

# Activity\_ Course 7 Salifort Motors project lab

December 12, 2024

## 1 Capstone project: Providing data-driven suggestions for HR

### 1.1 Description and deliverables

This capstone project is an opportunity for you to analyze a dataset and build predictive models that can provide insights to the Human Resources (HR) department of a large consulting firm.

Upon completion, you will have two artifacts that you would be able to present to future employers. One is a brief one-page summary of this project that you would present to external stakeholders as the data professional in Salifort Motors. The other is a complete code notebook provided here. Please consider your prior course work and select one way to achieve this given project question. Either use a regression model or machine learning model to predict whether or not an employee will leave the company. The exemplar following this activity shows both approaches, but you only need to do one.

In your deliverables, you will include the model evaluation (and interpretation if applicable), a data visualization(s) of your choice that is directly related to the question you ask, ethical considerations, and the resources you used to troubleshoot and find answers or solutions.

## 2 PACE stages

### 2.1 Pace: Plan

Consider the questions in your PACE Strategy Document to reflect on the Plan stage.

In this stage, consider the following:

#### 2.1.1 Understand the business scenario and problem

The HR department at Salifort Motors wants to take some initiatives to improve employee satisfaction levels at the company. They collected data from employees, but now they don't know what to do with it. They refer to you as a data analytics professional and ask you to provide data-driven suggestions based on your understanding of the data. They have the following question: what's likely to make the employee leave the company?

Your goals in this project are to analyze the data collected by the HR department and to build a model that predicts whether or not an employee will leave the company.

If you can predict employees likely to quit, it might be possible to identify factors that contribute to their leaving. Because it is time-consuming and expensive to find, interview, and hire new employees, increasing employee retention will be beneficial to the company.

### 2.1.2 Familiarize yourself with the HR dataset

The dataset that you'll be using in this lab contains 15,000 rows and 10 columns for the variables listed below.

**Note:** you don't need to download any data to complete this lab. For more information about the data, refer to its source on [Kaggle](#).

Variable	Description
satisfaction_level	Employee-reported job satisfaction level [0–1]
last_evaluation	Score of employee's last performance review [0–1]
number_project	Number of projects employee contributes to
average_monthly_hours	Average number of hours employee worked per month
time_spend_company	How long the employee has been with the company (years)
Work_accident	Whether or not the employee experienced an accident while at work
left	Whether or not the employee left the company
promotion_last_5years	Whether or not the employee was promoted in the last 5 years
Department	The employee's department
salary	The employee's salary (U.S. dollars)

## 2.2 Step 1. Imports

- Import packages
- Load dataset

### 2.2.1 Import packages

```
[2]: # Import packages
    ### YOUR CODE HERE ###
    import pandas as pd
    import numpy as np
```

```

import matplotlib.pyplot as plt
import seaborn as sns

from xgboost import XGBClassifier
from xgboost import XGBRegressor
from xgboost import plot_importance

from sklearn.linear_model import LogisticRegression
from sklearn.tree import DecisionTreeClassifier
from sklearn.ensemble import RandomForestClassifier

from sklearn.model_selection import GridSearchCV, train_test_split
from sklearn.metrics import accuracy_score, precision_score, recall_score, \
    f1_score, confusion_matrix, ConfusionMatrixDisplay, classification_report
from sklearn.metrics import roc_auc_score, roc_curve
from sklearn.tree import plot_tree

import pickle

```

### 2.2.2 Load dataset

Pandas is used to read a dataset called `HR_capstone_dataset.csv`. As shown in this cell, the dataset has been automatically loaded in for you. You do not need to download the .csv file, or provide more code, in order to access the dataset and proceed with this lab. Please continue with this activity by completing the following instructions.

```

[3]: # RUN THIS CELL TO IMPORT YOUR DATA.

# Load dataset into a dataframe
### YOUR CODE HERE ###
df = pd.read_csv("HR_capstone_dataset.csv")

# Display first few rows of the dataframe
### YOUR CODE HERE ###
df.head()

```

```

[3]:
  satisfaction_level  last_evaluation  number_project  average_monthly_hours \
0                0.38              0.53              2                157
1                0.80              0.86              5                262
2                0.11              0.88              7                272
3                0.72              0.87              5                223
4                0.37              0.52              2                159

  time_spend_company  Work_accident  left  promotion_last_5years  Department \
0                  3              0    1              0          sales

```

1	6	0	1	0	sales
2	4	0	1	0	sales
3	5	0	1	0	sales
4	3	0	1	0	sales

	salary
0	low
1	medium
2	medium
3	low
4	low

## 2.3 Step 2. Data Exploration (Initial EDA and data cleaning)

- Understand your variables
- Clean your dataset (missing data, redundant data, outliers)

### 2.3.1 Gather basic information about the data

```
[4]: # Gather basic information about the data
    ### YOUR CODE HERE ###
    df.describe()
```

```
[4]:
```

	satisfaction_level	last_evaluation	number_project \	
count	14999.000000	14999.000000	14999.000000	
mean	0.612834	0.716102	3.803054	
std	0.248631	0.171169	1.232592	
min	0.090000	0.360000	2.000000	
25%	0.440000	0.560000	3.000000	
50%	0.640000	0.720000	4.000000	
75%	0.820000	0.870000	5.000000	
max	1.000000	1.000000	7.000000	

	average_monthly_hours	time_spend_company	Work_accident	left \
count	14999.000000	14999.000000	14999.000000	14999.000000
mean	201.050337	3.498233	0.144610	0.238083
std	49.943099	1.460136	0.351719	0.425924
min	96.000000	2.000000	0.000000	0.000000
25%	156.000000	3.000000	0.000000	0.000000
50%	200.000000	3.000000	0.000000	0.000000
75%	245.000000	4.000000	0.000000	0.000000
max	310.000000	10.000000	1.000000	1.000000

	promotion_last_5years
count	14999.000000

```

mean          0.021268
std           0.144281
min           0.000000
25%           0.000000
50%           0.000000
75%           0.000000
max           1.000000

```

### 2.3.2 Gather descriptive statistics about the data

```

[5]: # Gather descriptive statistics about the data
    ## YOUR CODE HERE ##
    df.info()

```

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 14999 entries, 0 to 14998
Data columns (total 10 columns):
 #   Column                Non-Null Count  Dtype
---  -
 0   satisfaction_level     14999 non-null  float64
 1   last_evaluation        14999 non-null  float64
 2   number_project         14999 non-null  int64
 3   average_monthly_hours  14999 non-null  int64
 4   time_spend_company     14999 non-null  int64
 5   Work_accident          14999 non-null  int64
 6   left                   14999 non-null  int64
 7   promotion_last_5years  14999 non-null  int64
 8   Department             14999 non-null  object
 9   salary                 14999 non-null  object
dtypes: float64(2), int64(6), object(2)
memory usage: 1.1+ MB

```

### 2.3.3 Rename columns

As a data cleaning step, rename the columns as needed. Standardize the column names so that they are all in `snake_case`, correct any column names that are misspelled, and make column names more concise as needed.

