

# Luna Baalbaki - LSE Datathon 2022

Email: [l.baalbaki@lse.ac.uk](mailto:l.baalbaki@lse.ac.uk) (<mailto:l.baalbaki@lse.ac.uk>)

## Introduction:

IBM is struggling from a 16% employee attrition rate. But luckily they have great data that we can help to analyze. With the use of machine learning models and parameter optimization we can clearly point out to a few of the top reasons and signals that contribute to this high rate. We could also to help the HR department develop strategies to decrease the employee attrition rate.

```
In [ ]: 1 # Luna Baalbaki - LSE Datathon 2022
        2
        3 # Introduction:
        4 The aim behind this notebook is to predict employee turnover in the IBM
```

```
In [205]: 1 # Import libraries
          2 import seaborn as sns
          3 from sklearn.metrics import classification_report
          4 from sklearn.model_selection import train_test_split
          5 import pandas as pd
          6 import seaborn as sns
          7 import matplotlib as plt
          8 import matplotlib.pyplot as plt
          9 import numpy as np
         10 from numpy import mean, std, percentile
         11 from sklearn.preprocessing import LabelEncoder
         12 from sklearn.linear_model import LogisticRegression
         13 from sklearn.model_selection import cross_val_score, train_test_split,
         14 from sklearn.preprocessing import OneHotEncoder, OrdinalEncoder, Standa
         15 from sklearn.base import BaseEstimator, TransformerMixin
         16 from sklearn.exceptions import NotFittedError
         17 from sklearn.metrics import accuracy_score, mean_squared_error
         18 from sklearn.linear_model import LinearRegression
         19 from sklearn.svm import SVR
         20 from sklearn.tree import DecisionTreeRegressor
         21 from sklearn.ensemble import RandomForestRegressor, GradientBoostingReg
         22 from sklearn.feature_selection import RFECV
         23 from sklearn.compose import ColumnTransformer
         24 from sklearn.pipeline import Pipeline
```

```
In [206]: 1 df=pd.read_csv("/Users/user/Desktop/LSE/HR_IBM.csv")
```

In [207]: 1 df.head(10)

Out[207]:

	Age	Attrition	BusinessTravel	DailyRate	Department	DistanceFromHome	Education	Educatio
0	41	Yes	Travel_Rarely	1102	Sales	1	2	Life Sc
1	49	No	Travel_Frequently	279	Research & Development	8	1	Life Sc
2	37	Yes	Travel_Rarely	1373	Research & Development	2	2	
3	33	No	Travel_Frequently	1392	Research & Development	3	4	Life Sc
4	27	No	Travel_Rarely	591	Research & Development	2	1	N
5	32	No	Travel_Frequently	1005	Research & Development	2	2	Life Sc
6	59	No	Travel_Rarely	1324	Research & Development	3	3	N
7	30	No	Travel_Rarely	1358	Research & Development	24	1	Life Sc
8	38	No	Travel_Frequently	216	Research & Development	23	3	Life Sc
9	36	No	Travel_Rarely	1299	Research & Development	27	3	N

10 rows × 35 columns

## Data Cleaning

```
In [208]: 1 # Are there any duplicates in the dataset?
2 dups = df.duplicated()
3 print(dups.any())
4 print('Yay! No duplicates!')
```

False

Yay! No duplicates!

```
In [209]: 1 # Check for missing data
          2 df.isna().sum()
```

```
Out[209]: Age                                0
Attrition                                   0
BusinessTravel                             0
DailyRate                                  0
Department                                 0
DistanceFromHome                           0
Education                                  0
EducationField                             0
EmployeeCount                              0
EmployeeNumber                             0
EnvironmentSatisfaction                    0
Gender                                      0
HourlyRate                                 0
JobInvolvement                             0
JobLevel                                   0
JobRole                                    0
JobSatisfaction                            0
MaritalStatus                             0
MonthlyIncome                             0
MonthlyRate                                0
NumCompaniesWorked                         0
Over18                                     0
OverTime                                   0
PercentSalaryHike                          0
PerformanceRating                         0
RelationshipSatisfaction                   0
StandardHours                             0
StockOptionLevel                          0
TotalWorkingYears                         0
TrainingTimesLastYear                     0
WorkLifeBalance                           0
YearsAtCompany                            0
YearsInCurrentRole                        0
YearsSinceLastPromotion                   0
YearsWithCurrManager                      0
dtype: int64
```

```
In [210]: 1 # Display summary statistics
          2 df.describe()
```

Out[210]:

	Age	DailyRate	DistanceFromHome	Education	EmployeeCount	EmployeeNumb
<b>count</b>	1470.000000	1470.000000	1470.000000	1470.000000	1470.0	1470.000000
<b>mean</b>	36.923810	802.485714	9.192517	2.912925	1.0	1024.865300
<b>std</b>	9.135373	403.509100	8.106864	1.024165	0.0	602.024300
<b>min</b>	18.000000	102.000000	1.000000	1.000000	1.0	1.000000
<b>25%</b>	30.000000	465.000000	2.000000	2.000000	1.0	491.250000
<b>50%</b>	36.000000	802.000000	7.000000	3.000000	1.0	1020.500000
<b>75%</b>	43.000000	1157.000000	14.000000	4.000000	1.0	1555.750000
<b>max</b>	60.000000	1499.000000	29.000000	5.000000	1.0	2068.000000

8 rows × 26 columns

## Exploratory Data Analysis (EDA)

In this section, we will visualize the relationship between features and the target variable.

