

Predicting Rain in Sydney, Australia

Kristofer Schobert

Our Goal



Use a supervised learning classifier to answer:

“Will it rain tomorrow in Sydney, Australia?”



Significance of Weather Predictions

- Farmers
- Transit
- Safety
- Future of our planet

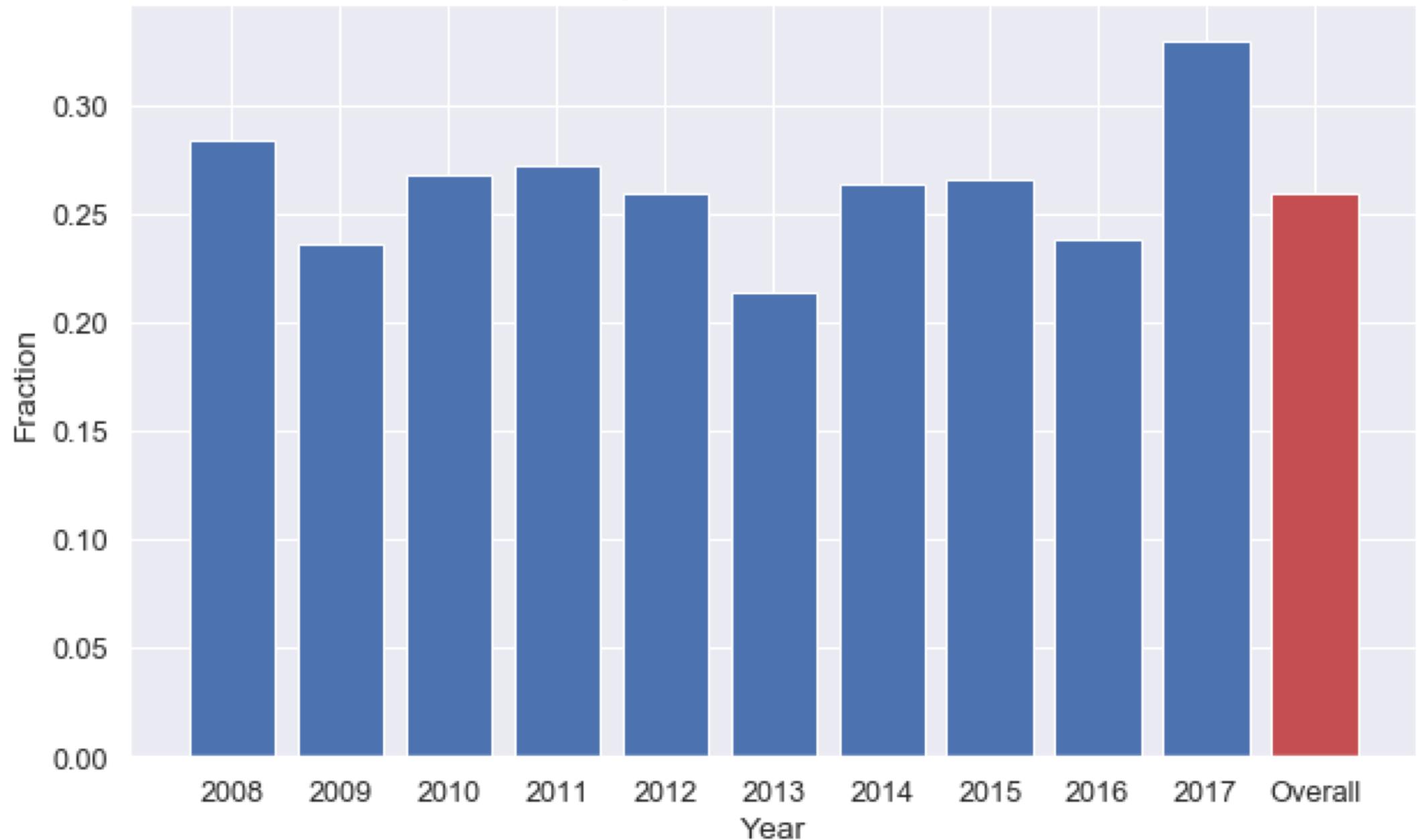


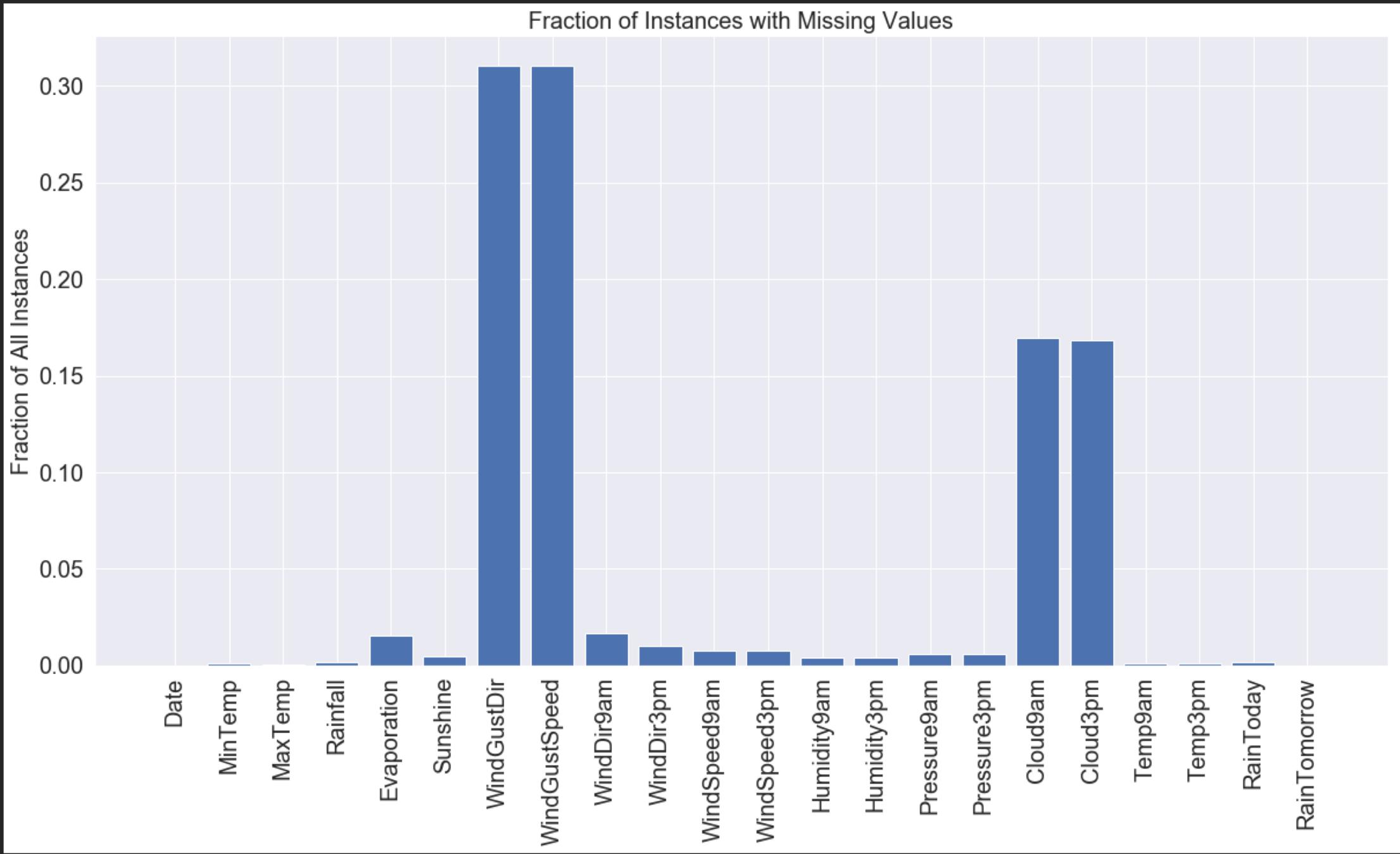
