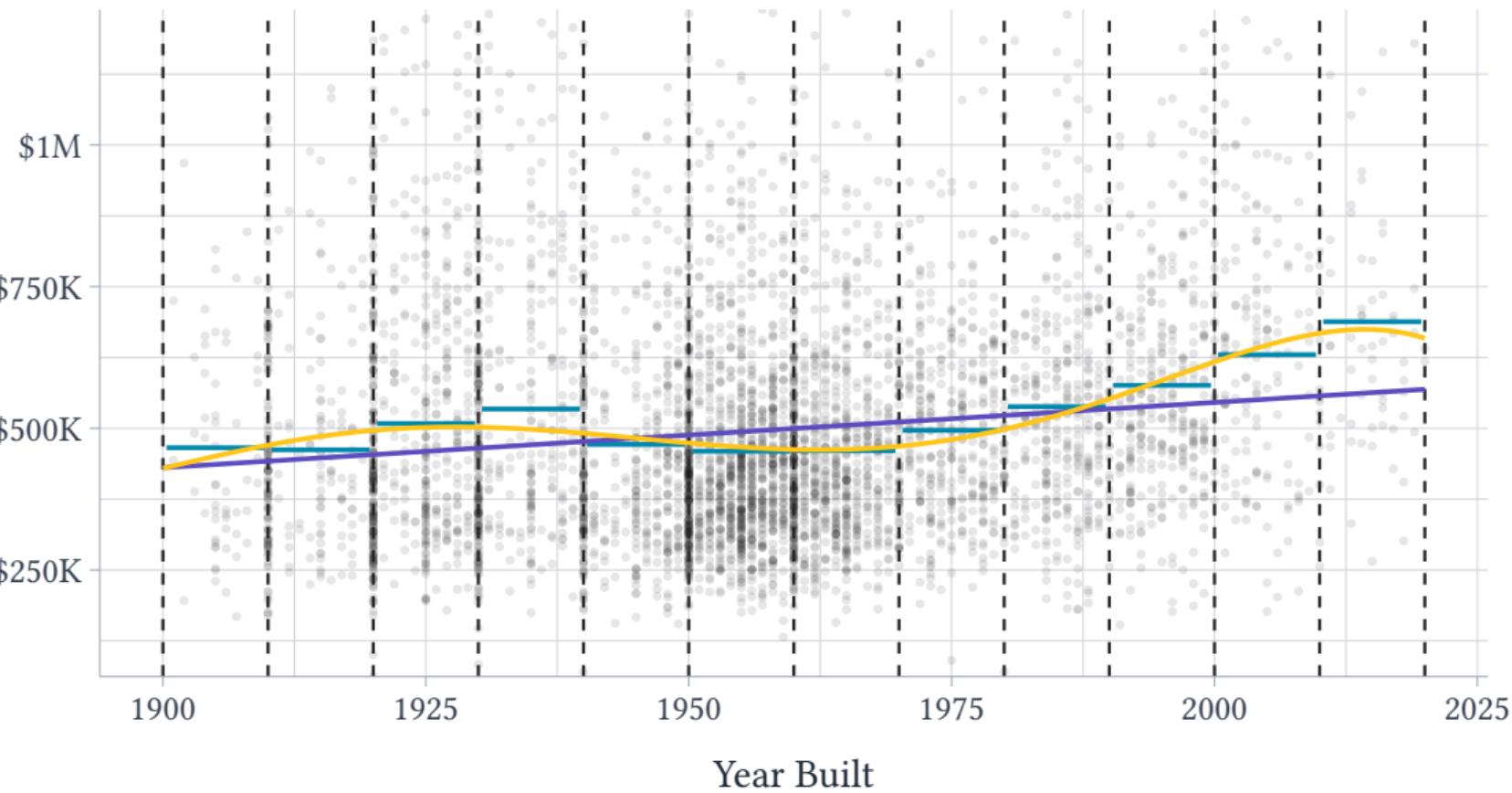
2000

2025

Year Built



Value of Property (\$100K)



## Pros and Cons

One advantage is this is relatively easy to explain and quite flexible:

- “The average property value for homes built from 1950 – 1959 is \$471,000”
- If we want more flexibility, we can increase the number of bins (if we have enough data)

## Pros and Cons

One advantage is this is relatively easy to explain and quite flexible:

- “The average property value for homes built from 1950 – 1959 is \$471,000”
- If we want more flexibility, we can increase the number of bins (if we have enough data)

One problem with this method is that the marginal effect estimates can be odd:

- The estimated function is flat (0 marginal effect) except at the “**knots**” where there is instant jump
- But, we are okay if we recognize this limitation

# Multivariate Regression – “All Else Equal”

## Regressors

Polynomials

Indicators

Multi-valued discrete variables

Interactions

Bins

Splines

log **transformations**

0/1 outcomes

## Previous approaches

$$y_i = \mu(X_{1i}) + W_i' \beta + u_i$$

We have discussed two ways to estimate  $\mu(X_{1i})$  flexibly, both had pros and cons

- Polynomials allow for flexibility, but could become quite sensitive / noisy
- Bins are flexible and simple to explain, but created artificial discontinuities and non-smoothness in our estimate

## Previous approaches

$$y_i = \mu(X_{1i}) + W'_i \beta + u_i$$

**Splines** are a way to try and blend the two approaches:

- Chop up the domain of  $X_1$  into a set of bins
- Within each bin, estimate a polynomial of a given order ( $p$ )
- Possibly, you can require the end points to “connect” ( $s$ )

