

# Regularization and the Bias-Variance Tradeoff

Park Effects: Estimate the park effect  $\alpha$  of each MLB ballpark, which represents the expected runs scored in one half-inning at that park above that of an average park, if an average offense faces an average defense

data all half-innings 2017-2019

variables  $i$  is the index of the  $i^{\text{th}}$  half-inning  
 $y_i$  Runs scored in half-inning  $i$

PARK( $i$ )

weather/day of the year

offensive team strength  
defensive team strength

ot( $i$ ) = offensive team at half-inning  $i$

dt( $i$ ) = defensive team

model

$$y_i = \beta_0 + \alpha_{\text{PARK}(i)} + \varepsilon_i, \quad \mathbb{E}\varepsilon_i = 0$$

model

$$y_i = \beta_0 + \alpha_{\text{PARK}(i)} + \beta_{\text{ot}(i)} + \gamma_{\text{dt}(i)} + \varepsilon_i$$

$$\beta_{ot(i)}$$

30 off. teams Phi, LAD, ...

$ot(i)$  = the off. team in half.ing  $i$

$$\beta_{LAD}, \beta_{Phi}, \dots$$

Equivalently, 
$$y_i = x_i^T \beta + \varepsilon_i$$

$X$  is a matrix whose  $i^{th}$  row is defined by

Row  $i$

$$x_i^T = \left[ 1, \underbrace{0's \text{ everywhere except } 1 \text{ at } park(i)}, \underbrace{0's \text{ everywhere except } 1 \text{ at } ot(i)}, \underbrace{0's \text{ everywhere except } 1 \text{ at } dt(i)} \right]$$

$$x_i^T = \left[ 1, \overset{\text{interest}}{\bullet}, \dots, \overset{\text{park } 1}{\bullet}, \dots, \overset{\text{park } 30}{\bullet}, \dots, \overset{\text{off } 1}{\bullet}, \dots, \overset{\text{off } 30}{\bullet}, \dots, \overset{\text{def } 1}{\bullet}, \dots, \overset{\text{def } 30}{\bullet} \right]$$

$$\beta = \begin{pmatrix} \beta_0 \\ \alpha_1 \\ \vdots \\ \alpha_{30} \\ \beta_1 \\ \vdots \\ \beta_{30} \\ \gamma_1 \\ \vdots \\ \gamma_{30} \end{pmatrix}$$

# Problem: Multicollinearity

When home team is on offense,  $\text{Park}(i) = \text{ot}(i)$

When road team is on offense,  $\text{Park}(i) = \text{dt}(i)$

So it is tough to disentangle

$\alpha_{\text{Park}(i)}$  from  $\beta_{\text{ot}(i)}$ ,  $\delta_{\text{dt}(i)}$

Are the runs scored in this half-inning due to the offensive home team being good or the park being easy?

$$y_i = \beta_0 + \underbrace{\alpha_{\text{Park}(i)}}_{\text{Park effect}} + \beta_{\text{ot}(i)} + \delta_{\text{dt}(i)} + \varepsilon_i$$

$y_i$ : runs scored in half-inning  $i$

To disentangle these effects, we need a huge number of instances of Road teams on offense to estimate  $\beta_{\text{ot}(i)}$  well, we need a huge number of instances of Home teams on offense

to estimate  $\gamma_{dt(i)}$  well, with a good  $\beta_{ot(i)}$  and  $\gamma_{ot(i)}$  we can estimate  $\alpha_{part(i)}$  well.

Our dataset of all half-innings from 2017-2019 consists of  $\approx 121,000$  half-innings (rows).

This may seem like a huge amount of data. But, with this multicollinearity, is it really?

