

Regularization & Ridge Regression

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
 $\text{ot}(i)$ is the offensive team-strength of half-inning i
 $\text{dt}(i)$ is the defensive team-strength of half-inning i
 Y_i is the Runs scored in half-inning i

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

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

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

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} .

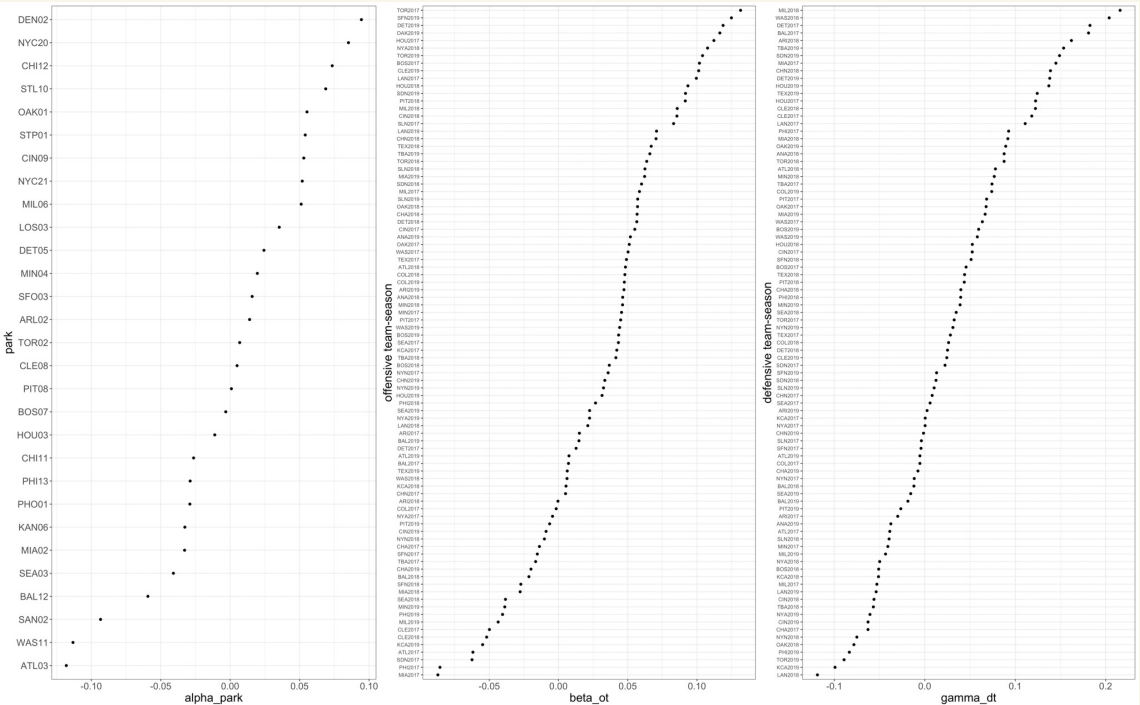
Our current dataset consists of 123,252 half-innings. This may seem like a lot of data, but due to our multicollinearity issue this actually isn't a huge amount of data. It's easy to overfit to noise and get the coefficients wrong. To demonstrate, we run a simulation study.

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 α, β, γ , generate fake historical data y (121,000 fake inning outcomes), and see how well we estimate the coefficients from this synthetic data.

* Suppose the "true" coefficients are



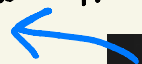
which are chosen to have a "reasonable" scale.

