

Regularization and the Bias Variance Tradeoff

Q (Park Effects) Estimate the park effect α 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.

→ Read my full analysis in Appendix of "Grid WAR" paper

training data all half innings from 2017–2019

Variables i indexes the i^{th} half inning in our dataset
 $\text{park}(i)$ is the ballpark of half-inning i
 $ot(i)$ is the offensive team-score of half-inning i
 $dt(i)$ is the defensive team-score of half-inning i
 y_i is the Runs scored in half-inning i

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

where ε_i is mean zero noise, $E \varepsilon_i = 0$

The park effects α and team quality coefficients β, γ are unknown parameters which need to be estimated from data.

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

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

$$x_i^T = \begin{bmatrix} 1 & \underbrace{\bullet}_{\text{1 at Park}(i)} & \underbrace{\bullet}_{0 \text{ else}} & \bullet & \underbrace{\bullet}_{\text{1 at } ot(i)} & \underbrace{\bullet}_{0 \text{ else}} & \bullet & \underbrace{\bullet}_{\text{1 at } dt(i)} & \underbrace{\bullet}_{0 \text{ else}} & \dots \end{bmatrix}$$

1 Part 1 Part 2 ... Part 30 ot_1 ot_2 ... ot_{30} dt_1 dt_2 ... dt_{30}

Problem : Multicollinearity

When home team is an offense, $park(i) = ot(i)$.
 When road team is an offense, $park(i) = dt(i)$.
 So, it is tough to disentangle $\alpha_{park(i)}$ from $\beta_{ot(i)}$ and $\beta_{dt(i)}$.

{ Are the runs scored in those half-innings due to the offensive home team being good or the park being easy?

To disentangle these effects, we need a huge number of instances of Road teams on offence to figure out $\beta_{ot(i)}$ well, and a huge number of instances of Home teams on offence to figure out $\delta_{dt(i)}$ well. Then, with β_{ot} and δ_{dt} good, we can figure disentangle α_{park} .

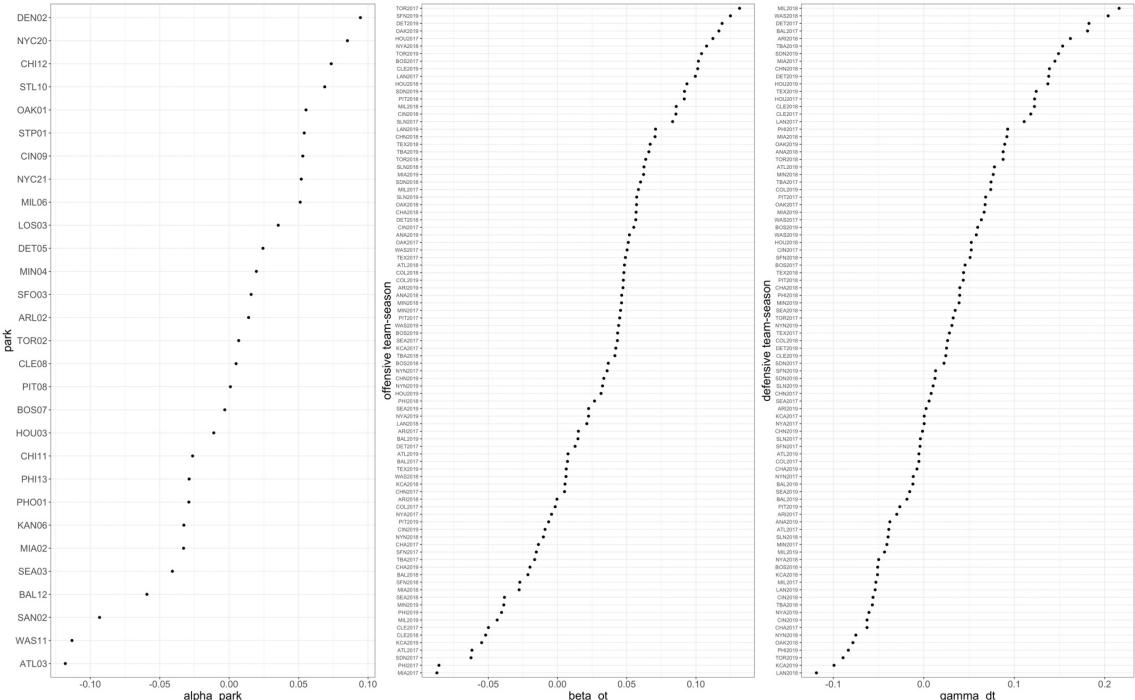
Our current dataset consists of 123,252 half-innings. This may seem like a lot of data, but due to our multicollinearity issue this actually isn't a huge amount of data. To demonstrate, we run a simulation study.