## Polynomial order and smoothness

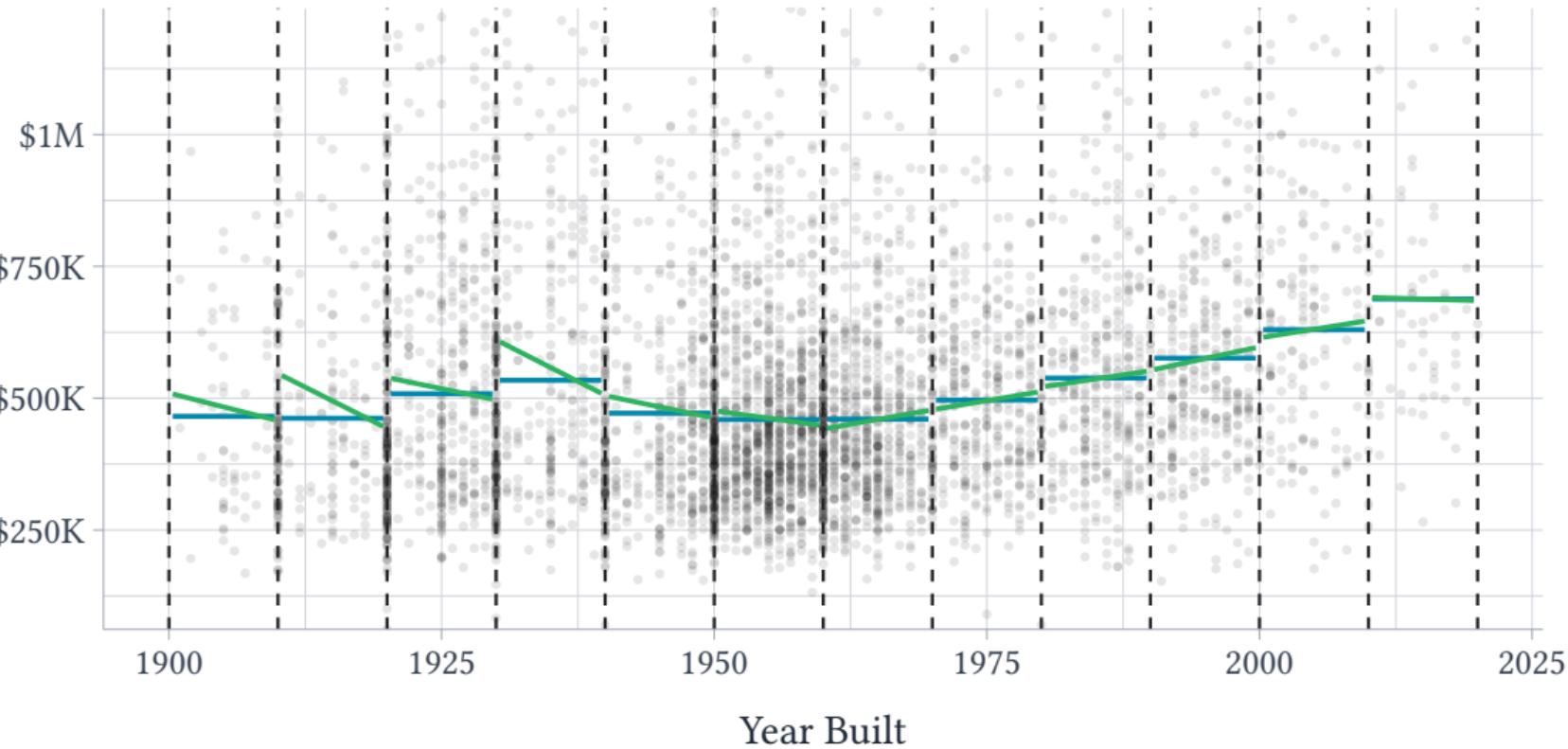
Our binned averages is a particular version of splines:

- Polynomial of order 0 ( $p = 0$ )
- Not required to connect ( $s = 0$ )

We could increase the polynomial order to  $p = 1$  to create linear functions within each bin

$p = 0, s = 0$ ;  $p = 1, s = 0$

Value of Property (\$100K)



