

# **Attrition ... Why and When!?**

Marty Fromuth

GW Bootcamp

Jun 2020– Jan 2021



# Overview

## The What

*What is the topic?*

*What are we trying to learn?*

*What is the data?*

## The Why

*Why did we choose this topic?*

*Why did we choose this data?*

## The How

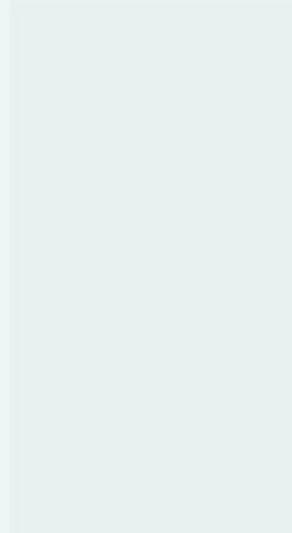
*How are we conducting data exploration?*

*How are we analyzing the data?*

*How are we storing the data?*

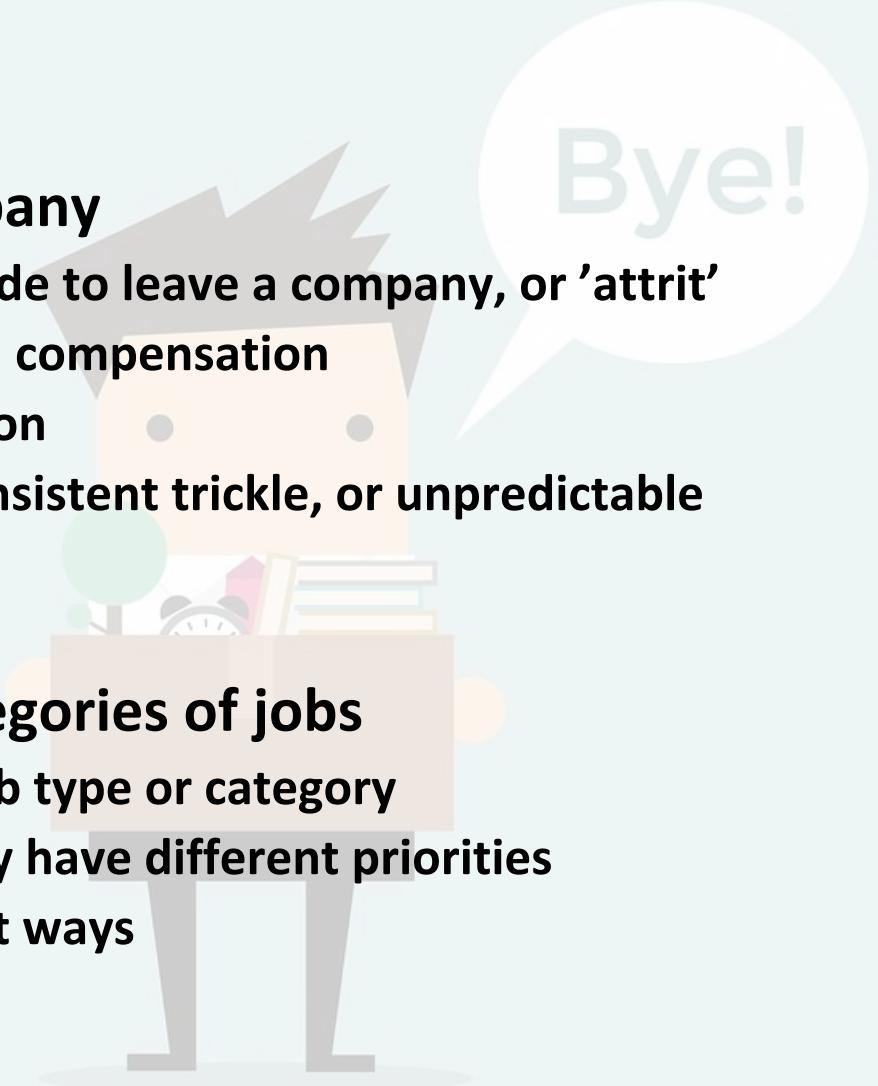
*How are we displaying our work?*

## The Results



# The What: What is the topic?

- Primary Topic: Attrition, aka people leaving a company
  - Companies want to know when and why employees decide to leave a company, or 'attrit'
  - Employers maintain information on job performance and compensation
  - Companies also put out anonymous surveys on satisfaction
  - A company's attrition 'cycle' can manifest in waves, a consistent trickle, or unpredictable cliffs
- Secondary Topic: Attrition in different types or categories of jobs
  - Most companies have smaller 'communities' based on job type or category
  - STEM, soft-skills and/or HR, and leadership positions may have different priorities
  - Like overall employees, attrition can manifest in different ways



# The What: What are we trying to learn?

- Primary questions:

- What are they key features that predict attrition?
- When and how many employees do we predict will leave?
- Do those answers change based on the job category?

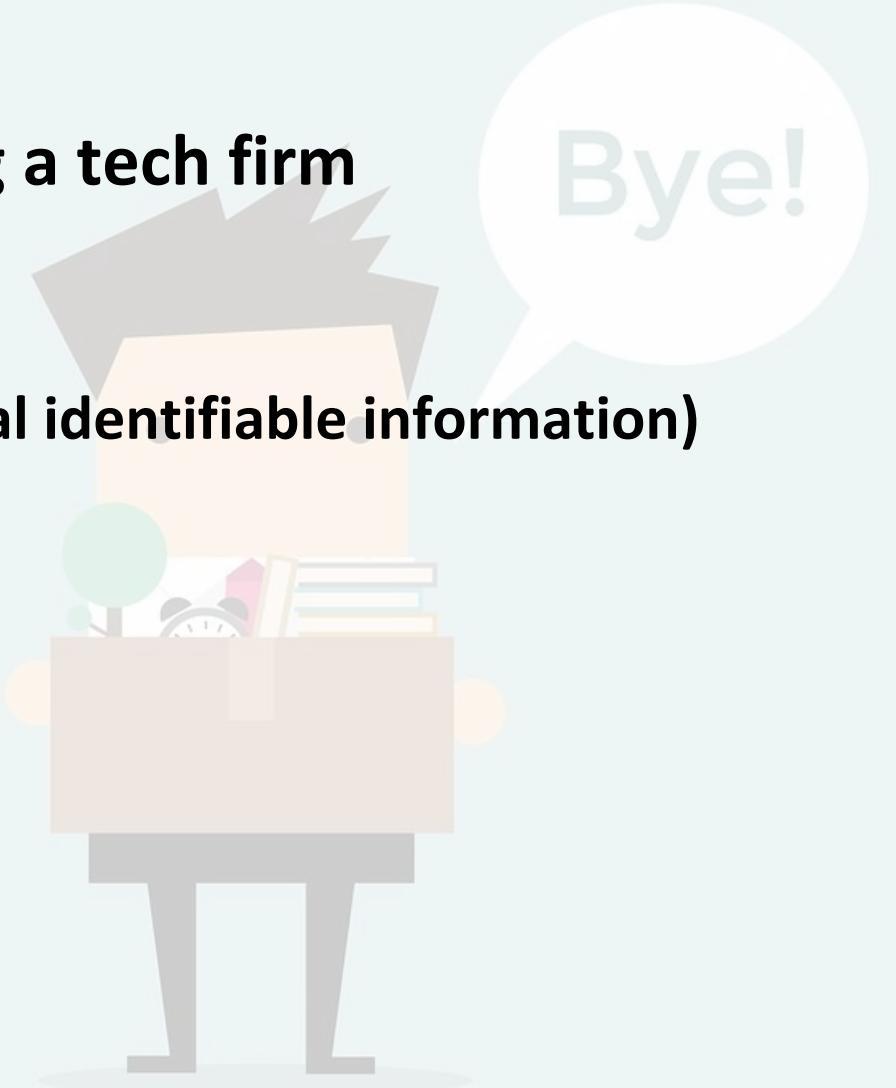
- Additional insights:

- What, if any, correlation is there between different features?
- Are there any major differences in features between the job categories?



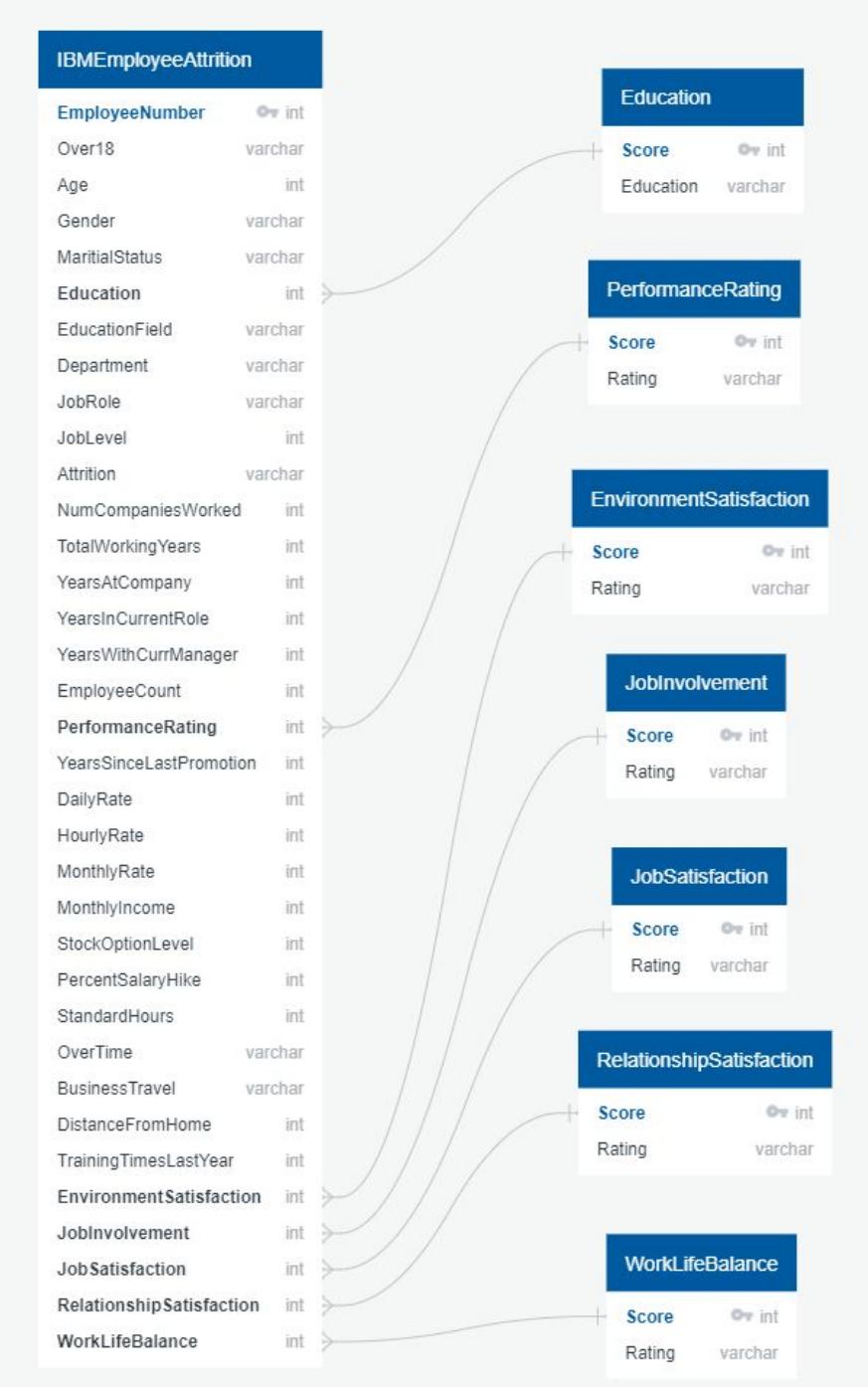
# The What: What is the data?