Q How much does multicollinearity affect our park effect estimates?
How well does OLS recover the park effects?

Simulation Study

Idea Pretend we knew the true coefficients, generate simulated data, and see how well we estimate the coefficients.

* Suppose the true coefficients are



which are chosen to have a "reasonable" scale.

* Then, assuming our model is true, let's generate the response y vector (Runs scored in a half inning) M times according to

$$y_i = \text{Round} \left(N_+ \left(x_i^T \beta, 1 \right) \right)$$

where x_i^T is the i^{th} half inning from our observed data matrix of all 123,252 half-innings from 2017 to 2019.

example snippet of simulated y

[, 1]	[, 1]
[1,]	1
[2,]	1
[3,]	1
[4,]	2
[5,]	1
[6,]	2
[7,]	1
[8,]	0
[9,]	1
[10,]	1

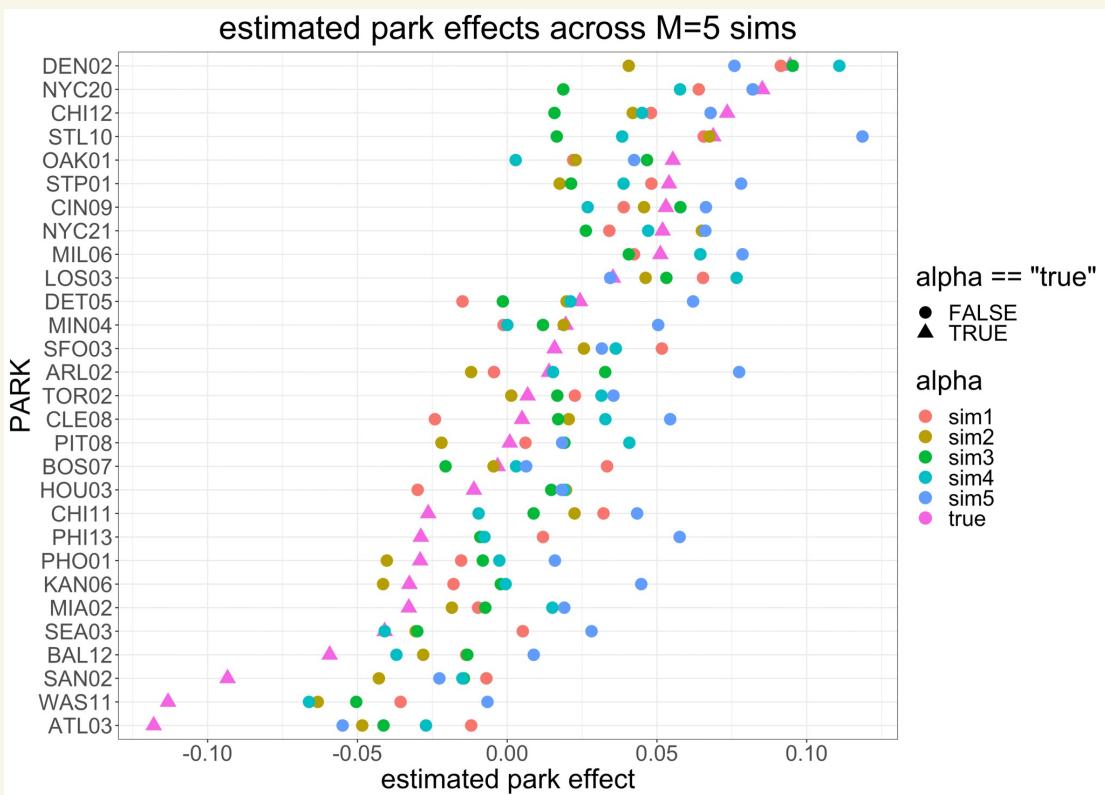
Here,

- N_+ means normal dist. conditional on it being ≥ 0
- "Round" because runs scored is an integer ≥ 0
- $EY_i \approx x_i^T \beta$
equivalently, $y_i \approx x_i^T \beta + \varepsilon_i$, $E\varepsilon_i = 0$