## Polynomial order and smoothness

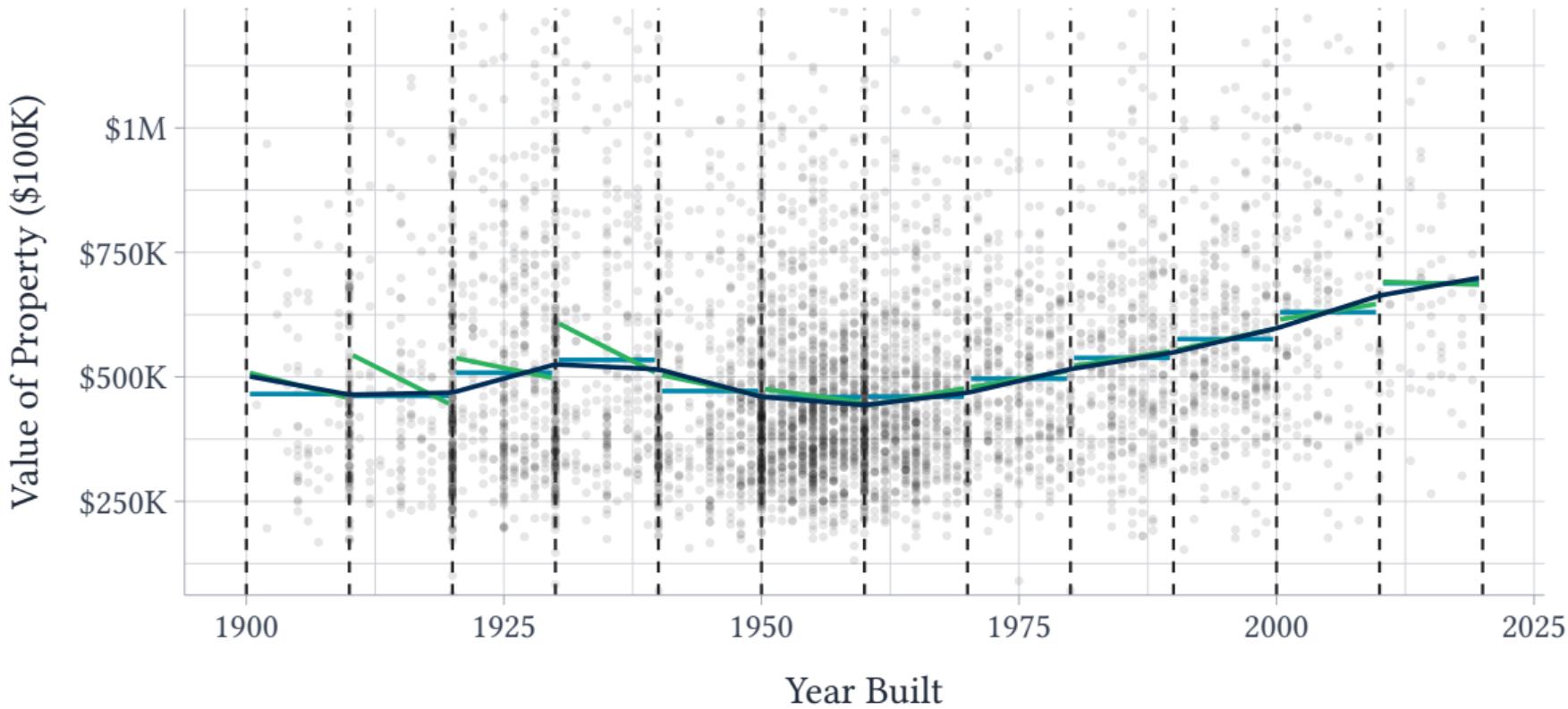
Our binned averages is a particular version of splines:

- Polynomial of order 0 ( $p = 0$ )
- Not required to connect ( $s = 0$ )

We could increase the polynomial order to  $p = 1$  to create linear functions within each bin

If we want the lines to connect, we can add a smoothness constraint  $s = 1$

$p = 0, s = 0$ ;  $p = 1, s = 0$ ;  $p = 1, s = 1$



## Polynomial order and smoothness

Our binned averages is a particular version of splines:

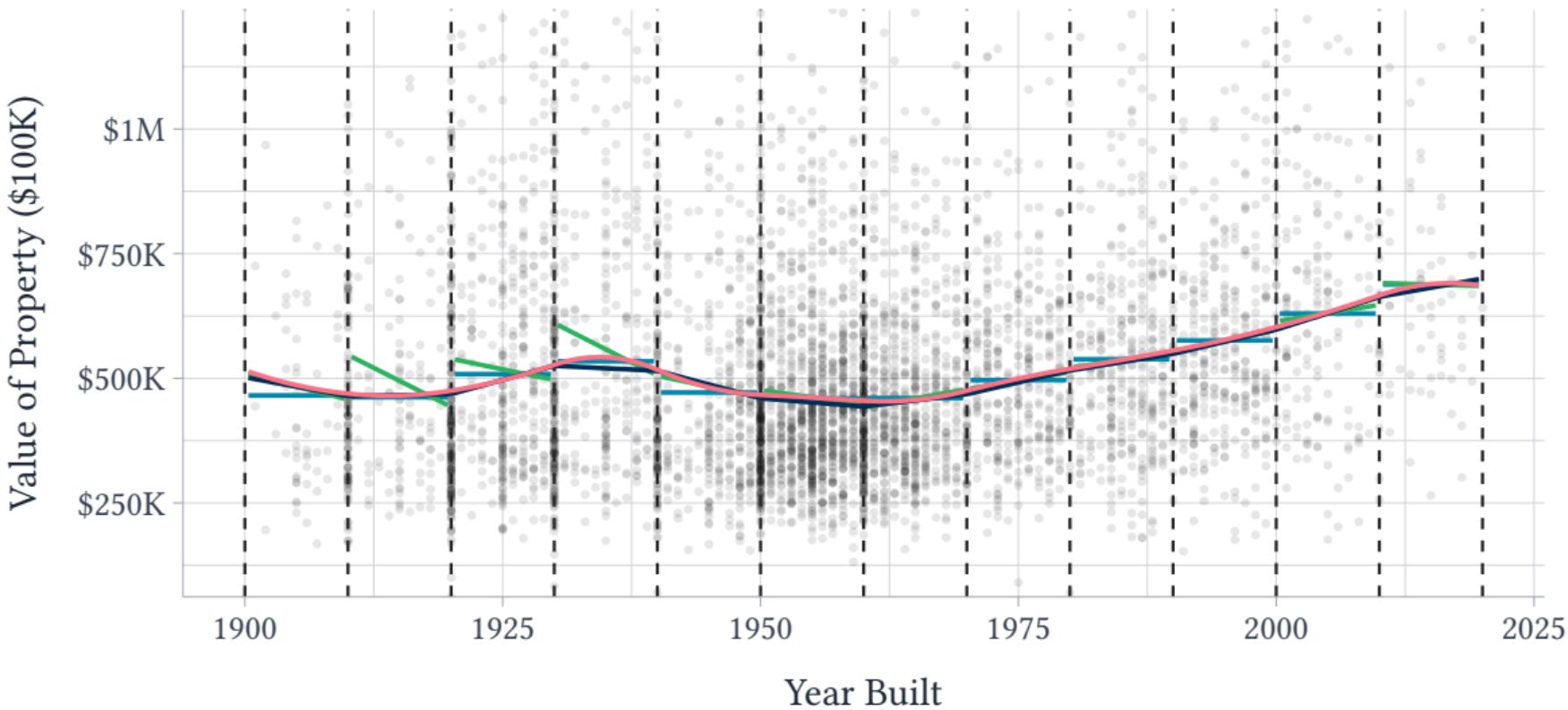
- Polynomial of order 0 ( $p = 0$ )
- Not required to connect ( $s = 0$ )

