# RUTGERS
UNIVERSITY | NEW BRUNSWICK

# Cricket Players (Cricketers) Data Analysis

**CS 439 - Introduction to Data Science**
**Course Project**
Professor: Gerard de Melo

**Name:** Smitkumar Patel
**NetID:** Shp109
**Group:** 06
**Recitation:** 02
**Email:** shp109@rutgers.edu

# DATA SETS

- **Cricket Players Data from ESPN:**
  The data set contains information about the cricket players all over the world that played Tests, ODI, List A, First-class, T20I, and T20. There are a total of 90308 observations (rows) and 176 variables(columns).
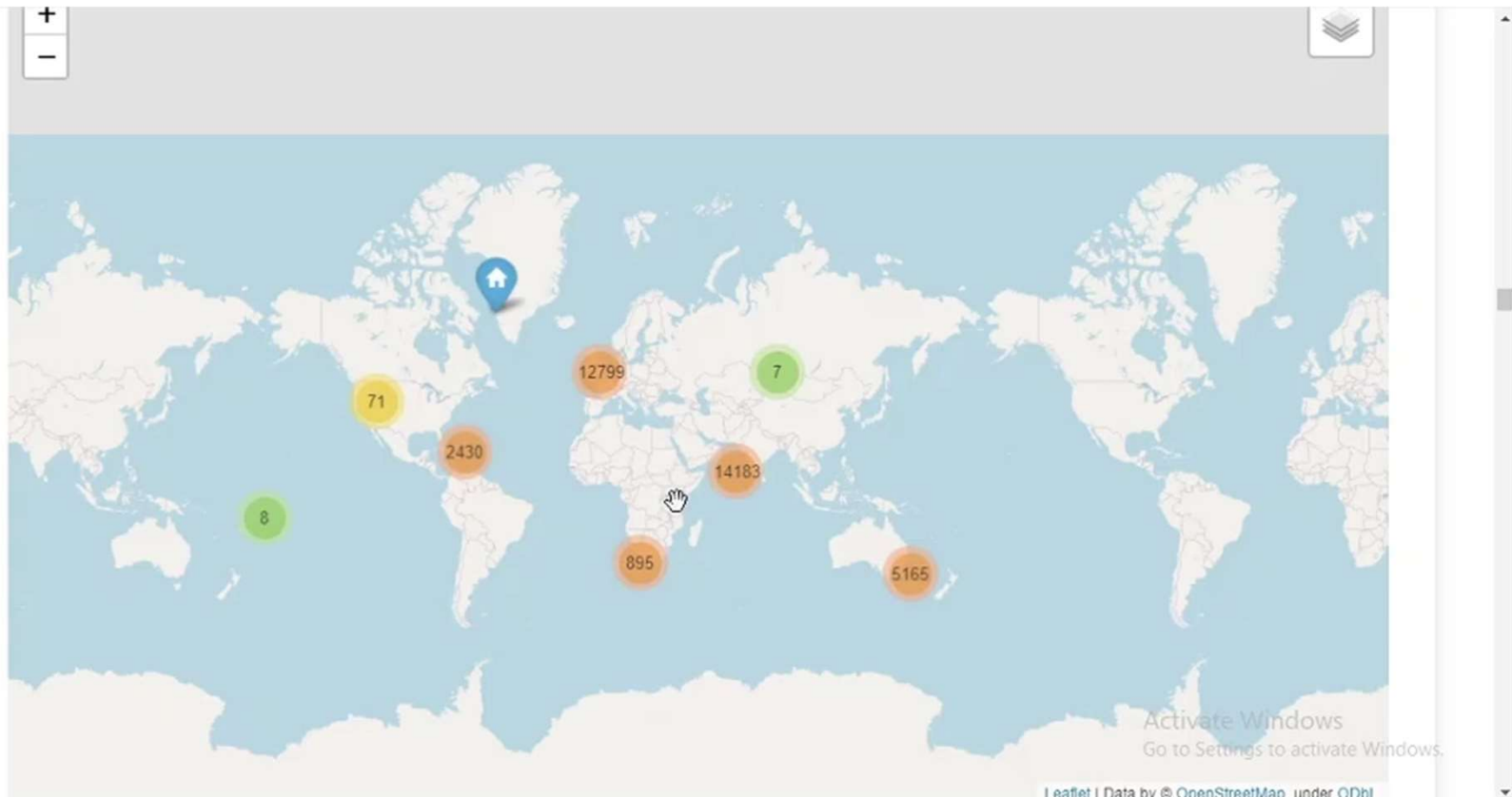  URL: https://data.world/raghav333/cricket-players-espn