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

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

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

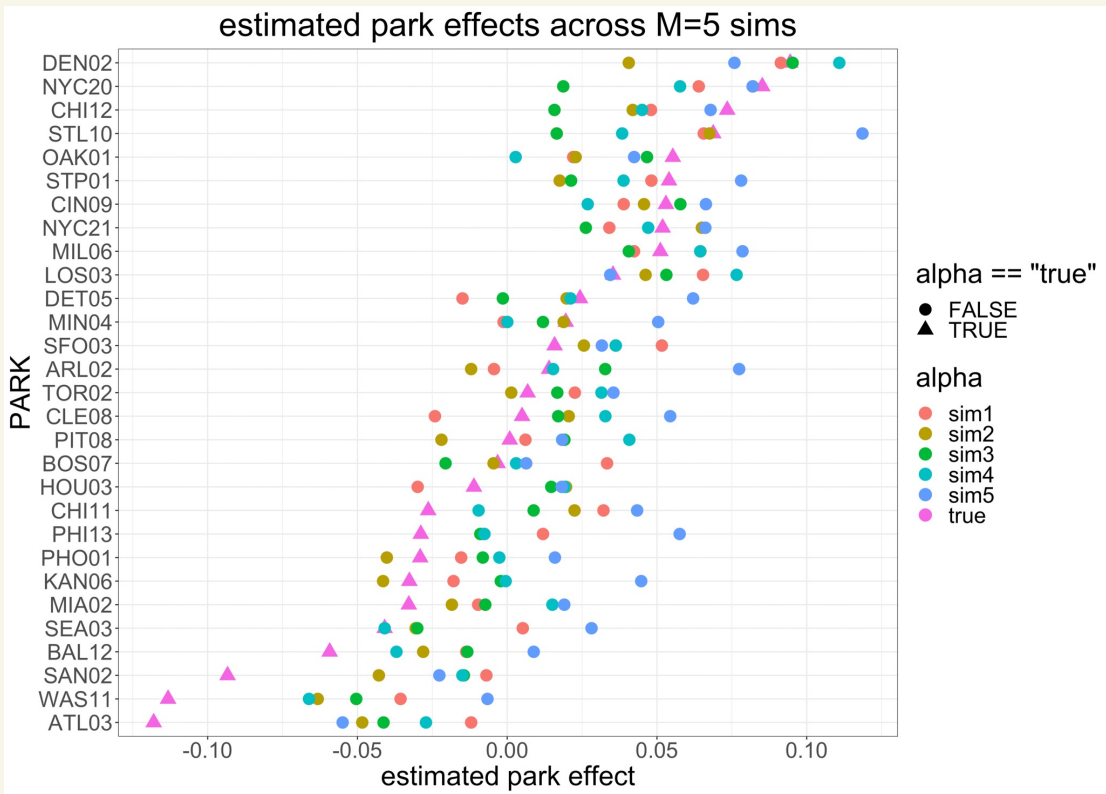
Here,

- \mathcal{N}_+ means normal dist. conditional on it being ≥ 0
- "Round" because runs scored is an integer ≥ 0
- $\mathbb{E} y_i \approx x_i^T \beta$
equivalently, $y_i \approx x_i^T \beta + \varepsilon_i$, $\mathbb{E} \varepsilon_i = 0$
so our original model assumptions hold here
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.



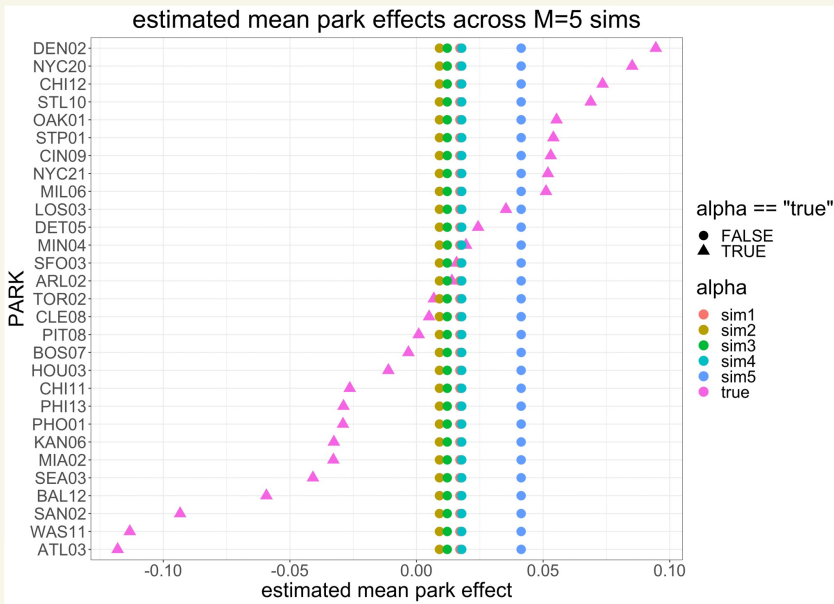
* Due to Randomness in the training dataset, from the noise in generating y , each simulation yields very different park effects estimates $\hat{\alpha}$, even though the "true" park effects are the same.

