

I. Introduction

II. Experimental Setting

III. Prediction Based on Different Models

IV. Summary

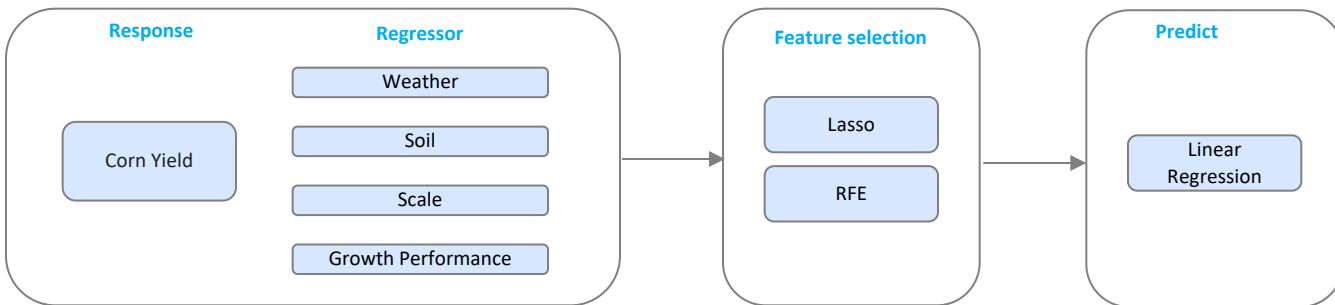
Using Linear Regression Method to Predict Corn Yield

Background and Method

Background

- Crop yield prediction is essential for global food security but is complex, we are focusing **on corn yield**.
- Our regressors represent four different aspects: **weather(weekly), soil, scale, and growth performance(weekly)**, and the response is **corn yield¹ (BU/A)**.
- We collected data from 1990 to 2018, with a total of **8,352 observations and 688 features**.
- Our goal is **to forecast corn yield performance** for the Corn Belt states Illinois, Iowa, and Indiana.