- **IPL Cricket Players Data:**
  The data set contains information about the cricket players that played IPL (Indian Premier League). There are a total of 497 observations (rows) and 7 variables(columns). URL: https://data.world/raghu543/ipl-data-till-2017/workspace/file?filename=Player.csv

**Note**: I left join 'Cricket Players Data from ESPN' and 'IPL Cricket Players Data'. That means all records from 'Cricket Players Data from ESPN', and matched records from the 'IPL Cricket Players Data' on 'Country', 'Birthdate', and 'Last_Name'. In addition, I created new column that contains binary labels 'Play' (player played IPL) and 'Not Play' (player did not play IPL).
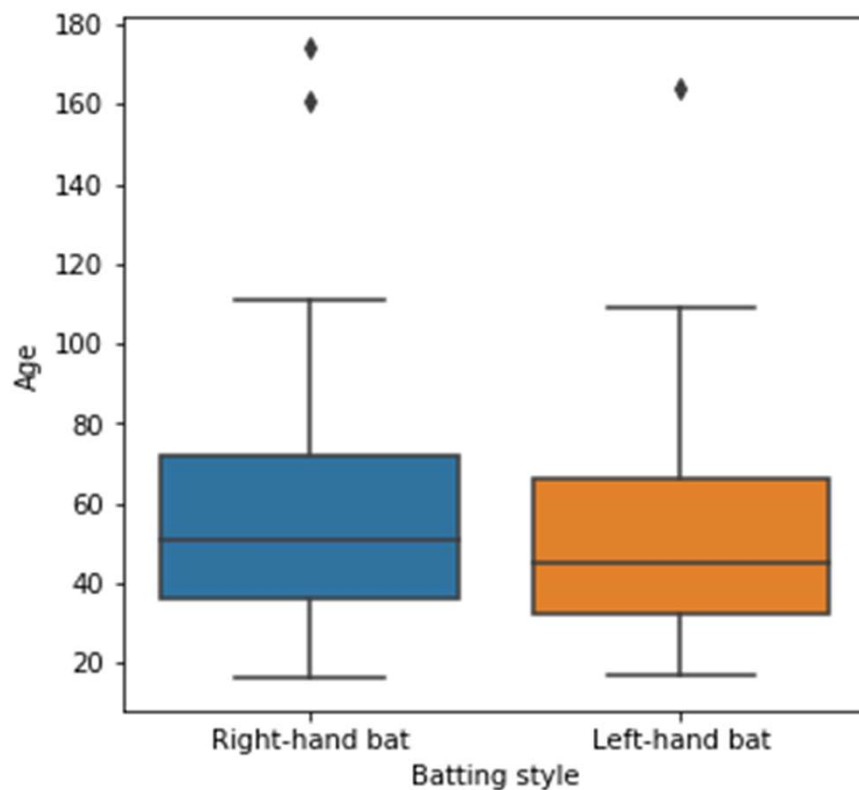
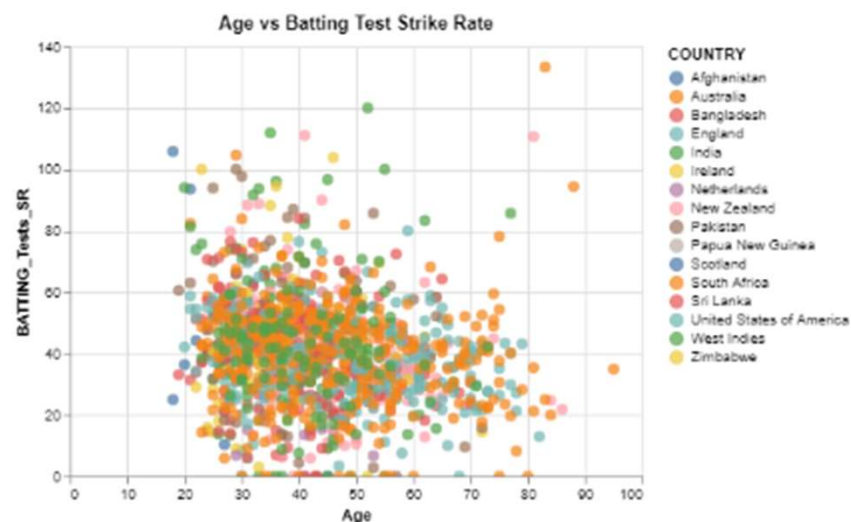# Number of players from country/city/county