We could increase the polynomial order to  $p = 1$  to create linear functions within each bin

If we want the lines to connect, we can add a smoothness constraint  $s = 1$

Or, we can set  $p = 2$  and  $s = 2$  to estimate quadratic functions and make them connect smoothly

$p = 0, s = 0$ ;  $p = 1, s = 0$ ;  $p = 1, s = 1$ ;  $p = 2, s = 2$



## Spline advantages and costs

From the previous figure, it seems like we've hit a really nice *sweet spot*

- Quite flexible model to estimate  $\mu(X_{1i})$  and can easily add covariates  $W_i$  as controls
- But, we can not really do this with *all* our covariates because splines add a lot of parameters (at least as many as bins you have)

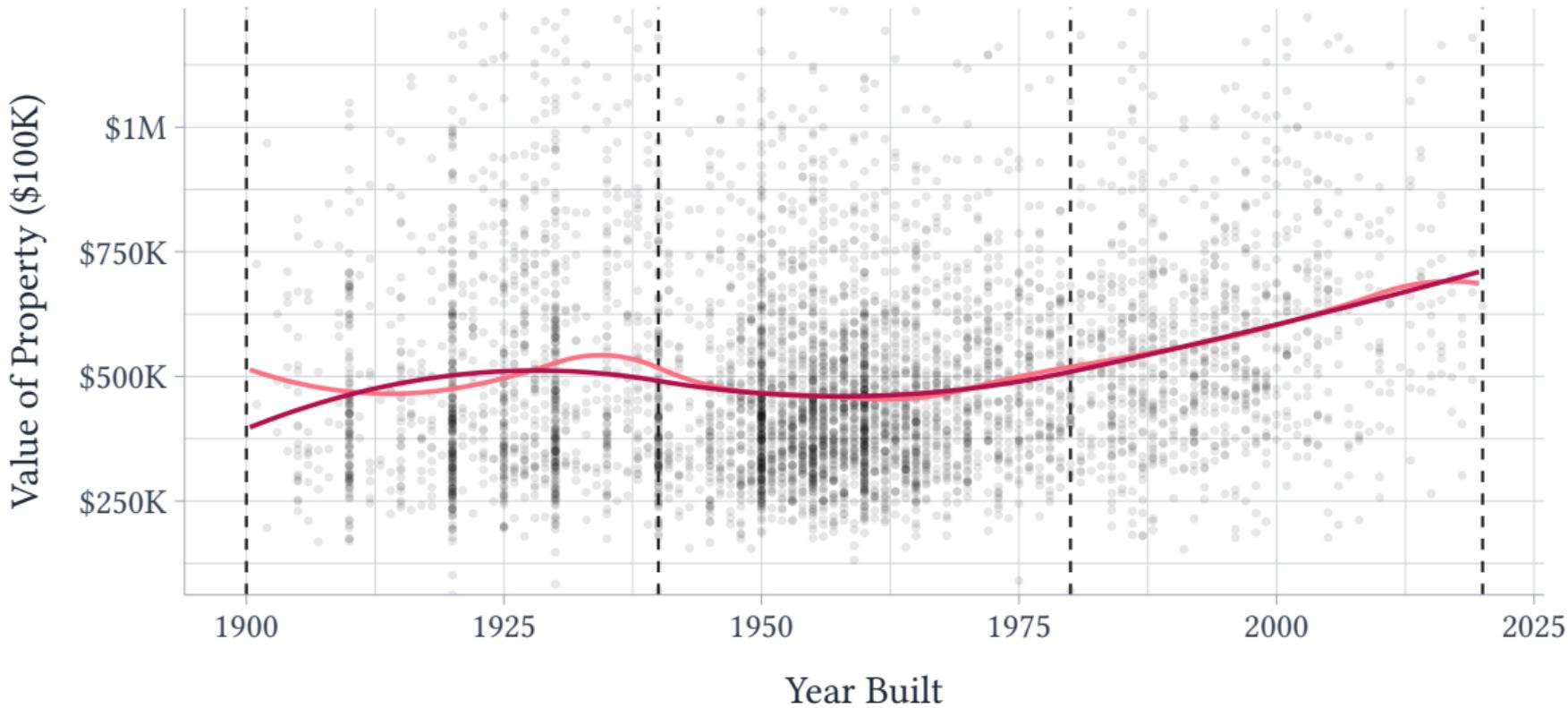
## Spline advantages and costs

From the previous figure, it seems like we've hit a really nice *sweet spot*

- Quite flexible model to estimate  $\mu(X_{1i})$  and can easily add covariates  $W_i$  as controls
- But, we can not really do this with *all* our covariates because splines add a lot of parameters (at least as many as bins you have)

