

Employee Attrition Prediction

May 12, 2021

0.1 Problem Statement

Employee attrition is downsizing in any organization where employees resign. Employees are valuable assets of any organization. It is necessary to know whether the employees are dissatisfied or whether there are other reasons for leaving their respective jobs. Nowadays, for better opportunities, employees are eager to move from one organization to another. But if they quit their jobs unexpectedly, it can result in a huge loss for the organization. A new hire will consume money and time, and newly hired employees will also take time to make the respective organization profitable.

Retaining skilled and hardworking employees is one of the most critical challenges many organizations face. Therefore, by improving employee satisfaction and providing a desirable working environment, we can certainly reduce this problem significantly

```
[1]: import numpy as np # linear algebra
import pandas as pd # data processing
import seaborn as sns
import matplotlib.pyplot as plt
%matplotlib inline

# Import statements required for Plotly
import plotly.offline as py
py.init_notebook_mode(connected=True)
import plotly.graph_objs as go
import plotly.tools as tls

from sklearn.ensemble import RandomForestClassifier, GradientBoostingClassifier
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import (accuracy_score, log_loss, classification_report)
from imblearn.over_sampling import SMOTE
import xgboost
```

```
[2]: attrition = pd.read_csv(r"C:
↳\Users\OSAGIE\Desktop\WA_Fn-UseC_-HR-Employee-Attrition.csv")
```

```
[3]: attrition.head()
```

```
[3]:   Age Attrition   BusinessTravel   DailyRate   Department \
0    41         Yes   Travel_Rarely    1102         Sales
```

1	49	No	Travel_Frequently	279	Research & Development
2	37	Yes	Travel_Rarely	1373	Research & Development
3	33	No	Travel_Frequently	1392	Research & Development
4	27	No	Travel_Rarely	591	Research & Development

	DistanceFromHome	Education	EducationField	EmployeeCount	EmployeeNumber	\
0	1	2	Life Sciences	1	1	
1	8	1	Life Sciences	1	2	
2	2	2	Other	1	4	
3	3	4	Life Sciences	1	5	
4	2	1	Medical	1	7	

	RelationshipSatisfaction	StandardHours	StockOptionLevel	\
0	...	1	80	0
1	...	4	80	1
2	...	2	80	0
3	...	3	80	0
4	...	4	80	1

	TotalWorkingYears	TrainingTimesLastYear	WorkLifeBalance	YearsAtCompany	\
0	8	0	1	6	
1	10	3	3	10	
2	7	3	3	0	
3	8	3	3	8	
4	6	3	3	2	

	YearsInCurrentRole	YearsSinceLastPromotion	YearsWithCurrManager
0	4	0	5
1	7	1	7
2	0	0	0
3	7	3	0
4	2	2	2