- Synthetic, anonymous HR data representing a tech firm
- Multiple fields/features in key areas:
  - Anonymized employee information (no personal identifiable information)
  - Job history
  - Education
  - Job role and department
  - Salary
  - Satisfaction scores
- 1 primary table with eight reference tables



# The What: What is the data?

- 1470 rows / unique employees
- Key = 'EmployeeNumber'
- Dropped = 'EmployeeCount', 'StandardHours', 'Over18', 'MonthlyIncome', 'HourlyRate', 'DailyRate', 'Department', 'TotalWorkingYears
- Created new column = 'JobCategory'
  - Leadership
  - Non-Tech
  - Tech
- Left join main table w/ satisfaction ratings



# The What: What is the data?

## **Combined:**

attrition\_combined\_text (Postgres)  
df\_attrition\_encoded (Python)

## **Non-Tech:**

attrition\_nontech\_text (Postgres)  
df\_attrition\_nontech\_encoded (Python)

## **Tech:**

attrition\_tech\_text (Postgres)  
df\_attrition\_tech\_encoded (Python)

## **Leadership:**

attrition\_ldrshp\_text (Postgres)  
df\_attrition\_ldrshp\_encoded (Python)

- **Created two main dataframes**

- Encoded dataframe – used for machine learning model building and inference, and feature analysis
- Visualization dataframe – used for the dashboard and to visualize data for the end user
  - Created tables in Postgres using SQL join & ‘INTO’
  - Stored these in new Postgres database

- **Created three sub-dataframes for each of the main dataframes**

- Leadership
- Tech
- Non-tech

- **Build & stored 8 dataframes in total**

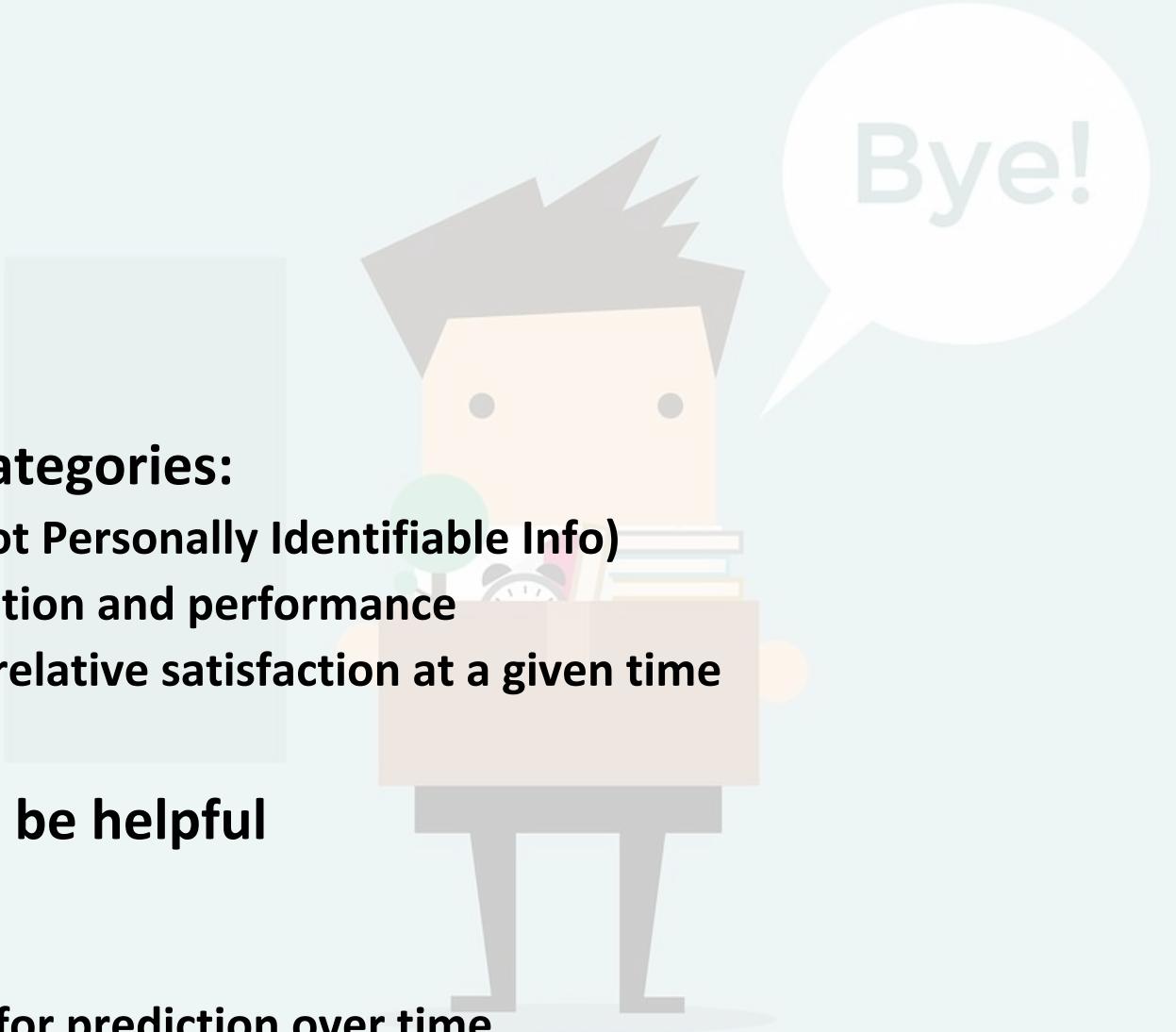
# The Why: Why did we choose this topic?

- Companies invest a lot in their employees ...
- Companies want to build programs to keep employees
  - Which feature is most common among those who leave → build a program for all employees related to that feature
  - Do those features change based on the job role → build programs tailored for job roles or job categories
- Companies understand attrition is part of life but want to be able to minimize risk of shortage
  - When will they leave → time recruiting efforts
  - How many can we expect to leave → drive size of recruiting
  - Which roles will they leave from → focus recruiting



# The Why: Why did we choose this data?

- Mixture of attrition results
  - People who stayed ...
  - People who left ...
- Includes multiple features in critical categories:
  - Personal background and information (not Personally Identifiable Info)
  - History and current employment information and performance
  - Survey results measuring an employees' relative satisfaction at a given time
- Missing additional features that could be helpful
  - Time/date of survey
  - Date of resignation
  - Multiple survey results to create models for prediction over time



# The How: How are we storing the data?

## *pgAdmin 4.24 with Postgres 12.4*

- Original tables
  - IBMEmployeeAttrition
  - Rating and/or Satisfaction score explanation
- New tables created from ETL process in Python
  - Added 'JobCategory' field
  - Incorporated text fields from satisfaction/rating tables
  - Created 3x sub-tables for each JobCategory type

## *Python w/ Jupyter Notebook*

- SQLAlchemy to upload original table from Postgres; used for ETL process
- Export new attrition dataframe to Postgres via SQLAlchemy
- Maintained encoded dataframes for data exploration, model build & training, and inference assessments