Our Dataset

- Daily Weather Observations
 - Several Stations
 - Several Australian Cities
 - 2007 - 2017
- Features
 - Twenty-two measurements of meteorology
- Outcome Variable
 - Rain Tomorrow

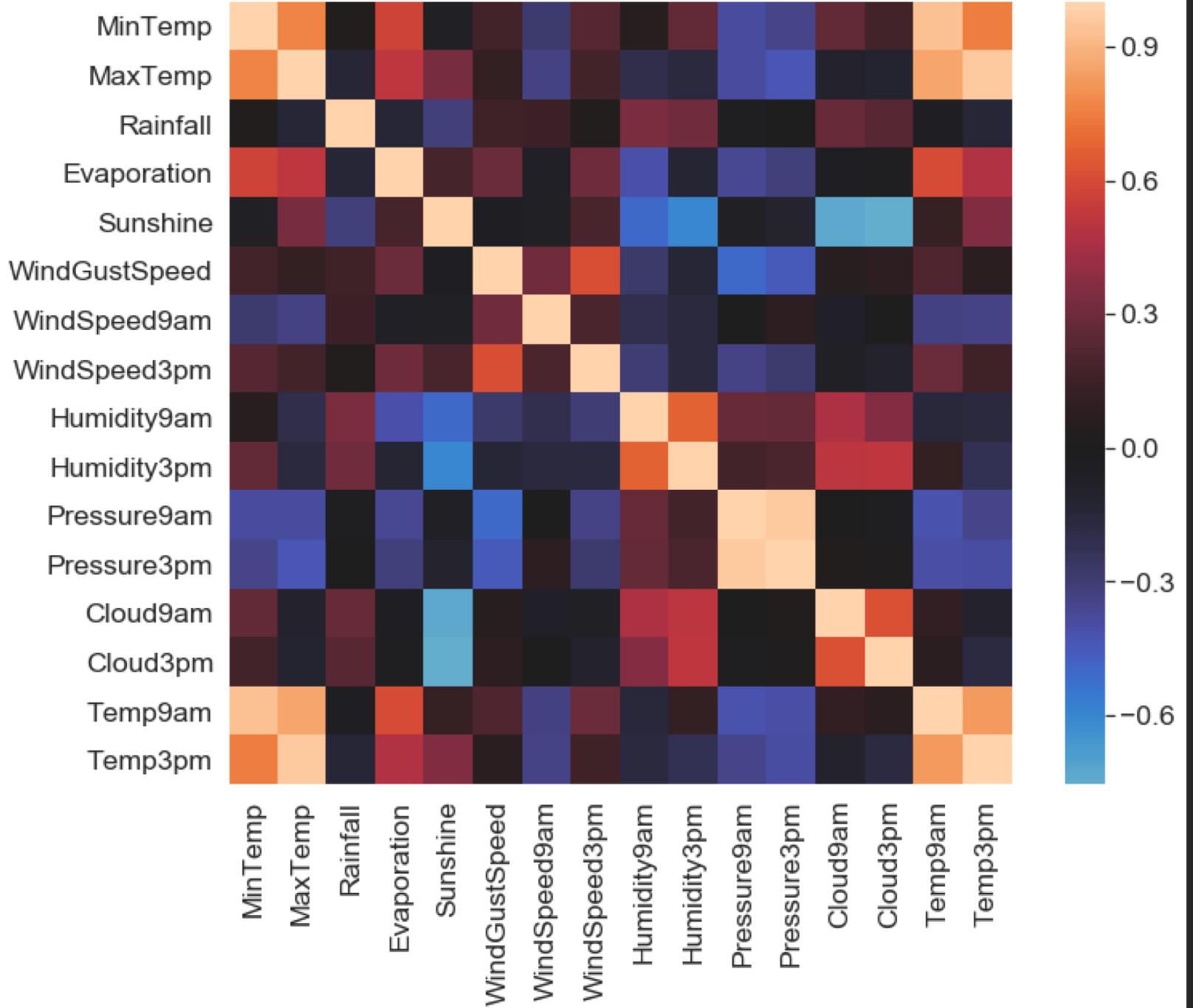


Minority Class Fraction per Year
Our Minority Class Is 'RainTomorrow' = 'Yes'