Method

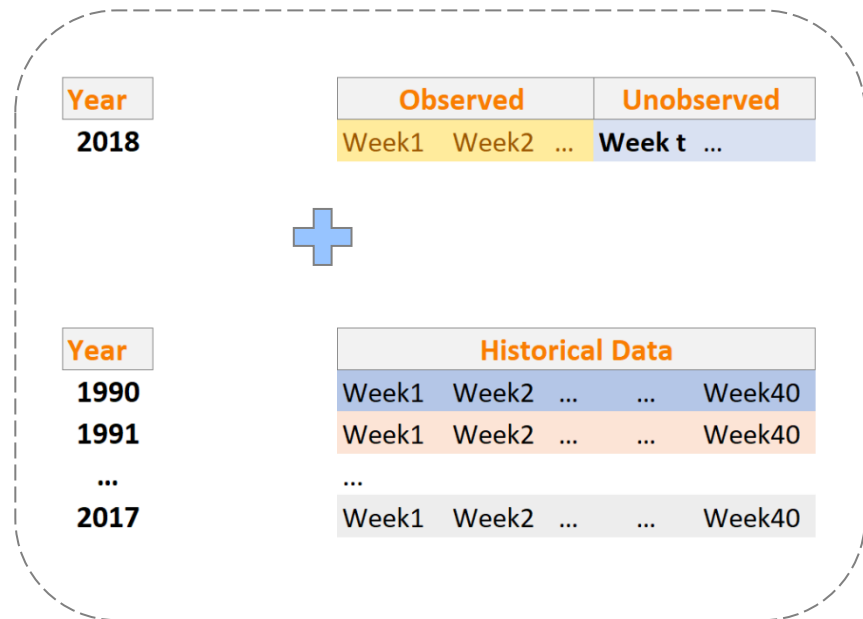


Using Historical Data as Potential Data Solution

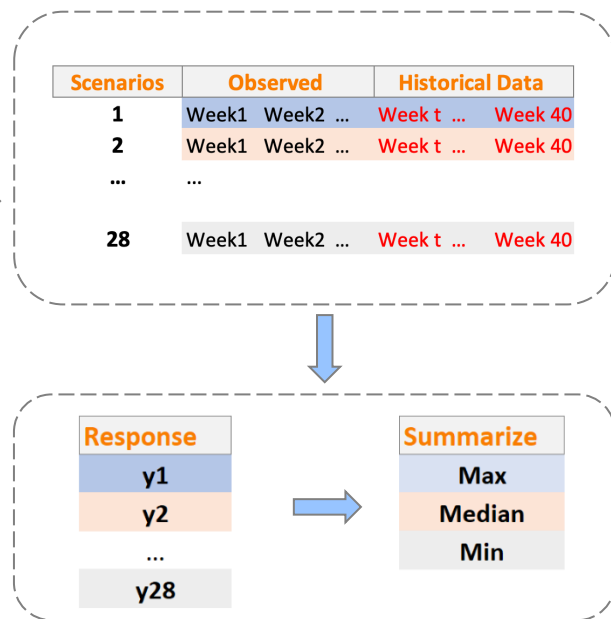
Experimental Setup

2018

Weather and Growth Performance Information at Week t



Weather and Growth Performance Information at Week t



Key Points

- As crops grow, we can collect more detailed information throughout the year, so we can **improve our forecasts on a weekly basis**, which could be beneficial to policy makers and farmers.
- We run all data from previous years as the **underlying data to construct prediction intervals**.
- We summarize the **maximum, median, and minimum** as significant result for our prediction.

Model I : Linear Regression Using Lasso as Feature Selection Method

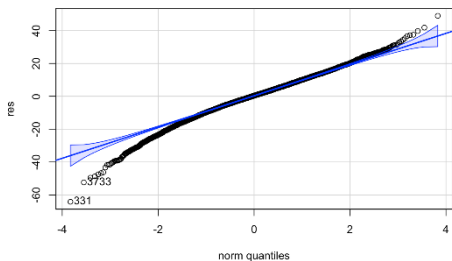
Function :

$$y = \beta_0 + \beta_w X_w + \beta_{soil} X_{soil} + \beta_{scale} X_{scale} + \beta_g X_g + \varepsilon$$

Model fitting Performance

2017-2018

QQ-Plot



Result

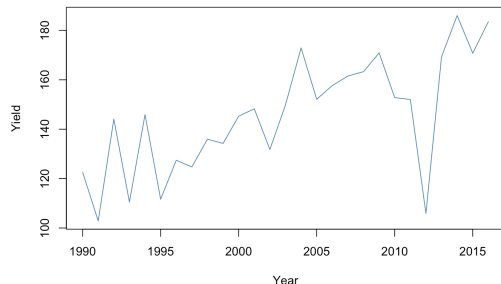
Residual standard error: 10.82 on 7247 degrees of freedom
Multiple R-squared: 0.8872, Adjusted R-squared: 0.8792
F-statistic: 110.1 on 518 and 7247 DF, p-value: < 2.2e-16

Explain

- The idea of the lasso is adding a penalty term to constrain the equation. The consequence of imposing such a penalty is to shrink the coefficient values towards zero, this would **set the less contributive regressor to have a zero coefficient** which achieves the purpose of variable elimination.
- We successfully reduce the variables **from 688 to 519 with 88% R-squared**.
- We select the data from 1990-2016 to train the model and apply the model to predict the corn yield in 2017. And train model based on the data from 1990-2017 to 2018.