* 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 idiosyncracies of our training set?

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

overall mean $\widehat{\alpha}$ estimated mean park effect
zero the constant value 0



* Constant values like zero, or overall mean — not too sensitive to the random idiosyncrasies of the training set, but are wrong for many parks

* OLS park effect estimates — very sensitive to the randomness of the training set, but are unbiased (on average, i.e. averaged over many training set generations, they are in the right spot)

* there is a tradeoff between sensitivity and unbiasedness

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

Idea Shrink the OLS estimates towards a constant value, like the overall mean or zero.

* The Bayesian way to shrink the coefficients is via a prior:

$$\left\{ \begin{array}{l} \beta_j \sim \mathcal{N}(\mu, \sigma^2) \quad \forall j \\ \text{or } \beta_j \sim \mathcal{N}(0, \sigma^2) \quad \forall j \\ \text{same for } \alpha_j, \gamma_j \end{array} \right.$$

* The drawback of fitting Bayesian models is that it can be extremely computationally intensive...

* Is there a quicker (and frequentist) way of accomplishing this?

Idea Make the coefficients less sensitive to noise by shrinking the OLS estimates towards zero by altering the loss function to be optimized.

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

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

$$= \operatorname{argmin}_{\beta} \sum_{i=1}^n (y_i - x_i^T \beta)^2$$

* 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),

$$\hat{\beta}^{(\text{Ridge})} = \operatorname{argmin}_{\beta} \underbrace{\sum_{i=1}^n (y_i - x_i^T \beta)^2}_{\substack{\text{Want } x_i^T \beta \text{ to be} \\ \text{close to } y}} + \underbrace{\lambda \sum_j \beta_j^2}_{\substack{\lambda > 0 \text{ is a} \\ \text{hyperparameter} \\ \text{Want } \beta \text{ to} \\ \text{be smaller, or} \\ \text{close to } 0}}$$

* 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_j .

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

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

in matrix notation.

Calculus: Set gradient equal to 0 and solve!

$$L(\beta) = (y - X\beta)^T (y - X\beta) + \lambda \beta^T \beta$$

$$= y^T y - 2\beta^T X^T y + \beta^T X^T X \beta + \lambda \beta^T \beta$$

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

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

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

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 $\frac{1}{\bullet + \lambda}$,

$(X^T X)^{-1}$ is like multiplying by $\frac{1}{\bullet}$

adding $\lambda > 0$ to the denominator
shrinks the estimates $\hat{\beta}$!

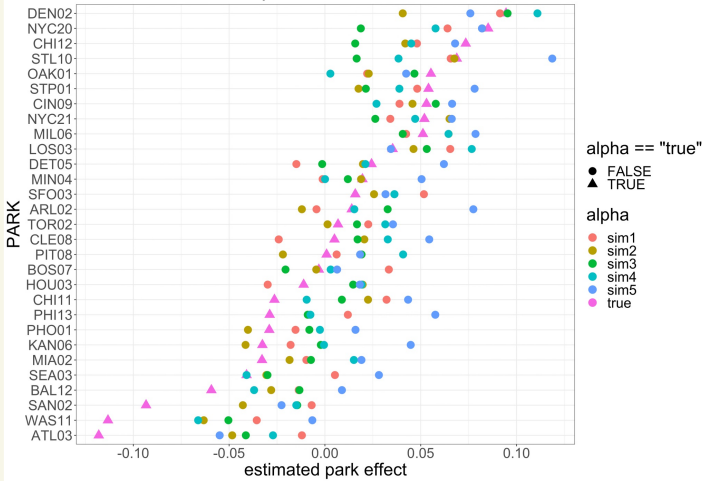
* Math HW: the Bayesian Regression Model $\begin{cases} y_i \stackrel{iid}{\sim} N(x_i^T \beta, \sigma^2) \\ \beta_i \stackrel{iid}{\sim} N(0, \frac{\sigma^2}{\lambda}) \end{cases}$

has maximum a posteriori (MAP) estimate, which is like a Bayesian version of MLE, equal to