so our original model assumption holds true even if we don't explicitly write ε_i here

* Then, let's use linear regression to estimate the coefficients $\hat{\beta}$ on each of our M simulated datasets (X, y) and see how well we recover the park effects!

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

We can do this because it's a simulation and we know the "true" park effects.



* Due to Randomness in the training dataset, from the noise in generating y , each simulation yields very different park effects estimates $\hat{\alpha}$, even though the "true" park effects are the same.

* The OLS (ordinary least squares
=ordinary linear regression)
coefficients $\hat{\alpha}_{OLS}$ change
quite significantly across different
simulations; they are quite
sensitive to the noise of the
training set

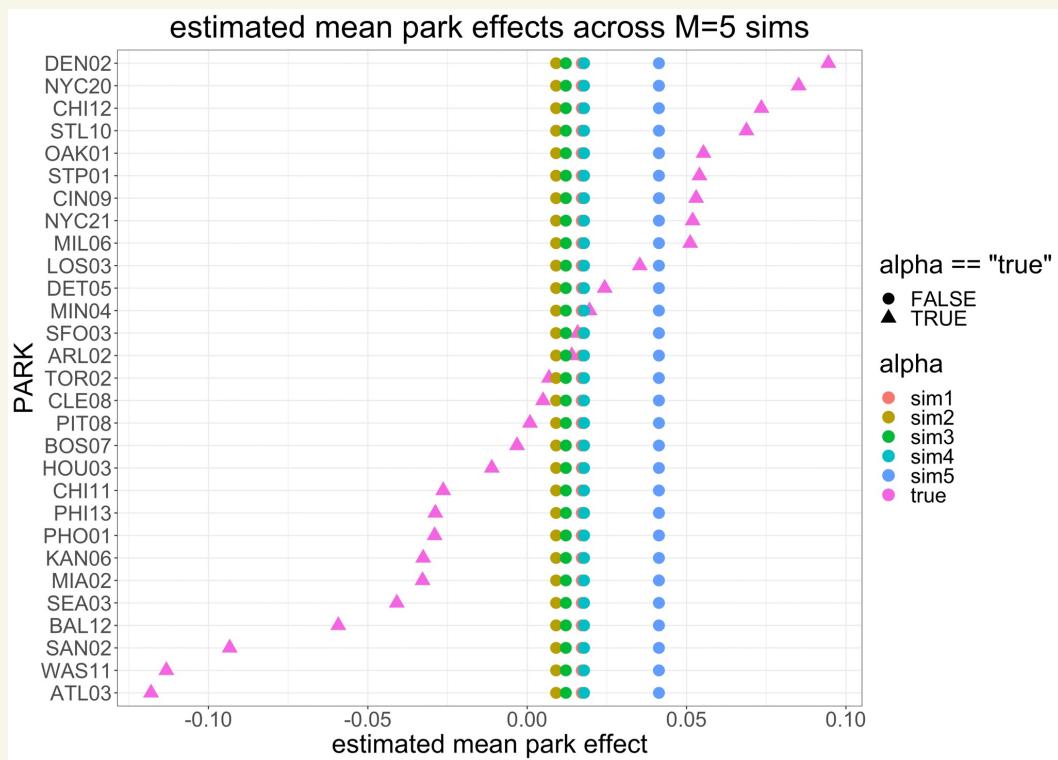
* How can we make the
coefficients less sensitive to the
random idiosyncrasies of our training set?