## Simulation Study

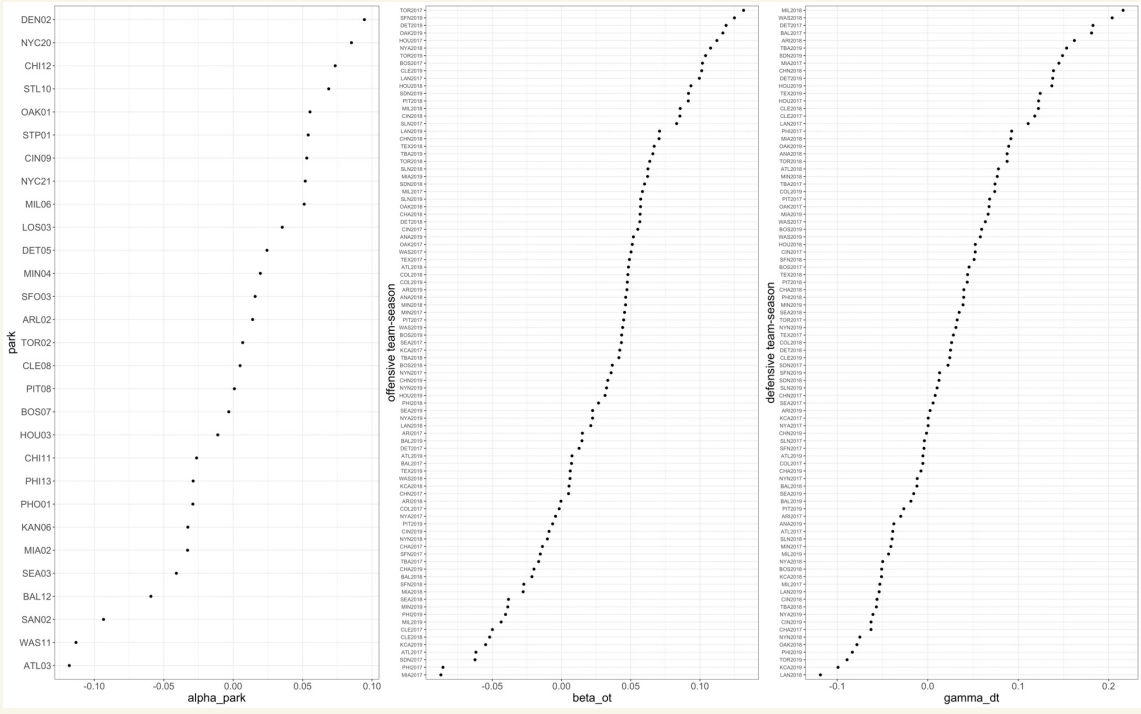
Idea Pretend we know  $\alpha, \gamma, \beta$

Generate fake historical data  $y$   
(121,000 fake innings)

Then estimate  $\alpha$  (and  $\beta, \gamma$ ) from this fake data.

Then compare  $\hat{\alpha}$  to "true"  $\alpha$  which we know.

\* Suppose the "true" coefficients are



Then, given these, lets generate  $y$  according to

$$y_i = \underbrace{\beta_0 + \alpha_{park(i)} + \beta_{ot(i)} + \gamma_{dt(i)}}_{x_i^T \beta} + \epsilon_i$$

$$y_i = \text{Round}(\mathcal{N}_+(x_i^T \beta, 1))$$

$$\mathcal{N}_+(\mathbf{x}_i^T \boldsymbol{\beta}, 1) = \mathbf{x}_i^T \boldsymbol{\beta} + \varepsilon_i$$

there exists some  $\varepsilon_i$  which satisfies the equation

↓

$$\mathbb{E} \varepsilon_i = 0$$

$$\text{VAR}(\varepsilon_i) = 1$$

$$y = X\beta + \varepsilon$$

if we're only interested in point estimate  $\hat{\beta}$ ,  
we only need to know  $\mathbb{E}\varepsilon = 0$   
and anything about the actual distribution of  $\varepsilon$   
is irrelevant.

$\hat{\beta}$  is unbiased

$$\mathbb{E}\hat{\beta}_{OLS} = \mathbb{E}[(X^T X)^{-1} X^T y] = (X^T X)^{-1} X^T \mathbb{E}[y]$$