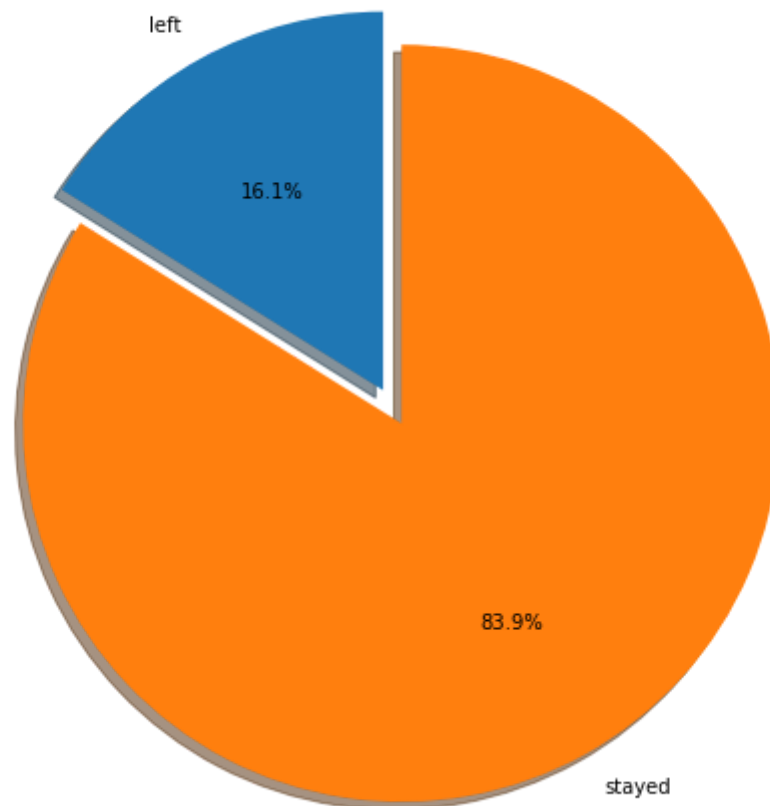
This would help us interpret which features affect customer churn through graphical representation.

```
In [211]: 1 # How many employees are churning?
          2 df['Attrition'].value_counts()
```

Out[211]: No 1233  
Yes 237  
Name: Attrition, dtype: int64

```
In [212]: 1 labels = 'left', 'stayed'
2 sizes = [df.Attrition[df['Attrition']=='Yes'].count(), df.Attrition[df[
3 explode = (0, 0.1)
4 fig1, ax1 = plt.subplots(figsize=(10, 8))
5 ax1.pie(sizes, explode=explode, labels=labels, autopct='%1.1f%%',
6         shadow=True, startangle=90)
7 ax1.axis('equal')
8 plt.title("Proportion of employees left and stayed", size = 20)
9 plt.show()
```

Proportion of employees left and stayed

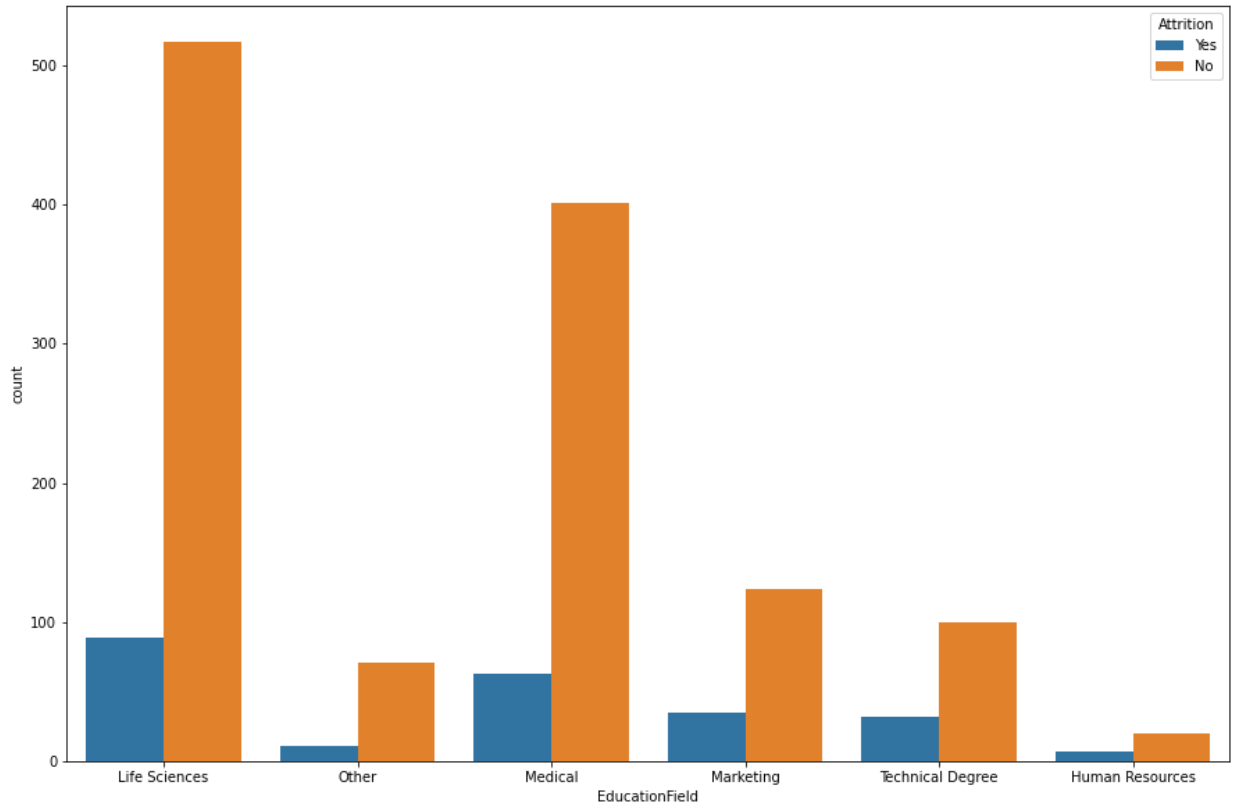


Note: 17.5% employees have decided to discontinue their services with the company, which is a major source of profit loss.