```

[6]: # Display all column names
    ## YOUR CODE HERE ##
    df.columns

```

```

[6]: Index(['satisfaction_level', 'last_evaluation', 'number_project',
          'average_monthly_hours', 'time_spend_company', 'Work_accident', 'left',
          'promotion_last_5years', 'Department', 'salary'],
          dtype='object')

```

```
[7]: # Rename columns as needed
### YOUR CODE HERE ###
df = df.rename(columns={'Work_accident': 'work_accident',
                        'average_montly_hours': 'average_monthly_hours',
                        'time_spend_company': 'tenure',
                        'Department': 'department'})

# Display all column names after the update
### YOUR CODE HERE ###
df.columns
```

```
[7]: Index(['satisfaction_level', 'last_evaluation', 'number_project',
          'average_monthly_hours', 'tenure', 'work_accident', 'left',
          'promotion_last_5years', 'department', 'salary'],
          dtype='object')
```

### 2.3.4 Check missing values

Check for any missing values in the data.

```
[8]: # Check for missing values
### YOUR CODE HERE ###
df.isna().sum()
```

```
[8]: satisfaction_level      0
last_evaluation             0
number_project              0
average_monthly_hours      0
tenure                     0
work_accident              0
left                       0
promotion_last_5years      0
department                 0
salary                    0
dtype: int64
```

### 2.3.5 Check duplicates

Check for any duplicate entries in the data.

```
[9]: # Check for duplicates
### YOUR CODE HERE ###
df.duplicated().sum()
```

```
[9]: 3008
```

```
[10]: # Inspect some rows containing duplicates as needed
      ### YOUR CODE HERE ###
      df[df.duplicated()].head()
```

```
[10]:
```

	satisfaction_level	last_evaluation	number_project	\
396	0.46	0.57	2	
866	0.41	0.46	2	
1317	0.37	0.51	2	
1368	0.41	0.52	2	
1461	0.42	0.53	2	

	average_monthly_hours	tenure	work_accident	left	\
396	139	3	0	1	
866	128	3	0	1	
1317	127	3	0	1	
1368	132	3	0	1	
1461	142	3	0	1	

	promotion_last_5years	department	salary
396	0	sales	low
866	0	accounting	low
1317	0	sales	medium
1368	0	RandD	low
1461	0	sales	low

```
[11]: # Drop duplicates and save resulting dataframe in a new variable as needed
      ### YOUR CODE HERE ###

      df1= df.drop_duplicates(keep='first')
      # Display first few rows of new dataframe as needed
      ### YOUR CODE HERE ###
      df1.head()
```

```
[11]:
```

	satisfaction_level	last_evaluation	number_project	average_monthly_hours	\
0	0.38	0.53	2	157	
1	0.80	0.86	5	262	
2	0.11	0.88	7	272	
3	0.72	0.87	5	223	
4	0.37	0.52	2	159	

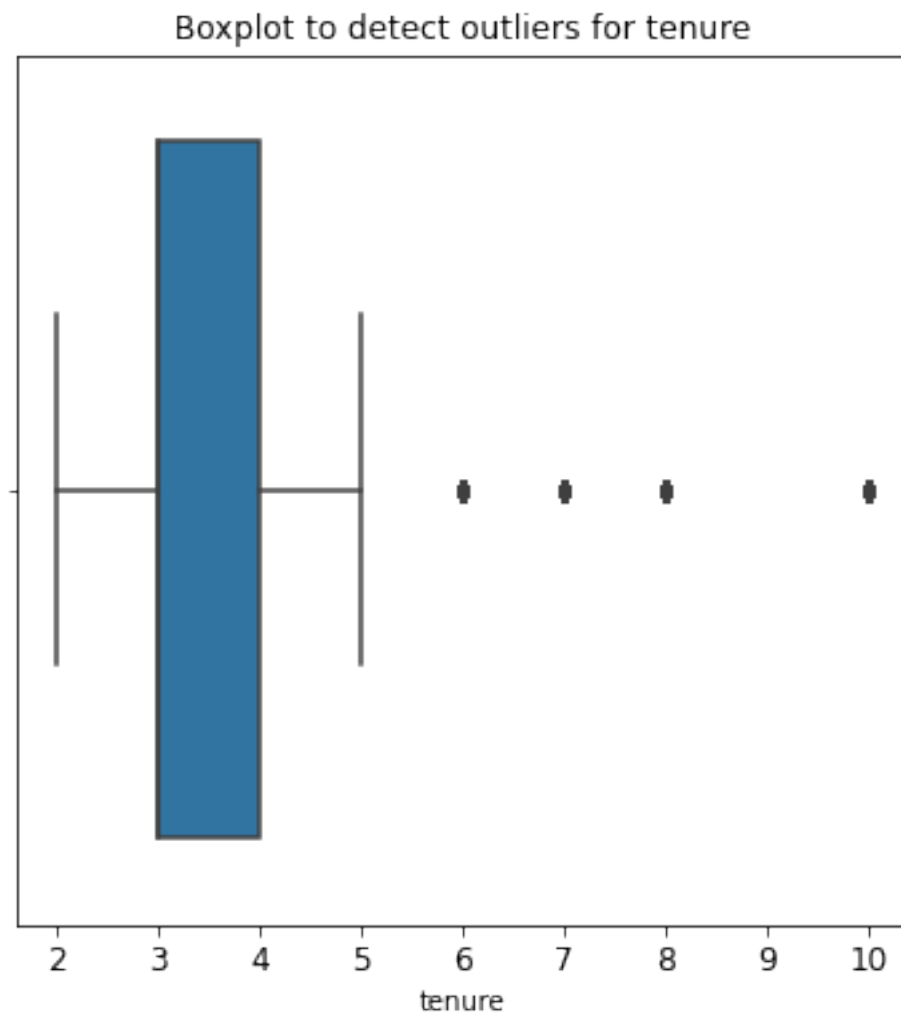
  

	tenure	work_accident	left	promotion_last_5years	department	salary
0	3	0	1	0	sales	low
1	6	0	1	0	sales	medium
2	4	0	1	0	sales	medium
3	5	0	1	0	sales	low
4	3	0	1	0	sales	low

### 2.3.6 Check outliers

Check for outliers in the data.

