



Factors Influencing Employee Retention

Griselda Arevalo

Coral Rosa-Falero

Jiemin Sheng

Bootcamp: GWU-VIRT-DATA-PT-08-2022-U-B-MW



Topic: Employee Retention Prediction

Purpose

- Build a predictive model determine the factors that influence employee retention.

Applications

- Improve selection and retention of employees.
- Understand what factors that influence employee retention

Tools

- Jupyter NB
 - pandas
 - sqlite3
 - Sklearn.metrics
- Tableau
- GitHub

Purpose

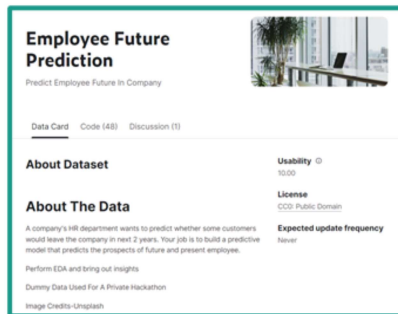
- Build a predictive model determine the factors that influence employee retention.

Applications

- Assist human resource departments improve selection of future employees and improve the retention of new employees.
- To better understand what variables influence the employees either stay or leave.

Data: Overview

csv file with 4,653 observations



Employee Future Prediction
Predict Employee Future In Company

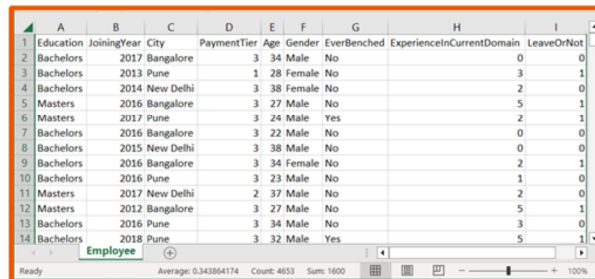
Data Card Code (48) Discussion (1)

About Dataset

About The Data
A company's HR department wants to predict whether some customers would leave the company in next 2 years. Your job is to build a predictive model that predicts the prospects of future and present employee.

Perform EDA and bring out insights
Dummy Data Used For A Private Hackathon
Image Credits-Ungliah

Usability 10.00
License CC0: Public Domain
Expected update frequency Never



	A	B	C	D	E	F	G	H	I
	Education	JoiningYear	City	PaymentTier	Age	Gender	EverBenched	ExperienceInCurrentDomain	LeaveOrNot
1	Bachelors	2017	Bangalore	3	34	Male	No	0	0
2	Bachelors	2013	Pune	1	28	Female	No	3	1
3	Bachelors	2014	New Delhi	3	38	Female	No	2	0
4	Masters	2016	Bangalore	3	27	Male	No	5	1
5	Masters	2017	Pune	3	24	Male	Yes	2	1
6	Bachelors	2016	Bangalore	3	22	Male	No	0	0
7	Bachelors	2015	New Delhi	3	38	Male	No	0	0
8	Bachelors	2016	Bangalore	3	34	Female	No	2	1
9	Bachelors	2016	Pune	3	23	Male	No	1	0
10	Masters	2017	New Delhi	2	37	Male	No	2	0
11	Masters	2012	Bangalore	3	27	Male	No	5	1
12	Bachelors	2016	Pune	3	34	Male	No	3	0
13	Bachelors	2018	Pune	3	32	Male	Yes	5	1

<https://www.kaggle.com/datasets/tejashvi14/employee-future-prediction>

1. Clean csv file from kaggle
2. Selected because high usability score, database completeness and cleanliness

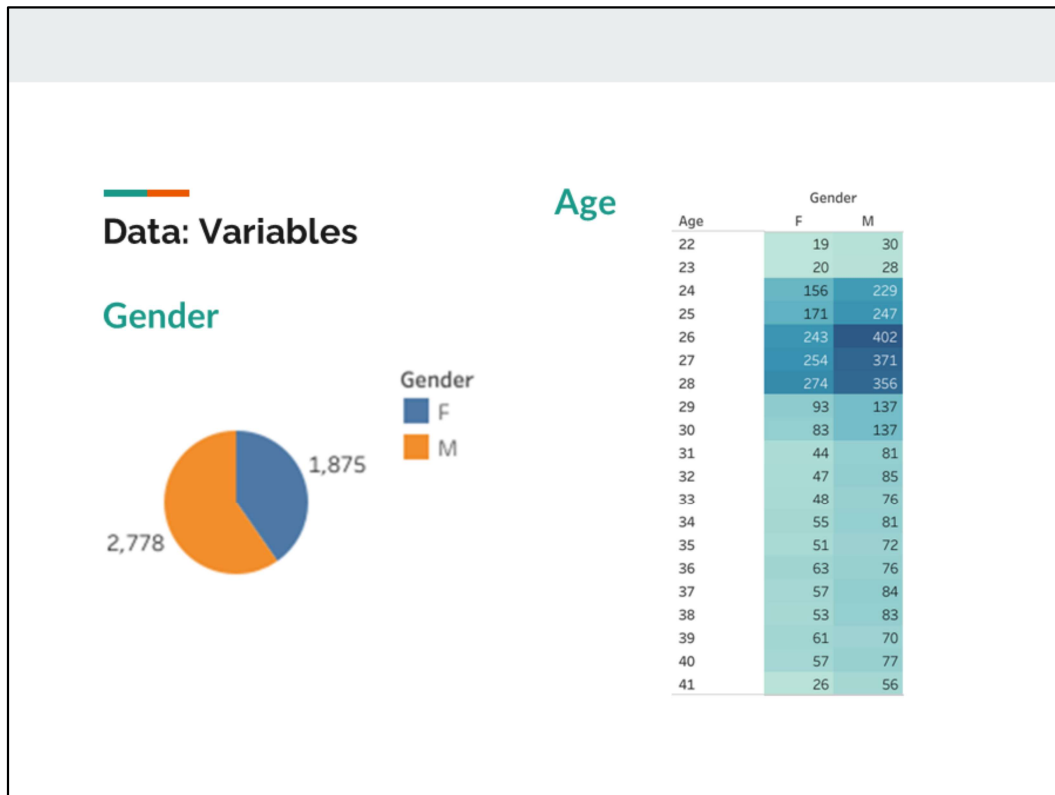
The firm in question does not appear to be stated but can be noted to be located in India based on the office locations in Indian cities.

Data

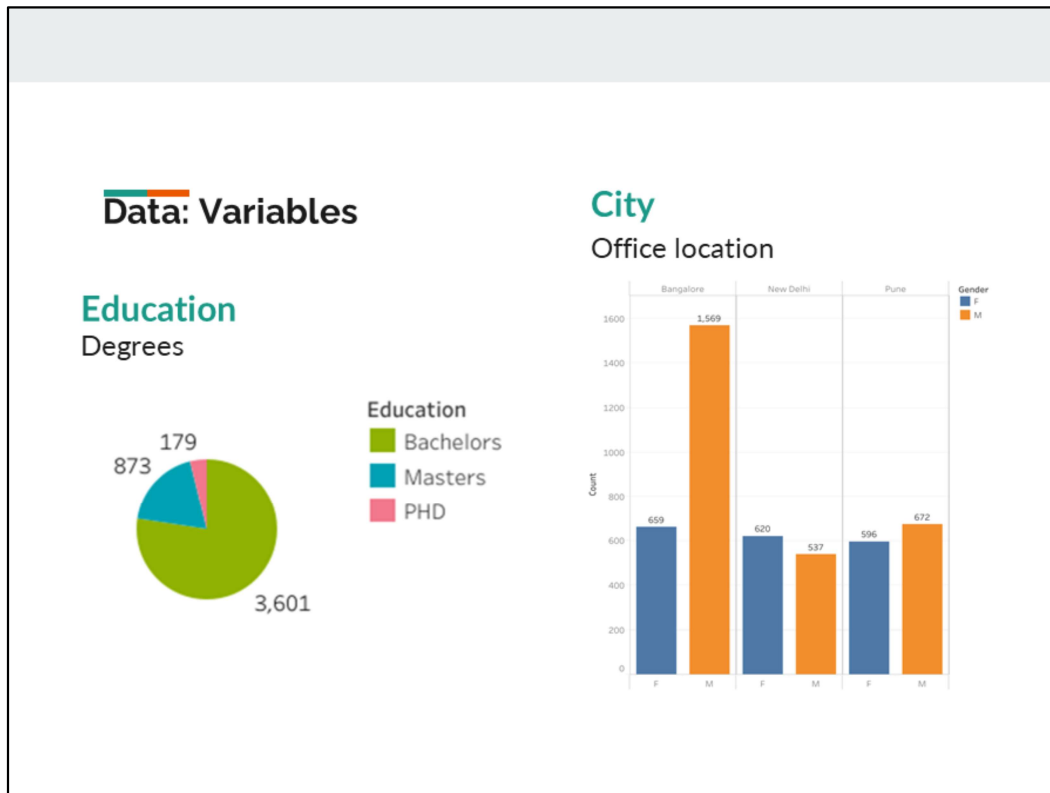
- Education (Degrees)
- Joining Year (Hiring Date)
- City (Office Location)
- Payment Tier (Salary Level)
- Age
- Gender
- Ever Benched (Productivity/Profitability)
- Experience In Current Domain (# Years)
- Leave Or Not (Retention)

```
employee_data_df.info()
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 4653 entries, 0 to 4652
Data columns (total 9 columns):
#   Column                      Non-Null Count  Dtype
---  -
0   Education                   4653 non-null   object
1   JoiningYear                 4653 non-null   int64
2   City                       4653 non-null   object
3   PaymentTier                 4653 non-null   int64
4   Age                        4653 non-null   int64
5   Gender                     4653 non-null   object
6   EverBenched                 4653 non-null   object
7   ExperienceInCurrentDomain    4653 non-null   int64
8   LeaveOrNot                  4653 non-null   int64
dtypes: int64(5), object(4)
memory usage: 327.3+ KB
```

1. There is a total of 4,653 observations.
2. Education level is a qualitative categorical variable consisting of three possible values: Bachelors, Masters, and Doctorate degrees.
3. The year of joining firm variable is a continuous quantitative variable varies from 2012 to 2018.
4. The city location of office is a qualitative categorical variable consisting of three cities: Bangalore, Pune, and New Delhi.
5. The salary tier is a qualitative categorical variable which include three tiers: 1 for the highest, 2 for the middle, and 3 for the lowest. Age is a continuous quantitative variable ranging from 22 to 41.
6. The gender variable is qualitative and categorical with either male or female.
7. The qualitative categorical variable of whether or not the employee was kept out of projects for longer than one month is either yes or no.
8. The experience length variable is continuous and quantitative and ranges from 0 to 7.
9. The dependent variable being predicted of whether or not the employee is expected to be leaving the firm within two years is qualitative categorical variable of either 1 or 0 for yes or no.



1. The **gender** variable is qualitative and categorical with either male or female
 - a. More Males than Females
2. **Age** is a continuous quantitative variable ranging from 22 to 41.
 - a. Most of the employees are within the ages of 24 and 28



1. **Education** level is a qualitative categorical variable consisting of three possible values: Bachelors, Masters, and Doctorate degrees
 - a. Great majority posses a BS degree
2. The **city** location of office is a qualitative categorical variable consisting of three cities: Bangalore, Pune, and New Delhi.
 - a. The bangalore office contains the majority of the employees.
 - b. New Delhi and Pune with comparable numbers of M/F. while in the Bangalore office the M employees are vet twice as many as F employees

===

Both images are set up to be filtered. Image can be replaced with link to tableau for live demo

Data: Variables

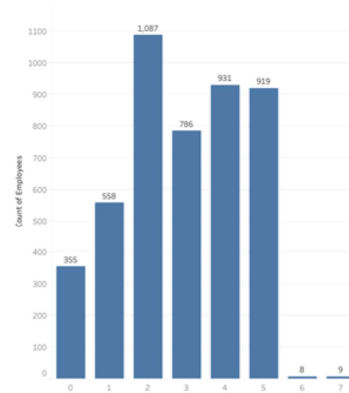
Join Year

Hiring Date

Join Year Distribution

Joining Year	Gender	
	F	M
2012	180	324
2013	253	416
2014	246	453
2015	440	341
2016	178	347
2017	450	658
2018	128	239

Experience In Current Domain

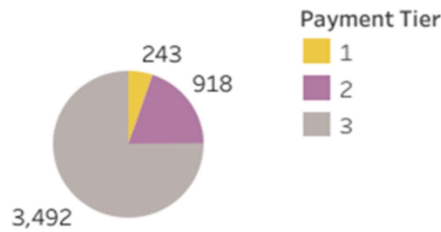


1. The **year of joining** firm variable is a continuous quantitative variable varies from 2012 to 2018.
 - a. Hiring increased from 2012 to 2015, dipped in 2016, doubled in 2017 and dropped again during 2018
2. The **experience** length variable is continuous and quantitative and ranges from 0 to 7y
 - a. The majority of employees have up to 5 years of experience in the current field.
 - b. A disproportionate minority has 5+ years experience
 - c. Most employees fall in the 2, 4 and 5y bracket and this distribution is consistent between genders

Data: Variables

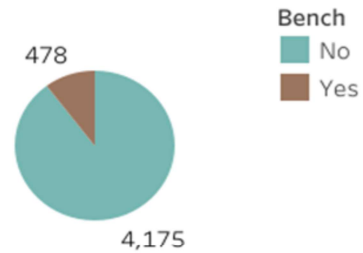
Payment Tier

Salary

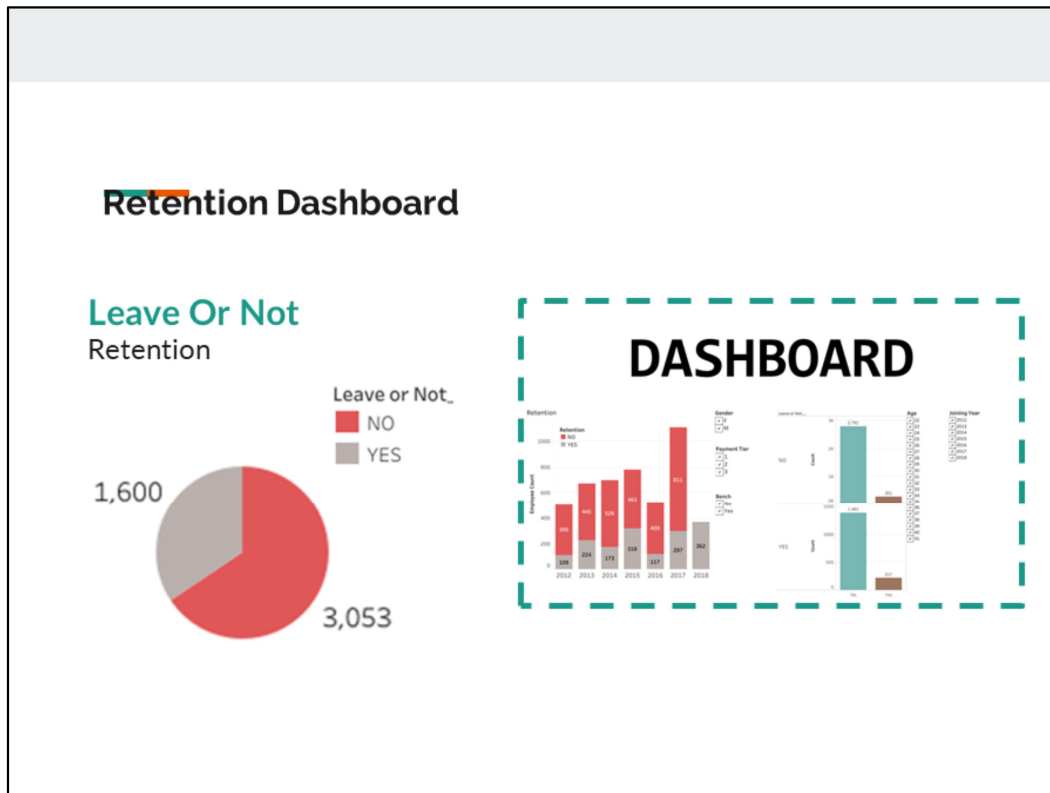


Ever Benched

Productivity/Profitability



1. The **salary** tier is a qualitative categorical variable which include three tiers: 1 for the highest, 2 for the middle, and 3 for the lowest
 - a. Lowest tier contains more employees than the other 2 tiers combined
2. The qualitative categorical variable of whether or not the employee was kept out of projects (**ever benched**) for longer than one month is either yes or no.
 - a. Benching similar between M/F
 - b. Highest benching was in 2015



1. The dependent variable being predicted of whether or not the employee is expected to be leaving the firm within two years (**retention**) is qualitative categorical variable of either 1 or 0 for yes or no
 - a. ???????

Click on dashboard to open the interactive view. Use the filters and discuss

Overview of Analysis

- Connect to sqlite3 database

- Assess the properties of the database

- Establish columns/target and remove unwanted columns
- Check the balance of the target values
- Establish the training set and train the model
- Resample the training data
- Calculate training accuracy score
- Calculate testing set accuracy score
- Calculate confusion matrix
- Evaluate the classification report
- Determine the importance rating of each feature

```
employee_data_df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 4653 entries, 0 to 4652
Data columns (total 9 columns):
#   Column              Non-Null Count  Dtype
---  ---
0   Education            4653 non-null  object
1   JoiningYear          4653 non-null  int64
2   City                 4653 non-null  object
3   PaymentTier          4653 non-null  int64
4   Age                  4653 non-null  int64
5   Gender               4653 non-null  object
6   EverBenched          4653 non-null  object
7   ExperienceInCurrentDomain 4653 non-null  int64
8   LeaveOrNot           4653 non-null  int64
dtypes: int64(5), object(4)
memory usage: 327.3+ KB
```

Overview of Analysis

- Connect to sqlite3 database
- Assess the properties of the database
- Establish columns/target and remove unwanted columns
- **Check the balance of the target values**
- **Establish the training set and train the model**
 - Calculate training accuracy score
 - Calculate testing set accuracy score
 - Calculate confusion matrix
 - Evaluate the classification report
 - Determine the importance rating of each feature

```
# Check the balance of our target values (1 = yes or 0 = no)
y.value_counts()

0    3053
1    1600
Name: LeaveOrNot, dtype: int64
```

```
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X,
                                                    y,
                                                    random_state=1,
                                                    stratify=y)


Counter(y_train)
Counter({1: 1200, 0: 2289})

X_train.shape
(3489, 14)
```

The observations for target 0 are 2x as those for target 1 therefore seasonably balanced but more experience is required to properly determine if this is a correct assumption

The train/test split is 75/25 and the y_train maintains the same pattern of 2 to 1 for 0 and 1.

Train set (3489, 14), test set (1164, 14) [3489 rows of data with 14 variables]



Overview of Analysis

- Connect to sqlite3 database
- Assess the properties of the database
- Establish columns/target and remove unwanted columns
- Check the balance of the target values
- Establish the training set and train the model
- **Calculate training accuracy score**
- **Calculate testing set accuracy score**
- **Calculate confusion matrix**
- **Evaluate the classification report**
- Determine the importance rating of each feature

```
# Calculated the balanced accuracy score: TRAINING
y_tr = brfc.predict(X_train)
balanced_accuracy_score(y_train, y_tr)

0.9137117737003058

# Calculated the balanced accuracy score: TESTING
y_pred = brfc.predict(X_test)
balanced_accuracy_score(y_test, y_pred)

0.7848625654450262
```

```
# Calculate confusion matrix.
cm = confusion_matrix(y_test, y_pred)
cm

array([[632, 132],
       [103, 297]], dtype=int64)
```

```
# classification report
from sklearn.metrics import classification_report
print(classification_report(y_test, y_pred))
```

	precision	recall	f1-score	support
0	0.86	0.83	0.84	764
1	0.69	0.74	0.72	400
accuracy			0.80	1164
macro avg	0.78	0.78	0.78	1164
weighted avg	0.80	0.80	0.80	1164

Accuracy score: how many predictions did the model got correct. *For binary classifications, balanced accuracy is equal to the arithmetic mean of **sensitivity** (true positive rate) and **specificity** (true negative rate)**

Accuracy score training: 91%,

Accuracy score test : 78%

=====

PRECISION: How reliable a positive classification is : Positive Predictive Value (PPV) for 0 = 86% and 1 = 69%

SENSITIVITY: How many items were correctly evaluated: 0 = 83% and 1 = 74%

SPECIFICITY: (true negative rate/ *recall of the negative class***

Overall, the model is better at predicting 0 than 1

F1 score combines precision and sensitivity, low value reflects major differences between precision and sensitivity. 80

Classification_report and classification_report imbalance yielded similar

results; unclear of the implications of that.