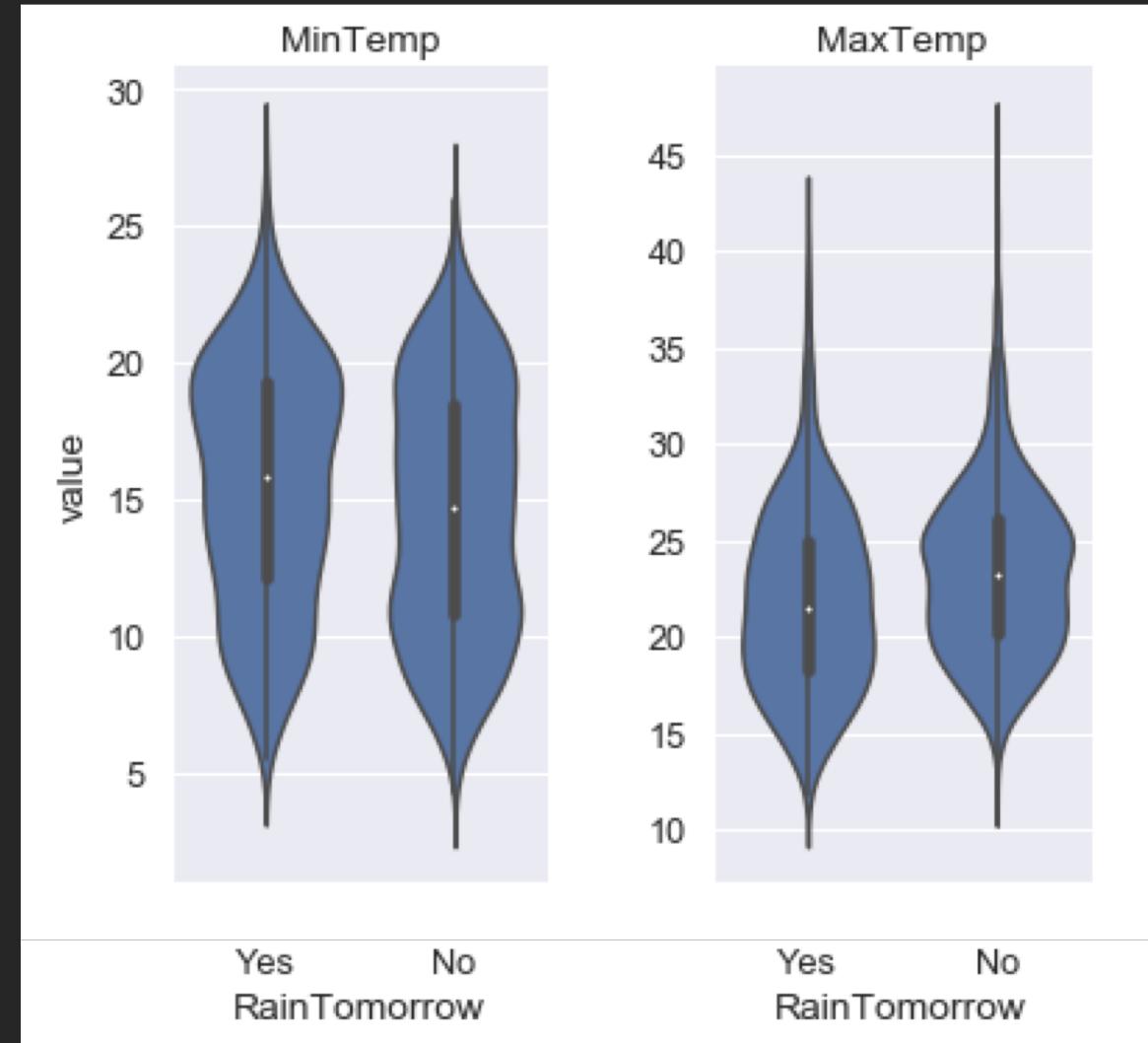


Feature Correlation Heatmap



Feature Creation

- TempRange
- Month

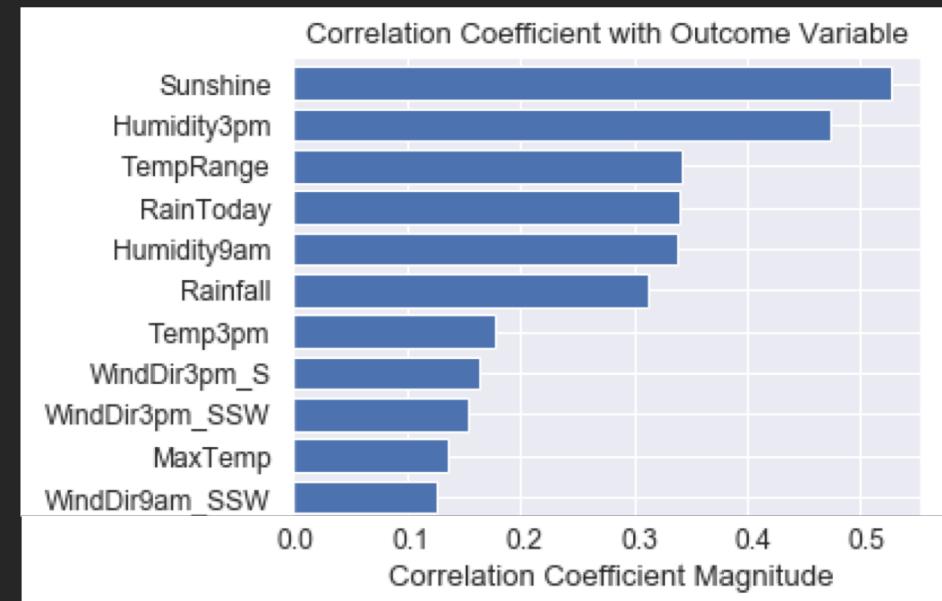


Method for Model Evaluation

- Classifiers
 - Nearest Neighbors
 - Random Forest
 - Logistic Regression
 - Support Vector Machine
 - Gradient Boosting

Method for Model Evaluation

- Vary model parameters
- Vary datasets
 - df_all
 - df_sign
 - df_pca
 - 5 components