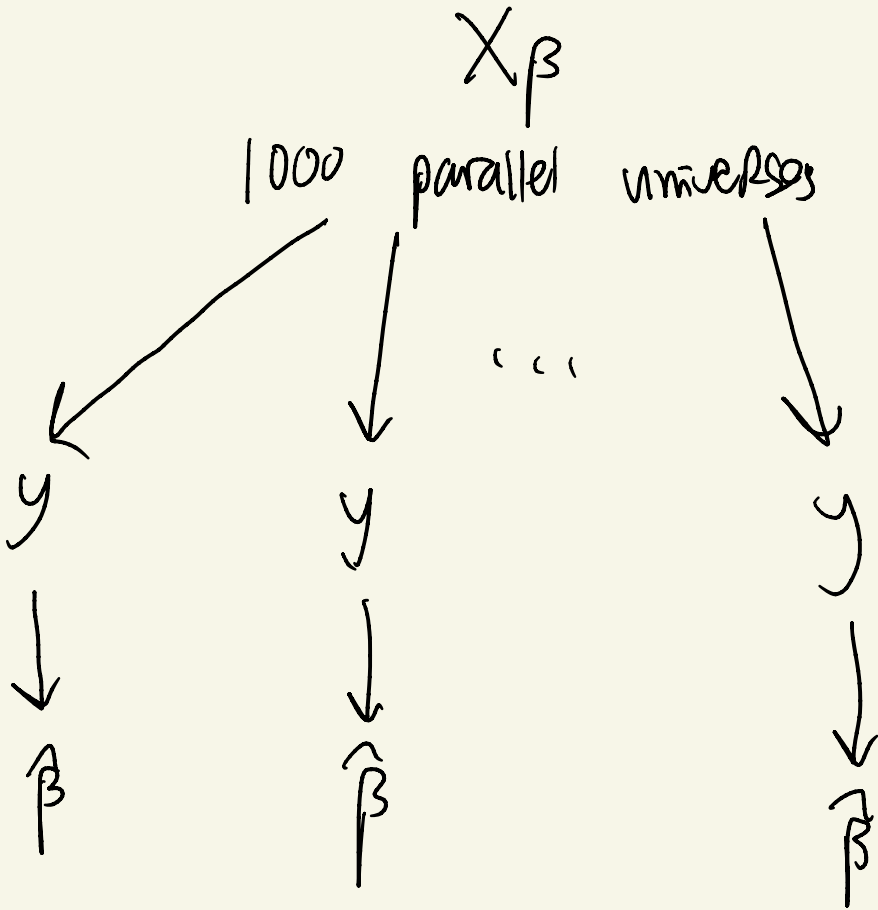
$$= (X^T X)^{-1} X^T \mathbb{E}[X\beta + \varepsilon]$$

$$= \boxed{(X^T X)^{-1} X^T X} \beta + (X^T X)^{-1} X^T \cancel{\mathbb{E}[\varepsilon]}$$

$$= \beta$$

$$\varepsilon \sim N(0, \sigma^2) \rightarrow \text{CI on } \beta$$
$$\hat{\beta} \pm 2\sigma$$

$$y = X\beta + \epsilon$$





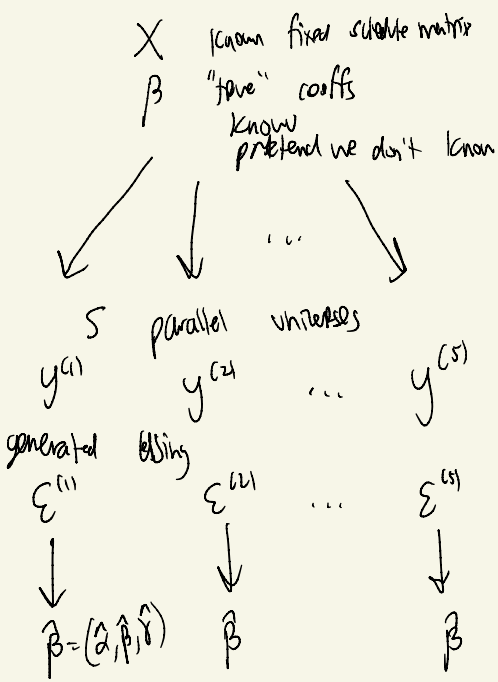
Where did  $x_i^T$  come from?  $\rightarrow$  used the actual  $X$  matrix (the actual schedule)

