



Brian McCabe

Baseball Classification

Purpose

Classify Baseball Players Feilding Positions

- Batting and **Fielding** Statistics

Comparing Classification Algorithms

- K-Nearest Neighbors Classifier
- SGD Classifier
- Gaussian Naïve Bayes
- Decision Tree Classifier
- Random Forest Classifier

Significance



What attributes are most important in predicting Player position?



Can organizations make a team around Players predicted positions?

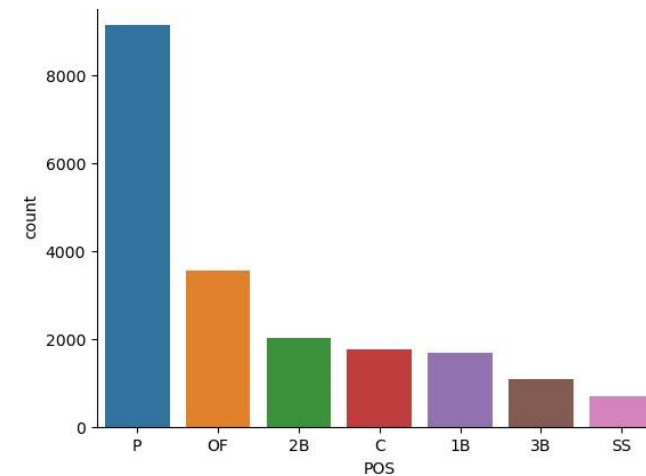
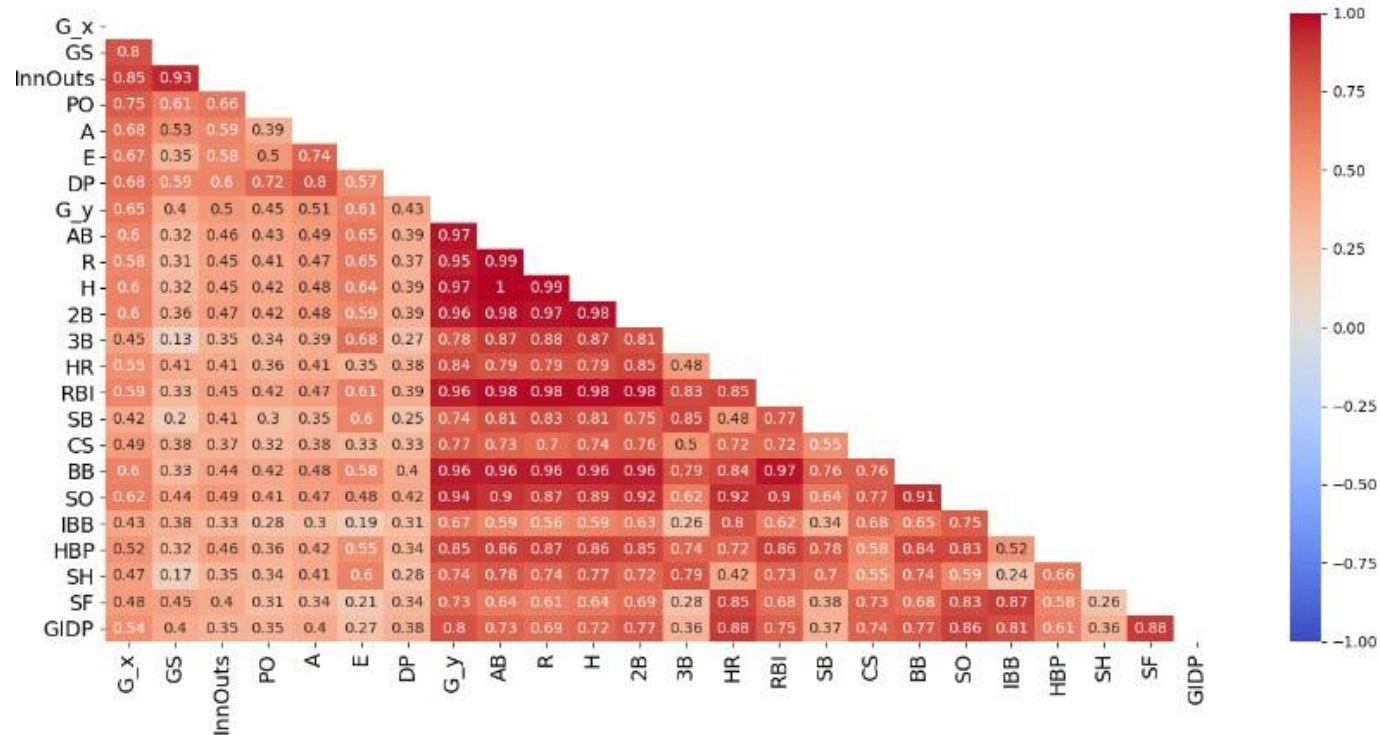