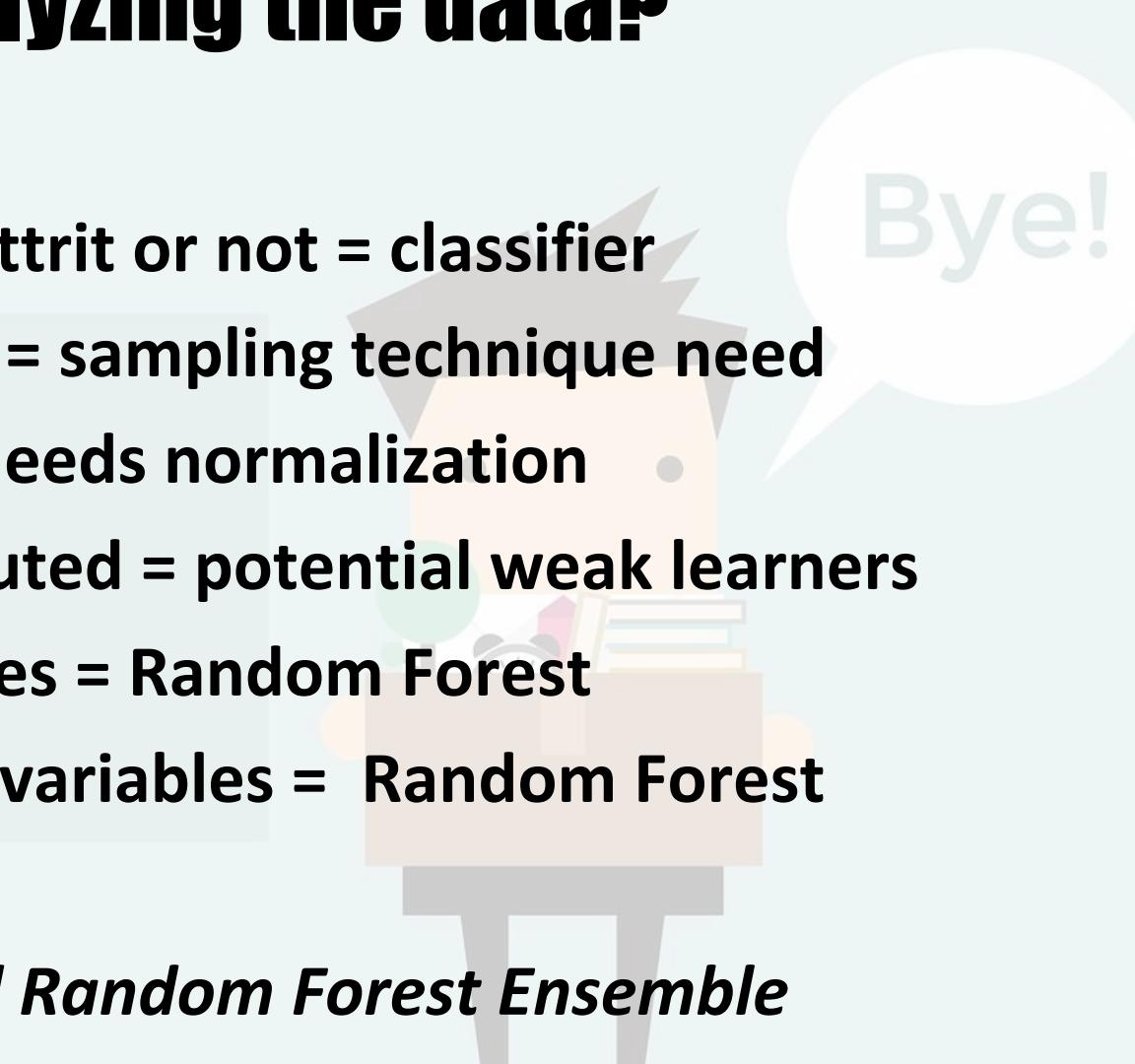
# The How: How are we conducting data exploration?

- Explore imbalance in the target
- Explore correlation between features
- Explore distribution of the features
- Explore distribution of features to one another



# The How: How are we analyzing the data?

- Want to predict if someone will attrit or not = classifier
- Highly imbalanced target dataset = sampling technique need
- Highly variable feature results = needs normalization
- Feature data is not evenly distributed = potential weak learners
- Large number of variables/features = Random Forest
- Rank the importance of features/variables = Random Forest



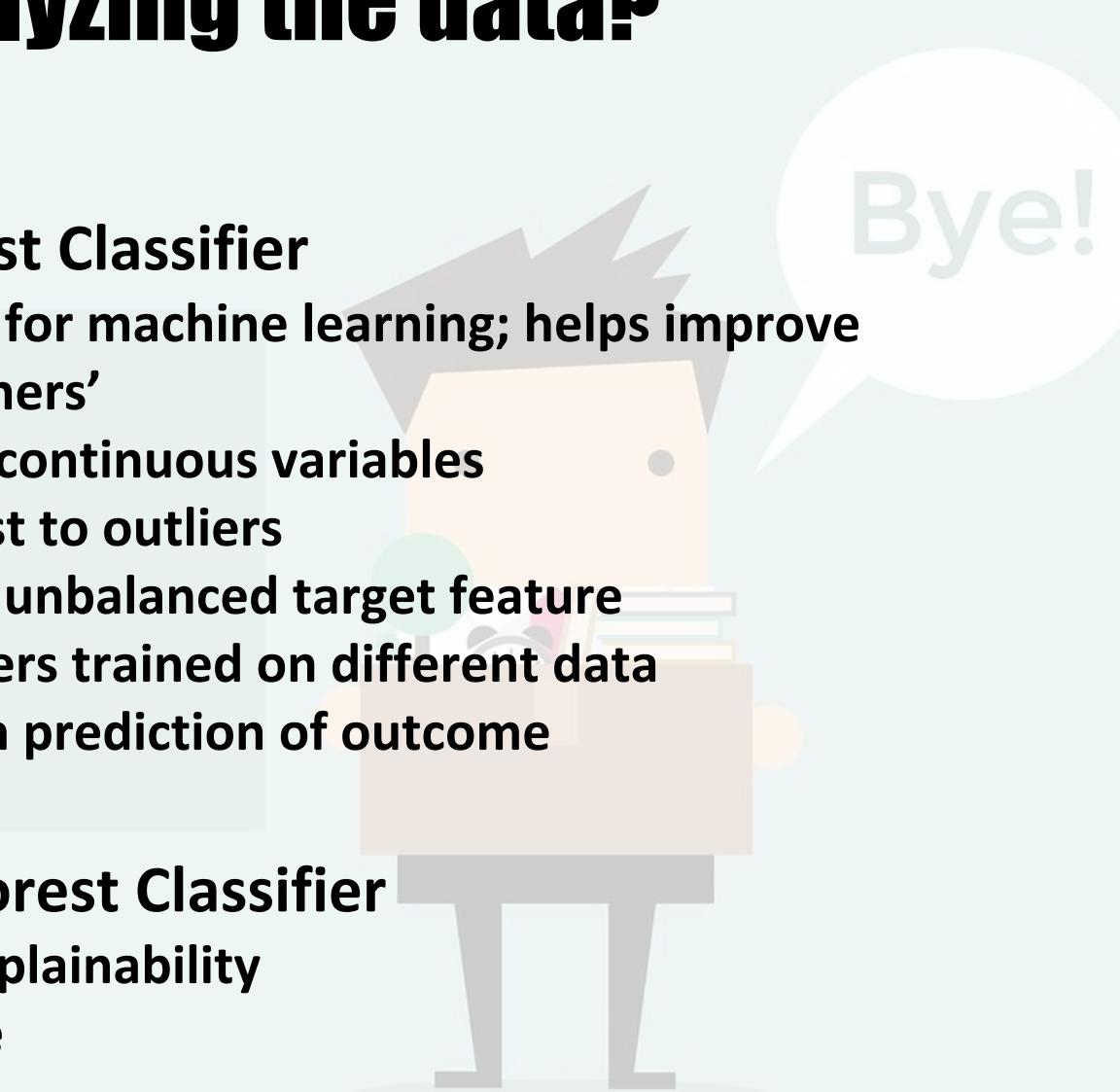
Bye!

*IDEAL ALGORITHM: Balanced Random Forest Ensemble*

*TARGET: Attrition ('Yes', 'No')*

# The How: How are we analyzing the data?

- **Strengths of Balanced Random Forest Classifier**
  - Based on bagging ensemble technique for machine learning; helps improve accuracy and robustness of ‘weak learners’
  - Good for categorical as well as well as continuous variables
  - Does not require feature scaling; robust to outliers
  - Incorporates techniques to balance an unbalanced target feature
  - Robust against overfitting; weak learners trained on different data
  - Good to rank importance of features in prediction of outcome
- **Weaknesses of Balanced Random Forest Classifier**
  - Complexity makes it challenging for explainability
  - Generally requires longer training time



Bye!

# The How: How are we analyzing the data?

- Train/Test Split:
  - 70 train // 30 test
  - Sklearn.model



# The How: How are we displaying our work?

- Main tool: JavaScript + HTML + CSS



# The How: How are we displaying our work?

Bye, Bye, Bye!

*Attrition: Why and When?*

Selection Bar: [Overview](#) // [Data Analysis](#) // [Machine Learning Model Assessments](#) // [Feature Analysis](#) // [Conclusions](#)

Overview of the project: The Why

Overview of the project: The How

Overview of the project: The Overall Results

Dataframe:

*Drop-down & update with different dataframes (combo, tech, non-tech, leadership)*

# The How: How are we displaying our work?

Bye, Bye, Bye!

*Attrition: Why and When?*

*Selection Bar: Overview // Data Analysis // Machine Learning Model Assessments // Feature Analysis // Conclusions*

**Overall Assessments of Data Analysis:**  
Text

**Target Imbalance Display:**  
*Visualization and discussion of imbalance in target in dataframes*

**Correlation Display:**  
*Visualization of Feature correlation to one another; highlight positive, negatives, and 'surprises'*

**Distribution Display:**  
*Visualization w/ drop-downs for each feature AND buttons different analysis based on JobCategory*

**Comparison Feature Display:**  
*Visualization w/ Comparison of Features*

# The How: How are we displaying our work?

Bye, Bye, Bye!

*Attrition: Why and When?*

*Selection Bar: Overview // Data Analysis // **Machine Learning Model Assessments** // Feature Analysis // Conclusions*

**Overall Assessments of the ML Model:**  
Text

**Machine Learning Assessment Approach:**  
Text

**Graphic Display:**

*Confusion Matrix, Accuracy, & Imbalance Report w/ buttons for the model used AND the sub-category*

# The How: How are we displaying our work?

Bye, Bye, Bye!

*Attrition: Why and When?*

*Selection Bar: Overview // Data Analysis // Machine Learning Model Assessments // Feature Analysis // Conclusions*

**Overall Assessments of the Feature Analysis & Comparison:**  
Text

**Graphic Display:**

*Comparison of each of the features importance from the classifier*

# The How: How are we displaying our work?

Bye, Bye, Bye!