References

*https://scikit-learn.org/stable/modules/model_evaluation.html#classification-report

**https://scikit-learn.org/stable/modules/generated/sklearn.metrics.classification_report.html

Contains additional calculations that can be explored as the future analysis

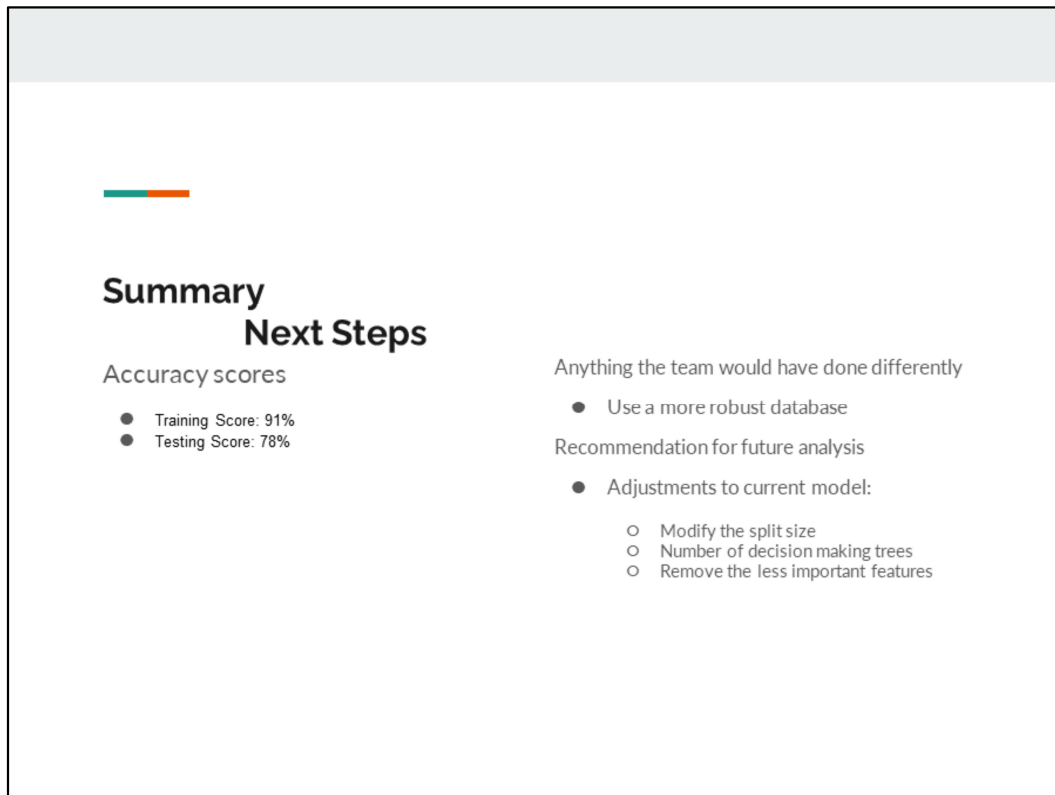


Overview of Analysis

- Connect to sqlite3 database
- Assess the properties of the database
- Establish columns/target and remove unwanted columns
- Check the balance of the target values
- Establish the training set and train the model
- Calculate training accuracy score
- Calculate testing set accuracy score
- Calculate confusion matrix
- Evaluate the classification report
- **Determine the importance rating of each feature**

```
# List the features sorted in descending order by feature importance
importances = brfc.feature_importances_
sorted(zip(brfc.feature_importances_, X.columns), reverse=True)
```

```
[(0.3086914210245113, 'JoiningYear'),
 (0.19190147118381745, 'Age'),
 (0.09890156348096112, 'ExperienceInCurrentDomain'),
 (0.09325636092904764, 'PaymentTier'),
 (0.05727410779116132, 'City_Pune'),
 (0.054031392148976856, 'Education_Masters'),
 (0.046093613205681915, 'Gender_Male'),
 (0.03876883320695893, 'Education_Bachelors'),
 (0.03204738082388843, 'Gender_Female'),
 (0.025990119766903574, 'City_Bangalore'),
 (0.025687810889074857, 'City_New_Delhi'),
 (0.009464044355872172, 'EverBenched_No'),
 (0.009311635638471467, 'EverBenched_Yes'),
 (0.008580245554672804, 'Education_PHD')]
```



Anything the team would have done differently

Use a more robust database with sufficient background information.

This database although clean and complete (no null values) seems to only focus on a very limited set of factors that seem more sided with what an Employer would consider when deciding to retain an employee and does not evaluate factors the Employees might be more included to consider when deciding to remain at a certain job including (1) training and development opportunities, (2) flexible telecommuting, (3) company culture/environment, (4) administrative/restructuring /leadership changes amongst others.

Some of the data is difficult to interpret due to lack of context, for example:

(1) the office location might be more significant if the database included residence information for each employee

(2) the salary tier is not informative because there are no details about the actual salary contained within the tiers