[5 rows x 35 columns]

```
[4]: f, axes = plt.subplots(3, 3, figsize=(10, 8),
                           sharex=False, sharey=False)

# Defining our colormap scheme
s = np.linspace(0, 3, 10)
cmap = sns.cubehelix_palette(start=0.0, light=1, as_cmap=True)

# Generate and plot
x = attrition['Age'].values
y = attrition['TotalWorkingYears'].values
sns.kdeplot(x, y, cmap=cmap, shade=True, cut=5, ax=axes[0,0])
axes[0,0].set( title = 'Age against Total working years')
```

```

cmap = sns.cubehelix_palette(start=0.333333333333, light=1, as_cmap=True)
# Generate and plot
x = attrition['Age'].values
y = attrition['DailyRate'].values
sns.kdeplot(x, y, cmap=cmap, shade=True, ax=axes[0,1])
axes[0,1].set( title = 'Age against Daily Rate')

cmap = sns.cubehelix_palette(start=0.666666666667, light=1, as_cmap=True)
# Generate and plot
x = attrition['YearsInCurrentRole'].values
y = attrition['Age'].values
sns.kdeplot(x, y, cmap=cmap, shade=True, ax=axes[0,2])
axes[0,2].set( title = 'Years in role against Age')

cmap = sns.cubehelix_palette(start=1.0, light=1, as_cmap=True)
# Generate and plot
x = attrition['DailyRate'].values
y = attrition['DistanceFromHome'].values
sns.kdeplot(x, y, cmap=cmap, shade=True, ax=axes[1,0])
axes[1,0].set( title = 'Daily Rate against DistancefromHome')

cmap = sns.cubehelix_palette(start=1.333333333333, light=1, as_cmap=True)
# Generate and plot
x = attrition['DailyRate'].values
y = attrition['JobSatisfaction'].values
sns.kdeplot(x, y, cmap=cmap, shade=True, ax=axes[1,1])
axes[1,1].set( title = 'Daily Rate against Job satisfaction')

cmap = sns.cubehelix_palette(start=1.666666666667, light=1, as_cmap=True)
# Generate and plot
x = attrition['YearsAtCompany'].values
y = attrition['JobSatisfaction'].values
sns.kdeplot(x, y, cmap=cmap, shade=True, ax=axes[1,2])
axes[1,2].set( title = 'Daily Rate against distance')

cmap = sns.cubehelix_palette(start=2.0, light=1, as_cmap=True)
# Generate and plot
x = attrition['YearsAtCompany'].values
y = attrition['DailyRate'].values
sns.kdeplot(x, y, cmap=cmap, shade=True, ax=axes[2,0])
axes[2,0].set( title = 'Years at company against Daily Rate')

cmap = sns.cubehelix_palette(start=2.333333333333, light=1, as_cmap=True)
# Generate and plot
x = attrition['RelationshipSatisfaction'].values
y = attrition['YearsWithCurrManager'].values

```

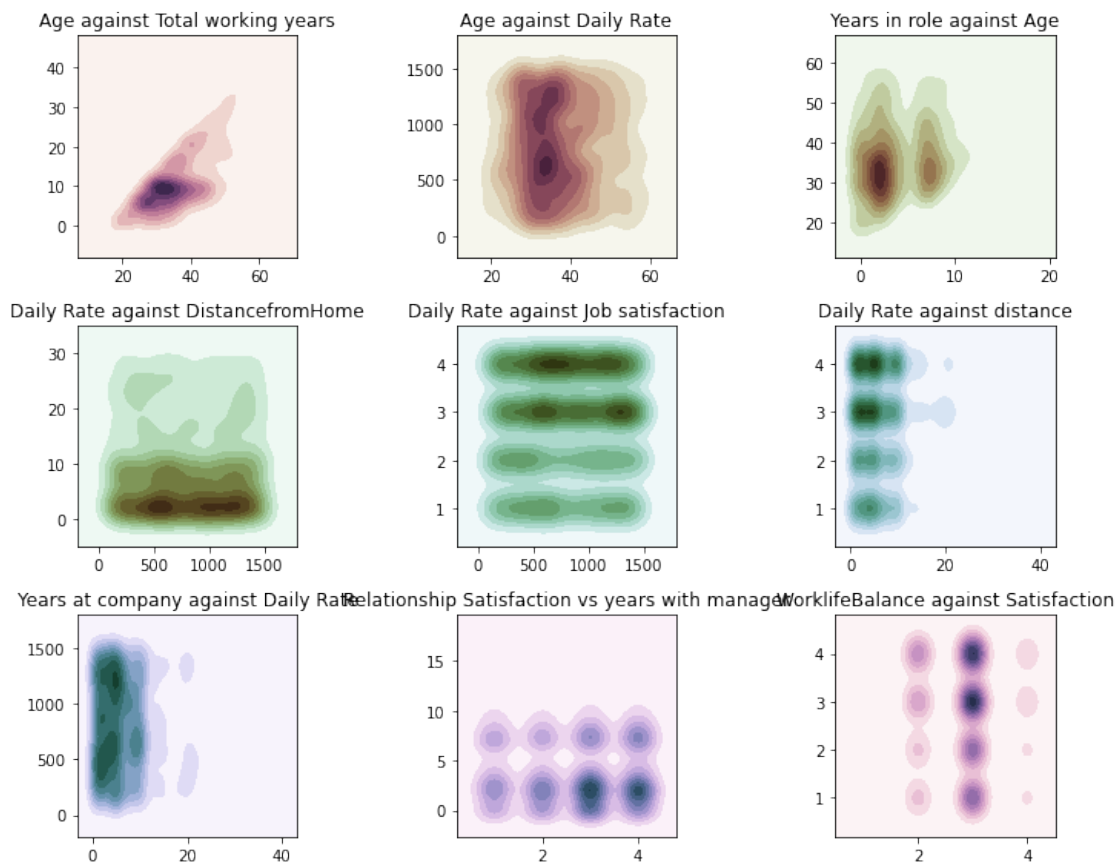
```

sns.kdeplot(x, y, cmap=cmap, shade=True, ax=axes[2,1])
axes[2,1].set( title = 'Relationship Satisfaction vs years with manager')

cmap = sns.cubehelix_palette(start=2.6666666666667, light=1, as_cmap=True)
# Generate and plot
x = attrition['WorkLifeBalance'].values
y = attrition['JobSatisfaction'].values
sns.kdeplot(x, y, cmap=cmap, shade=True, ax=axes[2,2])
axes[2,2].set( title = 'WorklifeBalance against Satisfaction')

f.tight_layout()