```
[12]: # Create a boxplot to visualize distribution of `tenure` and detect any outliers
      ### YOUR CODE HERE ###
      plt.figure(figsize=(6,6))
      plt.title('Boxplot to detect outliers for tenure',fontsize=12)
      plt.xticks(fontsize=12)
      plt.yticks(fontsize=12)
      sns.boxplot(x=df1['tenure'])
      plt.show()
```



```
[13]: # Determine the number of rows containing outliers
      ### YOUR CODE HERE ###
```



```

percentile25 = df1['tenure'].quantile(0.25)

# Compute the 75th percentile value in `tenure`
percentile75 = df1['tenure'].quantile(0.75)

# Compute the interquartile range in `tenure`
iqr = percentile75 - percentile25

# Define the upper limit and lower limit for non-outlier values in `tenure`
upper_limit = percentile75 + 1.5 * iqr
lower_limit = percentile25 - 1.5 * iqr
print("Lower limit:", lower_limit)
print("Upper limit:", upper_limit)

# Identify subset of data containing outliers in `tenure`
outliers = df1[(df1['tenure'] > upper_limit) | (df1['tenure'] < lower_limit)]

# Count how many rows in the data contain outliers in `tenure`
print("Number of rows in the data containing outliers in `tenure`:",
      len(outliers))

```

Lower limit: 1.5

Upper limit: 5.5

Number of rows in the data containing outliers in `tenure`: 824

Certain types of models are more sensitive to outliers than others. When you get to the stage of building your model, consider whether to remove outliers, based on the type of model you decide to use.

### 3 pAce: Analyze Stage

- Perform EDA (analyze relationships between variables)

#### 3.1 Step 2. Data Exploration (Continue EDA)

Begin by understanding how many employees left and what percentage of all employees this figure represents.

```

[14]: # Get numbers of people who left vs. stayed
      ### YOUR CODE HERE ###

      print(df1['left'].value_counts())

      # Get percentages of people who left vs. stayed
      ### YOUR CODE HERE ###
      print(df1['left'].value_counts(normalize=True))

```

```

0    10000
1     1991
Name: left, dtype: int64
0     0.833959
1     0.166041
Name: left, dtype: float64

```

### 3.1.1 Data visualizations

Now, examine variables that you're interested in, and create plots to visualize relationships between variables in the data.

```

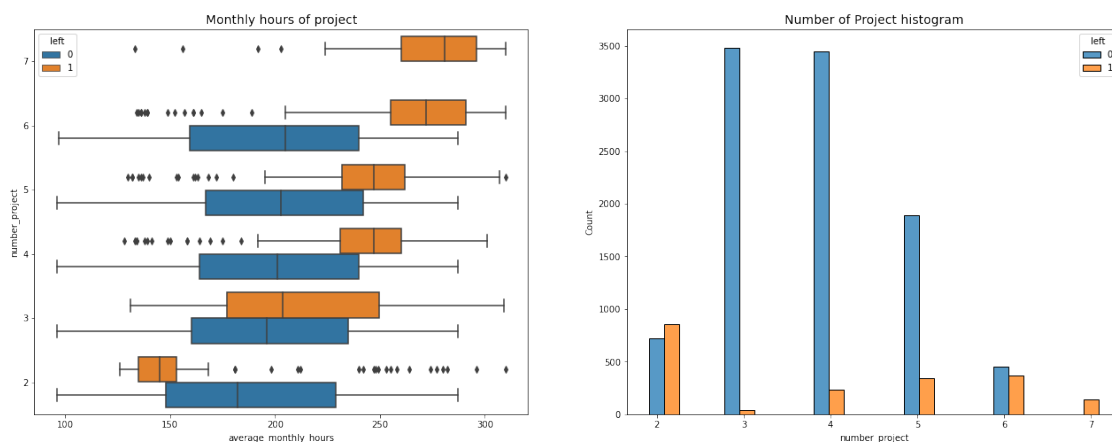
[17]: # Create a plot as needed
      ### YOUR CODE HERE ###
      fig, ax = plt.subplots(1,2, figsize=(22,8))

      sns.boxplot(data=df1, x='average_monthly_hours',y='number_project',hue='left',
                  ↪orient='h', ax=ax[0])
      ax[0].invert_yaxis()
      ax[0].set_title("Monthly hours of project",fontsize=14)

      tenure_stay = df1[df['left']==0]['number_project']
      tenure_left = df1[df['left']==1]['number_project']
      sns.histplot(data=df1,x='number_project',hue='left', multiple='dodge',
                  ↪shrink=2,ax=ax[1])
      ax[1].set_title("Number of Project histogram", fontsize='14')

      plt.show()

```



It might be natural that people who work on more projects would also work longer hours. This appears to be the case here, with the mean hours of each group (stayed and left) increasing with

number of projects worked. However, a few things stand out from this plot. 1. There are two groups of employees who left the company: (A) those who worked considerably less than their peers with the same number of projects, and (B) those who worked much more. Of those in group A, it's possible that they were fired. It's also possible that this group includes employees who had already given their notice and were assigned fewer hours because they were already on their way out the door. For those in group B, it's reasonable to infer that they probably quit. The folks in group B likely contributed a lot to the projects they worked in; they might have been the largest contributors to their projects. 2. Everyone with seven projects left the company, and the interquartile ranges of this group and those who left with six projects was ~255–295 hours/month—much more than any other group. 3. The optimal number of projects for employees to work on seems to be 3–4. The ratio of left/stayed is very small for these cohorts. 4. If you assume a work week of 40 hours and two weeks of vacation per year, then the average number of working hours per month of employees working Monday–Friday = 50 weeks \* 40 hours per week / 12 months = 166.67 hours per month. This means that, aside from the employees who worked on two projects, every group—even those who didn't leave the company—worked considerably more hours than this. It seems that employees here are overworked.