*Attrition: Why and When?*

*Selection Bar: Overview // Data Analysis // Machine Learning Model Assessments // Feature Analysis // Conclusions*

Overall Conclusions on Our Questions:

*Text*

# The Results: Data Exploration - Overall Findings

- Highly imbalanced dataset, especially within leadership dataframe
  - Target features imbalanced
  - Feature distribution is highly variable; leadership often older & paid more
- Data correlation is relatively even; most correlations expected
  - Positive:
    - Age → TotalWorkingYears / MonthlyIncome / JobLevel
    - MonthlyIncome → YearsAtCompany / TotalWorkingyears
    - PerformanceRating → PercentSalaryHike
    - JobLevel → YearsAtCompany / TotalWorkingYears / MonthlyIncome
  - Negative: JobCategory → JobLevel / MonthlyIncome / TotalWorkingYears
  - Unexpected: MaritalStatus / StockOptionLevel (negative)

## *Explore imbalance in the target*

```
Combined:  
No      1233  
Yes     237  
Name: Attrition, dtype: int64  
Tech:  
No      564  
Yes     118  
Name: Attrition, dtype: int64  
Non-Tech:  
No      359  
Yes     102  
Name: Attrition, dtype: int64  
Leadership:  
No      310  
Yes     17  
Name: Attrition, dtype: int64
```

## The Results: Data

- Target data for combined and three sub-categories are all imbalanced
  - Leadership most imbalanced
  - Non-tech least imbalanced
- Will require sampling technique to correct for imbalance in target data

# The Results: Data

- Multiple positive correlations:

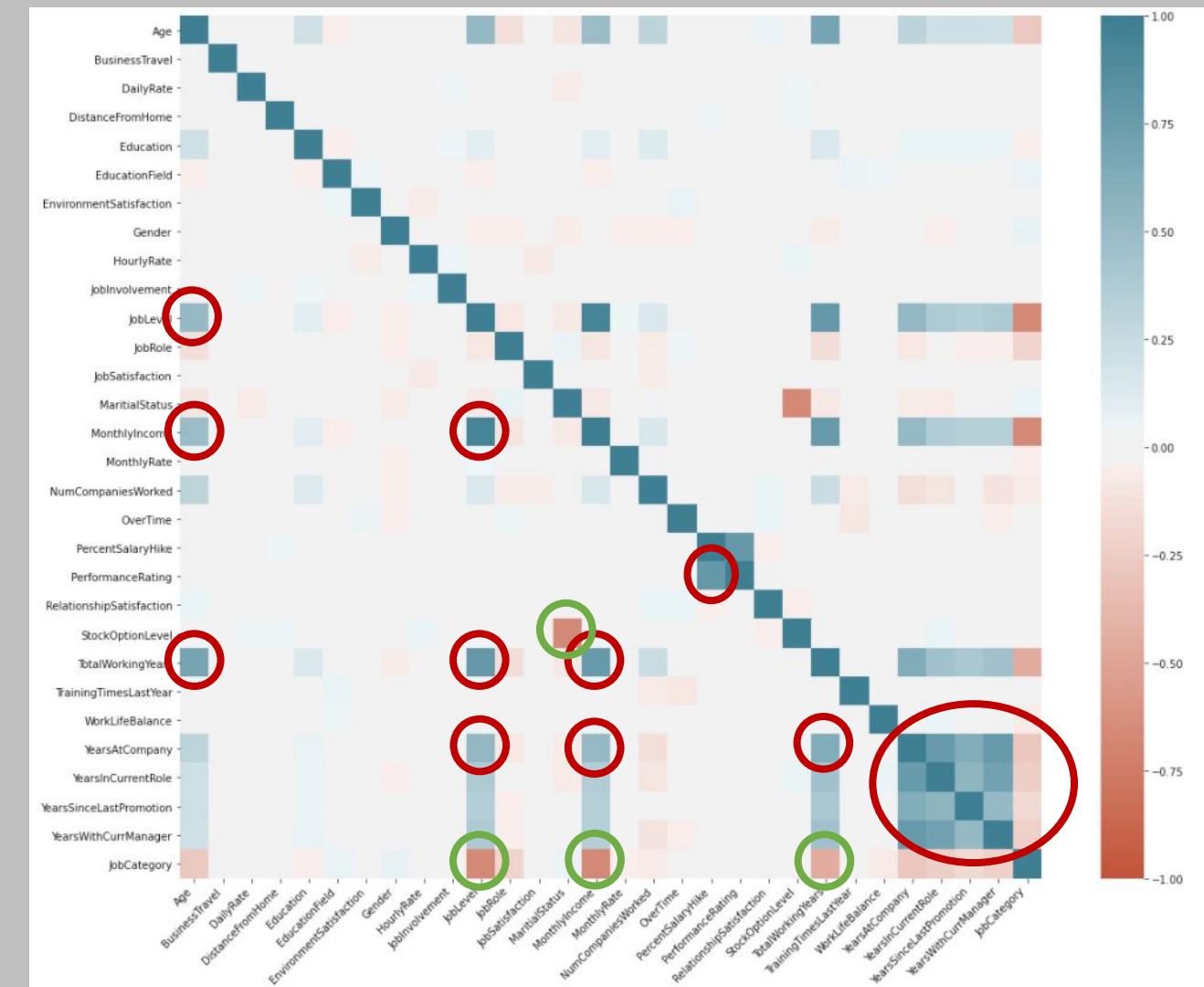
- Age: MonthlyIncome / JobLevel / TotalWorkingYears
- JobLevel: YearsAtCompany / TotalWorkingYears / MonthlyIncome
- MonthlyIncome: YearsAtCompany / TotalWorkingYears
- PercentSalaryHike: PerformanceRating
- TotalWorkingYears: YearsAtCompany

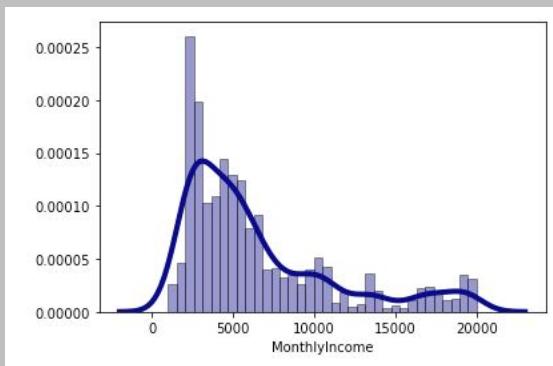
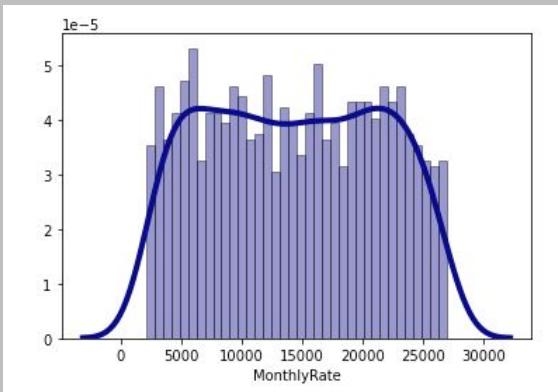
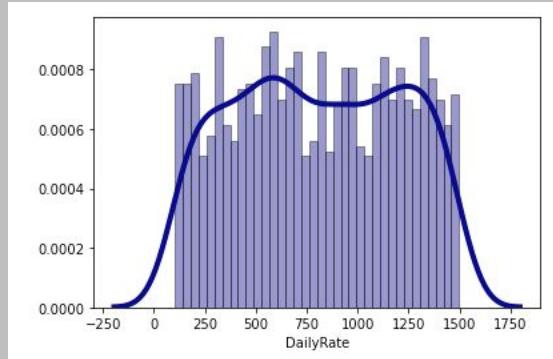
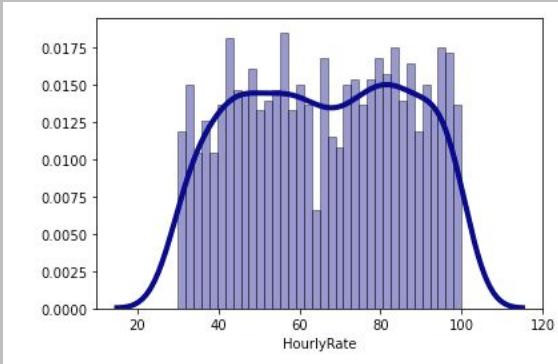
- Multiple negative correlations:

- JobLevel: JobCategory
- MaritalStatus: StockOptionLevel
- MonthlyIncome: JobCategory
- TotalWorkingYears: JobCategory

- Features representing similar/same information:

- HourlyRate / DailyRate / MonthlyRate / MonthlyIncome
- Age / TotalWorkingYears





# The Results: Data

- Four features provide the same info; built new dataframe to assess similarity between data
- Assessed distribution to determine which feature would remain, included in model
  - MonthlyIncome – least normally distributed
  - HourlyRate, DailyRate – normally distributed
  - MonthlyRate – most normally distributed

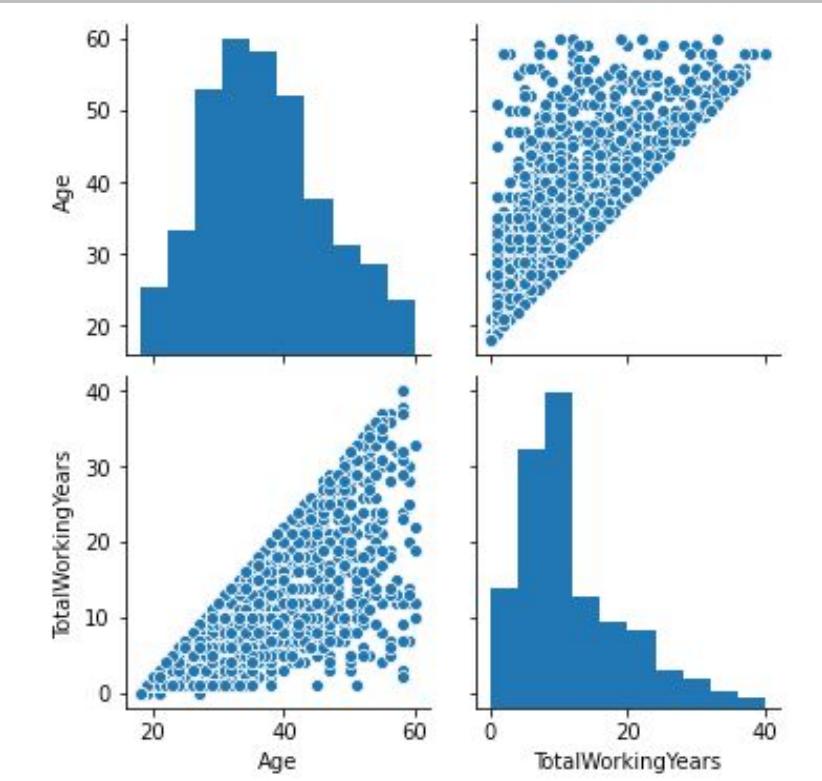
*Dropping HourlyRate, DailyRate; MonthlyRate & MonthlyIncome used in model development*