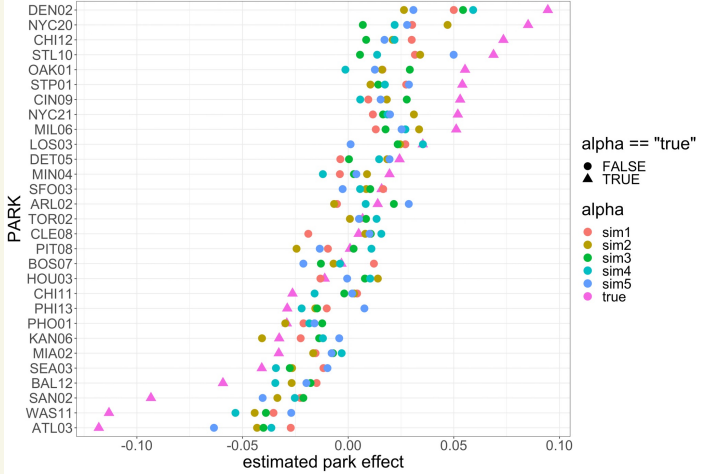
$$\hat{\beta} = \operatorname{argmax}_{\beta} P(\beta | X, y) = (X^T X + \lambda I)^{-1} X^T y$$

and so is equivalent to Ridge Regression.

estimated OLS park effects across M=5 sims



estimated Ridge park effects across M=5 sims



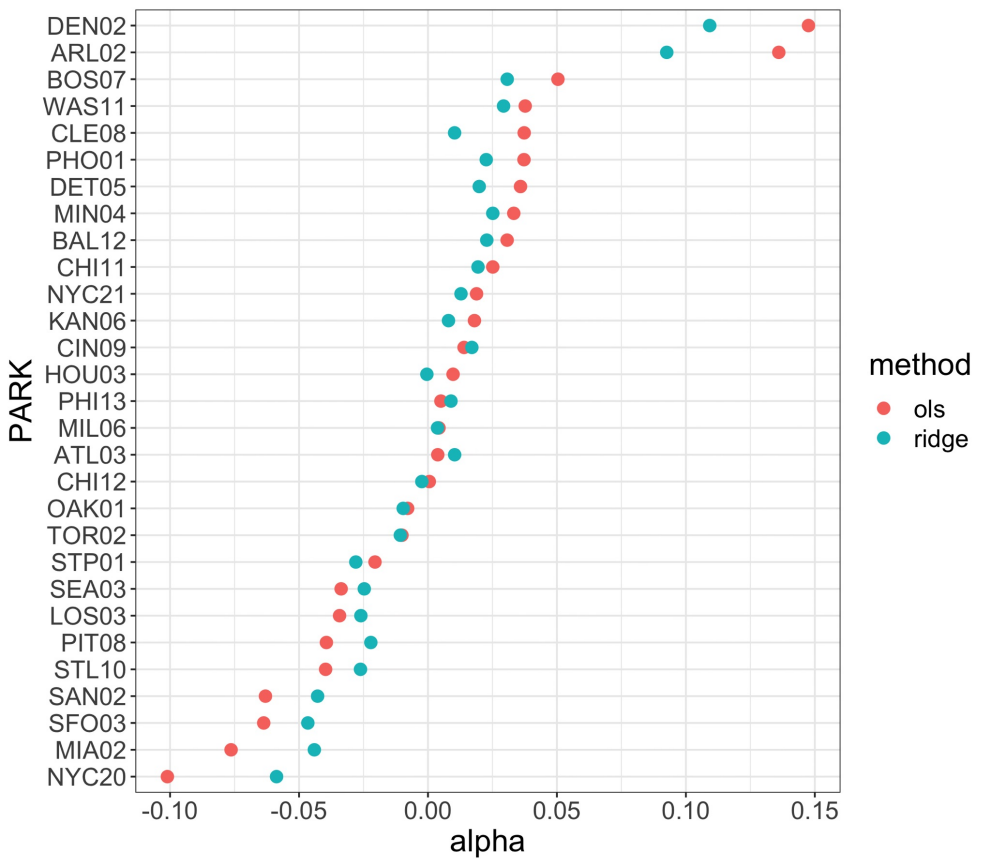
* Ridge regression park effect estimates indeed are more stable across simulations, i.e. are less sensitive to the noise of the training set!

```

> ## error
> err(beta.pk.df.sim)
[1] 0.03528335
> 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 towards 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)