```
[18]: # Create a plot as needed
      ### YOUR CODE HERE ###
      df1[df1['number_project']==7]['left'].value_counts()
```

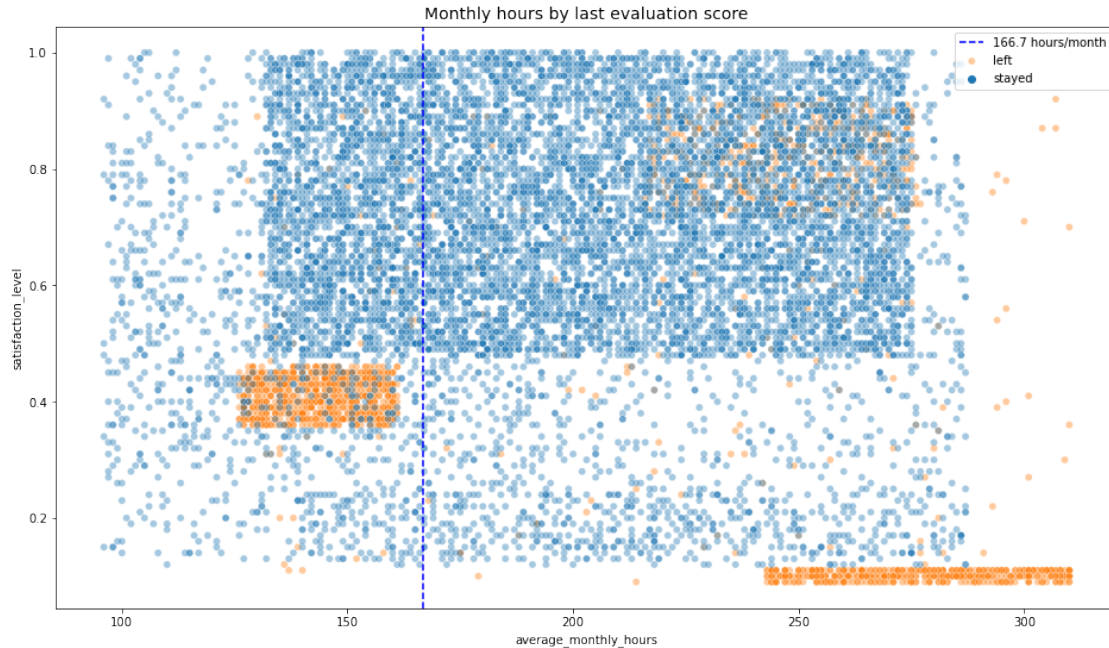
```
[18]: 1      145
      Name: left, dtype: int64
```

This confirms that all employees with 7 projects did leave.

Next, you could examine the average monthly hours versus the satisfaction levels.

```
[19]: plt.figure(figsize=(16,9))
      sns.
      ↪scatterplot(data=df1,x='average_monthly_hours',y='satisfaction_level',hue='left',alpha=0.
      ↪4)
      plt.axvline(x=166.7,color='blue',label='166.7 hours/month', ls='--')
      plt.legend(labels=['166.7 hours/month','left','stayed'])
      plt.title('Monthly hours by last evaluation score', fontsize='14')
```

```
[19]: Text(0.5, 1.0, 'Monthly hours by last evaluation score')
```



The scatterplot above shows that there was a sizeable group of employees who worked ~240–315 hours per month. 315 hours per month is over 75 hours per week for a whole year. It's likely this is related to their satisfaction levels being close to zero.

The plot also shows another group of people who left, those who had more normal working hours. Even so, their satisfaction was only around 0.4. It's difficult to speculate about why they might have left. It's possible they felt pressured to work more, considering so many of their peers worked more. And that pressure could have lowered their satisfaction levels.

Finally, there is a group who worked ~210–280 hours per month, and they had satisfaction levels ranging ~0.7–0.9.

Note the strange shape of the distributions here. This is indicative of data manipulation or synthetic data.

```
[20]: df1.groupby(['left'])['satisfaction_level'].agg([np.mean,np.median])
```

```
[20]:
```

	mean	median
left		
0	0.667365	0.69
1	0.440271	0.41

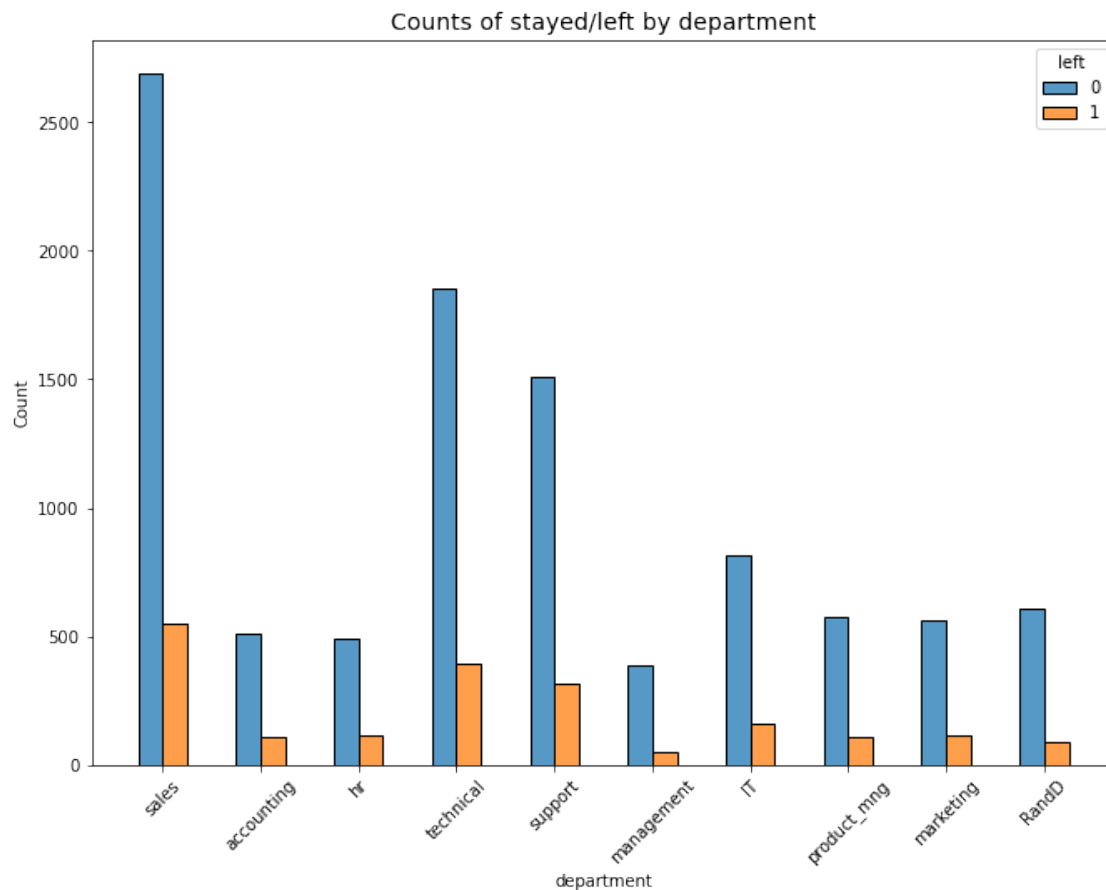
As expected, the mean and median satisfaction scores of employees who left are lower than those of employees who stayed. Interestingly, among employees who stayed, the mean satisfaction score appears to be slightly below the median score. This indicates that satisfaction levels among those who stayed might be skewed to the left.