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

overall mean

$$\bar{x}$$

estimated mean park effect



- * Less sensitive to the random noise of the training set, but the overall mean is in wrong for many parks

Q How can we blend the strengths of OLS with the strengths of the overall mean?

OLS — estimates are generally in the right neighborhood of the "true" park effects, but are quite sensitive to the random idiosyncrasies of the training set.

Overall Mean — not too sensitive to the random idiosyncrasies of the training set, but are in the wrong neighborhood of the "true" park effects.

Idea Shrink the OLS estimates towards the overall mean.

In other words, "combine" the two estimators in a smart way.

→ this can be pretty complicated if the overall mean is unknown; we do a version of this (Empirical Bayes) next lecture.

→ Overall mean is ≈ 0 , so we can just shrink towards zero

→ Since we want park effect relative to mean park effect, shrinking towards 0 is fine

Idea Shrink the OLS estimates towards zero, i.e. just make them smaller, which will make them less sensitive!

* In ordinary linear regression, we estimate the coefficients β by minimizing the Residual sum of squares,

$$\hat{\beta}^{(\text{OLS})} = \underset{\beta}{\operatorname{argmin}} \quad \text{RSS}(\beta)$$

$$= \underset{\beta}{\operatorname{argmin}} \quad \sum_{i=1}^n (y_i - x_i^\top \beta)^2$$

* In Ridge Regression we instead minimize the RSS with a penalty term that encourages the estimated coefficients $\hat{\beta}$ to be smaller (i.e., to lie closer to 0),

$$\hat{\beta}^{(\text{Ridge})} = \underset{\beta}{\operatorname{argmin}} \quad \sum_{i=1}^n (y_i - x_i^\top \beta)^2 + \lambda \sum_j \beta_j^2$$

* This technique of adding a penalty term to the loss function we are minimizing is called Regularization.

The hyperparameter $\lambda > 0$ describes by how much we are penalized for having large β_j .

λ is simply a number, which is tuned using cross-validation.

Large $\lambda \rightarrow$ large penalty for large β_j
 \rightarrow forces β_j to be smaller.

$\lambda = 0 \rightarrow$ equivalent to OLS
 \rightarrow no shrinkage of β .

$$\hat{\beta}^{(\text{Ridge})} = \underset{\beta}{\operatorname{argmin}} \sum_{i=1}^n (y_i - x_i^\top \beta)^2 + \lambda \sum_j \beta_j^2$$

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

in matrix notation.

Calculus: Set gradient equal to 0 and solve!

$$\begin{aligned} L(\beta) &= (y - X\beta)^\top (y - X\beta) + \lambda \beta^\top \beta \\ &= y^\top y - 2\beta^\top X^\top y + \beta^\top X^\top X\beta + \lambda \beta^\top \beta \end{aligned}$$

$$\nabla_{\beta} L(\beta) = -2X^\top y - 2X^\top X\beta + 2\lambda\beta = 0$$

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

$$\Rightarrow \boxed{\hat{\beta}^{(\text{ridge})} = (X^\top X + \lambda I)^{-1} X^\top y}$$