Where is  $\varepsilon_i$ ?  $\rightarrow$  implicit in the  $\mathcal{N}$

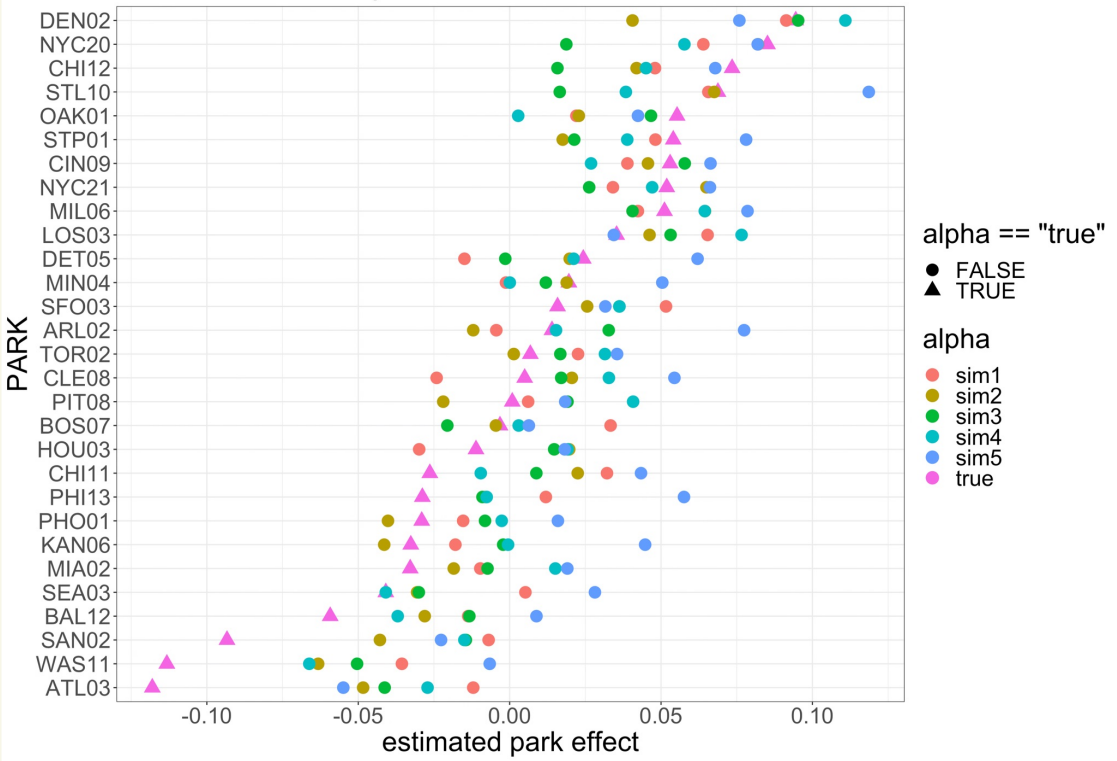
\* Now we have "true" known coefficients  $(\beta, \sigma, \alpha)$ , we have the observed schedule matrix  $X$ , the fake generated runs stored vector  $y$ .

\* Using our fake historical dataset  $(X, y)$ , let's use linear regression to estimate the coefficients,  $(\hat{\alpha}, \hat{\beta}, \hat{\sigma})$ .

$(X^T X)^{-1} X^T y$



estimated park effects across M=5 sims



\* Due to randomness in the training dataset, from the noise  $\epsilon$  in generating  $y$ ,

Each simulation yields very different park effect estimates  $\hat{\alpha}$  even though the "true" park effects are the same.

\* These OLS (Ordinary Least Squares) coefficients  $\hat{\alpha}$  are quite sensitive to the noise of the training set, even though we have 121,000 innings of data.

\* How can we make the estimated coefficients less sensitive to the Random idiosyncracies of our training dataset?

from data  $(X, y)$

Q What's the least sensitive, dumb, estimator you can think of?

zero

Idea Blend the strengths of OLS and zero:

- Constant values like zero are not sensitive to the random idiosyncrasies of training data, but are wrong (on average) for many parks
- OLS estimators are very sensitive to the randomness of the training set, especially in the presence of multicollinearity, but are unbiased (on average, they are right)

— tradeoff b/w sensitivity and unbiasedness

Idea Shrink the OLS estimates towards zero.

Make them smaller! And hence less sensitive.

$$\text{OLS: } \hat{\beta}^{(\text{OLS})} = \underset{\beta}{\text{argmin}} \text{RSS}(\beta) = \underset{\beta}{\text{argmin}} \sum_{i=1}^n (x_i^T \beta - y_i)^2$$

Now:

want  $x_i^T \beta$  close to  $y$       make  $\beta$  small

$$\underset{\beta}{\text{argmin}} \sum_{i=1}^n (y_i - x_i^T \beta)^2 + \underbrace{\lambda \sum_j \beta_j^2}_{\text{penalty term}} = L(\beta)$$

shrinkage

A larger  $\beta_j$  has larger  $\beta_j^2$   
so if we are minimizing  $L(\beta)$   
we want  $\beta_j$  to be smaller

\* This technique of adding a penalty term to our loss function is called Regularization.