```
[21]: df1['department'].value_counts()
```

```
[21]: sales          3239
      technical      2244
      support        1821
      IT             976
      RandD          694
      product_mng    686
      marketing      673
      accounting     621
      hr             601
      management     436
      Name: department, dtype: int64
```

```
[22]: plt.figure(figsize=(11,8))
      sns.
      ↪ histplot(data=df1,x='department',hue='left',discrete=1,hue_order=[0,1],multiple='dodge',shr
      ↪ 5)
      plt.xticks(rotation='45')
      plt.title('Counts of stayed/left by department',fontsize='14')
```

```
[22]: Text(0.5, 1.0, 'Counts of stayed/left by department')
```



There doesn't seem to be any department that differs significantly in its proportion of employees who left to those who stayed.

### 3.1.2 Insights

It appears that employees are leaving the company as a result of poor management. Leaving is tied to longer working hours, many projects, and generally lower satisfaction levels. It can be ungratifying to work long hours and not receive promotions or good evaluation scores. There's a sizeable group of employees at this company who are probably burned out. It also appears that if an employee has spent more than six years at the company, they tend not to leave.

## 4 paCe: Construct Stage

- Determine which models are most appropriate
- Construct the model
- Confirm model assumptions
- Evaluate model results to determine how well your model fits the data

## Recall model assumptions

**Logistic Regression model assumptions** - Outcome variable is categorical - Observations are independent of each other - No severe multicollinearity among X variables - No extreme outliers - Linear relationship between each X variable and the logit of the outcome variable - Sufficiently large sample size

[Double-click to enter your responses here.]

### 4.1 Step 3. Model Building, Step 4. Results and Evaluation

- Fit a model that predicts the outcome variable using two or more independent variables
- Check model assumptions
- Evaluate the model

#### 4.1.1 Identify the type of prediction task.

Your goal is to predict whether an employee leaves the company, which is a categorical outcome variable. So this task involves classification. More specifically, this involves binary classification, since the outcome variable left can be either 1 (indicating employee left) or 0 (indicating employee didn't leave).

#### 4.1.2 Logistic Regression Modeling

Before splitting the data, encode the non-numeric variables. There are two: **department** and **salary**.

department is a categorical variable, which means you can dummy it for modeling.

salary is categorical too, but it's ordinal. There's a hierarchy to the categories, so it's better not to dummy this column, but rather to convert the levels to numbers, 0-2.

```
[25]: df2=df1.copy()

df2['salary']= (df2['salary'].astype('category').cat.
    ↳set_categories(['low','medium','high']).cat.codes)

df2=pd.get_dummies(df2,drop_first=False)

df2.head()
```

```
[25]:
```

	satisfaction_level	last_evaluation	number_project	average_monthly_hours	\
0	0.38	0.53	2	157	
1	0.80	0.86	5	262	
2	0.11	0.88	7	272	
3	0.72	0.87	5	223	
4	0.37	0.52	2	159	

	tenure	work_accident	left	promotion_last_5years	salary	department_IT	\
0	3	0	1	0	0	0	
1	6	0	1	0	1	0	
2	4	0	1	0	1	0	
3	5	0	1	0	0	0	
4	3	0	1	0	0	0	

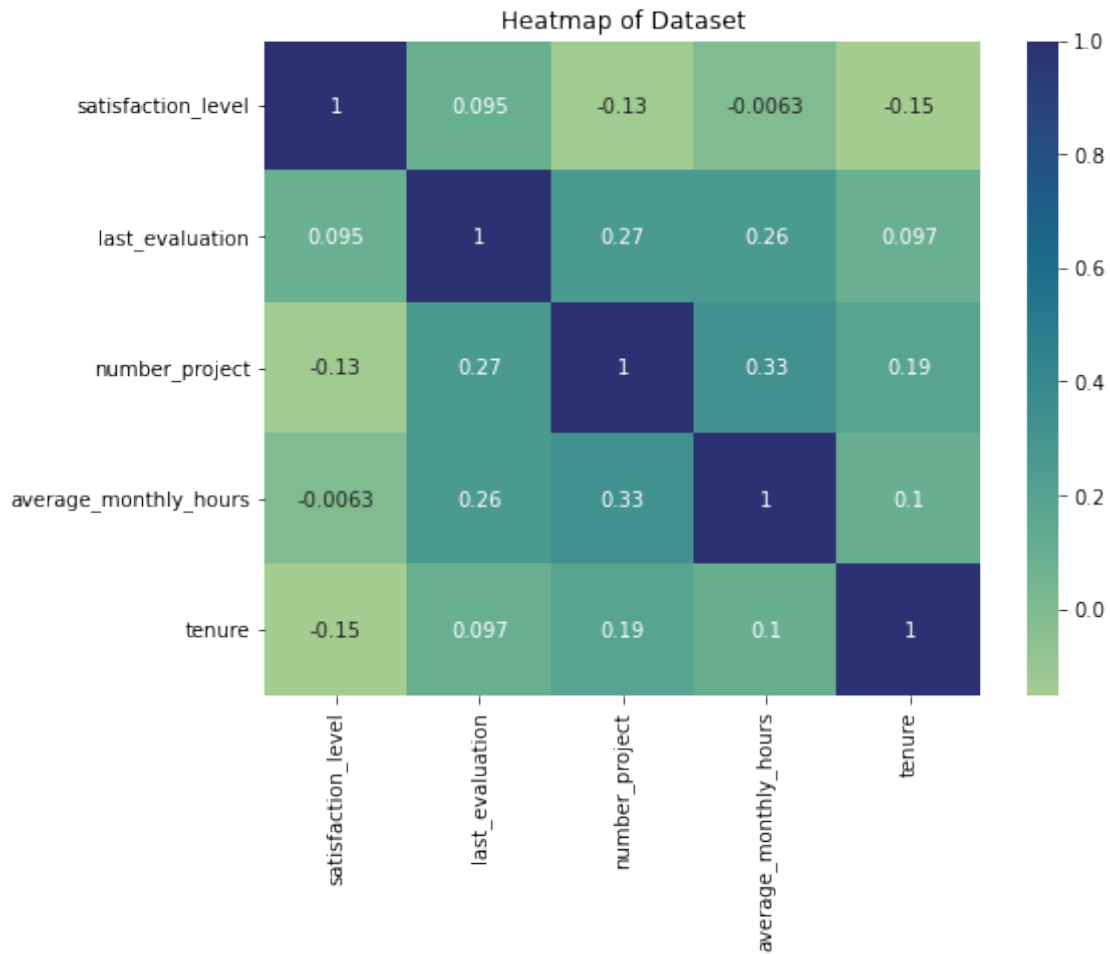
	department_RandD	department_accounting	department_hr	\
0	0	0	0	
1	0	0	0	
2	0	0	0	
3	0	0	0	
4	0	0	0	

	department_management	department_marketing	department_product_mng	\
0	0	0	0	
1	0	0	0	
2	0	0	0	
3	0	0	0	
4	0	0	0	

	department_sales	department_support	department_technical
0	1	0	0
1	1	0	0
2	1	0	0
3	1	0	0
4	1	0	0