Which position is the hardest to predict why are they difficult to predict?

Data Description

- 24 input variables
- Batting input variables
 - At Bats
 - Runs
 - Hits
 - Home Runs
 - Runs Batted In
 - Walks
 - Strike Outs
- Fielding input variables
 - Games Started
 - Put Outs
 - Assist
 - Errors
 - Inn Outs
 - Double Plays

- Out-Put variables: Player Position
 - Catcher (C)
 - Pitcher (P)
 - 1st Baseman (1B)
 - 2nd Baseman (2B)
 - 3rd Baseman (3B)
 - Short Stop (SS)
 - Out Fielder (OF)



BRIAN MCCABE

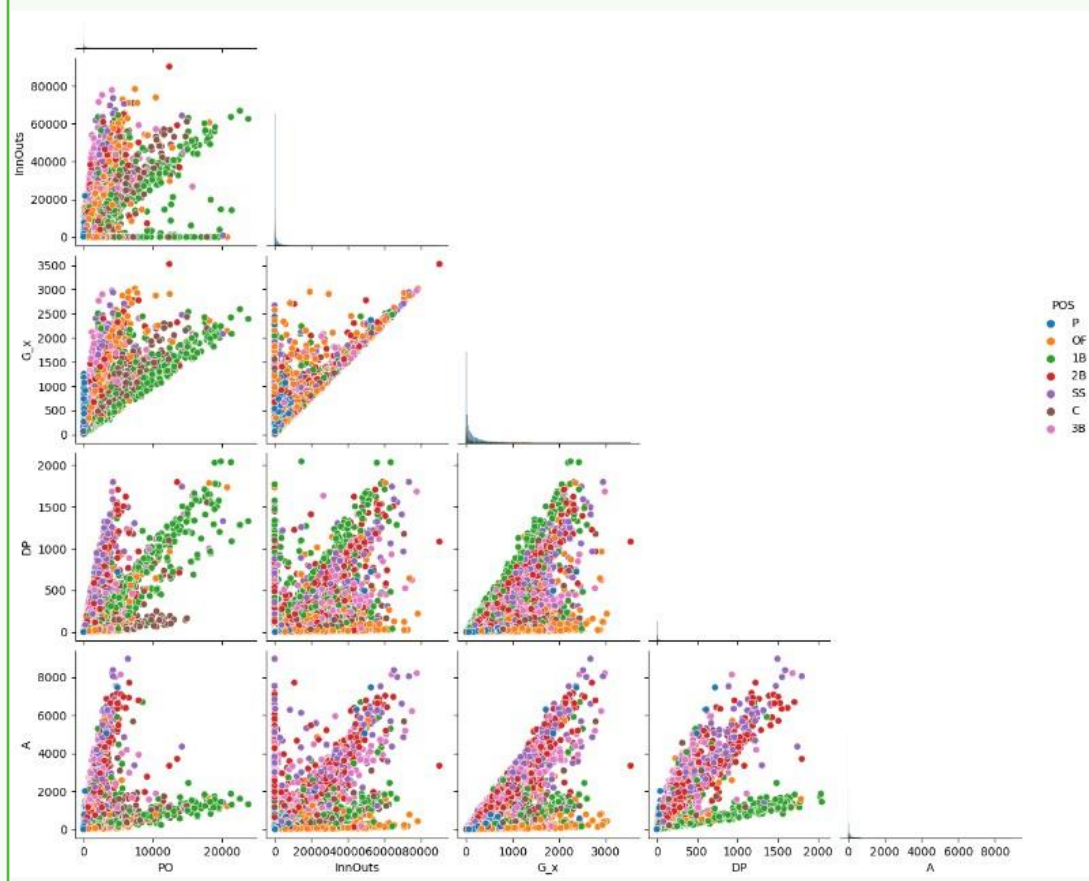
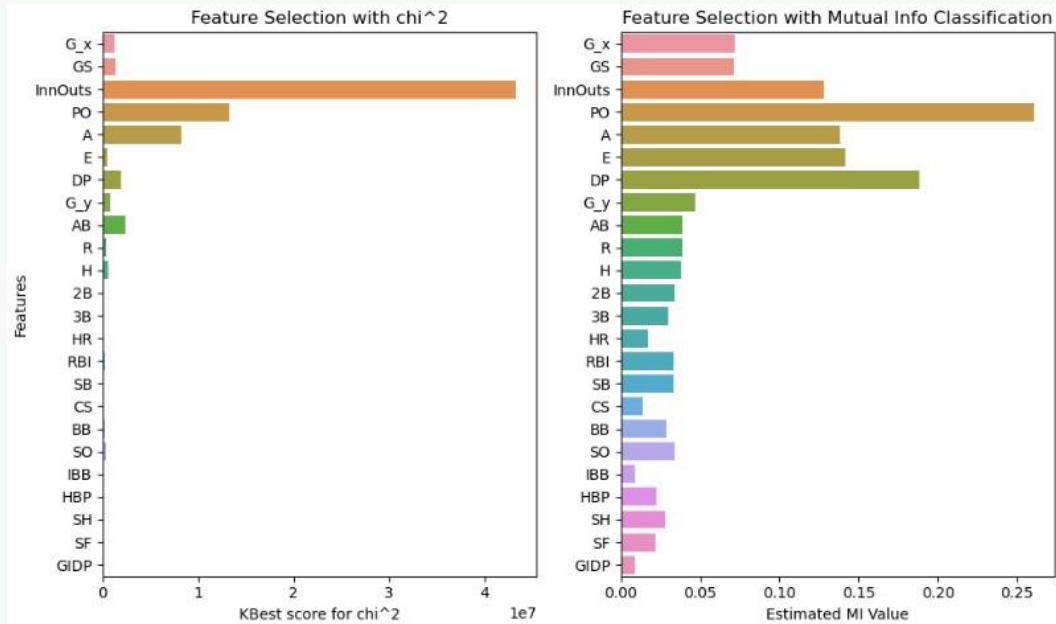
Exploratory