It's not clear yet how we should chose the bins...

knots every 10 years; knots every 40 years



# Choosing Bins

There is a *bias/variance* trade-off when selecting the number of bins

- More bins means we can better approximate  $\mu(X_{1i})$  (decreased bias)
- But, with fewer observations per bin, our estimates will be more noisy (increased variance)

More, it is not clear why the bins should be evenly-space

- It makes ‘intuitive’ sense, but maybe the bins should get smaller when the data is more ‘wiggly’ and larger where the data is less

# Choosing Bins

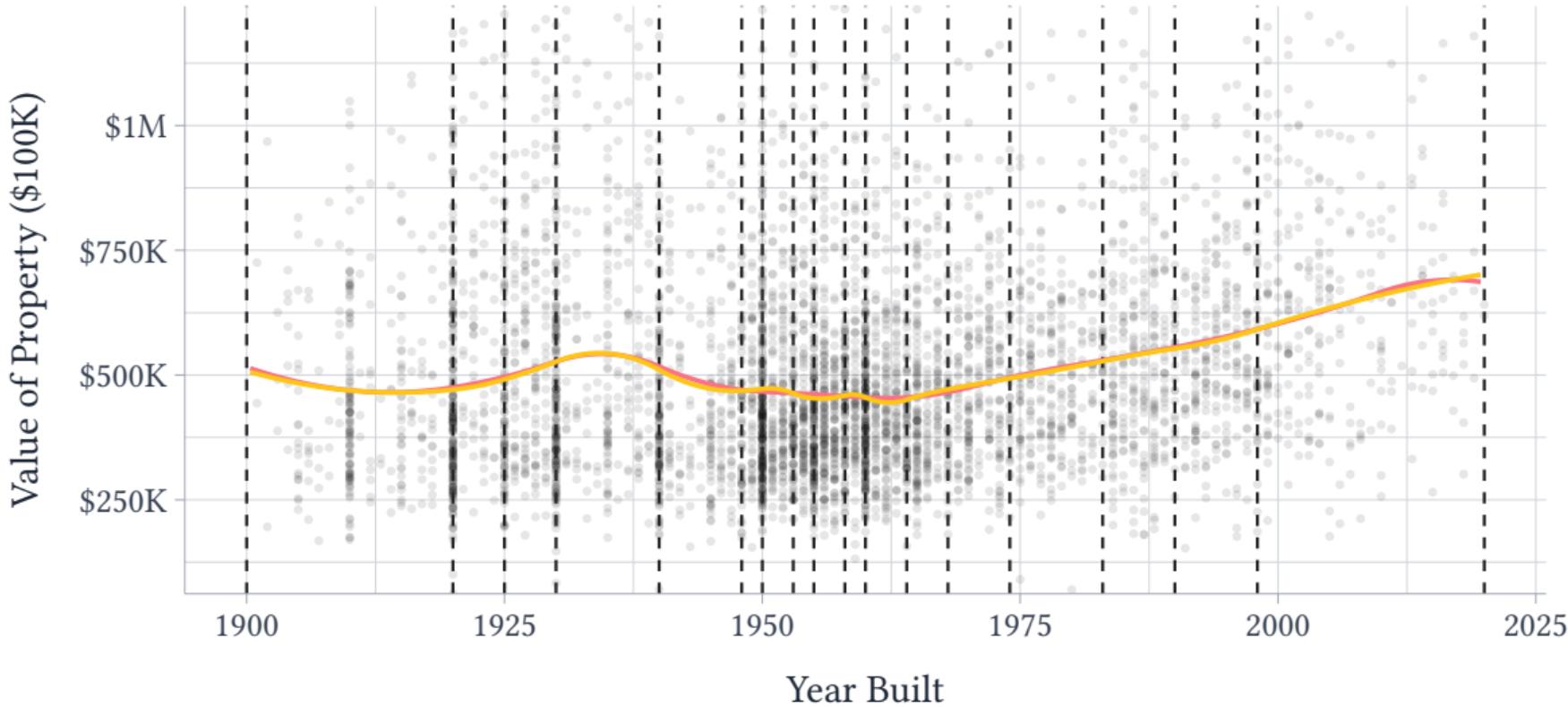
$$y_i = \mu(X_i) + W_i' \beta + u_i$$

One recent way of estimating this is using a ‘binscatter’ regression:

1. Chop  $X$  variable into  $J$  bins with an *equal number of observations into each bin*
2. Choose  $J$  *optimally* to minimize mean-squared prediction error (bias-variance trade-off)

Implemented in the `binsreg` package

knots every 10 years; binscatter selected



# Multivariate Regression – “All Else Equal”

## Regressors

Polynomials

Indicators

Multi-valued discrete variables

Interactions

Bins

Splines

log **transformations**

0/1 outcomes

## log-transformation

In economics, it is common to see log transformed outcomes:

$$\log(\text{wages}_i) = \beta_0 + \beta_1 \text{College Degree}_i + u_i$$

## log-transformation

In economics, it is common to see log transformed outcomes:

$$\log(\text{wages}_i) = \beta_0 + \beta_1 \text{College Degree}_i + u_i$$

This specification changes our interpretation of the slope coefficients:

“Having a college degree is associated with an increase in wages of  $\beta_1 * 100$  percent”

→ E.g. if  $\beta_1 = 0.02$ , then a college degree is associated with a 2% increase in wages.