Moving forward, we aim to identify the attributes and characteristics of churners, in hopes of predicting the likelihood of future churners. Subsequently, we would be able to design personalized marketing campaigns to ensure employee retention.

```
In [222]: 1 # Visualize the churn count by education field
          2
          3
          4 fig = plt.gcf()
          5 fig.set_size_inches(15,10)
          6 sns.countplot(x='EducationField', hue = 'Attrition',data = df)
```

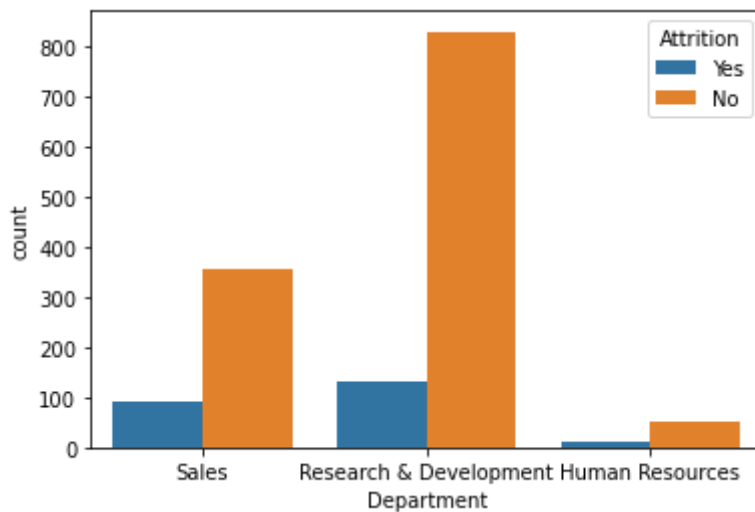
Out[222]: <matplotlib.axes.\_subplots.AxesSubplot at 0x7fd6aff8e1f0>



```
In [214]: 1 #People from Life Science Degrees are more likely to attrition
```

```
In [215]: 1 # Visualize the churn count for Department
          2 sns.countplot(x='Department', hue='Attrition', data= df)
          3
```

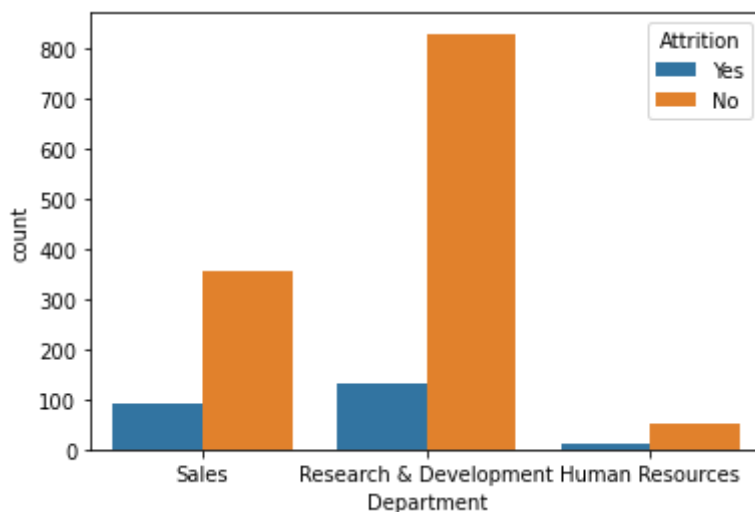
Out[215]: <matplotlib.axes.\_subplots.AxesSubplot at 0x7fd6b0812040>



```
In [216]: 1 #Employees from the Research Department left the company more than from
```

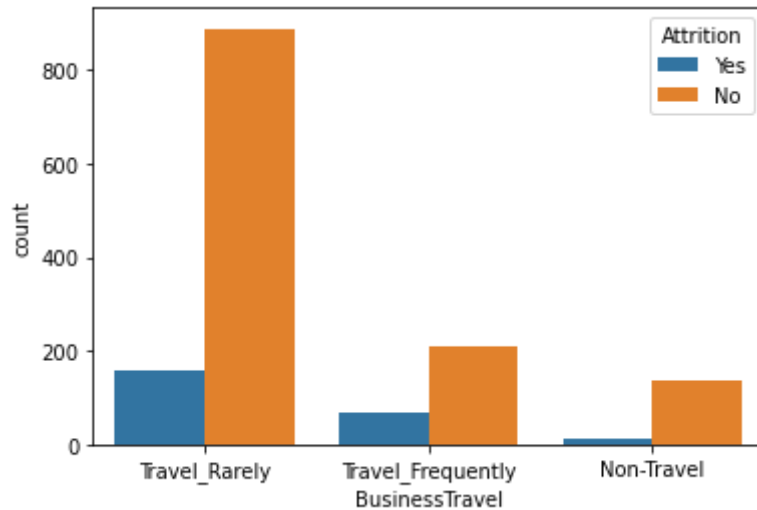
```
In [217]: 1 #People from Life Science Degrees are more likely to attrition
          2
          3 # Visualize the churn count for Department
          4 sns.countplot(x='Department', hue='Attrition', data= df)
          5
          6
          7 #Employees from the Research Department left the company more than from
```