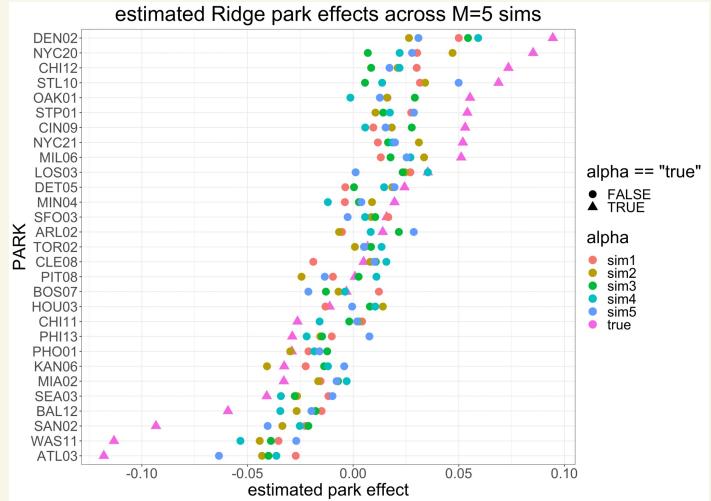
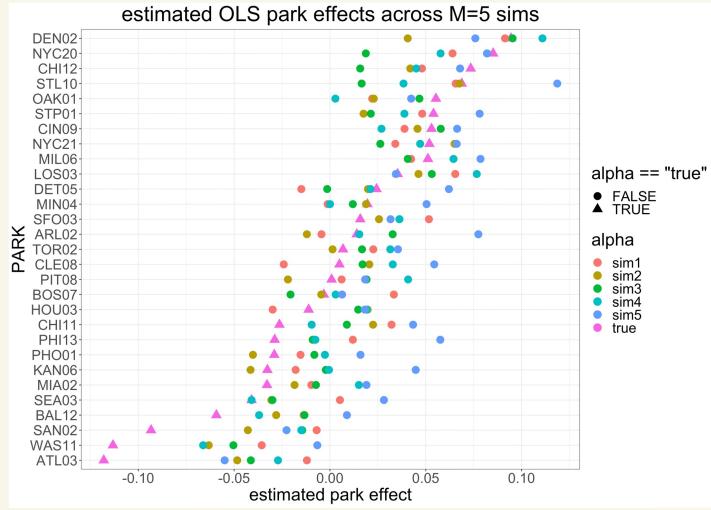
Solution always exists when $\lambda > 0$.

Ridge Regression — add matrix

$$\lambda I = \begin{pmatrix} \lambda & & 0 \\ & \ddots & \\ 0 & & \lambda \end{pmatrix} \text{ to } X^T X$$

prior to inverting. This is a "ridge" of λ 's.

$(X^T X + \lambda I)^{-1}$ is like multiplying by $\frac{1}{\sigma^2 + \lambda}$,
 $(X^T X)^{-1}$ is like multiplying by $\frac{1}{\sigma^2}$
adding $\lambda > 0$ to the denominator
shinks the estimates $\hat{\beta}_0$!