## Derivation of log-transformation interpretation

Compare two individuals: unit 1 with and unit 0 without a college degree. Then, we have

$$\log(\text{wages}_1) - \log(\text{wages}_0) = \beta_0 + \beta_1 - \beta_0$$

$$\implies \log(\text{wages}_1 / \text{wages}_0) = \beta_1$$

$$\implies \log\left(1 + \frac{\text{wages}_1 - \text{wages}_0}{\text{wages}_0}\right) = \beta_1$$

## Derivation of log-transformation interpretation

Compare two individuals: unit 1 with and unit 0 without a college degree. Then, we have

$$\log(\text{wages}_1) - \log(\text{wages}_0) = \beta_0 + \beta_1 - \beta_0$$

$$\implies \log(\text{wages}_1 / \text{wages}_0) = \beta_1$$

$$\implies \log\left(1 + \frac{\text{wages}_1 - \text{wages}_0}{\text{wages}_0}\right) = \beta_1$$

If you recall, exponentiating gets rid of the the log

$$\frac{\text{wages}_1 - \text{wages}_0}{\text{wages}_0} = \exp(\beta_1) - 1$$

## Derivation of log-transformation interpretation

$$\frac{\text{wages}_1 - \text{wages}_0}{\text{wages}_0} = \exp(\beta_1) - 1$$

The left-hand side is our percent-change formula from high-school science class

## Derivation of log-transformation interpretation

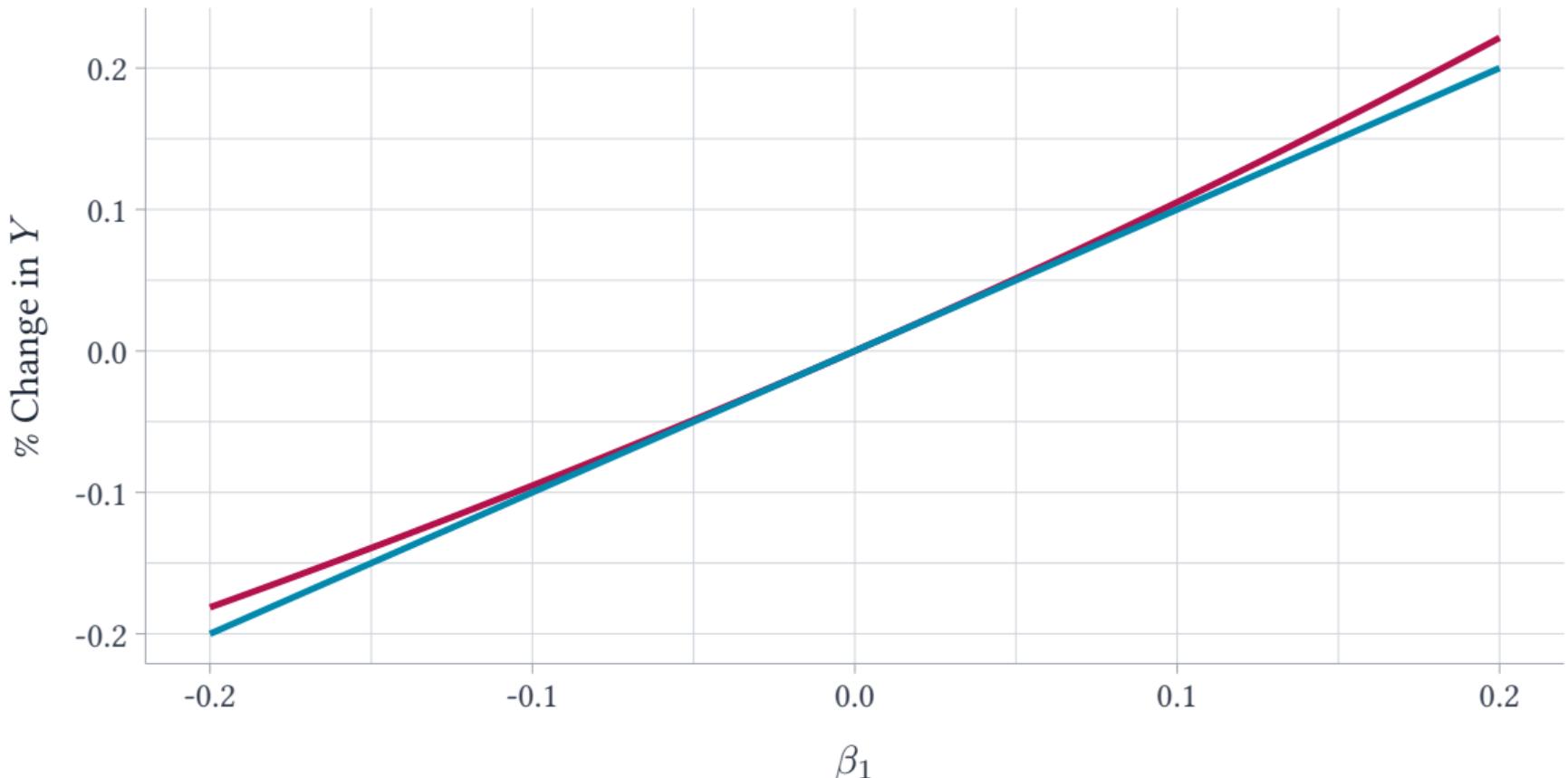
$$\frac{\text{wages}_1 - \text{wages}_0}{\text{wages}_0} = \exp(\beta_1) - 1$$