Corn Yield Trend

1990-2016



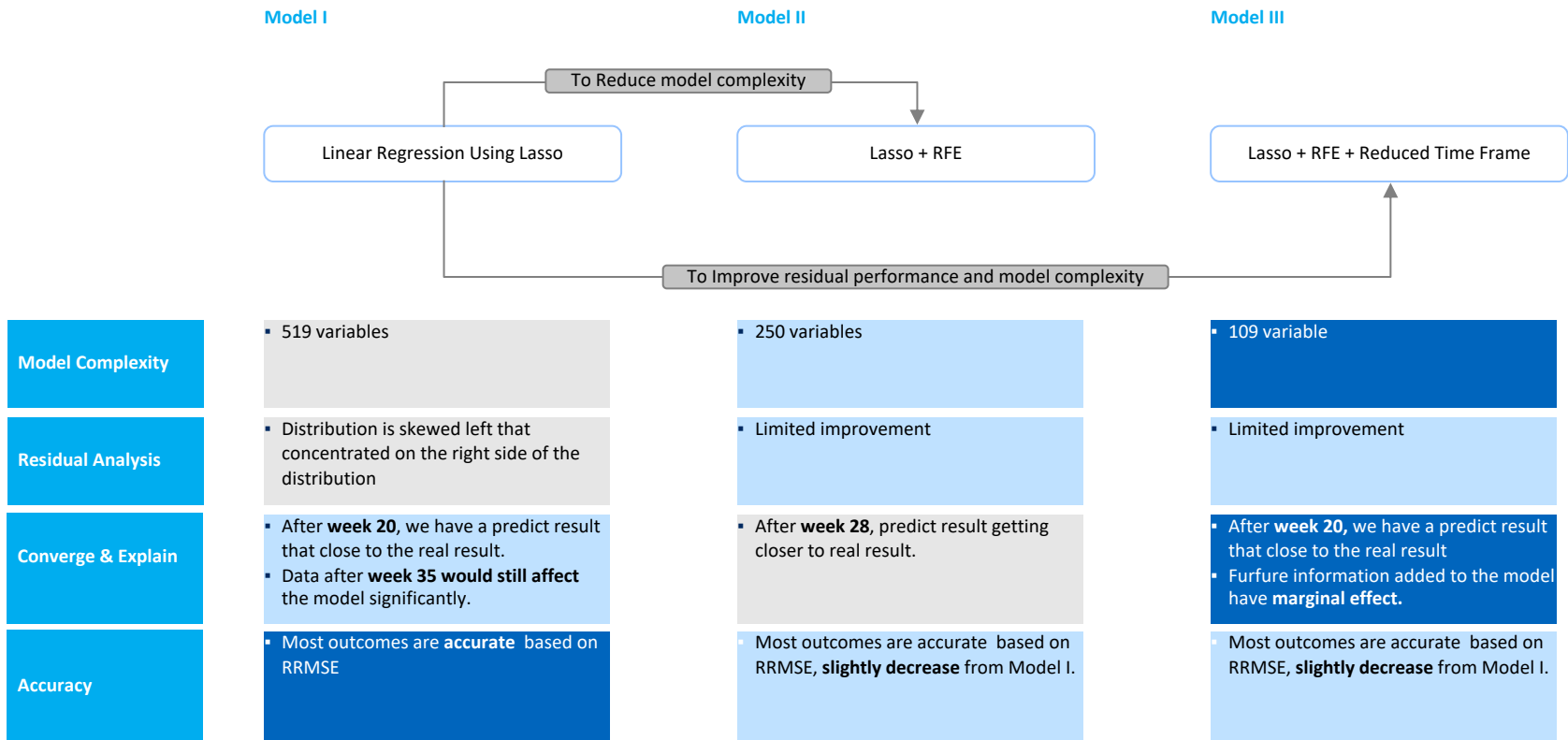
Bias and Improvements

- Model is **complex** with 519 variables: RFE
- **Residual not fit perfectly** : applied interaction terms; get ride of Genetic Modified Seed effects **by reduce time frame**