```
[29]: plt.figure(figsize=(8,6))
sns.heatmap(df2[['satisfaction_level', 'last_evaluation', 'number_project', 'average_monthly_hours', 'tenure']].corr(), annot=True, cmap="crest")
plt.title('Heatmap of Dataset')
```

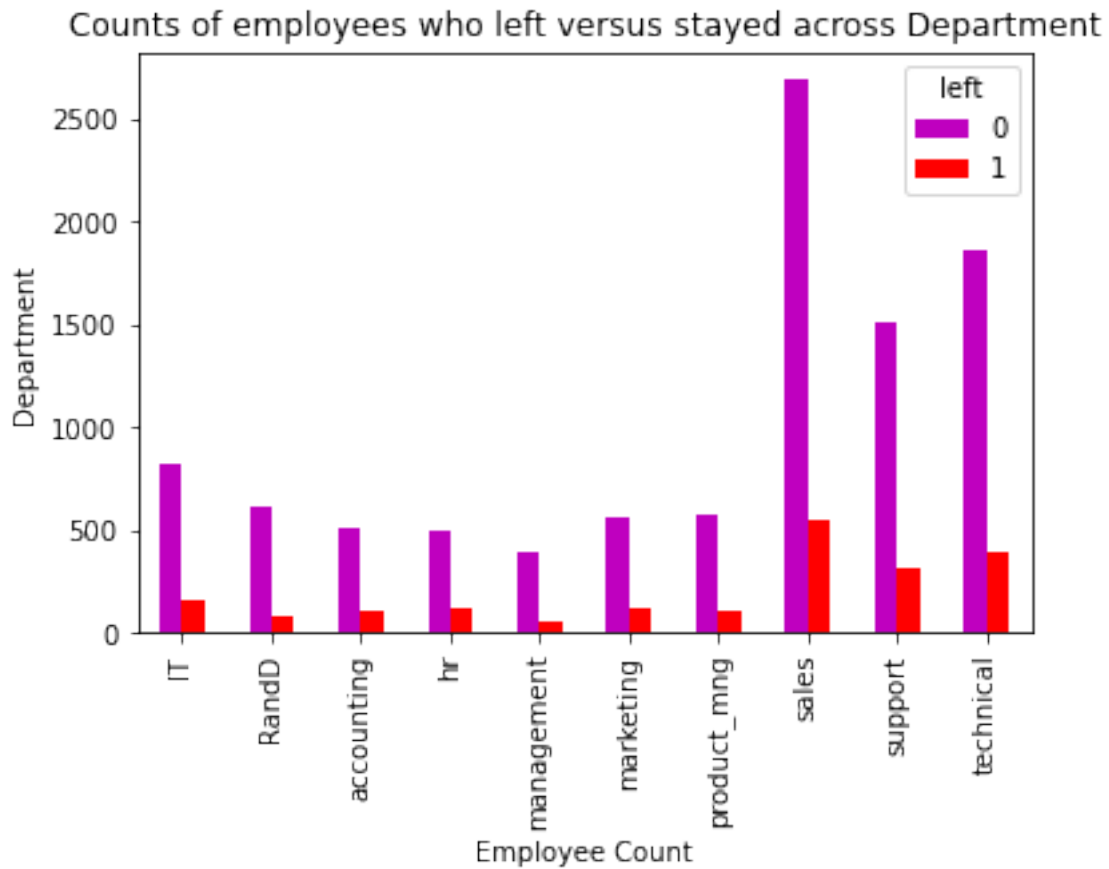
```
[29]: Text(0.5, 1.0, 'Heatmap of Dataset')
```



```
[30]: pd.crosstab(df1['department'],df1['left']).plot(kind='bar',color='mr')
plt.title("Counts of employees who left versus stayed across Department")
plt.xlabel("Employee Count")
plt.ylabel("Department")
```

```
[30]: Text(0, 0.5, 'Department')
```





```
[31]: # Select rows without outliers in `tenure` and save resulting dataframe in a
      ↳ new variable
dflogreg = df2[(df2['tenure'] >= lower_limit) & (df2['tenure'] <= upper_limit)]

# Display first few rows of new dataframe
dflogreg.head()
```

```
[31]:
```

	satisfaction_level	last_evaluation	number_project	average_monthly_hours	\
0	0.38	0.53	2	157	
2	0.11	0.88	7	272	
3	0.72	0.87	5	223	
4	0.37	0.52	2	159	
5	0.41	0.50	2	153	

	tenure	work_accident	left	promotion_last_5years	salary	department_IT	\
0	3	0	1	0	0	0	
2	4	0	1	0	1	0	
3	5	0	1	0	0	0	
4	3	0	1	0	0	0	