The left-hand side is our percent-change formula from high-school science class

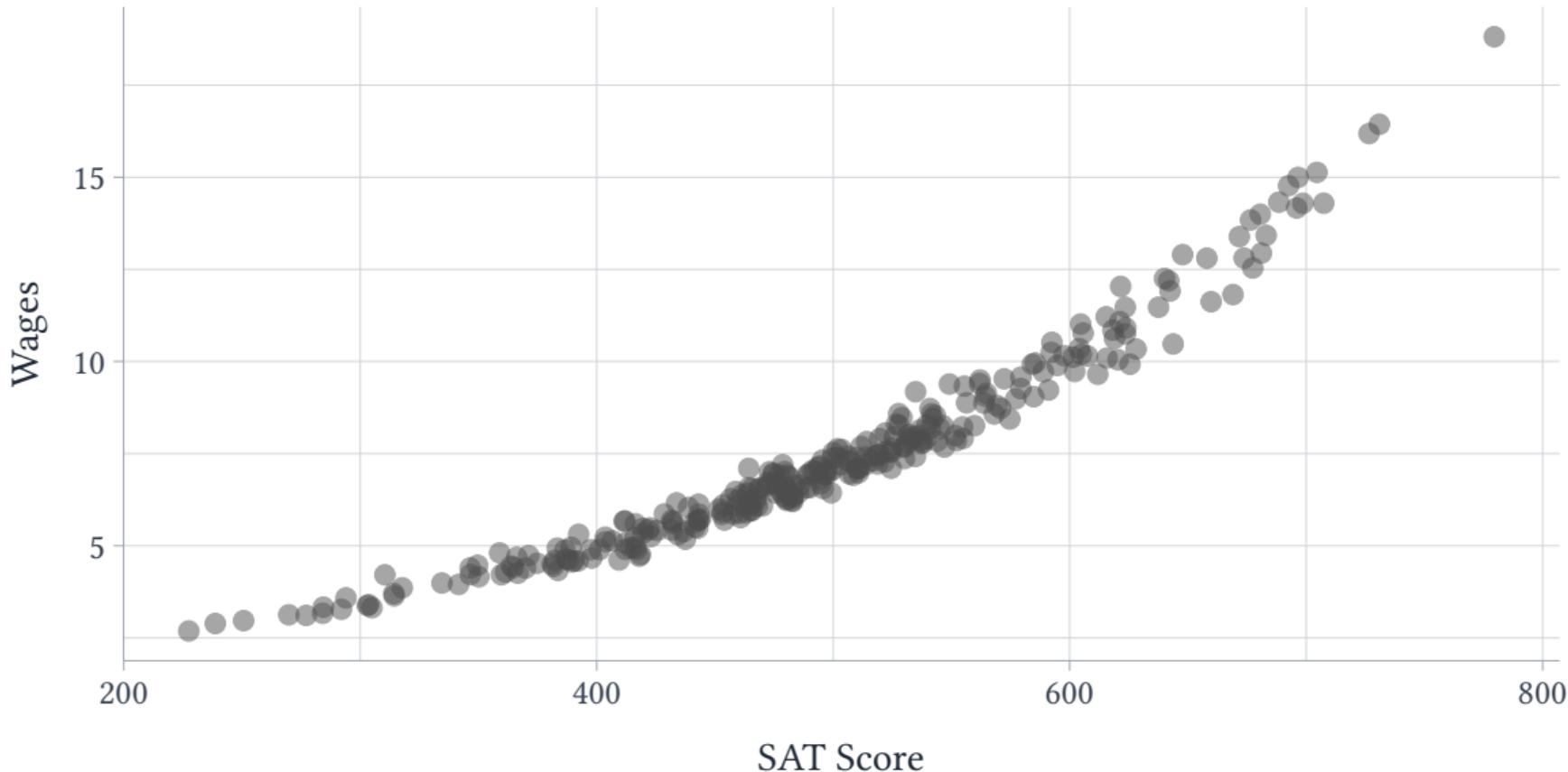
In this case, the more *precise* answer is that having a college degree is associated with an  $\exp(\beta_1) - 1$  percent change in wages

→ But for  $-0.10 < \beta_1 < 0.10$ ,  $\exp(\beta_1) - 1$  is approximately equal to  $\beta_1$  so it's simpler to use the latter