Improvements for Model I Using Different Strategies

Model Progress

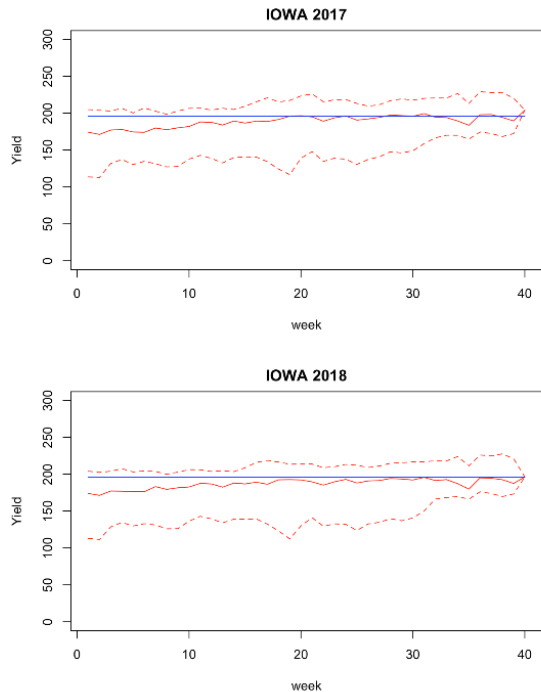
Improvement needed Good Great



Iowa: Prediction Interval based on Weekly Updated Data under Different Model

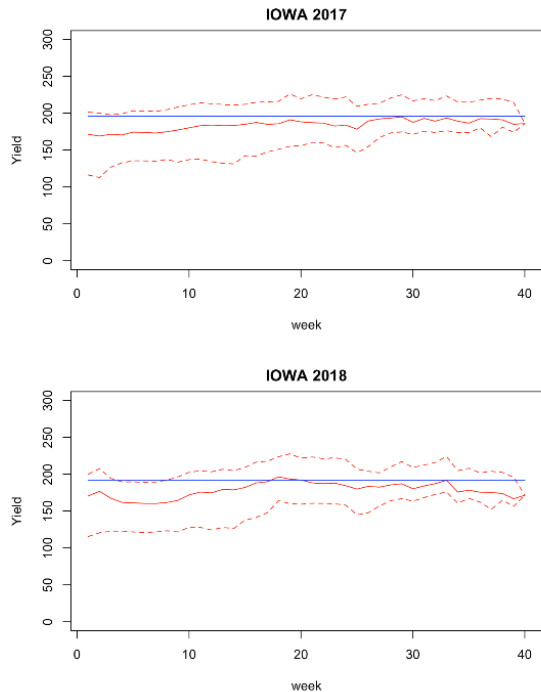
Model I

Yield¹:BU/A ; 2017 - 18



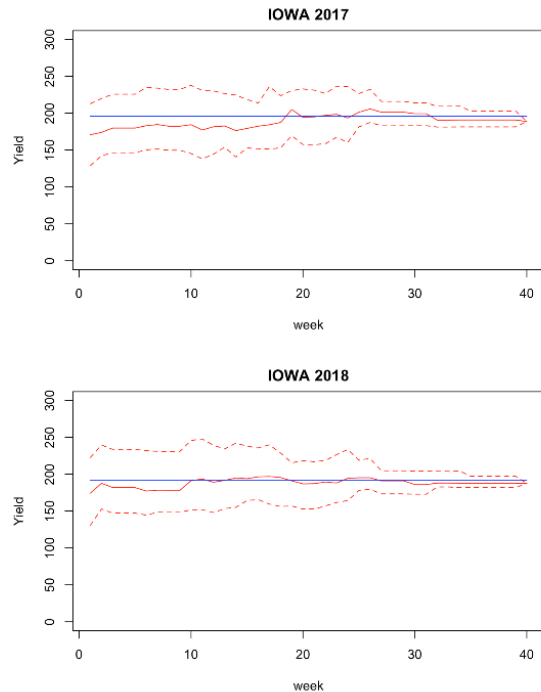
Model II

Yield:BU/A ; 2017 - 18



Model III

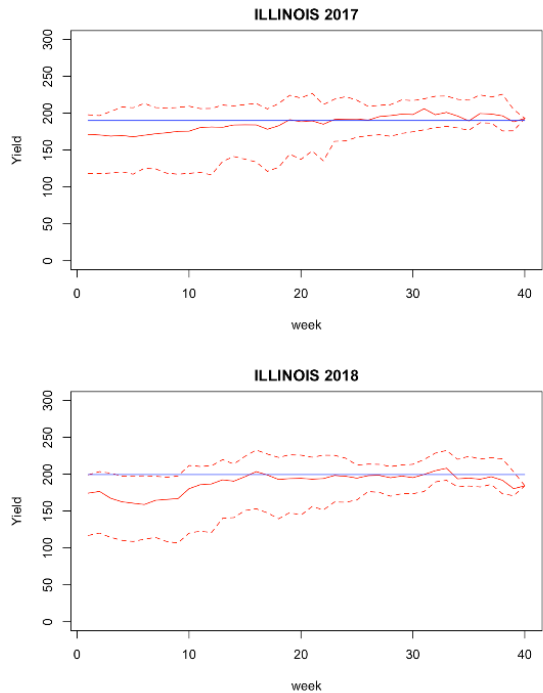
Yield:BU/A ; 2017 - 18



Illinois: Prediction Interval based on Weekly Updated Data under Different Model

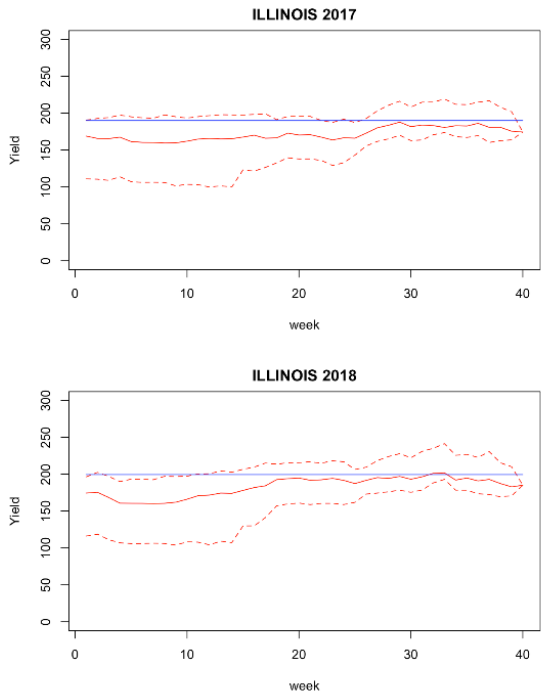
Model I

Yield:BU/A ; 2017 - 18



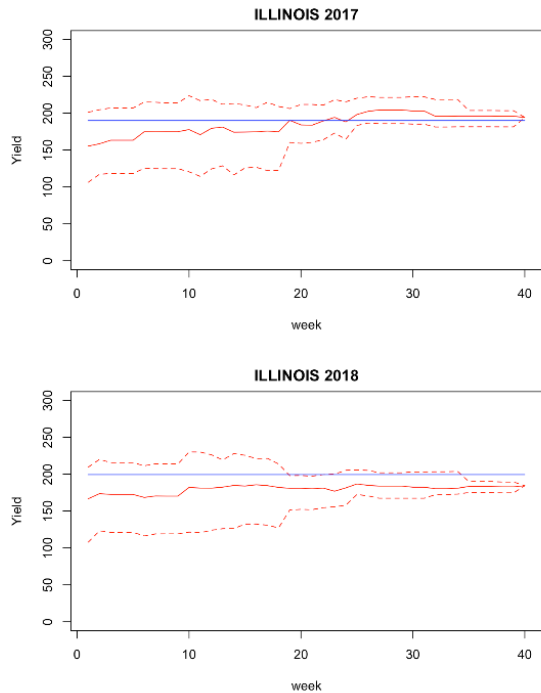
Model II

Yield:BU/A ; 2017 - 18



Model III