# The Results: Data

- Age and TotalWorkingYears appear to highlight similar info, but not the same
  - Used describe and distribution analysis
  - Both skewed left, but still independent

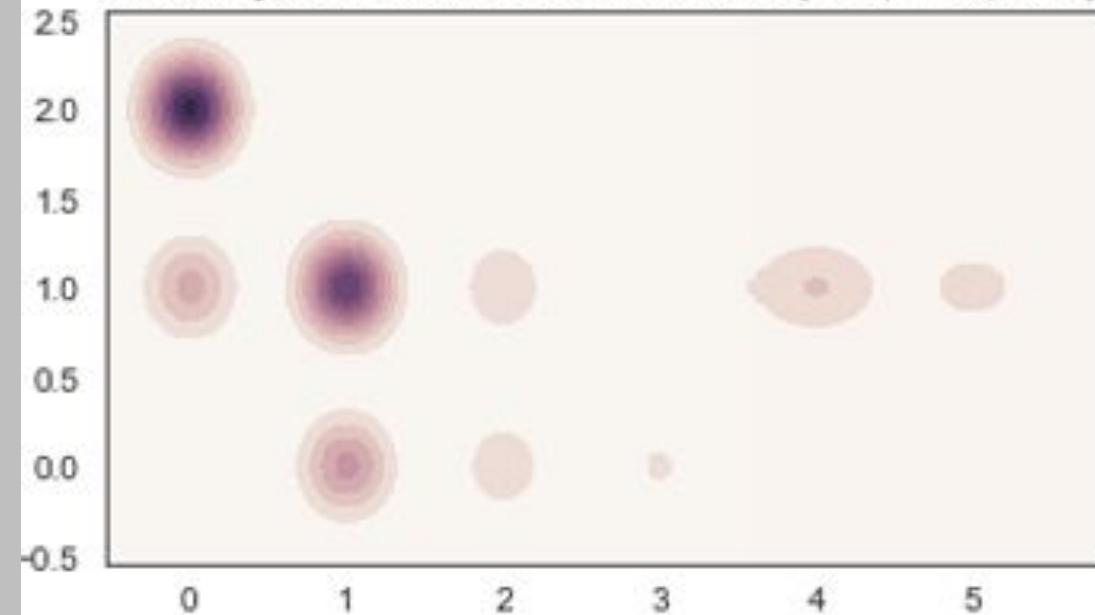
*Maintain both features for the model development and inference*

	Age	TotalWorkingYears	Difference
count	1470.000000	1470.000000	1470.000000
mean	36.923810	11.279592	25.644218
std	9.135373	7.780782	6.875481
min	18.000000	0.000000	18.000000
25%	30.000000	6.000000	20.000000
50%	36.000000	10.000000	24.000000
75%	43.000000	15.000000	30.000000
max	60.000000	40.000000	56.000000



# The Results: Data

Stock Option Level vs. Marital Status (2=S, 1=M, 0=D)



- Single employees tend to have lower stock levels
- Married employees appear to have higher levels
- Divorced employees fall between the two above

*Younger employees are often single vs. married employees; some employees who are divorced get remarried later in life*

# The Results: Analysis - Overall Findings

- Ideal model use is to predict employees who do NOT attrit ('no')
  - Model prediction of employees who attrition ('yes') is not robust
  - F1 scores for categories and overall employees no higher than .51
  - Largest deficiency in 'yes' for attrition is precisions; no higher than .39
- Leadership imbalance too large; not able to accurately predict
  - Poor ability to predict employees who attrit/leave the company
  - F1 scores no higher than .06 for 'yes' to attrition target feature
  - Model appears to be overfitting for predicting 'no' for attrition
- Feature importance in attrition does vary based on job category
  - Tech more sensitive to age, TotalWorkingYears, YearsInCurrentRole; less to YearsWithCurrManager, JobRole
  - Non-tech more sensitive to StockOptionLevel, NumberofCompaniesWorked, DistanceFromHome; less to YearsWithCurrManager
  - Leadership more sensitive to YearsSinceLastPromotion, Education, JobLevel, RelationshipSatisfaction, PerformanceRating; less to Age, MonthlyIncome

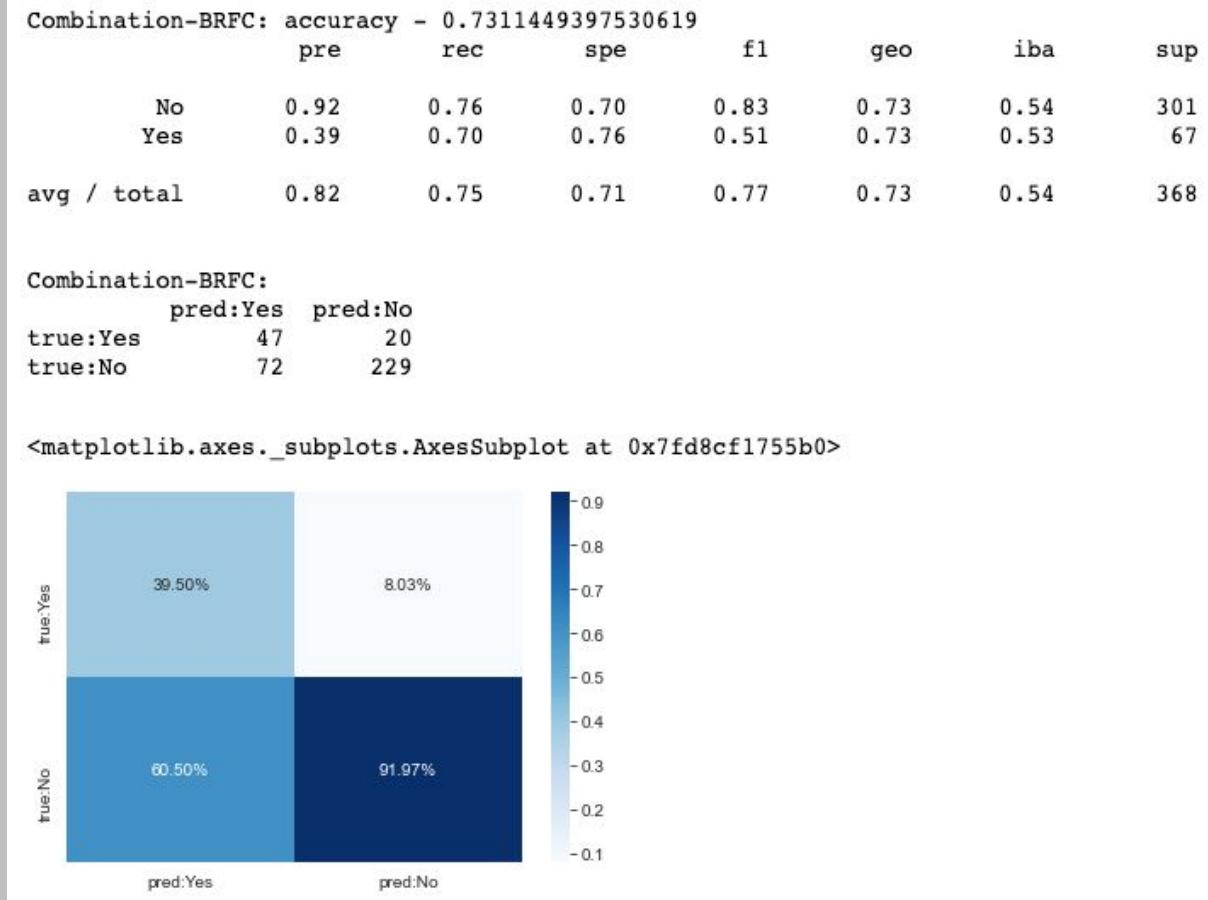


Bye!

# The Results: Analysis - Combo

- Precision – Where all the ‘yes/no’s right?
  - ‘Yes’: Worse than a coin toss
  - ‘No’: Fairly accurate
- Recall – Did we get all the ‘yes/no’s?
  - ‘Yes’: Fairly accurate
  - ‘No’: Fairly accurate
- Overfitting –
  - ‘Yes’: Potentially overfit
  - ‘No’: Not likely overfit

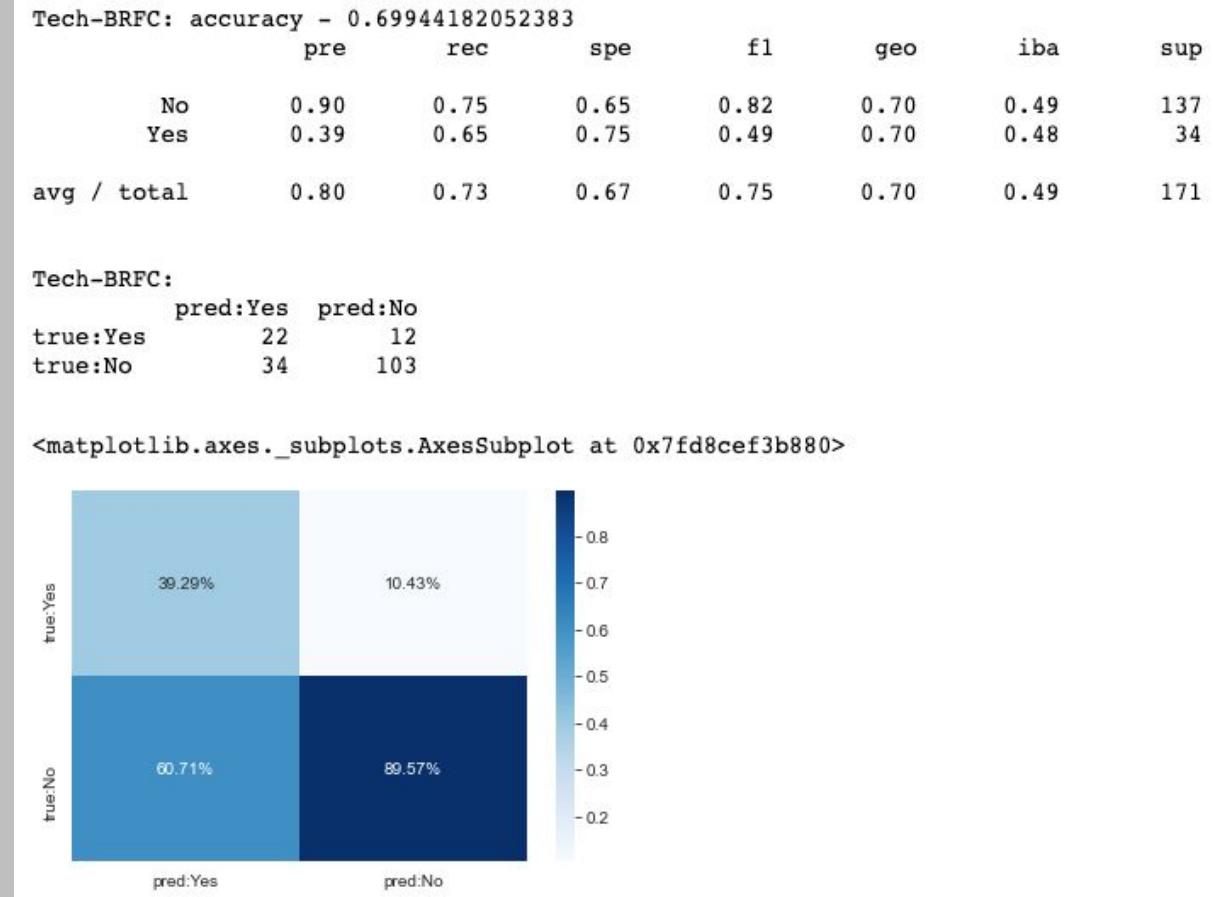
*Overall, no better than a coin toss at predicting employees who left, but far better at predicting those who stayed*