Out[217]: <matplotlib.axes.\_subplots.AxesSubplot at 0x7fd6b24449d0>



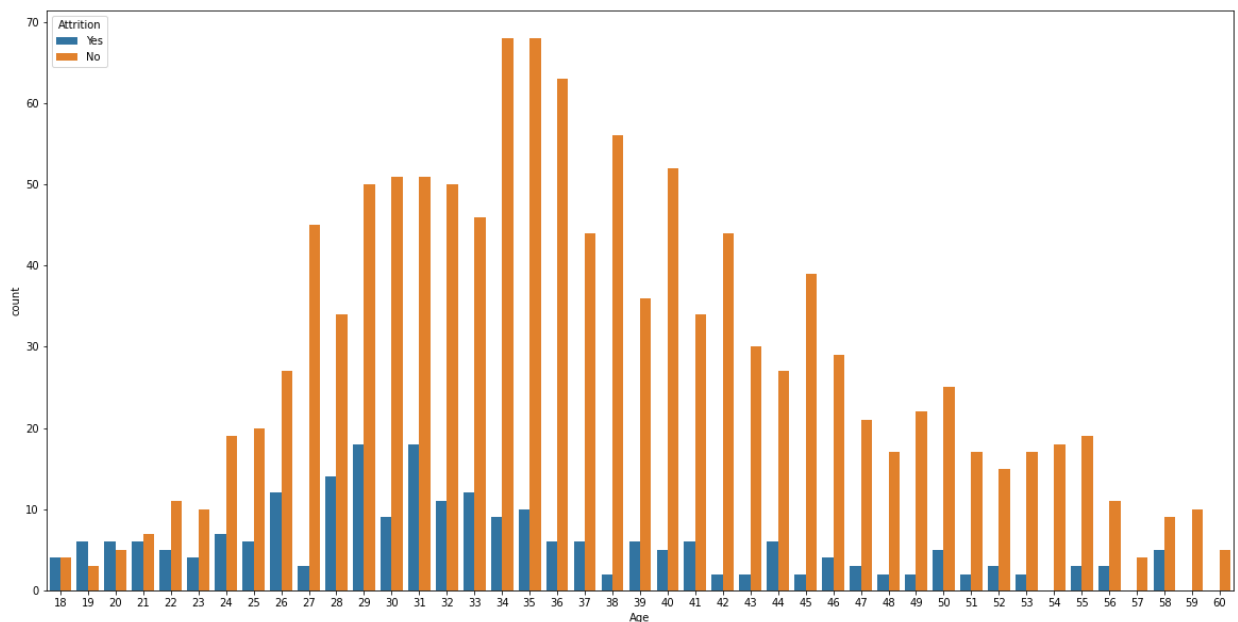
```
In [218]: 1 # Visualize the churn count for Business Travel
          2 sns.countplot(x='BusinessTravel', hue='Attrition', data= df)
```

Out[218]: <matplotlib.axes.\_subplots.AxesSubplot at 0x7fd6b113a400>



```
In [219]: 1 # Visualize the churn count for age
          2
          3
          4 fig = plt.gcf()
          5 fig.set_size_inches(20,10)
          6 sns.countplot(x='Age', hue='Attrition', data= df)
```

Out[219]: <matplotlib.axes.\_subplots.AxesSubplot at 0x7fd6aff90820>



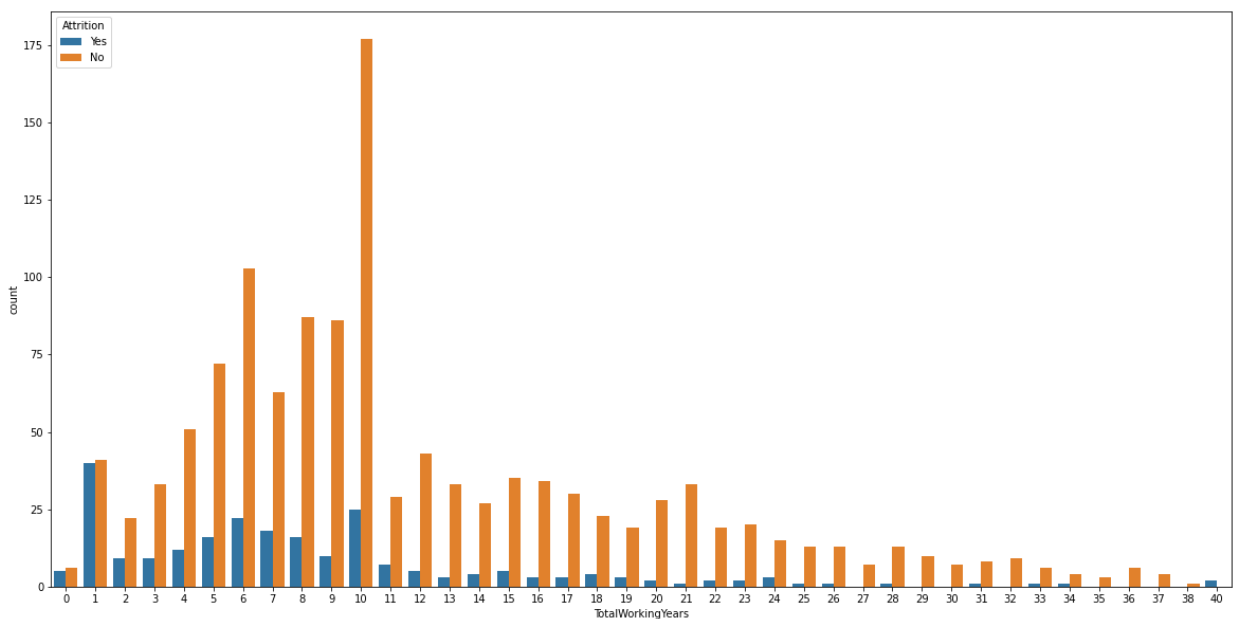


```
In [ ]: 1 import numpy as np; np.random.seed(0)
2 import matplotlib.pyplot as plt
3 import pandas as pd
4
5 x = np.random.rand(120)
6 df = pd.DataFrame({"x":x})
7
8 bins= [10,20,30,40]
9 plt.hist(df.values, bins=bins, edgecolor="k")
10 plt.xticks(bins)
11
12 plt.show()
```

```
In [178]: 1 #People in their early 30s tend to churn more
```

```
In [179]: 1 # Visualize the churn count for TotalWorkingYears
2
3
4 fig = plt.gcf()
5 fig.set_size_inches(20,10)
6 sns.countplot(x='TotalWorkingYears', hue='Attrition', data= df)
```