10/12/2022

5

Exploratory

BRIAN MCCABE



10/12/2022

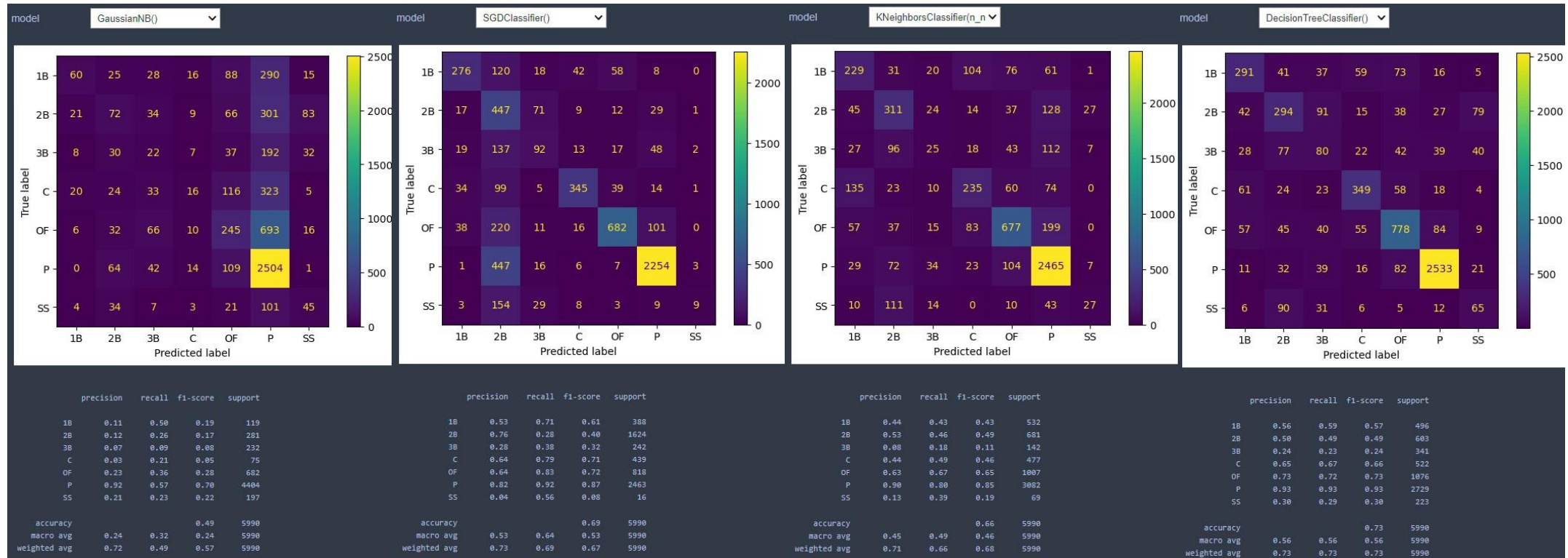
6

Data Preparation

```
aggregation_functions_f = {'POS': 'first', 'G': 'sum', 'GS': 'sum', 'InnOuts': 'sum',  
                           'PO': 'sum', 'A': 'sum', 'E': 'sum', 'DP': 'sum'}  
combined_f = fielding.groupby(fielding['playerID']).aggregate(aggregation_functions_f)  
aggregation_functions_b = {'G': 'sum', 'AB': 'sum', 'R': 'sum', 'H': 'sum',  
                           '2B': 'sum', '3B': 'sum', 'HR': 'sum', 'RBI': 'sum', 'SB': 'sum',  