# The Results: Analysis - Tech

- Precision – Where all the ‘yes/no’s right?
  - ‘Yes’: Worse than a coin toss
  - ‘No’: Fairly accurate
- Recall – Did we get all the ‘yes/no’s?
  - ‘Yes’: Fairly accurate
  - ‘No’: Fairly accurate
- Overfitting –
  - ‘Yes’: Potentially overfit
  - ‘No’: Not likely overfit

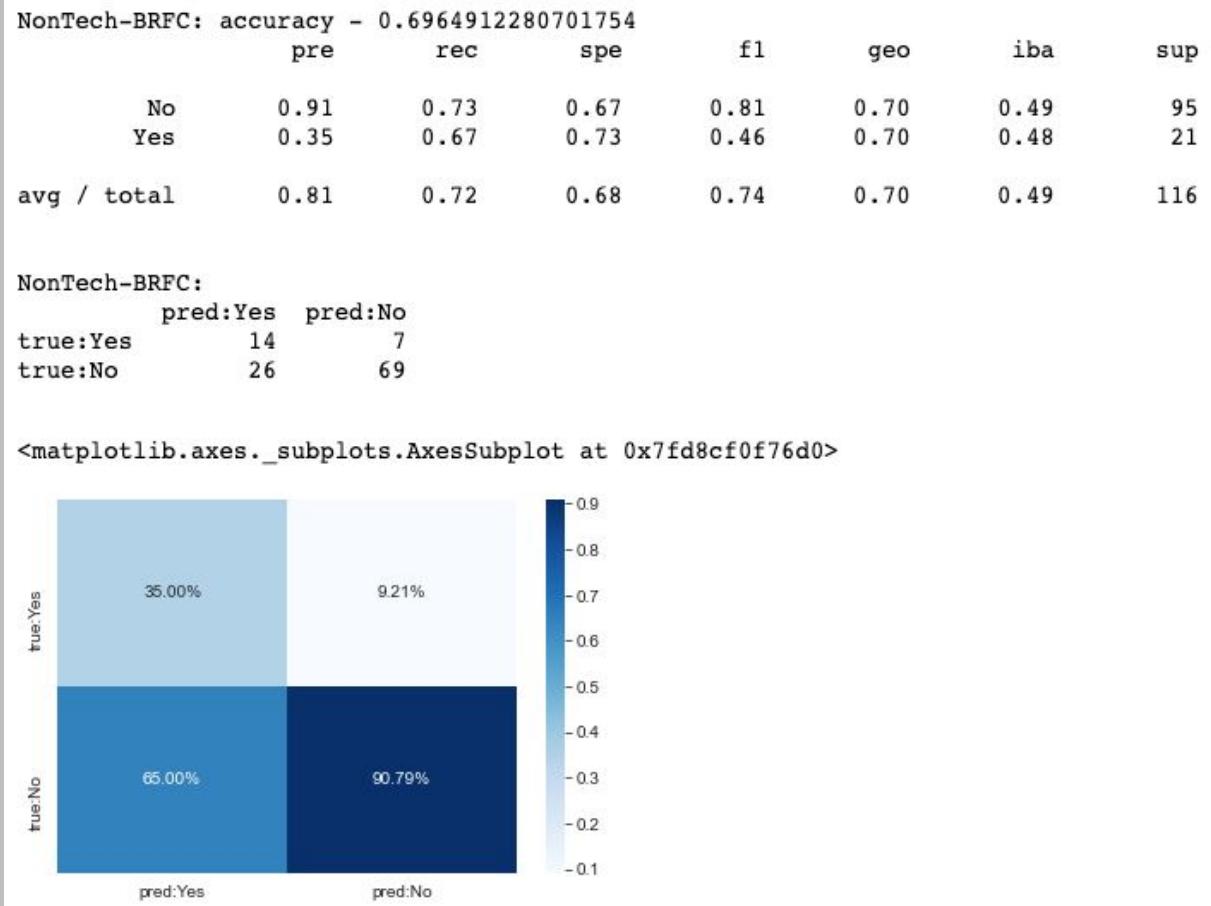
*Overall, worse than a coin toss at predicting employees who left, but far better at predicting those who stayed*



# The Results: Analysis - Non-Tech

- Precision – Where all the ‘yes/no’s right?
  - ‘Yes’: Worse than a coin toss
  - ‘No’: Fairly accurate
- Recall – Did we get all the ‘yes/no’s?
  - ‘Yes’: Fairly accurate
  - ‘No’: Fairly accurate
- Overfitting –
  - ‘Yes’: Potentially overfit
  - ‘No’: Not likely overfit

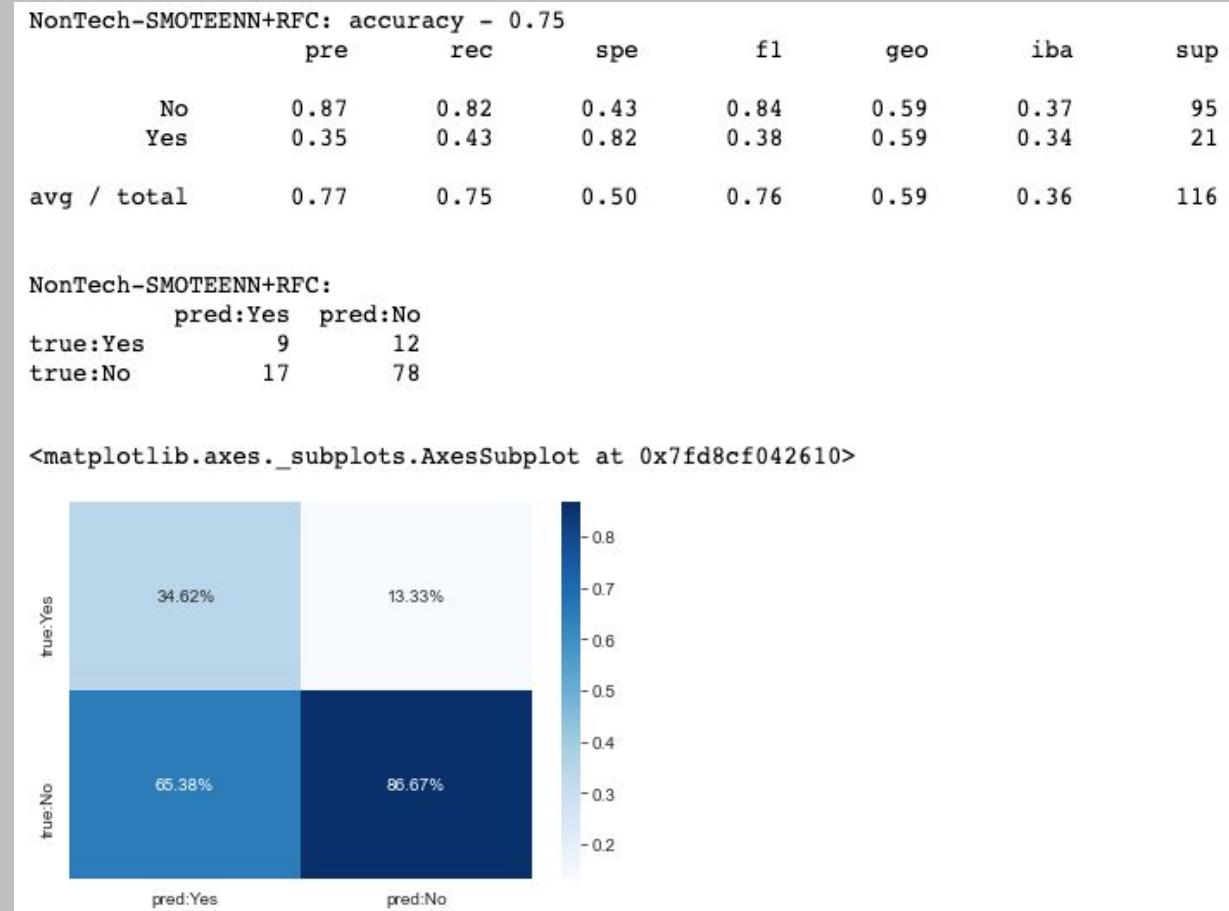
*Overall, worse than a coin toss at predicting employees who left, but far better at predicting those who stayed*



# The Results: Analysis - Non-Tech

- Precision – Where all the ‘yes/no’s right?
  - ‘Yes’: Worse than a coin toss
  - ‘No’: Fairly accurate
- Recall – Did we get all the ‘yes/no’s?
  - ‘Yes’: Fairly accurate
  - ‘No’: Worse than a coin toss
- Overfitting –
  - ‘Yes’: Potentially overfit
  - ‘No’: Not likely overfit

