

## Lab 16: Bias-Variance Tradeoff

*Instructor: Jonathan Pipping**Author: Ryan Brill*

## 16.1 Park Effects Simulation Study

### 16.1.1 Motivation

It would be great to compute the bias and variance

$$\begin{aligned}\text{Bias}(\hat{f}) &= (f - \mathbb{E}_{\mathcal{D}}[\hat{f}]) \\ \text{Var}(\hat{f}) &= \mathbb{E}_{\mathcal{D}}[\hat{f}^2] - (\mathbb{E}_{\mathcal{D}}[\hat{f}])^2\end{aligned}$$

for real-world estimators  $\hat{f}$  of real-world sports datasets  $\mathcal{D}$ , but this is impossible because the "true" function  $f$  is unknown and  $\mathbb{E}_{\mathcal{D}}$  is an expectation over the randomness of drawing dataset  $\mathcal{D}$ , which is impossible to calculate. To understand the nature of this tradeoff, we turn to a **simulation study**, using park effects as our example.

### 16.1.2 Simulation Setup

In this simulation study, we assume that the park effects, team offensive quality, and team defensive quality have known coefficients. You will use these coefficients to simulate  $y$ , the runs scored in each half-inning in the dataset. Then you will model your data using OLS and Ridge Regression, using bias and variance to compare the two models.

### 16.1.3 Your Task

1. Generate a "true" parameter vector  $\beta$  in  $\mathbf{R}$  according to the following distributions:

$$\begin{aligned}\beta_0 &= 0.4 \\ \beta_j^{(\text{park})} &\stackrel{\text{iid}}{\sim} \mathcal{N}(0.04, 0.065) \\ \beta_k^{(\text{off})} &\stackrel{\text{iid}}{\sim} \mathcal{N}(0.02, 0.045) \\ \beta_k^{(\text{def})} &\stackrel{\text{iid}}{\sim} \mathcal{N}(0.03, 0.07)\end{aligned}$$

2. Assemble your park effects data matrix  $\mathbf{X}$  associated with the model

$$y_i = \beta_0 + \beta_{\text{park}(i)} + \beta_{\text{off}(i)} + \beta_{\text{def}(i)} + \epsilon_i$$

where each  $i$  is a half-inning.

3. Use a for loop to simulate a "true" outcome vector  $y^{(m)}$  for  $m = 1, \dots, 100$  datasets according to the following formula:

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

where  $\mathcal{N}_+$  is the normal distribution truncated to be positive.

4. Your goal is to recover the park effects  $\beta^{(\text{park})}$  from the simulated data  $(X, y^{(m)})$ . Do this by fitting OLS and Ridge Regression models, which will give estimates  $\hat{\beta}_{\text{OLS}}^{(m, \text{park})}$  and  $\hat{\beta}_{\text{Ridge}}^{(m, \text{park})}$ .
5. Then, estimate  $\mathbb{E}_{\mathcal{D}}[\hat{\beta}_{\text{OLS}}^{(m, \text{park})}]$  and  $\mathbb{E}_{\mathcal{D}}[\hat{\beta}_{\text{Ridge}}^{(m, \text{park})}]$  by averaging your estimates  $\beta^{(m)}$  over the 100 datasets. Then, estimate the average bias for each model using the following formula:

$$\|\beta^{(\text{park})} - \mathbb{E}_{\mathcal{D}}[\hat{\beta}_{\text{model}}^{(m, \text{park})}]\|_2$$

where  $\|\cdot\|_2$  is the Euclidean norm.

6. Estimate the variance similarly by averaging over the 100 datasets.
7. Compare OLS to Ridge Regression via estimated bias and variance. Visualize your results.