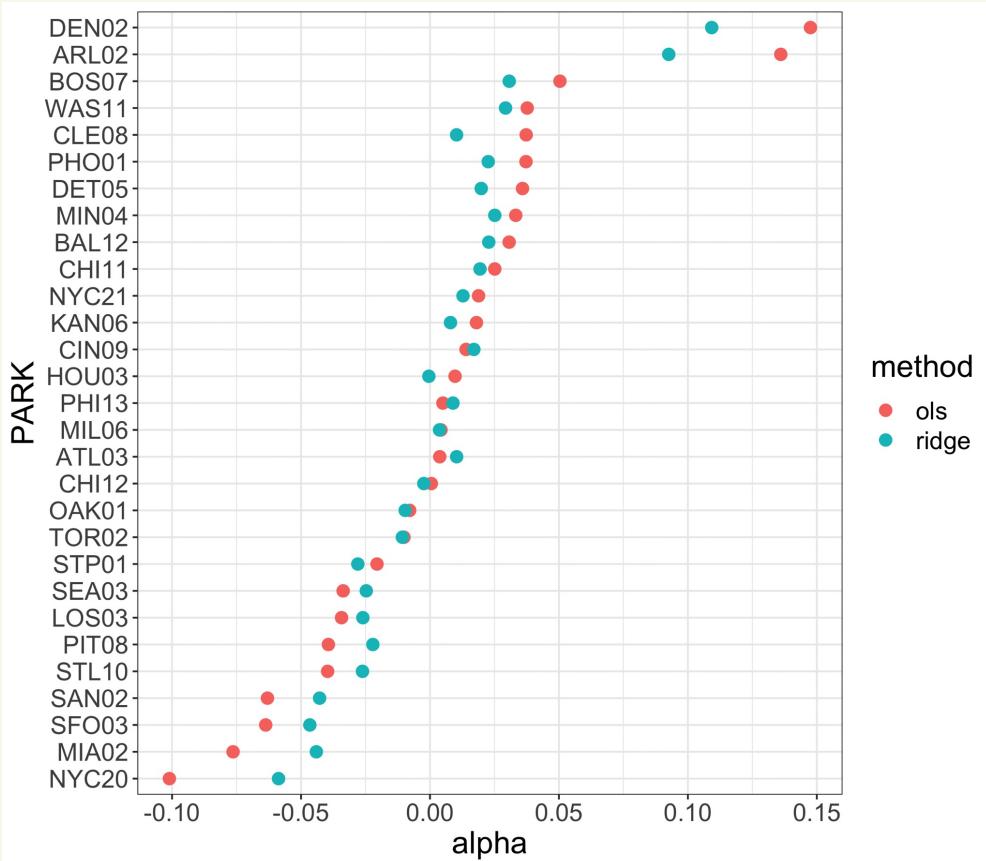


* Ridge regression park effect estimates indeed are more stable across simulations, i.e. are less sensitive to the noise of the training set!

```
> ### error
> err<-beta_pk$df$sim
[1] 0.95528335
> err<-beta_pk$df$sim_ridge
[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$sim_ridge %>% filter( abs(beta_pk$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$sim_ridge %>% filter( abs(beta_pk$true) >= 0.05 )
[1] 0.05359246
```

* Shrinking outliers isn't always a great idea; OLS outperforms on outliers

* Park effects on Real MLB data, 2017-2019



- * We see that Ridge indeed shrinks the park effects toward zero!
- * On this Real data, it turns out that the Ridge shrunken park effects are everywhere better than OLS since OLS overfits...
(based on out-of-sample predictive performance)

Q How do we quantify the sensitivity of an estimator to the random idiosyncrasies of a training dataset?

Model Suppose $y_i = f(x_i) + \varepsilon_i$ for some "true" underlying function f and noise ε_i with $E\varepsilon_i = 0$.

Goal is to estimate f with \hat{f}

e.g.
 [OLS
 Ridge
 overall mean
 etc.

training dataset $D = \{(x_i, y_i)\}_{i=1}^n$
 $\hat{f} = \hat{f}(x; D)$

Want our estimator \hat{f} to be as "close" to true f as possible, which we can measure from data as the smallest **out-of-sample Mean Squared Error** which uses datapoints (X, Y) not in the training dataset,

$$MSE(f, \hat{f}) := E[Y - \hat{f}(x; D)]^2.$$

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

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

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

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

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

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

since $Y = f(x) + \varepsilon$

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

using $f = f(x)$ as shorthand

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

$$+ \mathbb{E}\hat{f}^2$$

since $f(x)$ is deterministic and not random
and $\hat{f}(x)$ is independent of ε

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

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

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

$$\Rightarrow \text{Bias}(\hat{f})^2 + \text{Var}(\hat{f}) + \sigma_\varepsilon^2.$$

Bias Variance Tradeoff

$$\text{MSE}(x; D) = [\text{Bias } \hat{f}(x; D)]^2 + \text{Var}(\hat{f}(x; D)) + \sigma_\varepsilon^2$$

* $\sigma_\varepsilon^2 = \mathbb{E}\varepsilon^2$ is irreducible error,
the noise inherent to the problem

(e.g. for pure effects, σ_ε^2 is the inherent noise if scoring a certain number of runs in a half inning)

* $\text{Var}(\hat{f}) = \mathbb{E}\hat{f}^2 - (\mathbb{E}\hat{f})^2$
is how variable \hat{f} is depending on the training set, i.e. how much it responds to randomness in the training set

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

is how close \hat{f} is to f
on average (averaged over $N \rightarrow \infty$
training sets if the same size)

* Bias-Variance Tradeoff for Park Effects:

1. Overall Mean $\begin{cases} \text{very low variance} \\ \text{very high bias} \end{cases}$

2. OLS $\begin{cases} \text{low bias} \\ \text{high variance} \end{cases}$

3. Ridge — introduces bias with the
penalty term $+\lambda \sum_j \beta_j^2$
in order to lower
variance!

Takeaway To make better predictions, sometimes it helps to introduce some bias to lower the variance!