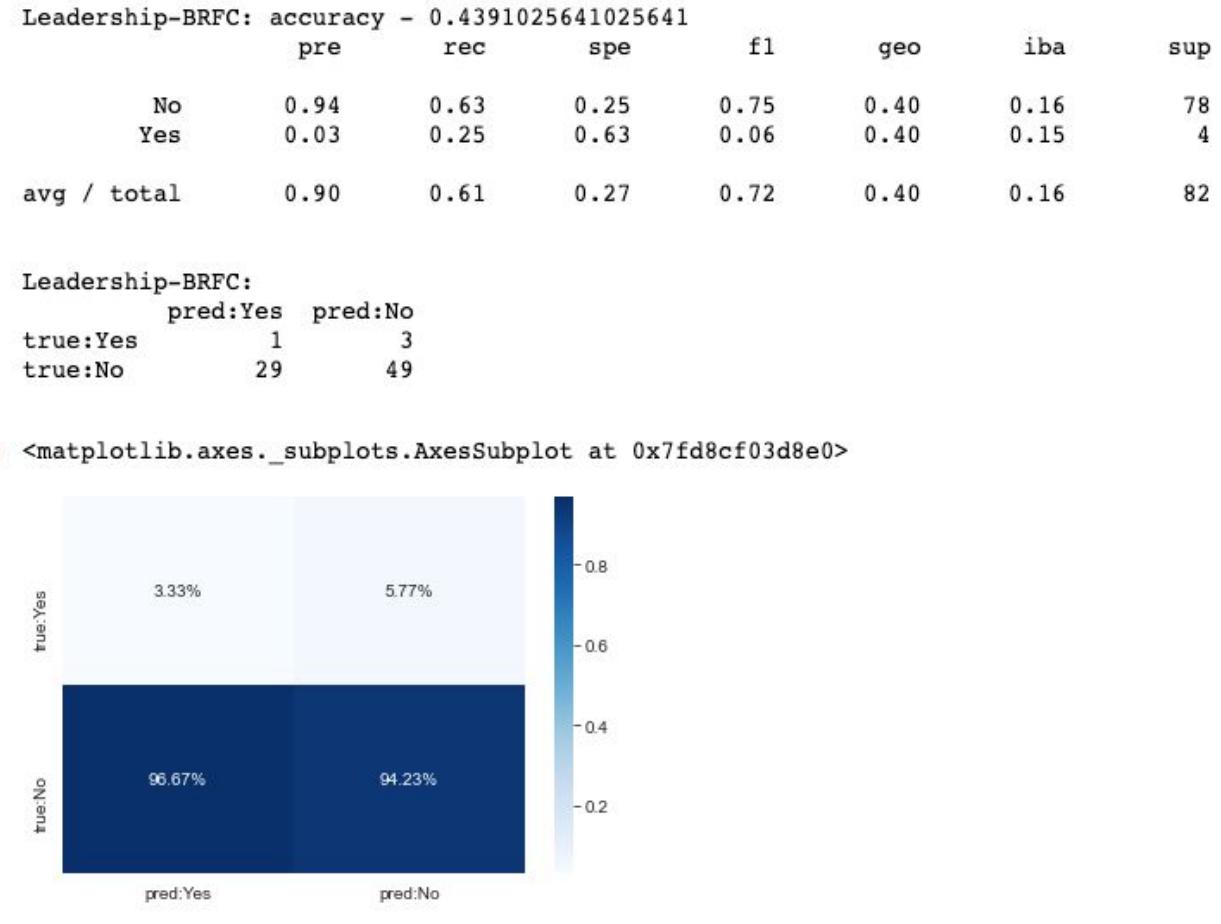
*Overall, although accuracy is better, the F1 score is no better than a coin toss at predicting people who leave; ability to predict those who stay is less accurate than the BRFC model*



# The Results: Analysis - Ldrshp

- Precision – Where all the ‘yes/no’s right?
  - ‘Yes’: Completely unreliable
  - ‘No’: Fairly accurate
- Recall – Did we get all the ‘yes/no’s?
  - ‘Yes’: Only right  $\frac{1}{4}$  of the time
  - ‘No’: Slightly better than a coin toss
- Overfitting –
  - ‘Yes’: Likely overfit
  - ‘No’: Possibly overfit

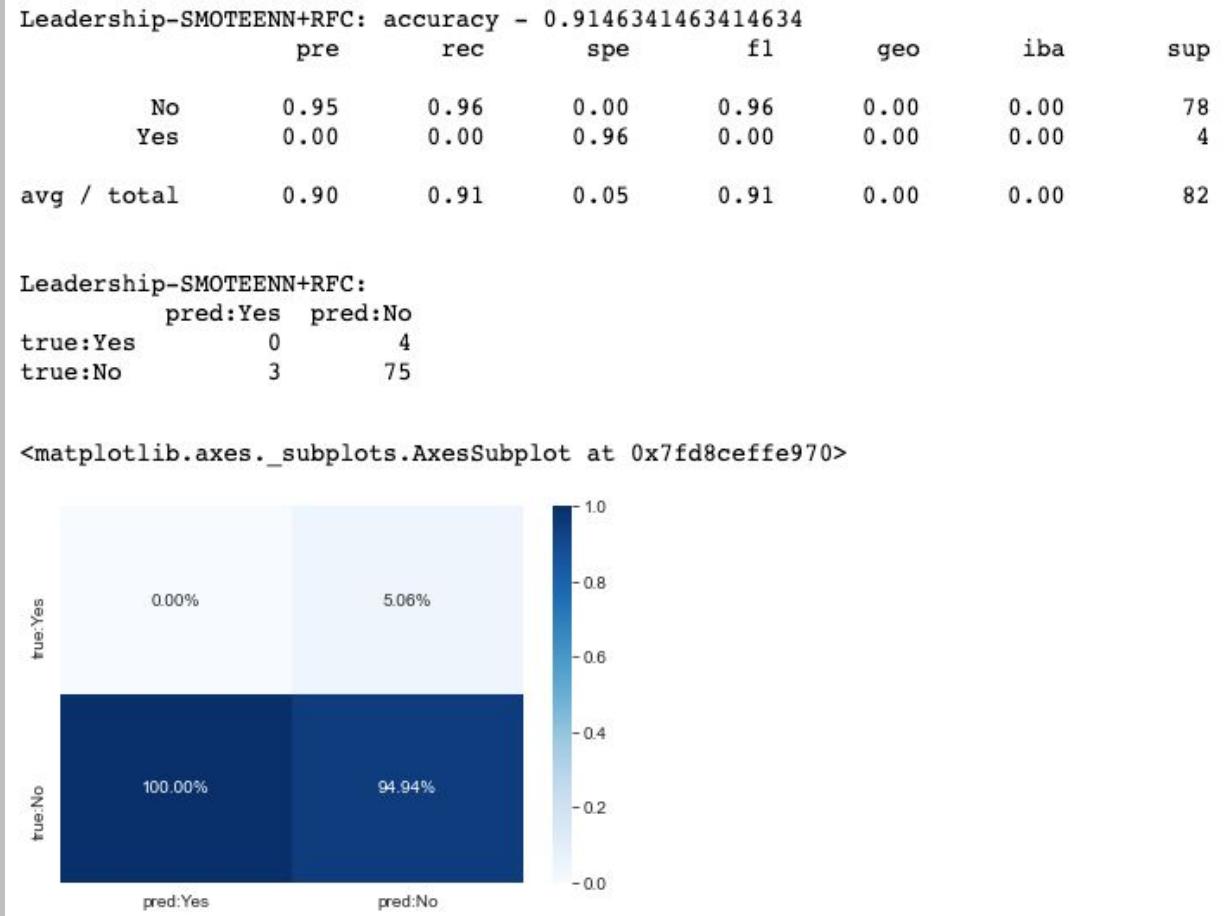
*Overall, completely unable to predict employees who left, but far better at predicting those who stayed, aka did not attrit*



# The Results: Analysis - Ldrshp

- Precision – Where all the ‘yes/no’s right?
  - ‘Yes’: Completely unreliable
  - ‘No’: Fairly accurate
- Recall – Did we get all the ‘yes/no’s?
  - ‘Yes’: Completely unreliable
  - ‘No’: Fairly reliable
- Overfitting –
  - ‘Yes’: Likely overfit
  - ‘No’: Unlikely overfit

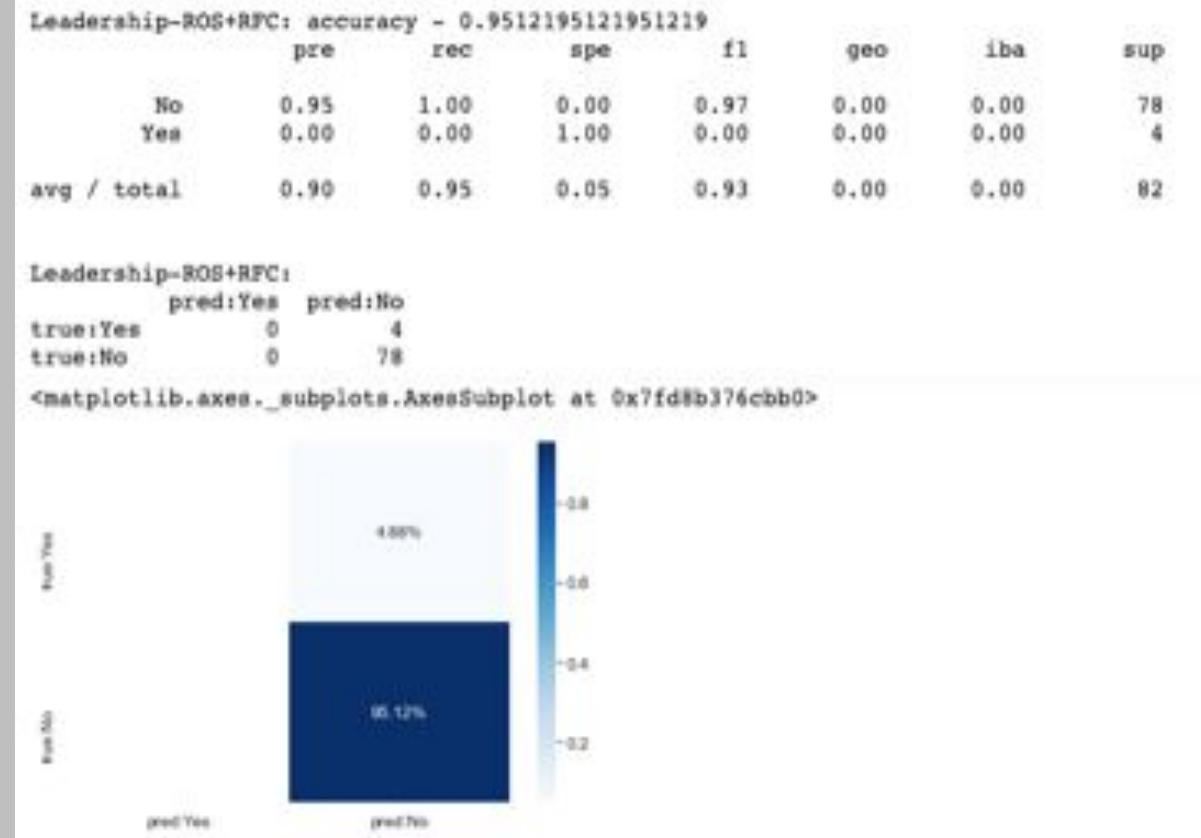
*Overall, far less able to predict employees who left as opposed to BRFC; although predicted employees who stayed better, it is still an unreliable model*



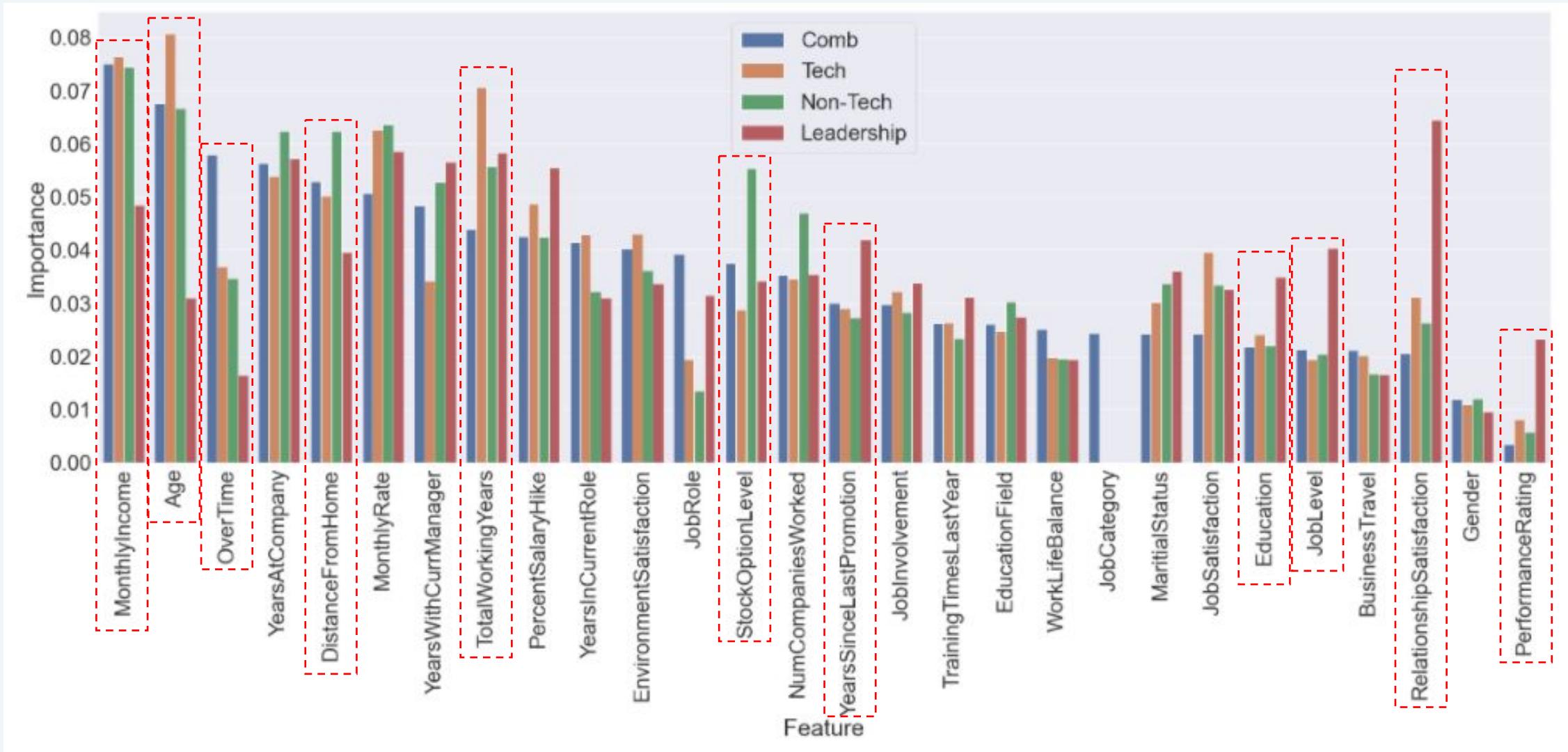
# The Results: Analysis - Ldrshp

- Precision – Where all the ‘yes/no’s right?
  - ‘Yes’: Completely unreliable
  - ‘No’: Fairly accurate
- Recall – Did we get all the ‘yes/no’s?
  - ‘Yes’: Completely unreliable
  - ‘No’: Fairly reliable
- Overfitting –
  - ‘Yes’: Likely overfit
  - ‘No’: Unlikely overfit

*Overall, far less able to predict employees who left as opposed to BRFC; although predicted employees who stayed better, it is still an unreliable model*

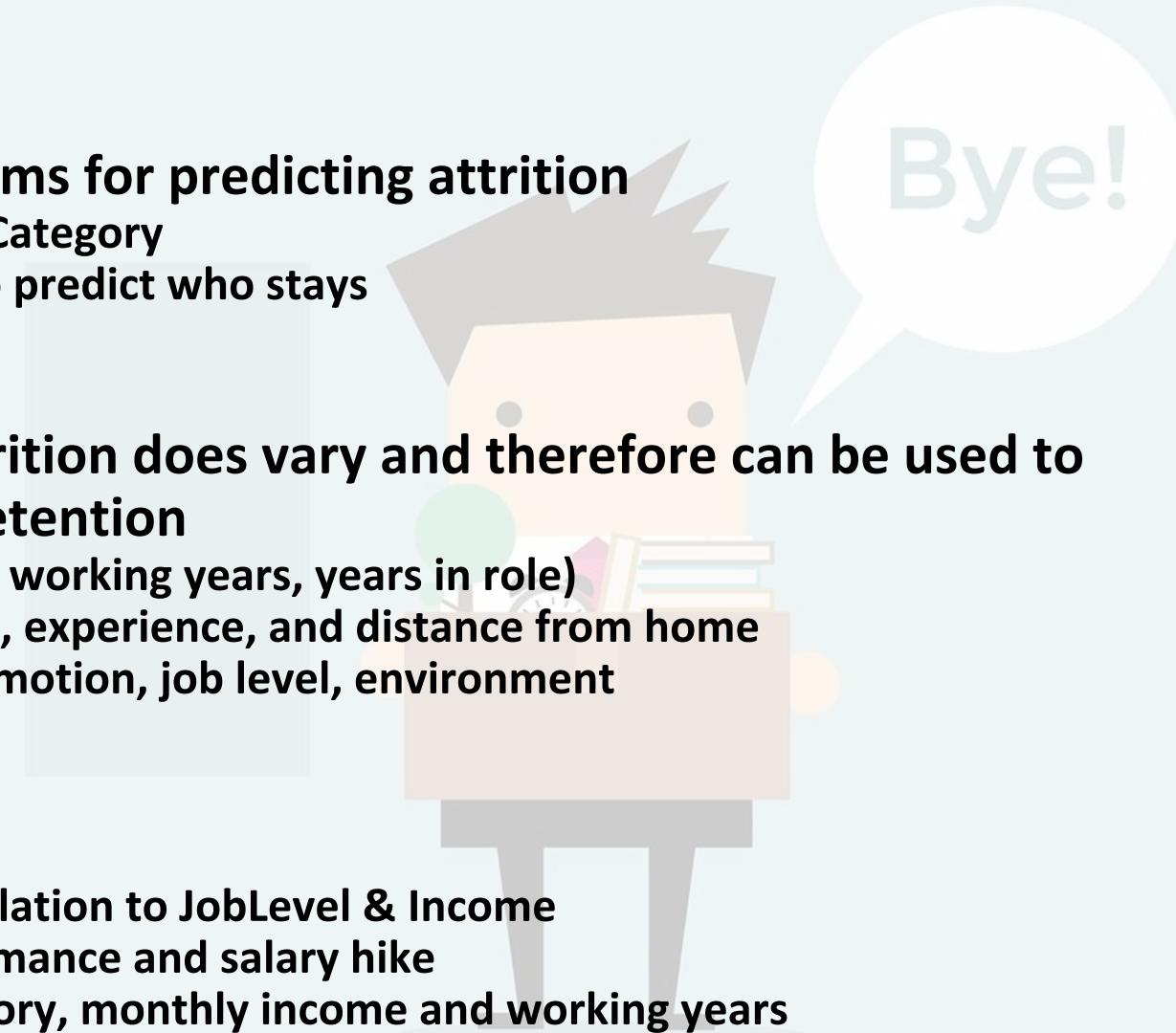


# The Results: Analysis - Feature Comparison



# Overall Conclusions

- Imbalanced targetset, creates problems for predicting attrition
  - This is especially for the Leadership JobCategory
  - Much more accurate and appropriate to predict who stays
- Feature importance in predicting attrition does vary and therefore can be used to build out targeted mechanisms for retention
  - Tech roles - more sensitive to time (age, working years, years in role)
  - Non-Tech roles - more sensitive to stock, experience, and distance from home
  - Leadership roles - more sensitive to promotion, job level, environment
- We saw unique correlations:
  - Time based features have positive correlation to JobLevel & Income
  - Positive correlation between job performance and salary hike
  - Negative correlation between job category, monthly income and working years



Bye!