                           'CS': 'sum', 'BB': 'sum', 'SO': 'sum', 'IBB': 'sum', 'HBP': 'sum',  
                           'SH': 'sum', 'SF': 'sum', 'GIDP': 'sum'}  
combined_b = batting.groupby(fielding['playerID']).aggregate(aggregation_functions_b)  
merged_df = pd.merge(combined_f, combined_b, on='playerID', how='outer')  
merged_df = merged_df.fillna(0)
```

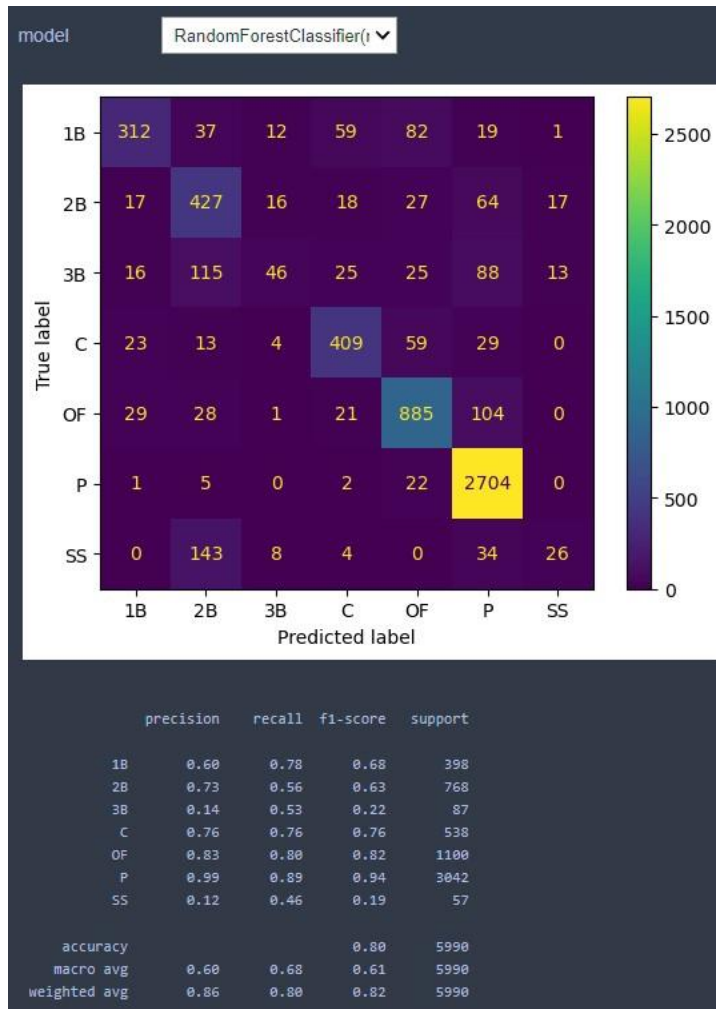
```
y = merged_df['POS']  
X = merged_df.drop('POS', axis=1)  
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.30)
```