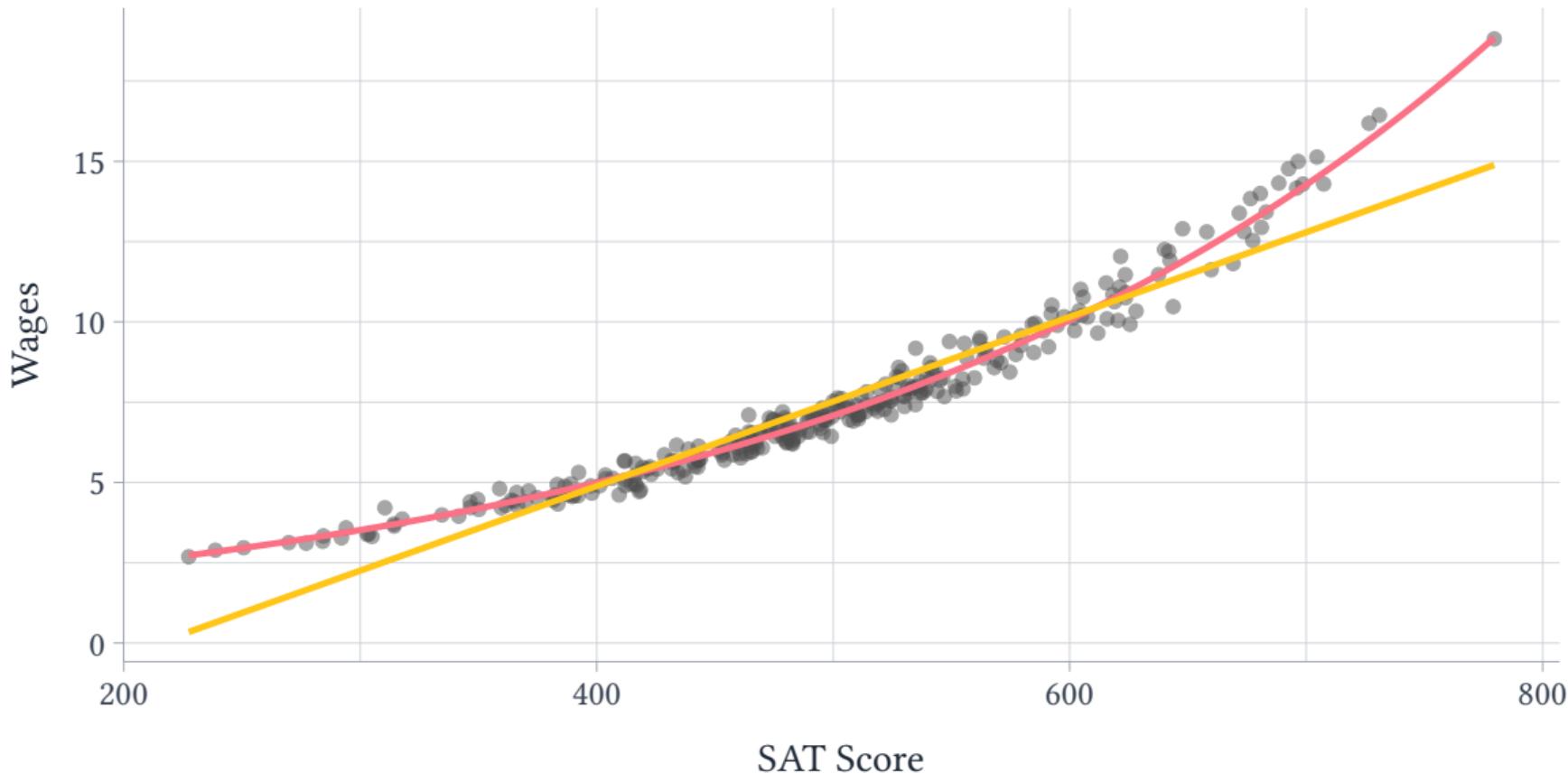
## Comparison of $\exp(\beta_1) - 1$ and $\beta_1$



# Data on SAT score and wages



## True log CEF vs. Linear Approximation



## When to use log transformations

You should take the log of an outcome variable when you think a 1 unit change in  $X$  is related to a % change in  $Y$ .

## When to use log transformations

You should take the log of an outcome variable when you think a 1 unit change in  $X$  is related to a % change in  $Y$ .

$\log(Y) \sim X$  is called fitting an ‘exponential’ relationship. These are common in:

1. Financial markets where compounding returns imply  $Y_t = Y_0 e^{rt}$
2. Epidemiology where disease growth rate is exponential (it is not actually, but early growth rate is approximately)

## When to use log transformations

You should take the log of an outcome variable when you think a 1 unit change in  $X$  is related to a % change in  $Y$ .

$\log(Y) \sim X$  is called fitting an ‘exponential’ relationship. These are common in:

1. Financial markets where compounding returns imply  $Y_t = Y_0 e^{rt}$
2. Epidemiology where disease growth rate is exponential (it is not actually, but early growth rate is approximately)
3. Settings with skewed distributions (e.g. home prices, GDP, population)
  - Skewness makes a ‘unit’ change in  $X$  difficult to think about

## log-log transformations

Alternatively, you may see log transformations of both variables:

$$\log(Y_i) = \beta_0 + \beta_1 \log(X_i) + u_i$$

The interpretation is now simpler: a 1% change in  $X_1$  is associated with a  $\beta_1$  % change in  $Y$

# Multivariate Regression – “All Else Equal”

## Regressors

Polynomials

Indicators

Multi-valued discrete variables

Interactions

Bins

Splines

log **transformations**

**0/1 outcomes**