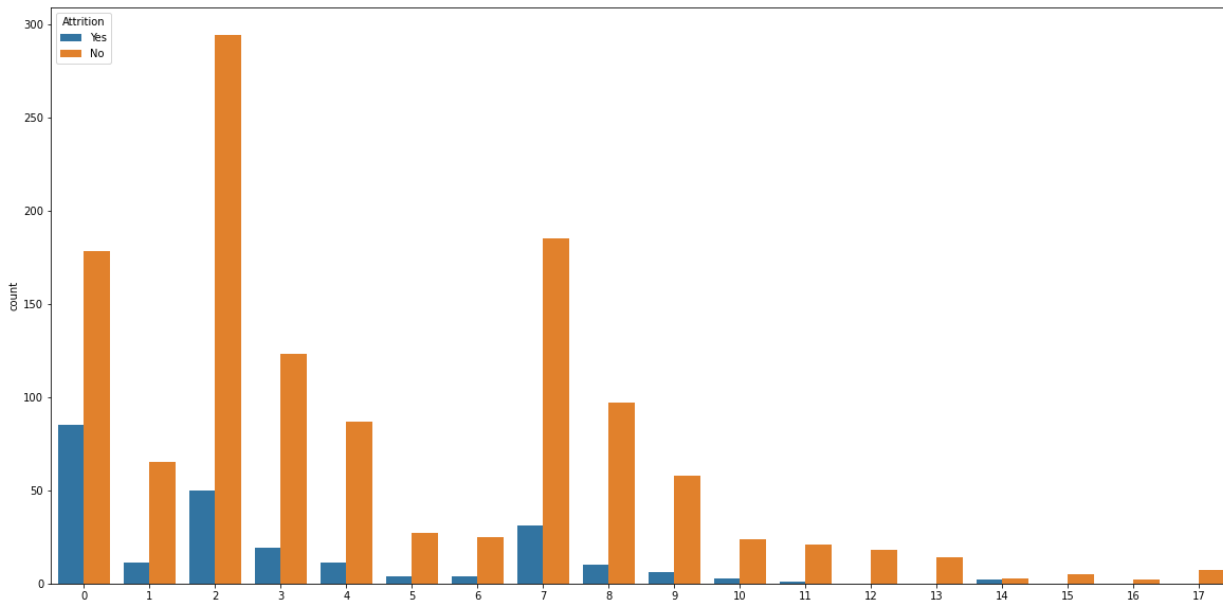
```
Out[179]: <matplotlib.axes._subplots.AxesSubplot at 0x7fd6b44aba00>
```



```
In [180]: 1 #People with less working years tend to attrition more
```

```
In [181]: 1 # Visualize the churn count for YearsWithCurrManager
2
3
4 fig = plt.gcf()
5 fig.set_size_inches(20,10)
6 sns.countplot(x='YearsWithCurrManager', hue='Attrition', data= df)
```

Out[181]: <matplotlib.axes.\_subplots.AxesSubplot at 0x7fd6b1b55f70>



## Feature Engineering

```
In [112]: 1 ### One-hot Encoding
```

```
In [182]: 1 #Split teh data into train and test
2 y = df['Attrition'].reset_index(drop=True)
3 X = df.drop(columns='Attrition')
4 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
5
```

```

In [184]: 1 # Defining a function for OneHotEncoding that retains feature names (S
2 # Reference: https://github.com/gdiepen/PythonScripts/blob/master/data
3
4 class DataFrameOneHotEncoder(BaseEstimator, TransformerMixin):
5     """Specialized version of OneHotEncoder that plays nice with pandas
6     will automatically set the feature/column names after fit/transform
7     """
8
9     def __init__(
10         self,
11         categories="auto",
12         drop=None,
13         sparse=None,
14         dtype=np.float64,
15         handle_unknown="error",
16         col_overrule_params={},
17     ):
18         """Create DataFrameOneHotEncoder that can be fitted to and tran
19         and that will set up the column/feature names automatically to
20         original_column_name[categorical_value]
21         If you provide the same arguments as you would for the sklearn
22         OneHotEncoder, these parameters will apply for all of the colum
23         to have specific overrides for some of the columns, provide the
24         argument col_overrule_params.
25
26         For example:
27         DataFrameOneHotEncoder(col_overrule_params={"col2":{"drop"
28         will create a OneHotEncoder for each of the columns with default
29         uses a drop=first argument for columns with the name col2
30         Args:
31             categories 'auto' or a list of array-like, default='auto'
32                 'auto' : Determine categories automatically from the t
33                 list : categories[i] holds the categories expected in
34                 The passed categories should not mix strings and numer
35                 within a single feature, and should be sorted in case
36                 values.
37             drop: {'first', 'if_binary'} or a array-like of shape (n_f
38                 default=None
39                 See OneHotEncoder documentation
40             sparse: Ignored, since we always will work with dense data
41             dtype: number type, default=float
42                 Desired dtype of output.
43             handle_unknown: {'error', 'ignore'}, default='error'
44                 Whether to raise an error or ignore if an unknown cate
45                 is present during transform (default is to raise). Whe
46                 is set to 'ignore' and an unknown category is encounte
47                 transform, the resulting one-hot encoded columns for th
48                 be all zeros. In the inverse transform, an unknown cate
49                 denoted as None.
50             col_overrule_params: dict of {column_name: dict_params} wh
51                 are exactly the options cateogires,drop,sparse,dtype,h
52                 For the column given by the key, these values will ove
53                 parameters
54         """
55         self.categories = categories
56         self.drop = drop