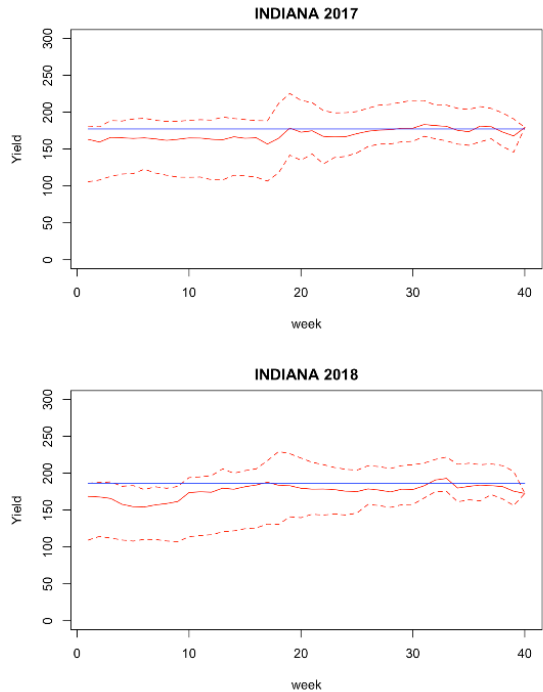
Yield:BU/A ; 2017 - 18



Indiana: Prediction Interval based on Weekly Updated Data under Different Model

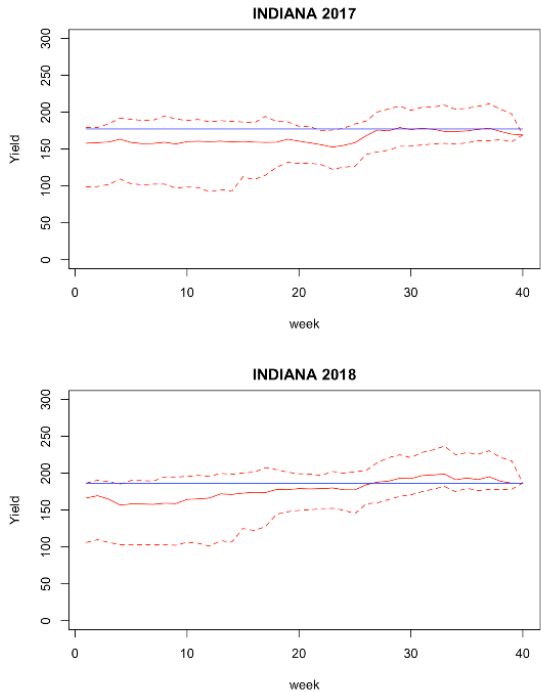
Model I

Yield:BU/A ; 2017 - 18



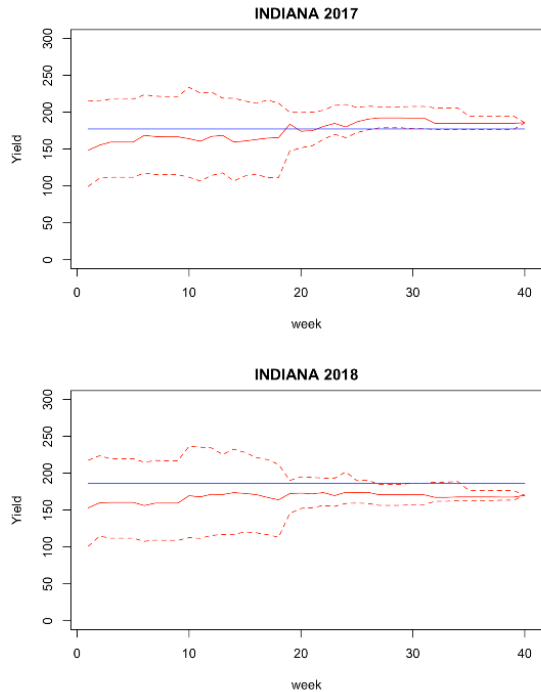
Model II

Yield:BU/A ; 2017 - 18



Model III

Yield:BU/A ; 2017 - 18

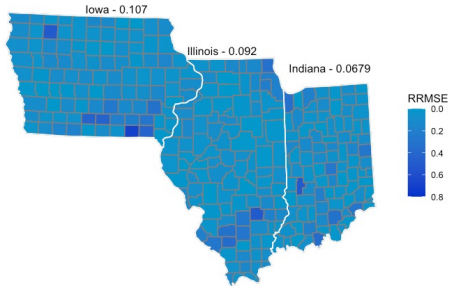


Prediction Accuracy at County Level for Different Models

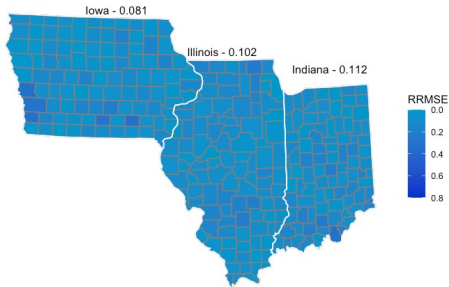
Model I

2017 - 18

Predict Accuracy Based on RRMSE for 2017 - Lasso



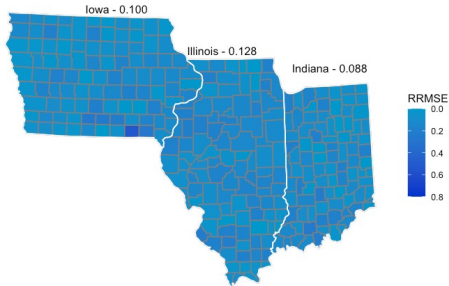
Predict Accuracy Based on RRMSE for 2018 - Lasso



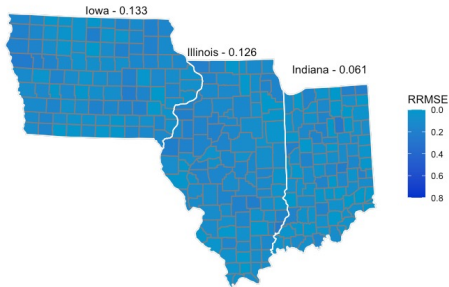
Model II