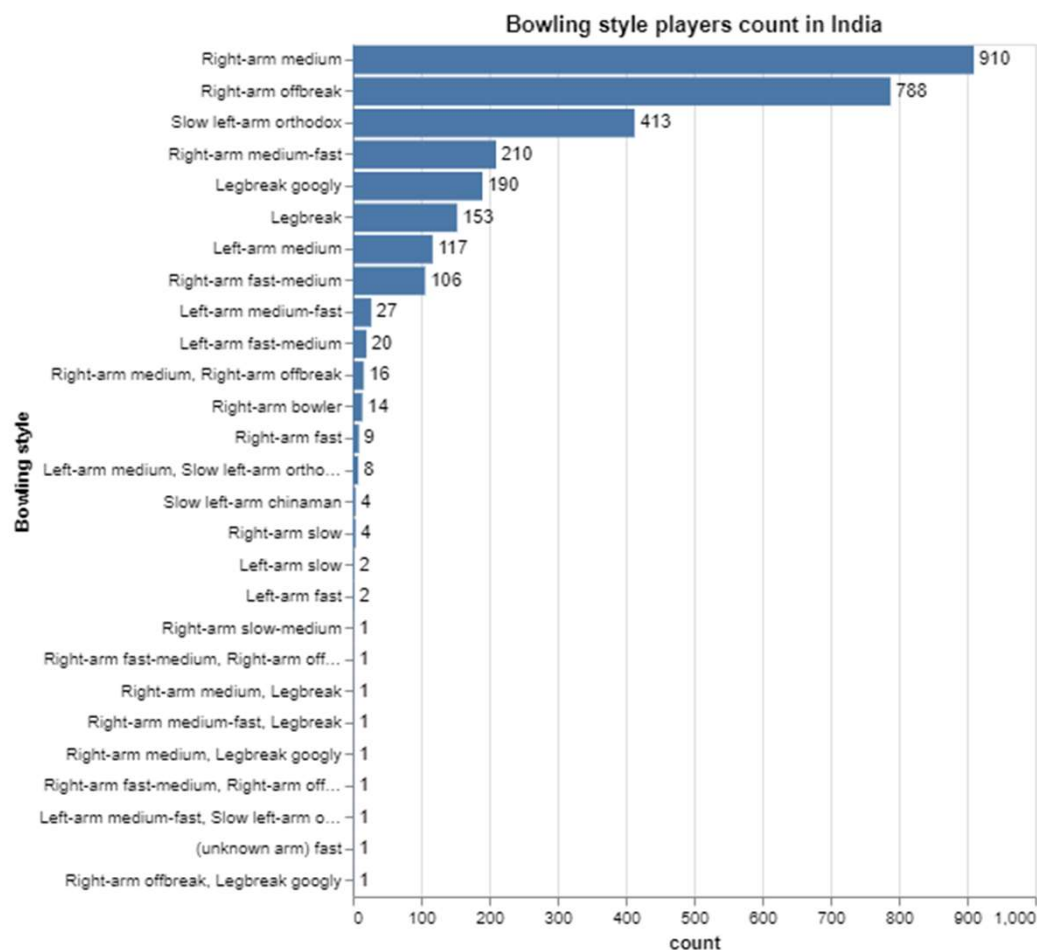# Top 10 batsmen and bowlers in Test, ODI, and T20
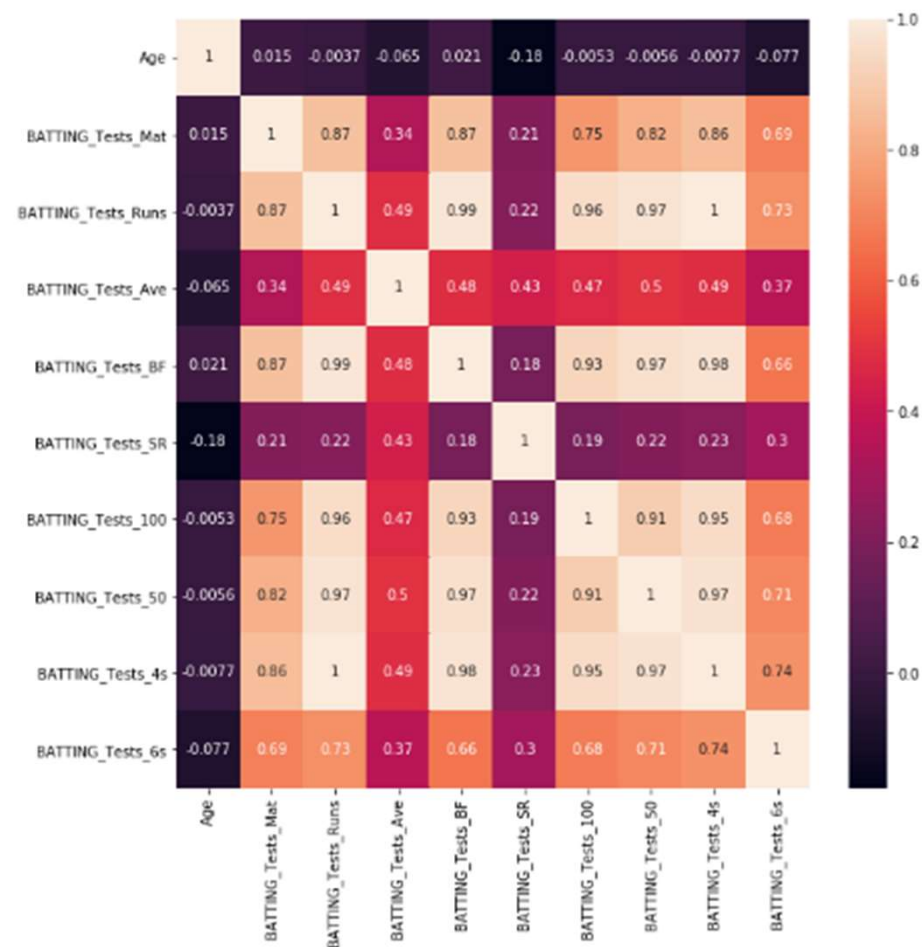
## Age of players for various Batting style:

## Age vs Batting Test Strike Rate (Test Run/Ball faced in Test):

# Bowling style players count in India



Bowling style players count in India

| Bowling style | count |
|---|---|
| Right-arm medium | 910 |
| Right-arm offbreak | 788 |
| Slow left-arm orthodox | 413 |
| Right-arm medium-fast | 210 |
| Legbreak googly | 190 |
| Legbreak | 153 |
| Left-arm medium | 117 |
| Right-arm fast-medium | 106 |
| Left-arm medium-fast | 27 |
| Left-arm fast-medium | 20 |
| Right-arm medium, Right-arm offbreak | 16 |
| Right-arm bowler | 14 |
| Right-arm fast | 9 |
| Left-arm medium, Slow left-arm ortho... | 8 |
| Slow left-arm chinaman | 4 |
| Right-arm slow | 4 |
| Left-arm slow | 2 |
| Left-arm fast | 2 |
| Right-arm slow-medium | 1 |
| Right-arm fast-medium, Right-arm off... | 1 |
| Right-arm medium, Legbreak | 1 |
| Right-arm medium-fast, Legbreak | 1 |
| Right-arm medium, Legbreak googly | 1 |
| Right-arm fast-medium, Right-arm off... | 1 |
| Left-arm medium-fast, Slow left-arm o... | 1 |
| (unknown arm) fast | 1 |
| Right-arm offbreak, Legbreak googly | 1 |