```

```

57     self.sparse = sparse
58     self.dtype = dtype
59     self.handle_unknown = handle_unknown
60     self.col_ouerrule_params = col_ouerrule_params
61     pass
62
63     def fit(self, X, y=None):
64         """Fit a separate OneHotEncoder for each of the columns in the
65         Args:
66             X: dataframe
67             y: None, ignored. This parameter exists only for compatibility
68                 Pipeline
69         Returns
70             self
71         Raises
72             TypeError if X is not of type DataFrame
73         """
74         if type(X) != pd.DataFrame:
75             raise TypeError(f"X should be of type dataframe, not {type(X)}")
76
77         self.onehotencoders_ = []
78         self.column_names_ = []
79
80         for c in X.columns:
81             # Construct the OHE parameters using the arguments
82             ohe_params = {
83                 "categories": self.categories,
84                 "drop": self.drop,
85                 "sparse": False,
86                 "dtype": self.dtype,
87                 "handle_unknown": self.handle_unknown,
88             }
89             # and update it with potential ouerrule parameters for the
90             ohe_params.update(self.col_ouerrule_params.get(c, {}))
91
92             # Regardless of how we got the parameters, make sure we always
93             # sparsity to False
94             ohe_params["sparse"] = False
95
96             # Now create, fit, and store the onehotencoder for current
97             ohe = OneHotEncoder(**ohe_params)
98             self.onehotencoders_.append(ohe.fit(X.loc[:, [c]]))
99
100             # Get the feature names and replace each x0_ with empty and
101             # surround the categorical value with [] and prefix it with
102             # column name
103             feature_names = ohe.get_feature_names()
104             feature_names = [x.replace("x0_", "") for x in feature_names]
105             feature_names = [f"{c}[{x}]" for x in feature_names]
106
107             self.column_names_.append(feature_names)
108
109         return self
110
111     def transform(self, X):
112         """Transform X using the one-hot-encoding per column
113         Args:

```

```

114         X: Dataframe that is to be one hot encoded
115     Returns:
116         Dataframe with onehotencoded data
117     Raises
118         NotFittedError if the transformer is not yet fitted
119         TypeError if X is not of type DataFrame
120     """
121     if type(X) != pd.DataFrame:
122         raise TypeError(f"X should be of type dataframe, not {type(X)}")
123
124     if not hasattr(self, "onehotencoders_"):
125         raise NotFittedError(f"{type(self).__name__} is not fitted")
126
127     all_df = []
128
129     for i, c in enumerate(X.columns):
130         ohe = self.onehotencoders_[i]
131
132         transformed_col = ohe.transform(X.loc[:, [c]])
133
134         df_col = pd.DataFrame(transformed_col, columns=self.column_names_[i])
135         all_df.append(df_col)
136
137     return pd.concat(all_df, axis=1)

```

In [185]:

```

1  # Let's use the new function for OneHotEncoding
2  # Reference: https://www.guidodiepen.nl/2021/02/keeping-column-names-when-one-hot-encoding/
3
4
5  df = X_train[['BusinessTravel', 'Department', 'EducationField', 'Gender']]
6  df_ohe = DataFrameOneHotEncoder().fit_transform(df)
7  ohe_done = pd.DataFrame(df_ohe)
8  # Merge both back into one dataframe
9  X_train.reset_index(inplace=True)
10 X_train = pd.concat([X_train, ohe_done], axis=1)
11 # Drop old columns
12 X_train = X_train.drop(['BusinessTravel', 'Department', 'EducationField'], axis=1)
13 X_train.reset_index(inplace=True)
14 X_train.drop('index', axis=1, inplace=True)
15 X_train.drop('level_0', axis=1, inplace=True)
16

```

```
In [186]: 1 X_train.head()
```

Out[186]:

	JobRole[Manufacturing Director]	JobRole[Research Director]	JobRole[Research Scientist]	JobRole[Sales Executive]	JobRole[Sales Representative]	MaritalS
	1.0	0.0	0.0	0.0	0.0	
	0.0	0.0	0.0	0.0	0.0	
	0.0	0.0	1.0	0.0	0.0	
	0.0	0.0	0.0	0.0	0.0	
	0.0	1.0	0.0	0.0	0.0	

```
In [187]: 1 # Let's not forget to encode the Target Variable using LabelEncoder!
          2 # But to avoid data leakage, let's split the dataset
```

```
In [188]: 1 # Now we encode the target variable
          2 from sklearn import preprocessing
          3 le = preprocessing.LabelEncoder()
          4 y_train = pd.DataFrame(le.fit_transform(y_train))
          5 y_test = pd.DataFrame(le.transform(y_test))
```

```
In [189]: 1 # Reset indices for easier interpretation
          2 pd.options.mode.chained_assignment = None # Disables harmless 'Setting
          3 X_train.reset_index(inplace=True)
          4 X_test.reset_index(inplace=True)
          5 X_train.drop('index', axis=1, inplace=True)
          6 X_test.drop('index', axis=1, inplace=True)
```

```
In [190]: 1 X_train.drop('Over18',axis='columns', inplace=True)
```

```
In [191]: 1 X_train.drop('OverTime',axis='columns', inplace=True)
```

```
In [192]: 1 # How does our revised train dataset look?
          2 X_train.head()
```

Out[192]:

	Age	DailyRate	DistanceFromHome	Education	EmployeeCount	EmployeeNumber	EnvironmentS
0	37	370	10	4	1	1809	
1	40	543	1	4	1	2012	
2	32	977	2	3	1	1671	
3	28	791	1	4	1	1286	
4	34	419	7	4	1	28	