Method for Model Evaluation

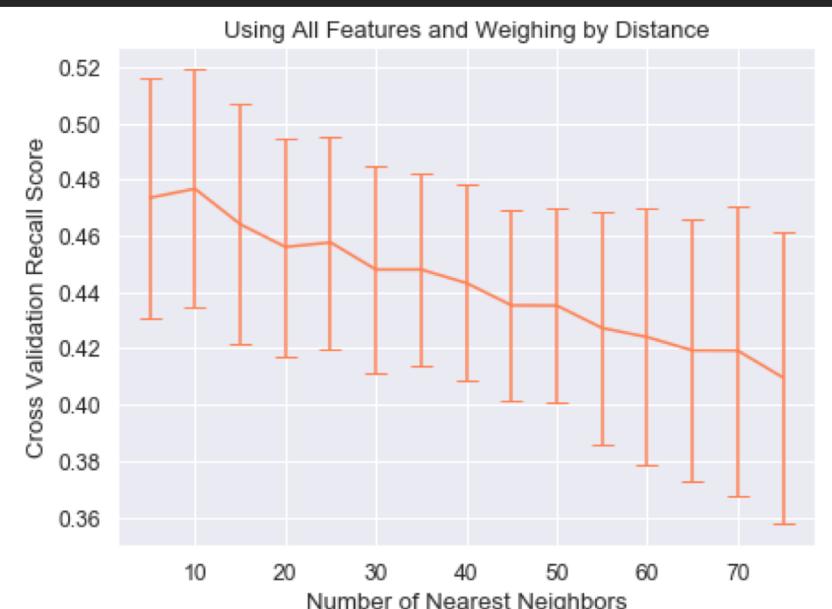
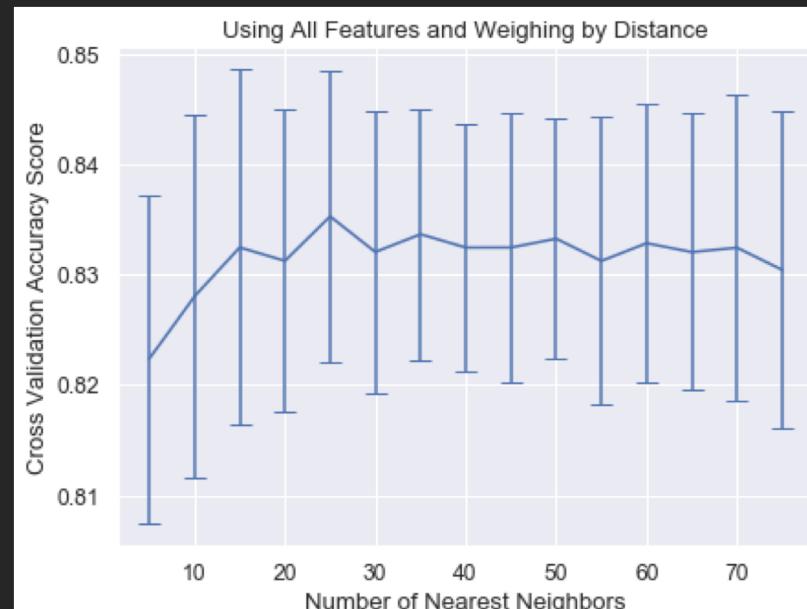
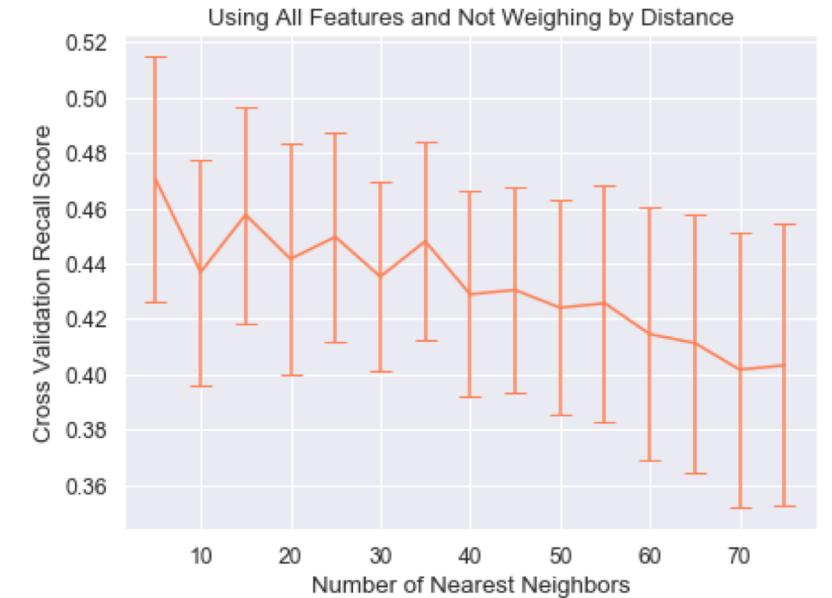
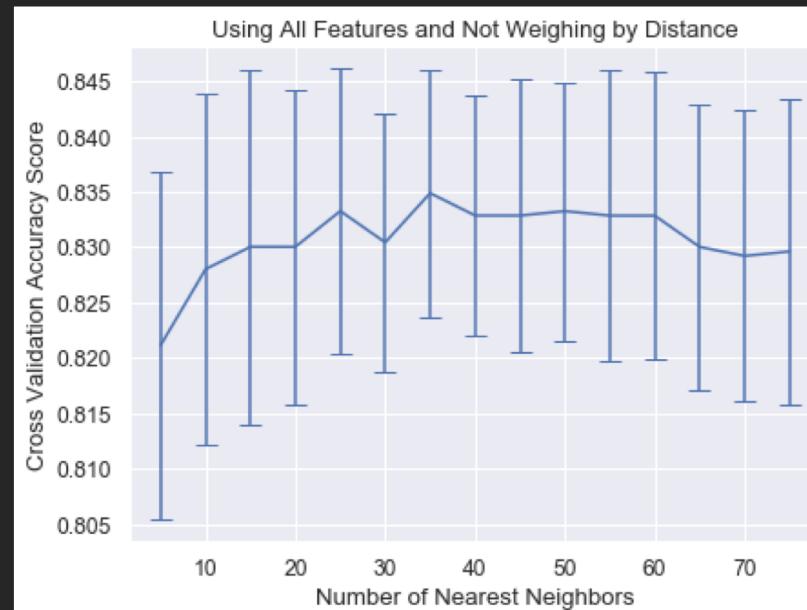
- Winning Criterion
 - Best Mean Cross Validation **Accuracy** Score
 - If tie: Best Mean Cross Validation **Recall** Score

Comparing Model Parameters

How many neighbors?