- Combined records that had same player ID
- Aggregate Function
- Merged batting and Fielding on player ID
- Filled NANs with Zeros
- 70/30 Train Test Split



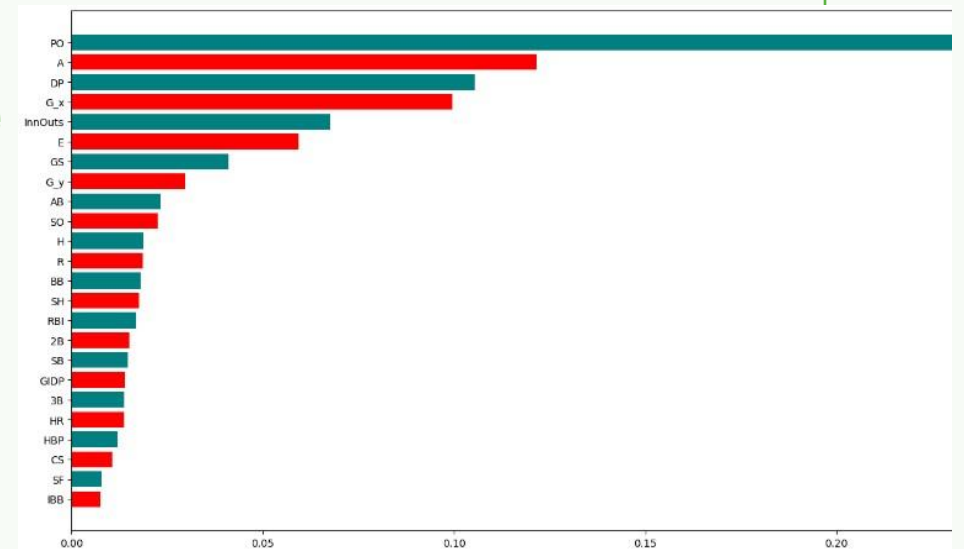
Model Analysis

Model	Data	Accuracy	Precision	Recall	F1	Train Acc
GaussianNB()	features_df	0.494825	0.495	0.703	0.495	0.493202
SGDClassifier()	features_df	0.610518	0.611	0.781	0.611	0.610634
KNeighborsClassifier(n_neighbors=4)	features_df	0.662604	0.663	0.814	0.663	0.780163
DecisionTreeClassifier()	features_df	0.729048	0.729	0.854	0.729	0.999356
RandomForestClassifier()	features_df	0.796995	0.797	0.893	0.797	0.999356
RandomForestClassifier(max_features='log2', min_samples_split=5, n_estimators=450, n_jobs=-1, random_state=4)	features_df	0.802838	0.803	0.896	0.803	0.98626



Conclusion

- With 80% accuracy the Random Forest Classifier performed the best
- Feature Importance
 - Put outs
 - Assist
 - Double Plays
 - Games
 - Inn Outs
- Hardest Positions to predict
 - Short Stop
 - 3rd Baseman





Questions



Thank You