5 rows x 52 columns

```
In [193]: 1 X_train.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1176 entries, 0 to 1175
Data columns (total 52 columns):
 #   Column                                          Non-Null Count  Dtype
---  -
 0   Age                                           1176 non-null   int64
 1   DailyRate                                     1176 non-null   int64
 2   DistanceFromHome                             1176 non-null   int64
 3   Education                                     1176 non-null   int64
 4   EmployeeCount                               1176 non-null   int64
 5   EmployeeNumber                             1176 non-null   int64
 6   EnvironmentSatisfaction                     1176 non-null   int64
 7   HourlyRate                                   1176 non-null   int64
 8   JobInvolvement                              1176 non-null   int64
 9   JobLevel                                     1176 non-null   int64
10   JobSatisfaction                             1176 non-null   int64
11   MonthlyIncome                               1176 non-null   int64
12   MonthlyRate                                 1176 non-null   int64
13   NumCompaniesWorked                         1176 non-null   int64
14   PercentSalaryHike                           1176 non-null   int64
15   PerformanceRating                           1176 non-null   int64
16   RelationshipSatisfaction                     1176 non-null   int64
17   StandardHours                              1176 non-null   int64
18   StockOptionLevel                           1176 non-null   int64
19   TotalWorkingYears                           1176 non-null   int64
20   TrainingTimesLastYear                       1176 non-null   int64
21   WorkLifeBalance                             1176 non-null   int64
22   YearsAtCompany                             1176 non-null   int64
23   YearsInCurrentRole                         1176 non-null   int64
24   YearsSinceLastPromotion                     1176 non-null   int64
25   YearsWithCurrManager                       1176 non-null   int64
26   BusinessTravel[Non-Travel]                  1176 non-null   float64
27   BusinessTravel[Travel_Frequently]           1176 non-null   float64
28   BusinessTravel[Travel_Rarely]               1176 non-null   float64
29   Department[Human Resources]                 1176 non-null   float64
30   Department[Research & Development]          1176 non-null   float64
31   Department[Sales]                           1176 non-null   float64
32   EducationField[Human Resources]              1176 non-null   float64
33   EducationField[Life Sciences]               1176 non-null   float64
34   EducationField[Marketing]                   1176 non-null   float64
35   EducationField[Medical]                     1176 non-null   float64
36   EducationField[Other]                       1176 non-null   float64
37   EducationField[Technical Degree]             1176 non-null   float64
38   Gender[Female]                              1176 non-null   float64
39   Gender[Male]                                1176 non-null   float64
40   JobRole[Healthcare Representative]           1176 non-null   float64
41   JobRole[Human Resources]                    1176 non-null   float64
42   JobRole[Laboratory Technician]              1176 non-null   float64
43   JobRole[Manager]                            1176 non-null   float64
44   JobRole[Manufacturing Director]              1176 non-null   float64
45   JobRole[Research Director]                  1176 non-null   float64
46   JobRole[Research Scientist]                 1176 non-null   float64
47   JobRole[Sales Executive]                    1176 non-null   float64
48   JobRole[Sales Representative]                1176 non-null   float64
49   MaritalStatus[Divorced]                     1176 non-null   float64
```

```
50 MaritalStatus[Married]          1176 non-null    float64
51 MaritalStatus[Single]          1176 non-null    float64
dtypes: float64(26), int64(26)
memory usage: 477.9 KB
```

## Model Optimization and Selection

```
In [194]: 1 #Support vetctor Machine
2 svr = SVR(kernel='rbf')
3 svr.fit(X_train, y_train.values.ravel())
4
5 svr_t_predictions = svr.predict(X_train)
6 svr_train_rmse = mean_squared_error(y_train, svr_t_predictions, squared=
7 print(f"RMSE for testing: {svr_train_rmse:.1f}")
8
```

RMSE for testing: 0.4

```
In [195]: 1 #Linear Regression
2 lr = LinearRegression()
3 lr.fit(X_train, y_train.values.ravel())
4
5 lr_tr_predictions = lr.predict(X_train)
6 lr_train_rmse = mean_squared_error(y_train, lr_tr_predictions, squared=
7 print(f"RMSE for testing: {lr_train_rmse:.1f}")
8
9
```

RMSE for testing: 0.3

```
In [196]: 1 #random Forest Regressor
2 rfr = RandomForestRegressor(n_estimators=10, random_state=1111)
3 rfr.fit(X_train, y_train.values.ravel())
4
5 rfr_tr_predictions = rfr.predict(X_train)
6 rfr_train_rmse = mean_squared_error(y_train, rfr_tr_predictions, square
7 print(f"RMSE for testing: {rfr_train_rmse:.1f}")
8
9
```

RMSE for testing: 0.2

## Note

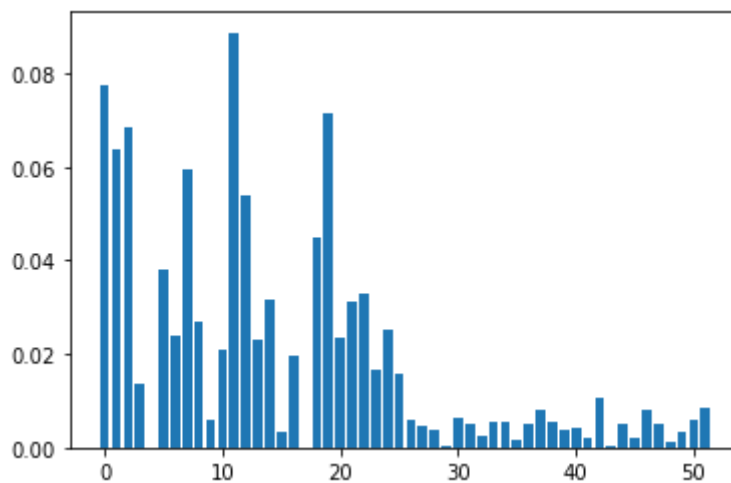
Random Forrest Regresosr Model has the least error among the other Machine Learning models. Thus, we go with Random Forrest Regressor