```
5      3      0      1      0      0      0
```

```

department_RandD  department_accounting  department_hr  \
0                0                0                0
2                0                0                0
3                0                0                0
4                0                0                0
5                0                0                0

```

```

department_management  department_marketing  department_product_mng  \
0                0                0                0
2                0                0                0
3                0                0                0
4                0                0                0
5                0                0                0

```

```

department_sales  department_support  department_technical
0                1                0                0
2                1                0                0
3                1                0                0
4                1                0                0
5                1                0                0

```

```
[32]: y = dflogreg['left']
      y.head()
```

```
[32]: 0    1
      2    1
      3    1
      4    1
      5    1
      Name: left, dtype: int64
```

```
[33]: X = dflogreg.drop('left',axis=1)
```

```
[35]: X_train,X_test,y_train,y_test = train_test_split(X,y,test_size=0.
      ↪25,stratify=y,random_state=42)
```

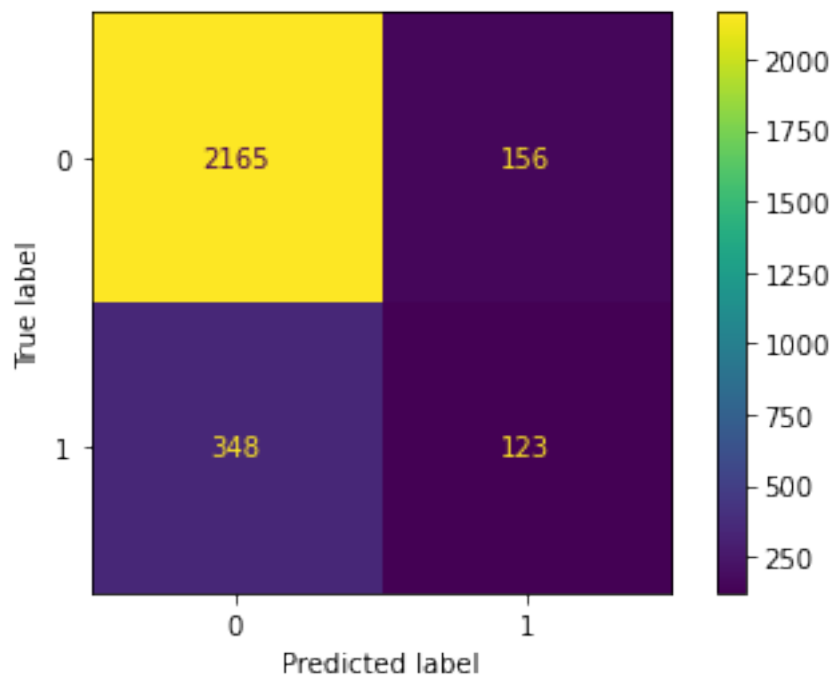
```
[37]: logclf = LogisticRegression(random_state=42,max_iter=500).fit(X_train,y_train)
```

```
[38]: y_pred = logclf.predict(X_test)
```

```
[39]: logcm = confusion_matrix(y_test, y_pred, labels= logclf.classes_)
      logdisp = ConfusionMatrixDisplay(confusion_matrix=logcm,display_labels=logclf.
      ↪classes_)
      logdisp.plot(values_format='')

```

```
[39]: <sklearn.metrics._plot.confusion_matrix.ConfusionMatrixDisplay at 0x751884f3ca10>
```



The upper-left quadrant displays the number of true negatives. The upper-right quadrant displays the number of false positives. The bottom-left quadrant displays the number of false negatives. The bottom-right quadrant displays the number of true positives.

True negatives: The number of people who did not leave that the model accurately predicted did not leave.

False positives: The number of people who did not leave the model inaccurately predicted as leaving.

False negatives: The number of people who left that the model inaccurately predicted did not leave

True positives: The number of people who left the model accurately predicted as leaving

A perfect model would yield all true negatives and true positives, and no false negatives or false positives.

```
[40]: dflogreg['left'].value_counts(normalize=True)
```

```
[40]: 0    0.831468
      1    0.168532
      Name: left, dtype: float64
```

There is an approximately 83%-17% split. So the data is not perfectly balanced, but it is not too imbalanced. If it was more severely imbalanced, you might want to resample the data to make

it more balanced. In this case, you can use this data without modifying the class balance and continue evaluating the model.

```
[41]: target_names=['Predicted wouldnt leave','Predicted would leave']
print(classification_report(y_test, y_pred, target_names=target_names))
```

	precision	recall	f1-score	support
Predicted wouldnt leave	0.86	0.93	0.90	2321
Predicted would leave	0.44	0.26	0.33	471
accuracy			0.82	2792
macro avg	0.65	0.60	0.61	2792
weighted avg	0.79	0.82	0.80	2792

The classification report above shows that the logistic regression model achieved a precision of 79%, recall of 82%, f1-score of 80% (all weighted averages), and accuracy of 82%. However, if it's most important to predict employees who leave, then the scores are significantly lower.

## 5 Tree Based Model

```
[42]: y= df2['left']
y.head()
```

```
[42]: 0    1
      1    1
      2    1
      3    1
      4    1
      Name: left, dtype: int64
```

```
[44]: X= df2.drop('left',axis=1)
X.head()
```

```
[44]: satisfaction_level  last_evaluation  number_project  average_monthly_hours  \
0                0.38           0.53                2                157
1                0.80           0.86                5                262
2                0.11           0.88                7                272
3                0.72           0.87                5                223
4                0.37           0.52                2                159

      tenure  work_accident  promotion_last_5years  salary  department_IT  \
0          3              0                    0        0                0
1          6              0                    0        1                0
2          4              0                    0        1                0
3          5              0                    0        0                0
```

4	3	0	0	0	0
---	---	---	---	---	---

	department_RandD	department_accounting	department_hr	\
0	0	0	0	
1	0	0	0	
2	0	0	0	
3	0	0	0	
4	0	0	0	

	department_management	department_marketing	department_product_mng	\
0	0	0	0	
1	0	0	0	
2	0	0	0	
3	0	0	0	
4	0	0	0	

	department_sales	department_support	department_technical
0	1	0	0
1	1	0	0
2	1	0	0
3	1	0	0
4	1	0	0