Weigh by distance or not?

KNN with df_all

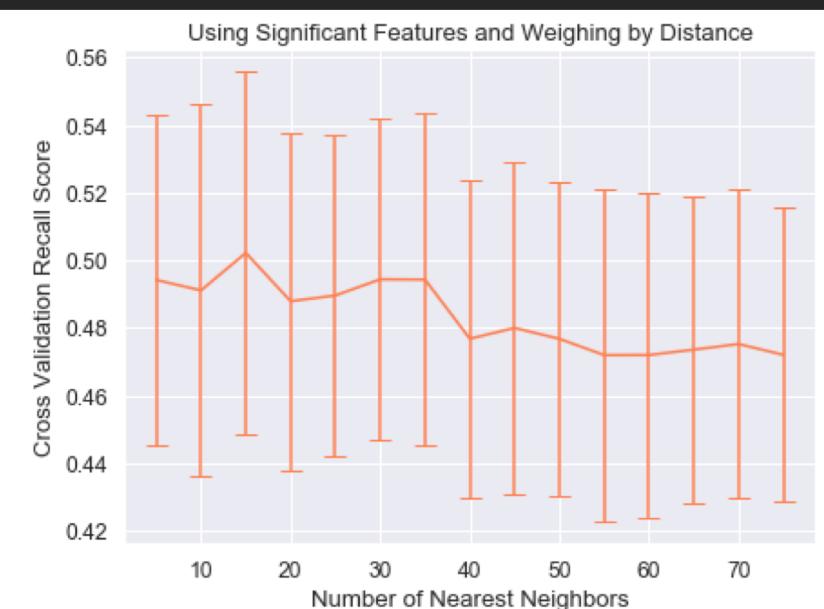
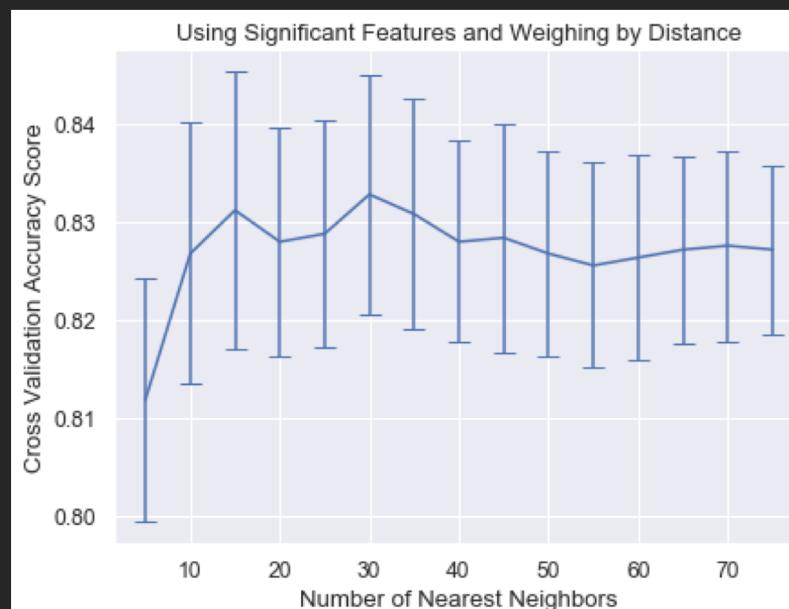
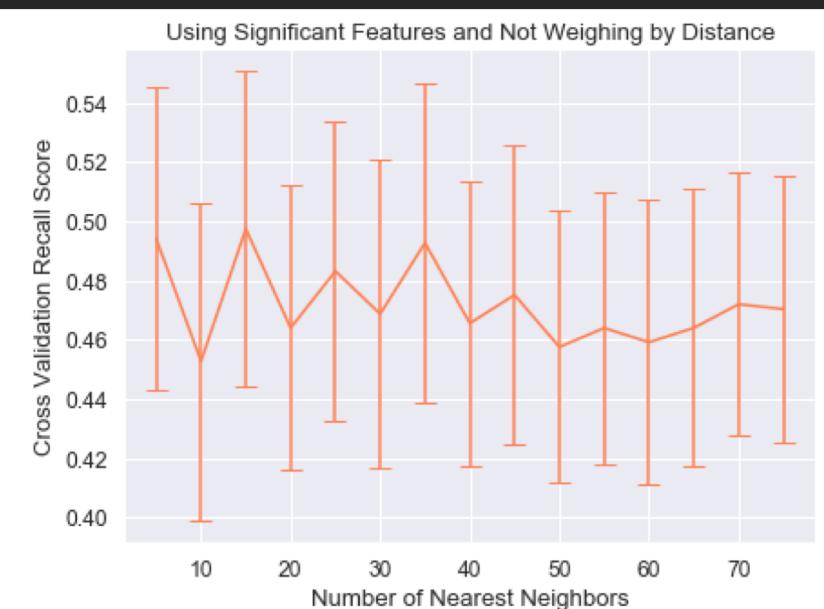
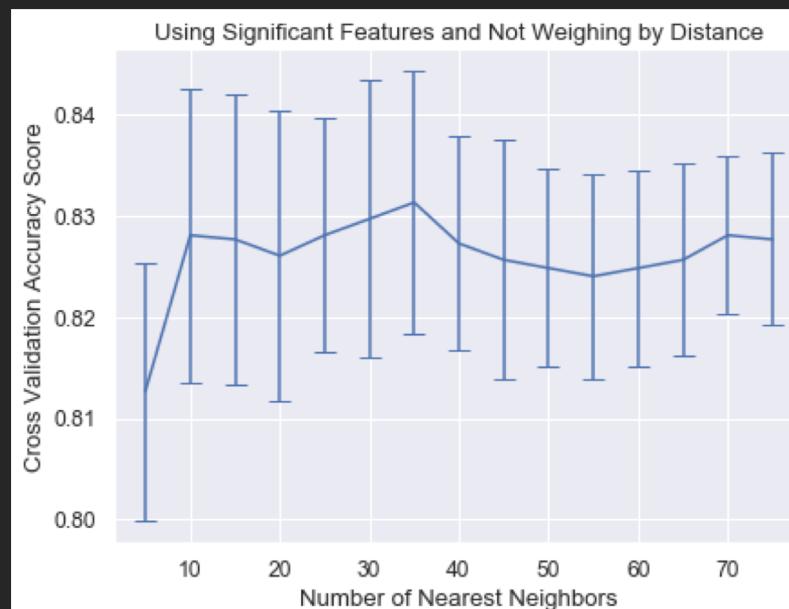


Comparing Model Parameters

Current winner:

- df_all
- 25 neighbors
- with weights
- Accuracy = 0.835
- Recall = 0.46

KNN with df_sign

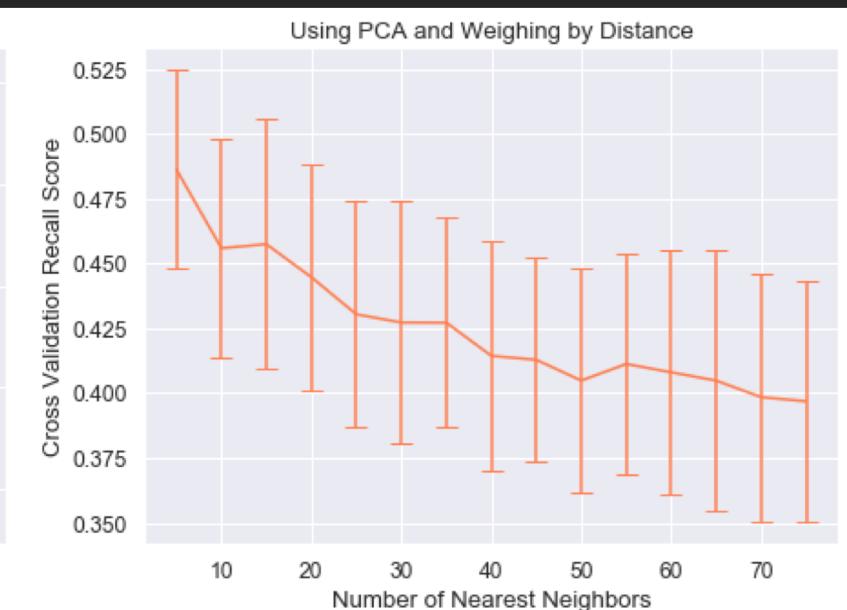
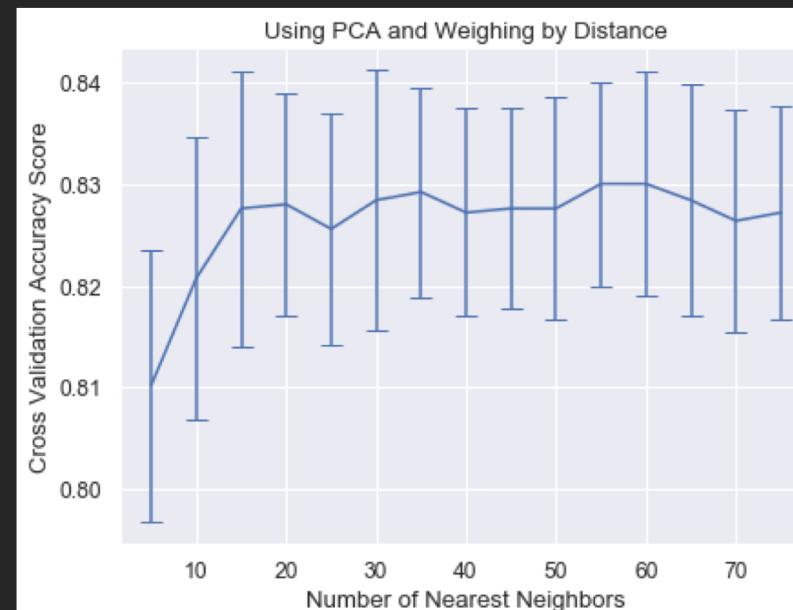
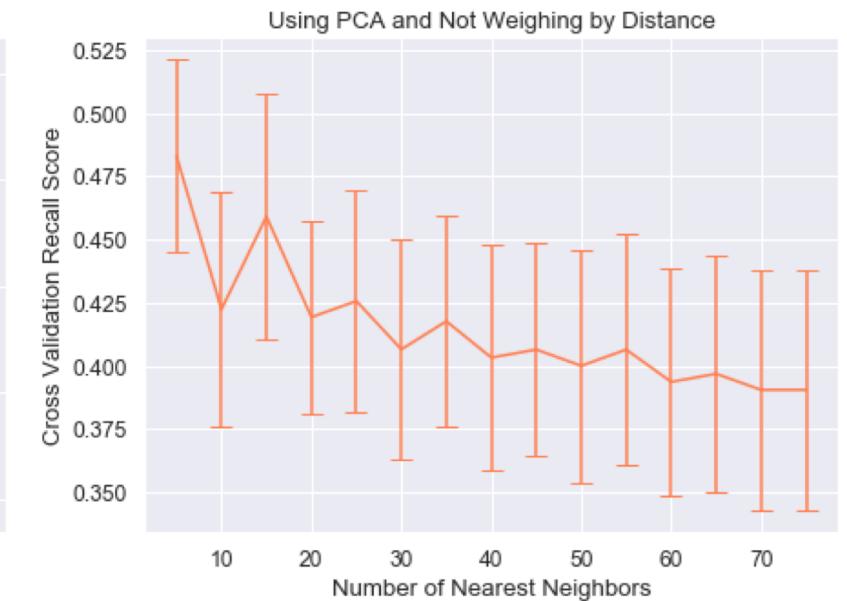
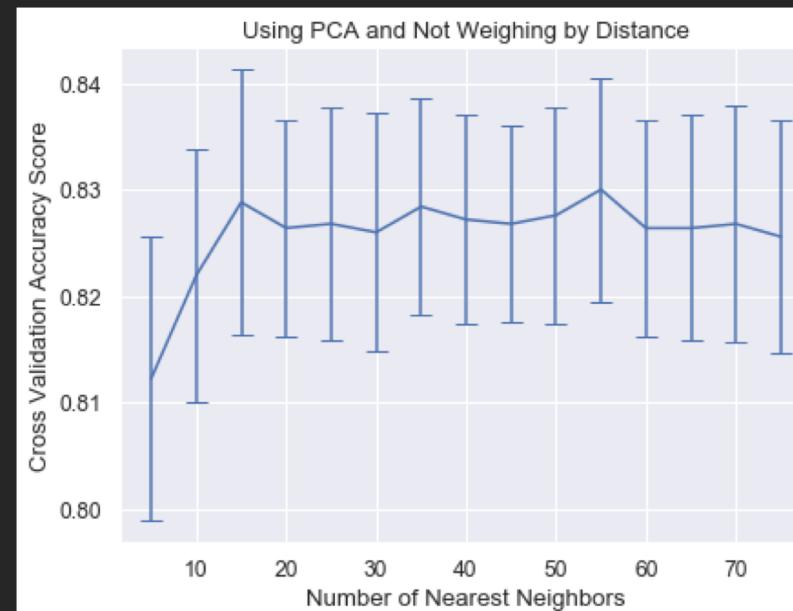


Comparing Model Parameters

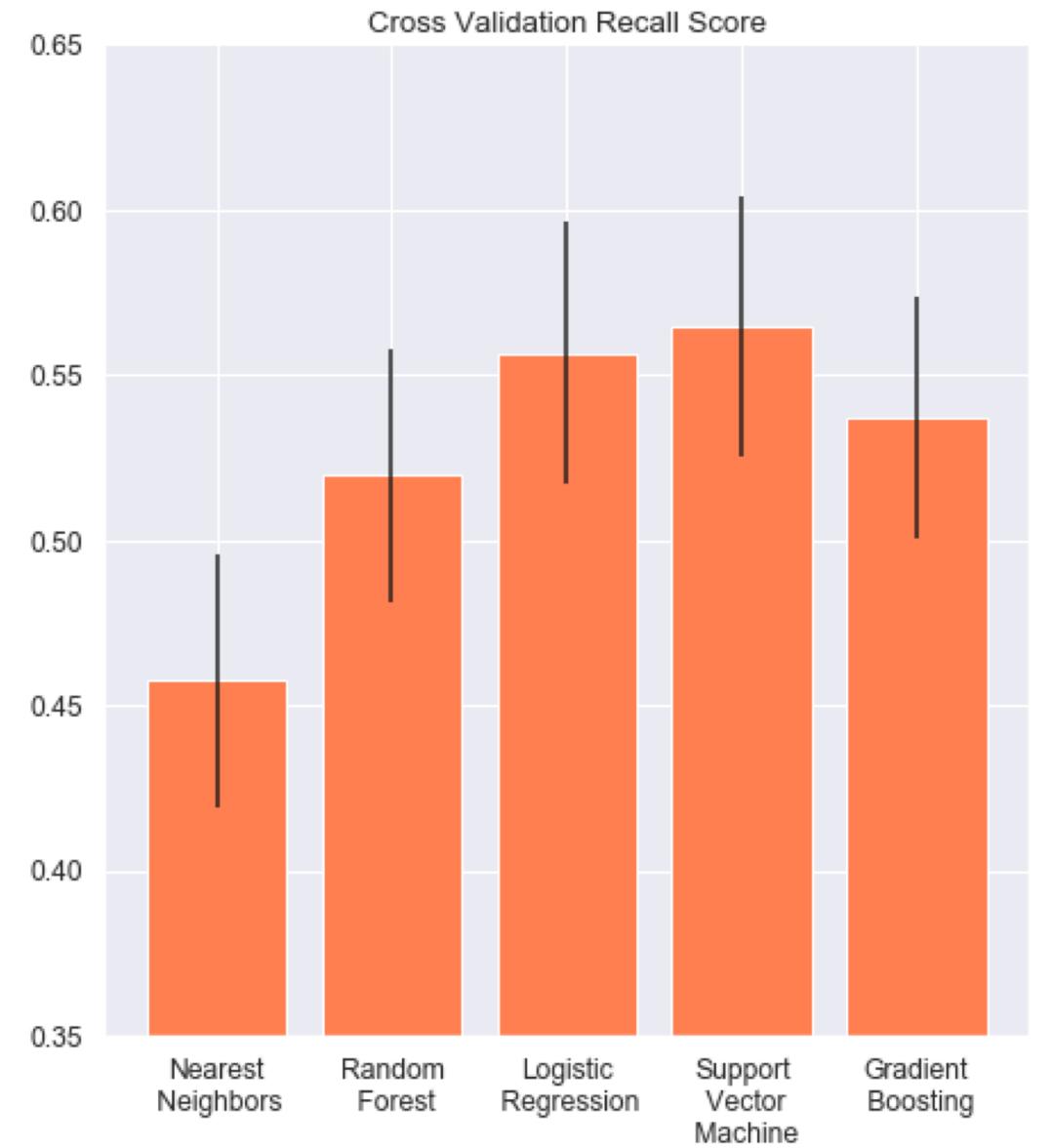
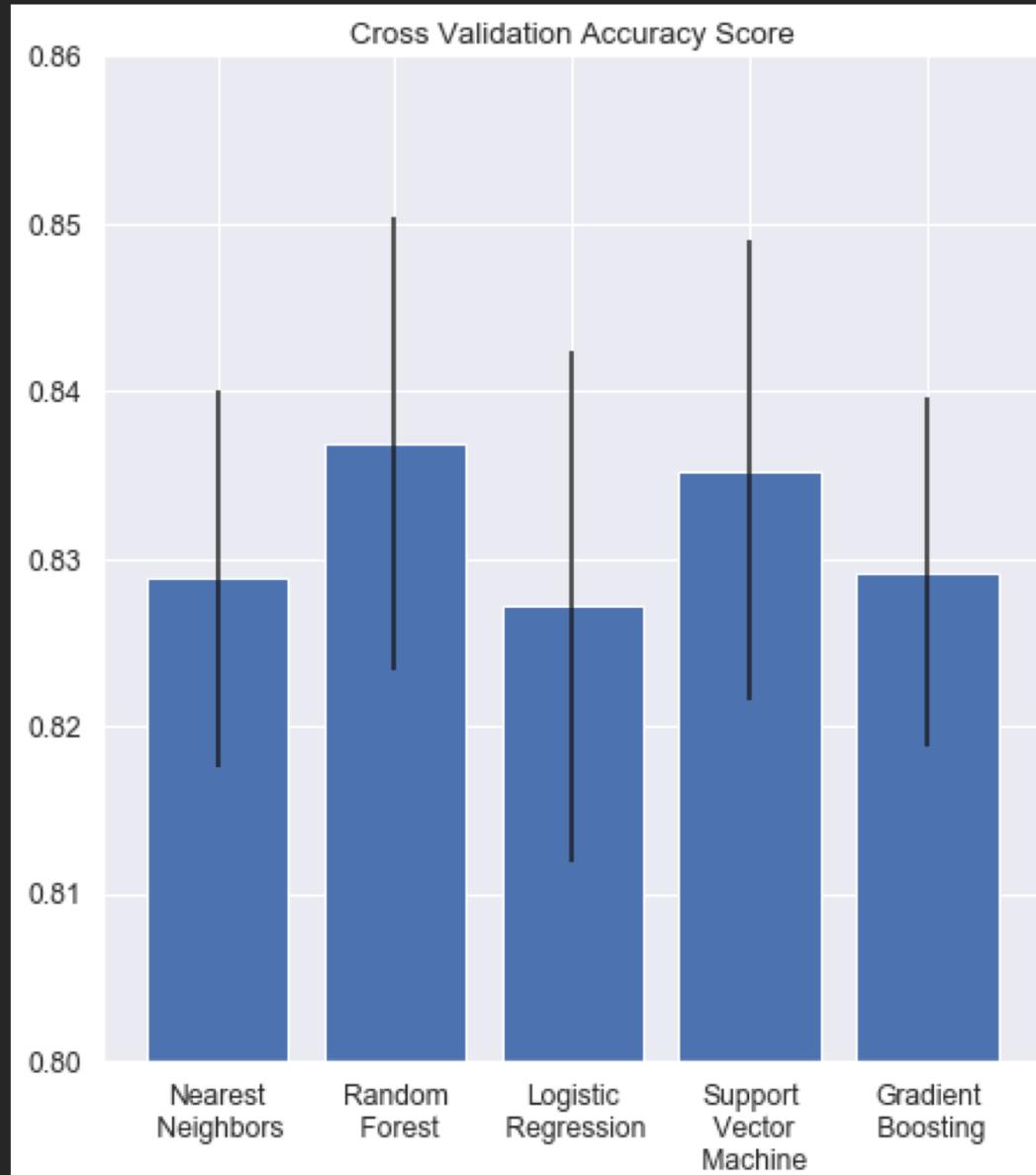
Current winner:

- df_sign
- 30 neighbors
- with weights
- Accuracy = 0.834
- Recall = 0.49

KNN with df_pca



Best of the Five Models



Our Winner: Support Vector Machine

Prediction of Testing Data (2011 and 2014)

Confusion Matrix

`[[462 39]`

`[77 100]]`

Scores

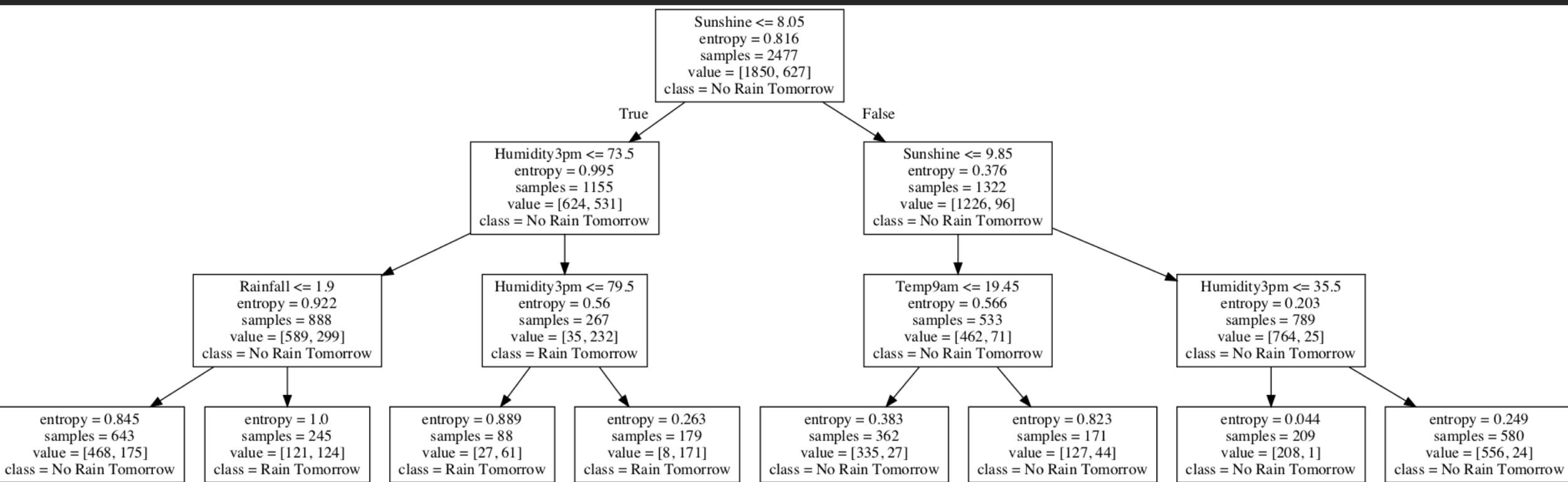
Accuracy: 82.9%

Recall: 56.5%

Specificity: 92.2%

Model Shortcomings – Feature Importance

Decision Tree



Random Forest Feature Importance