2017 - 18

Predict Accuracy Based on RRMSE for 2017 - RFE250



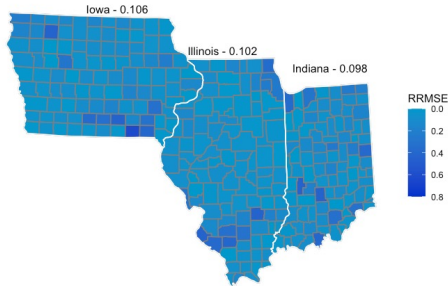
Predict Accuracy Based on RRMSE for 2018 - RFE250



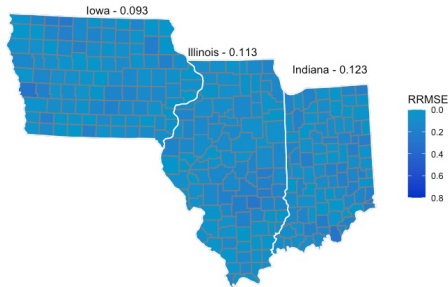
Model III

2017 - 18

Predict Accuracy Based on RRMSE for 2017 - RFE109



Predict Accuracy Based on RRMSE for 2018 - RFE109



Conclusion: Model III Would Be the Optimize Solution for Prediction

Model Comparison

| Model | Model Fitting | | | RRMSE | | |
|----------|---------------|-----------|---------------|----------|---------|------|
| | Variables | R-squared | Adj R-squared | Illinois | Indiana | Iowa |
| I- 17 | 519 | 89% | 88% | 9% | 7% | 11% |
| II- 17 | 250 | 84% | 83% | 13% | 9% | 10% |
| III - 17 | 109 | 75% | 74% | 10% | 10% | 11% |
| I- 18 | 519 | 89% | 88% | 10% | 11% | 8% |
| II- 18 | 250 | 84% | 84% | 13% | 6% | 13% |
| III - 18 | 109 | 76% | 75% | 11% | 12% | 9% |

- Model I performed the **best in both model fitting and prediction accuracy**, but the **complexity** needs to be improved.
- The complexity of Model II is lower, but the prediction **accuracy dropped** based on RRMSE.
- Model III performed **best in terms of model complexity**, and prediction result **outperforming Model II** but slightly weaker than Model I.
- After comparison, all models have their own biases and strengths. We choose **Model III as the optimal solution** for our project.

Thank you