This is  $L_2$  Regularization  $(\beta_j^2)$

$$|\beta_j| \rightarrow L_1$$

\* The hyperparameter  $\lambda > 0$  describes by how much we are penalized for having larger  $\beta_j$ .

$\lambda$  is simply a number which is tuned to optimize predictive performance.

\* Ridge Regression = OLS +  $L_2$  Regularization

$$\hat{\beta}^{(\text{Ridge})} = \underset{\beta}{\text{Argmin}} \sum_{i=1}^n (y_i - x_i^T \beta)^2 + \lambda \sum_{j=1}^p \beta_j^2$$

$$= \underset{\beta}{\operatorname{argmin}} \quad (y - X\beta)^T (y - X\beta) + \lambda \beta^T \beta$$

$$\beta^T \beta = (\beta_1 \dots \beta_p) \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_p \end{pmatrix} = \beta_1^2 + \dots + \beta_p^2$$

Calculus: gradient = 0 and solve

$$\nabla_{\beta} L(\beta) = \nabla_{\beta} \left[ (y - X\beta)^T (y - X\beta) + \lambda \beta^T \beta \right]$$

$$= \nabla_{\beta} \left[ y^T y - 2\beta^T X^T y + \beta^T X^T X \beta \right] + \lambda \nabla_{\beta} [\beta^T \beta]$$

$$= \left[ -2X^T y + 2X^T X \beta \right] + \lambda \left[ 2\beta \right]$$

$$\frac{\partial}{\partial \beta_j} [\beta^T \beta] = \frac{\partial}{\partial \beta_j} (\beta_1^2 + \dots + \beta_p^2) = 2\beta_j$$

$$\nabla_{\beta} [\beta^T \beta] = \left( \frac{\partial}{\partial \beta_1} [\beta^T \beta], \dots, \frac{\partial}{\partial \beta_p} [\beta^T \beta] \right)$$

$$= 0$$

$$\Rightarrow \left( \underbrace{2X^T X}_{p \times p} + 2\lambda \underbrace{I}_{1 \times 1} \right) \beta = 2X^T y$$

$$2\lambda \beta \rightarrow 2\lambda I \beta$$

$$I = \begin{pmatrix} 1 & & 0 \\ & \ddots & \\ 0 & & 1 \end{pmatrix}$$

$$\Rightarrow (X^T X + \lambda I) \cdot \beta = X^T y$$

$$\underbrace{(X^T X + \lambda I)^{-1} (X^T X + \lambda I)}_I \beta = (X^T X + \lambda I)^{-1} X^T y$$

$$\Rightarrow \hat{\beta}^{(Ridge)} = (X^T X + \lambda I)^{-1} X^T y$$

$$\hat{\beta}^{(OLS)} = (X^T X)^{-1} X^T y$$

$$\lambda I = \begin{pmatrix} \lambda & & 0 \\ & \ddots & \\ 0 & & \lambda \end{pmatrix} \quad \text{Ridge}$$



pretend  $x^T x$  is a number

$(X^T X)^{-1}$  is like  $\frac{1}{x^T x}$

$(X^T X + \lambda I)^{-1}$  is like  $\frac{1}{x^T x + \lambda}$

if  $x^T x \approx 0$ ,  $\frac{1}{x^T x}$  is large and unstable

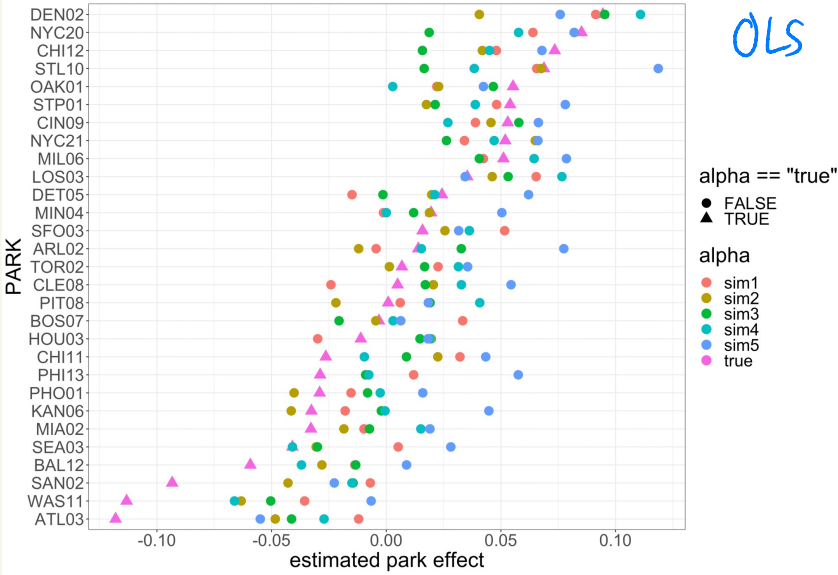
$\frac{1}{x^T x + \lambda}$  is more stable

In the presence of multicollinearity,

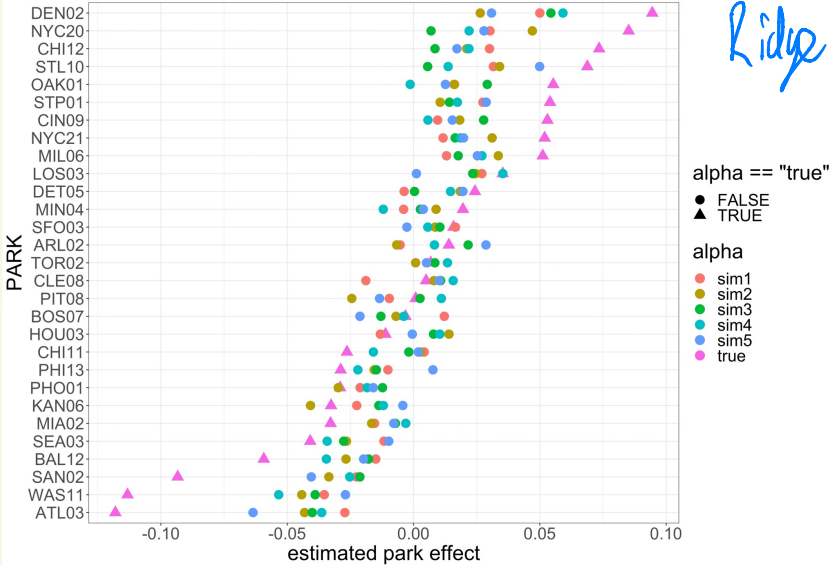
$(X^T X)^{-1}$  is like dividing by zero in some way.

Let's look at  $\hat{\beta}^{(Ridge)}$  and  $\hat{\beta}^{(OLS)}$   
 $\lambda = 0.3$