# Test Batting Data Correlations

# Machine learning algorithms to classify IPL status to be able to know player will play IPL or not
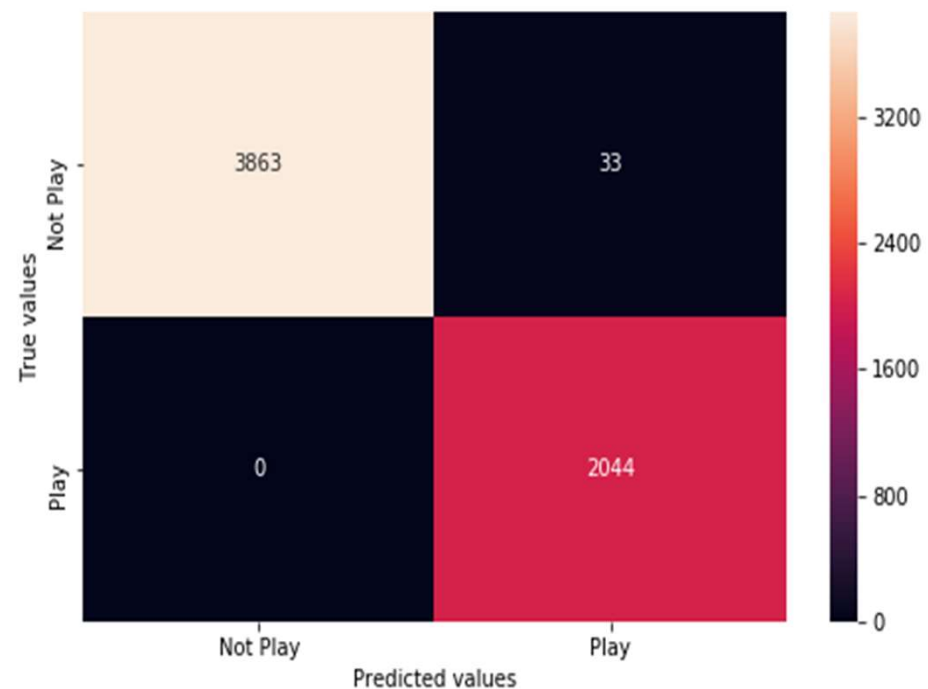
**Note:** All the accuracy and area under roc curve are acquire by 10-fold cross validation on training dataset.

| Model | Average Accuracy | Average Area under Roc Curve | Note |
|---|---|---|---|
| LogisticRegression | 0.966788161 | 0.990200503 | |
| Support Vector Classification | -- | -- | Took too much time to fit training data |
| K-nearest neighbors | 0.981100088 | 0.994323267 | |
| Decision Tree | 0.989266064 | 0.994361751 | |
| Random Forest | 0.99040264 | 0.999867555 | |
| Stochastic Gradient Descent | 0.899986776 | 0.964317684 | |
| Gaussian Processes classification | -- | -- | Took too much time to fit training data |
| Naive Bayes | 0.821736747 | 0.937844178 | |
| Adaboost | 0.9989056 | 1 | |
| Gradient Boosting | 0.979584458 | 0.995913181 | |
| Histogram-Based Gradient Boosting | 0.99360147 | 0.999849177 | |
| Neural network models | 0.939638222 | 0.977137591 | |

# Voting Classifier model on Test Data

I used voting classifier to use majority vote to predict the IPL status because Random Forest, Adaboost, and Histogram-Based Gradient Boosting model have better accuracy and area under ROC curve.

# Acknowledgments