```



0.2 In this correlation plot, I will be using the Plotly library in Python to produce an interactive Pearson correlation matrix via the Heatmap function as follows:

```

[5]: # Define a dictionary for the target mapping
target_map = {'Yes':1, 'No':0}
# Use the pandas apply method to numerically encode our attrition target_
→variable

```

```

attrition["Attrition_numerical"] = attrition["Attrition"].apply(lambda x:
↳target_map[x])

# creating a list of only numerical values
numerical = [u'Age', u'DailyRate', u'DistanceFromHome',
             u'Education', u'EmployeeNumber', u'EnvironmentSatisfaction',
             u'HourlyRate', u'JobInvolvement', u'JobLevel', u'JobSatisfaction',
             u'MonthlyIncome', u'MonthlyRate', u'NumCompaniesWorked',
             u'PercentSalaryHike', u'PerformanceRating',
↳u'RelationshipSatisfaction',
             u'StockOptionLevel', u'TotalWorkingYears',
             u'TrainingTimesLastYear', u'WorkLifeBalance', u'YearsAtCompany',
             u'YearsInCurrentRole',
↳u'YearsSinceLastPromotion', u'YearsWithCurrManager']

data = [
    go.Heatmap(
        z= attrition[numerical].astype(float).corr().values, # Generating the
↳Pearson correlation
        x=attrition[numerical].columns.values,
        y=attrition[numerical].columns.values,
        colorscale='Viridis',
        reversescale = False,
#         text = True ,
        opacity = 1.0

    )
]

layout = go.Layout(
    title='Pearson Correlation of numerical features',
    xaxis = dict(ticks='', nticks=36),
    yaxis = dict(ticks=''),
    width = 900, height = 700,

)

fig = go.Figure(data=data, layout=layout)
py.iplot(fig, filename='labelled-heatmap')

```

From the correlation plot, we can see that a lot of our columns appear to be poorly correlated to each other. Generally, when building a predictive model, it would be better to train a model with features that are not too correlated with each other so that we don't need to deal with redundant features.

```
[7]: attrition = attrition.drop(['Attrition_numerical'], axis=1)
```

```
# Empty list to store columns with categorical data
categorical = []
for col, value in attrition.iteritems():
    if value.dtype == 'object':
        categorical.append(col)

# Store the numerical columns in a list numerical
numerical = attrition.columns.difference(categorical)
```

```
[8]: numerical
```

```
[8]: Index(['Age', 'DailyRate', 'DistanceFromHome', 'Education', 'EmployeeCount',
        'EmployeeNumber', 'EnvironmentSatisfaction', 'HourlyRate',
        'JobInvolvement', 'JobLevel', 'JobSatisfaction', 'MonthlyIncome',
        'MonthlyRate', 'NumCompaniesWorked', 'PercentSalaryHike',
        'PerformanceRating', 'RelationshipSatisfaction', 'StandardHours',
        'StockOptionLevel', 'TotalWorkingYears', 'TrainingTimesLastYear',
        'WorkLifeBalance', 'YearsAtCompany', 'YearsInCurrentRole',
        'YearsSinceLastPromotion', 'YearsWithCurrManager'],
        dtype='object')
```