estimated OLS park effects across M=5 sims



estimated Ridge park effects across M=5 sims

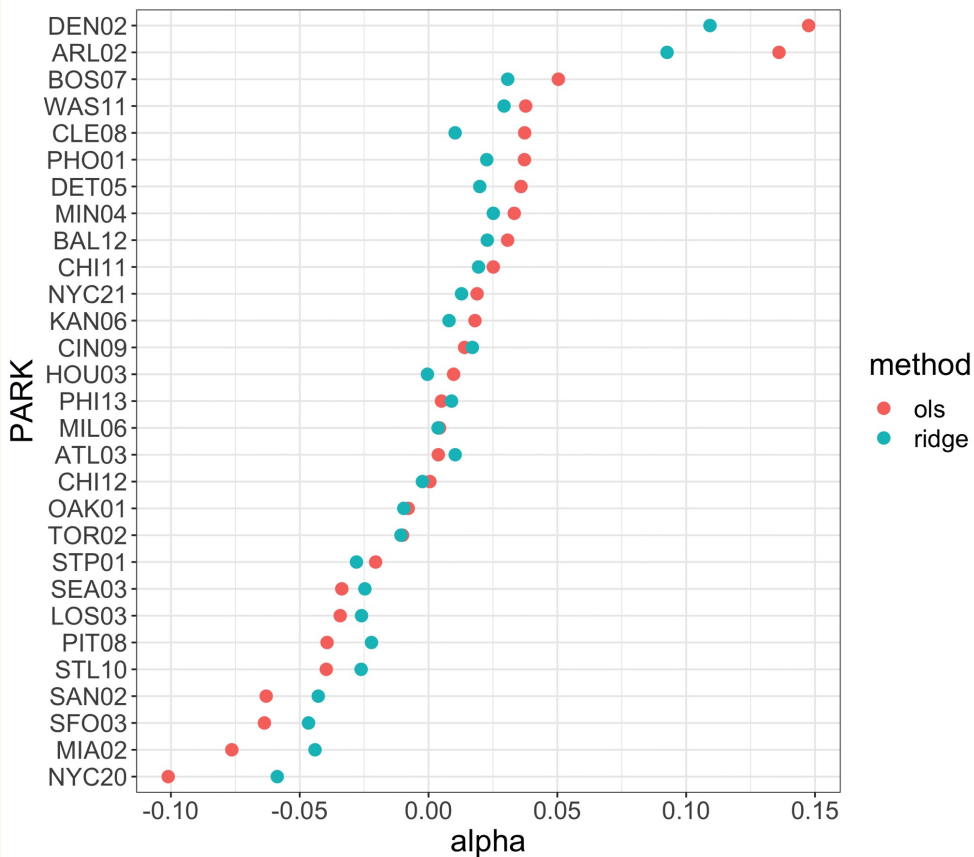


Ridge park effects are more stable across simulations, even though they are unbiased!

```

> ### error
> err(beta.pk.df.sim)
[1] 0.03528335
> err(beta.pk.df.sim_ride)
[1] 0.03804942
> ### error on non-outliers
> err(beta.pk.df.sim %>% filter( abs(beta.pk.true) < 0.05) )
[1] 0.02533202
> err(beta.pk.df.sim_ride %>% filter( abs(beta.pk.true) < 0.05) )
[1] 0.01690153
> ### error on outliers
> err(beta.pk.df.sim %>% filter( abs(beta.pk.true) >= 0.05) )
[1] 0.04406852
> err(beta.pk.df.sim_ride %>% filter( abs(beta.pk.true) >= 0.05) )
[1] 0.05359246

```



How do we quantify the sensitivity of an estimator to the Random Idiosyncrasies of a training dataset?

Model

$$y_i = f(x_i) + \varepsilon_i$$

"true" underlying function  $f$   
noise  $\varepsilon_i$        $\mathbb{E}\varepsilon_i = 0$

$$\mathbb{E}(y_i | x_i) = f(x_i)$$

Goal of ML: estimate  $f$  (obtain  $\hat{f}$ )  
as best as possible  
from a training dataset

$$\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$$

$$\hat{f} = \hat{f}(x; \mathcal{D})$$

Want our estimator  $\hat{f}$  to be as "close"  
to the true  $f$  as possible;  
on average we want  $\hat{f}$  to be  
as close to  $f$  as possible,


$$\text{MSE}(f, \hat{f}) := \mathbb{E} \left[ (f(x) - \hat{f}(x; \mathcal{D}))^2 \right]$$