## Parameter Optimization



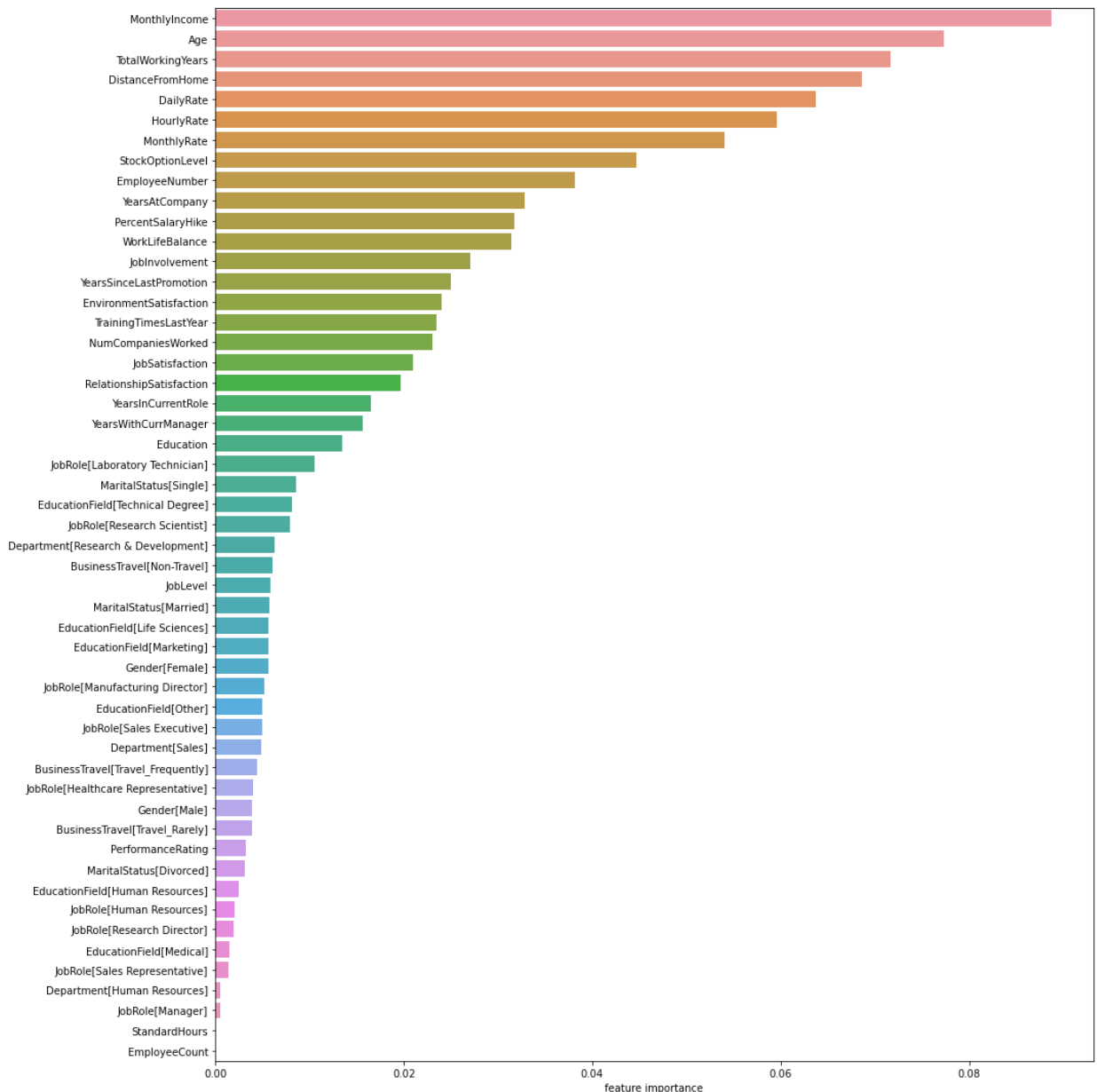
```
In [197]: 1 #Paraemter Optimization
2 # Optimize the parameters using GridSearchCV
3 parameters = {'n_estimators': [4, 6, 9],
4               'max_features': ['log2', 'sqrt', 'auto'],
5               'max_depth': [2, 3, 5, 10],
6               'min_samples_split': [2, 3, 5],
7               'min_samples_leaf': [1, 5, 8]
8             }
9 rfr_cv = GridSearchCV(estimator= rfr, param_grid = parameters, scoring
10 rfr_cv = rfr_cv.fit(X_train, y_train.values.ravel())
```

```
In [198]: 1 #Feature Selection
2
3 from matplotlib import pyplot
4 # get importance
5 importance = rfr.feature_importances_
6
7 ## summarize feature importance (commented to save space)
8 #for i,v in enumerate(importance):
9     #print('Feature: %0d, Score: %.5f' % (i,v))
10
11 # plot feature importance
12 pyplot.bar([x for x in range(len(importance))], importance)
13 pyplot.show()
```



## Feature Selection

```
In [200]: 1 ranking = np.argsort(-rfr.feature_importances_)
2 f, ax = plt.subplots(figsize=(15, 15))
3 sns.barplot(x=rfr.feature_importances_[ranking], y=X_train.columns.values)
4 ax.set_xlabel("feature importance")
5 plt.tight_layout()
6 plt.show()
```



## Note:

Based on the model feature importance, the top 3 features are (Monthly Income, Age and Total Working years). These are the features that affect the willingness of the employee to leave the company or not.

This seems logical as an employee's salary is one of the biggest reasons for an employee to leave their job or not.

The second affecting variable in employee attrition is the age of an employee. So as it shows that age plays a significant role in the attrition of the employee. So here we can inform the HR department to put some strategies depending on ages and especially the employees in their late 20s and early 30s as we have seen earlier in the explanatory analysis that they tend to leave the company more than others. SO the HR can target their strategy more towards employees of these ages and focus on them more so that they won't leave the company.

The third affecting variable in employee attrition is total working years. The longer an employee stays in the company the less likely they tend to leave the company. Here the HR department could work on strategies so that the employee would stay longer in the company.

In [ ]:

1