After identifying which of our features contain categorical data, we can start to digitally encode the data. To do this, I'll use Pandas' `get_dummies` method in Python which creates dummy variables encoded from the categorical variables:

```
[9]: attrition_cat = attrition[categorical]
attrition_cat = attrition_cat.drop(['Attrition'], axis=1) # Dropping the target_
    ↪ column
attrition_cat = pd.get_dummies(attrition_cat)
attrition_cat.head(3)
attrition_num = attrition[numerical]
attrition_final = pd.concat([attrition_num, attrition_cat], axis=1)
```

One last step we need to remember is to generate our target variable. The target, in this case, is given by the Attrition column which contains categorical variables therefore requires numeric coding. We digitally encode it by creating a dictionary with the given mapping as 1: Yes and 0: No

```
[10]: target_map = {'Yes':1, 'No':0}
# Use the pandas apply method to numerically encode our attrition target_
    ↪ variable
target = attrition["Attrition"].apply(lambda x: target_map[x])
```

But before implementing Machine Learning for prediction of Employee Attrition prediction we need to split the data into a training set and test set:

```
[13]: from sklearn.model_selection import train_test_split
      from sklearn.model_selection import StratifiedShuffleSplit

      # Split data into train and test sets as well as for validation and testing
      train, test, target_train, target_val = train_test_split(attrition_final,
                                                                target,
                                                                train_size= 0.80,
                                                                random_state=0);

      #train, test, target_train, target_val =
      ↪StratifiedShuffleSplit(attrition_final, target, random_state=0)
```

Now let's train the Random forest classification model for the task of Employee Attrition prediction using Machine Learning and Python:

```
[14]: oversampler=SMOTE(random_state=0)
      smote_train, smote_target = oversampler.fit_sample(train,target_train)

      seed = 0  # We set our random seed to zero for reproducibility
      # Random Forest parameters
      rf_params = {
          'n_jobs': -1,
          'n_estimators': 1000,
          # 'warm_start': True,
          'max_features': 0.3,
          'max_depth': 4,
          'min_samples_leaf': 2,
          'max_features' : 'sqrt',
          'random_state' : seed,
          'verbose': 0
      }

      rf = RandomForestClassifier(**rf_params)
      rf.fit(smote_train, smote_target)
      rf_predictions = rf.predict(test)
      print("Accuracy score: {}".format(accuracy_score(target_val, rf_predictions)))
      print("="*80)
      print(classification_report(target_val, rf_predictions))
```

Accuracy score: 0.8537414965986394

```
=====
```

	precision	recall	f1-score	support
0	0.90	0.93	0.91	245
1	0.57	0.49	0.53	49
accuracy			0.85	294
macro avg	0.74	0.71	0.72	294
weighted avg	0.85	0.85	0.85	294

As observed, our Random Forest returns around 88% accuracy for its predictions and at first glance, this may seem like a fairly good mode

Sklearn's Random Forest classifier also contains a very handy attribute for analyzing feature importance which tells us which features in our dataset have received the most importance by the Random Forest algorithm. Let's visualize the features taken into account by our machine learning model for employee attrition:

```
[15]: trace = go.Scatter(
    y = rf.feature_importances_,
    x = attrition_final.columns.values,
    mode='markers',
    marker=dict(
        sizemode = 'diameter',
        sizeref = 1,
        size = 13,
        #size= rf.feature_importances_,
        #color = np.random.randn(500), #set color equal to a variable
        color = rf.feature_importances_,
        colorscale='Portland',
        showscale=True
    ),
    text = attrition_final.columns.values
)
data = [trace]

layout= go.Layout(
    autosize= True,
    title= 'Random Forest Feature Importance',
    hovermode= 'closest',
    xaxis= dict(
        ticklen= 5,
        showgrid=False,
        zeroline=False,
        showline=False
    ),
    yaxis=dict(
        title= 'Feature Importance',
        showgrid=False,
        zeroline=False,
        ticklen= 5,
        gridwidth= 2
    ),
    showlegend= False
)
fig = go.Figure(data=data, layout=layout)
```



```
py.ipplot(fig,filename='scatter2010')
```

```
[ ]:
```