averaging over the randomness  
in the training set  $\mathcal{D}$

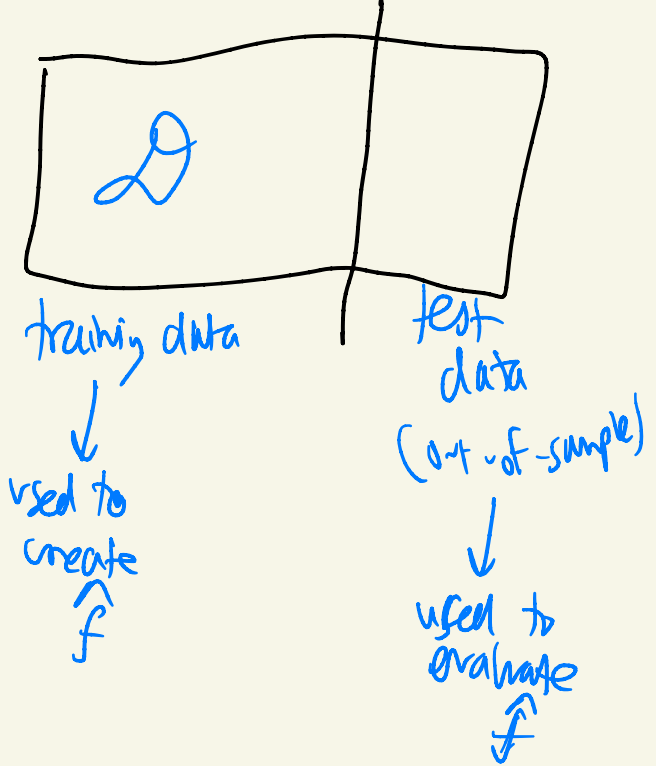
We don't actually observe  $f$ , so

$$\text{MSE} = \mathbb{E} \left[ (Y(x) - \hat{f}(x; \mathcal{D}))^2 \right]$$



out-of-sample  
testing data

full dataset



$$\begin{aligned} \text{MSE}(x; \mathcal{D}) &= \mathbb{E} \left[ (y - \hat{f}(x; \mathcal{D}))^2 \right] \\ &= \mathbb{E} (y - \hat{f})^2 && \hat{f} = \hat{f}(x; \mathcal{D}) \\ &&& x \mapsto \hat{f}(x) \\ &= \mathbb{E} (y^2 - 2y\hat{f} + \hat{f}^2) \\ &= \mathbb{E} y^2 - 2 \mathbb{E} (y\hat{f}) + \mathbb{E} (\hat{f}^2) \\ & y = f + \epsilon \end{aligned}$$

$$= \mathbb{E} (f + \varepsilon)^2 - 2 \mathbb{E} [(f + \varepsilon) \hat{f}] + \mathbb{E} \hat{f}^2$$

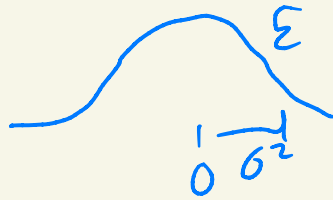
$$= \mathbb{E} [f^2 + 2f\varepsilon + \varepsilon^2] - 2 \mathbb{E} [f\hat{f} + \varepsilon\hat{f}] + \mathbb{E} \hat{f}^2$$

$x \mapsto f(x)$  is (unknown) fixed/constant

$$= f^2 + \cancel{2f \mathbb{E} \varepsilon} + \mathbb{E} \varepsilon^2$$

$$- 2f \mathbb{E} \hat{f} - \cancel{2 \mathbb{E}(\hat{f}) \mathbb{E}(\varepsilon)} + \mathbb{E} \hat{f}^2$$

$$\mathbb{E} \varepsilon = 0$$



$$\mathbb{E} \varepsilon = 0$$

$$\mathbb{E} \varepsilon^2 = \sigma^2$$

$$= f^2 - 2f \mathbb{E} \hat{f} + \mathbb{E} \hat{f}^2 + \mathbb{E} \varepsilon^2$$

$$= \underline{f^2} - \underline{2f \mathbb{E} \hat{f}} + \underline{(\mathbb{E} \hat{f})^2} + \mathbb{E} \hat{f}^2 - (\mathbb{E} \hat{f})^2 + \mathbb{E} \varepsilon^2$$

$$= \underbrace{(f - \mathbb{E} \hat{f})^2} + \left[ \mathbb{E} \hat{f}^2 - (\mathbb{E} \hat{f})^2 \right] + \mathbb{E} \varepsilon^2$$

$$y = f + \varepsilon$$

$$\text{MSE}(x; \mathcal{D}) = \mathbb{E} \left[ (Y - \hat{f}(x; \mathcal{D}))^2 \right]$$

$$= (f - \mathbb{E} \hat{f})^2 + \left[ \mathbb{E} \hat{f}^2 - (\mathbb{E} \hat{f})^2 \right] + \mathbb{E} \varepsilon^2$$

out-of-sample

$$\text{MSE}(\hat{f}) = \text{Bias}(\hat{f})^2 + \text{VAR}(\hat{f}) + \text{IRREDUCIBLE ERROR}$$

the Bias-Variance Tradeoff

- OLS: unbiased  $\mathbb{E} \hat{\beta}^{(\text{OLS})} = \beta$   
higher variance

↓  
Can't actually compute this



- Ridge: biased  
lower variance (more stable)

• variance  $\leftarrow$  sensitivity to the training set

overfitting = too much variance  
= memorizes noise in the training set, rather than getting the underlying trend

$$\text{Bias}(\hat{f}) = \mathbb{E} (f - \hat{f})^2$$

\*  $f$  in real life is the "truest" underlying function

$$\hat{f} = \text{OLS} \rightarrow X\hat{\beta} + \epsilon \quad \text{FALSE}$$