```
[45]: X_train,X_test,y_train,y_test = train_test_split(X,y,test_size=0.
      ↪25,random_state=0)
```

### 5.0.1 Decision Tree

```
[46]: tree = DecisionTreeClassifier(random_state=0)
cv = {'max_depth':[4,6,8,None],
      'min_samples_leaf':[2,5,1],
      'min_samples_split':[2,4,6]}

scoring={'accuracy','precision','recall','f1','roc_auc'}

tree1 = GridSearchCV(tree, cv, scoring=scoring, cv=4,refit='roc_auc')
```

```
[47]: %%time
tree1.fit(X_train,y_train)
```

```
CPU times: user 2.61 s, sys: 329 ms, total: 2.94 s
Wall time: 2.94 s
```

```
[47]: GridSearchCV(cv=4, error_score=nan,
                  estimator=DecisionTreeClassifier(ccp_alpha=0.0, class_weight=None,
                                                    criterion='gini', max_depth=None,
```

```

max_features=None,
max_leaf_nodes=None,
min_impurity_decrease=0.0,
min_impurity_split=None,
min_samples_leaf=1,
min_samples_split=2,
min_weight_fraction_leaf=0.0,
presort='deprecated',
random_state=0, splitter='best'),
iid='deprecated', n_jobs=None,
param_grid={'max_depth': [4, 6, 8, None],
            'min_samples_leaf': [2, 5, 1],
            'min_samples_split': [2, 4, 6]},
pre_dispatch='2*n_jobs', refit='roc_auc', return_train_score=False,
scoring={'precision', 'f1', 'accuracy', 'roc_auc', 'recall'},
verbose=0)

```

```
[48]: tree1.best_params_
```

```
[48]: {'max_depth': 6, 'min_samples_leaf': 1, 'min_samples_split': 4}
```

```
[49]: tree1.best_score_
```

```
[49]: 0.9698382034497465
```

```

[50]: def make_results(model_name:str, model_object, metric:str):
    """
    Arguments:
        model_name (string): what you want the model to be called in the output_
    ↪table
        model_object: a fit GridSearchCV object
        metric (string): precision, recall, f1, accuracy, or auc

    Returns a pandas df with the F1, recall, precision, accuracy, and auc scores
    for the model with the best mean 'metric' score across all validation folds.
    ↪
    """

    # Create dictionary that maps input metric to actual metric name in_
    ↪GridSearchCV
    metric_dict = {'auc': 'mean_test_roc_auc',
                  'precision': 'mean_test_precision',
                  'recall': 'mean_test_recall',
                  'f1': 'mean_test_f1',
                  'accuracy': 'mean_test_accuracy'
                  }

```

```

# Get all the results from the CV and put them in a df
cv_results = pd.DataFrame(model_object.cv_results_)

# Isolate the row of the df with the max(metric) score
best_estimator_results = cv_results.iloc[cv_results[metric_dict[metric]].
↳idxmax(), :]

# Extract Accuracy, precision, recall, and f1 score from that row
auc = best_estimator_results.mean_test_roc_auc
f1 = best_estimator_results.mean_test_f1
recall = best_estimator_results.mean_test_recall
precision = best_estimator_results.mean_test_precision
accuracy = best_estimator_results.mean_test_accuracy
# Create table of results
table = pd.DataFrame()
table = pd.DataFrame({'model': [model_name],
                        'precision': [precision],
                        'recall': [recall],
                        'F1': [f1],
                        'accuracy': [accuracy],
                        'auc': [auc]
                      })

return table

```

```

[51]: treeresult=make_results('decision tree cv', tree1 , 'auc')
treeresult

```

```

[51]:
      model  precision  recall    F1  accuracy    auc
0  decision tree cv   0.973531  0.915069  0.943315   0.98143  0.969838

```

All of these scores from the decision tree model are strong indicators of good model performance. Recall that decision trees can be vulnerable to overfitting, and random forests avoid overfitting by incorporating multiple trees to make predictions. You could construct a random forest model next.

## 6 pacE: Execute Stage

- Interpret model performance and results
- Share actionable steps with stakeholders

## Recall evaluation metrics

- **AUC** is the area under the ROC curve; it's also considered the probability that the model ranks a random positive example more highly than a random negative example.
- **Precision** measures the proportion of data points predicted as True that are actually True, in other words, the proportion of positive predictions that are true positives.

- **Recall** measures the proportion of data points that are predicted as True, out of all the data points that are actually True. In other words, it measures the proportion of positives that are correctly classified.
- **Accuracy** measures the proportion of data points that are correctly classified.
- **F1-score** is an aggregation of precision and recall.

## 6.1 Step 4. Results and Evaluation

- Interpret model
- Evaluate model performance using metrics
- Prepare results, visualizations, and actionable steps to share with stakeholders

### 6.1.1 Summary of model results

#### Logistic Regression

The logistic regression model achieved precision of 80%, recall of 83%, f1-score of 80% (all weighted averages), and accuracy of 83%, on the test set.

#### Tree-based Machine Learning

After conducting feature engineering, the decision tree model achieved AUC of 93.8%, precision of 87.0%, recall of 90.4%, f1-score of 88.7%, and accuracy of 96.2%, on the test set. The random forest modestly outperformed the decision tree model.

### 6.1.2 Conclusion, Recommendations, Next Steps

The models and the feature importances extracted from the models confirm that employees at the company are overworked.

To retain employees, the following recommendations could be presented to the stakeholders:

- Cap the number of projects that employees can work on.
- Consider promoting employees who have been with the company for at least four years, or conduct further investigation about why four-year tenured employees are so dissatisfied.
- Either reward employees for working longer hours, or don't require them to do so.
- If employees aren't familiar with the company's overtime pay policies, inform them about this. If the expectations around workload and time off aren't explicit, make them clear.
- Hold company-wide and within-team discussions to understand and address the company work culture, across the board and in specific contexts.
- High evaluation scores should not be reserved for employees who work 200+ hours per month. Consider a proportionate scale for rewarding employees who contribute more/put in more effort.

#### Next Steps

It may be justified to still have some concern about data leakage. It could be prudent to consider how predictions change when `last_evaluation` is removed from the data. It's possible that evaluations aren't performed very frequently, in which case it would be useful to be able to predict



employee retention without this feature. It's also possible that the evaluation score determines whether an employee leaves or stays, in which case it could be useful to pivot and try to predict performance score. The same could be said for satisfaction score.

[ ]: