# Activity_ Course 7 Salifort Motors project lab

April 10, 2025

# 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 actiivty 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) |

### Reflect on these questions as you complete the plan stage.

- Who are your stakeholders for this project?
- What are you trying to solve or accomplish?
- What are your initial observations when you explore the data?
- What resources do you find yourself using as you complete this stage? (Make sure to include the links.)
- Do you have any ethical considerations in this stage?

1. Who are your stakeholders for this project?

Primary Stakeholders: * HR Department at Salifort Motors – They want to understand why employees leave and how to reduce attrition. * Company Executives and Management – Interested in reducing hiring/training costs and improving team performance.

Secondary Stakeholders: * Employees – As changes in HR policies could affect their work environ-

ment and satisfaction. * Recruitment Team – Will benefit from predictive insights in planning.

2. What are you trying to solve or accomplish?

Identify key factors that influence employee attrition and build a predictive model to classify employess as likely to stay or leave. * reduce attrition * Improve job satisfaction * Increase employee retention through actionable HR strategies

3. What are your initial observations when you explore the data?

- The dataset includes 10 features such as satisfaction level, evaluation scores, and workload.
- There are some potentially duplicated rows, which may affect model accuracy.
- Some features show non-linear relationships with attrition (e.g., number of projects and average monthly hours).
- Class imbalance might be present (fewer people left than stayed).

4. What resources do you find yourself using as you complete this stage? (Make sure to include the links.)

- Kaggle
- Seaborn Boxplot
- Scikit_learn for Classification Models.
- ROC and AUC explanation
- Google Advanced Data Analytics Course Resources

5. Do you have any ethical considerations in this stage?

- Data Privacy - Ensure no personally identifiable information is used.
- Bias - Be aware of biases in data.
- Model Impact - Predicitions should not be used punitively but to proactively improve conditions.
- Transparency- HR and employees should be informed how predicitions will be used to support, not penalize them.

## 2.2 Step 1. Imports

- Import packages
- Load dataset

### 2.2.1 Import packages

```python
[1]: # Import packages
     #For data manipulation
     import pandas as pd
     import numpy as np

     #For data Visualization
     import matplotlib.pyplot as plt
     import seaborn as sns
```

```
#For displaying all of the columns in dataframes
pd.set_option('display.max_columns', None)

#For data modeling
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

#For metrics and helpful functions
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

#For saving models
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.

```
[2]: # RUN THIS CELL TO IMPORT YOUR DATA.
     # Load dataset into a dataframe
     df0 = pd.read_csv("HR_capstone_dataset.csv")

     # Display first few rows of the dataframe
     df0.head(5)
```

```
[2]:    satisfaction_level  last_evaluation  number_project  average_montly_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

```
[3]: # Gather basic information about the data
     df0.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_montly_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
```

### 2.3.2  Gather descriptive statistics about the data

```
[4]: # Gather descriptive statistics about the data
     df0.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_montly_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
```

### 2.3.3 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.

```
[5]: # Display all column names
     df0.columns
```

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

```
[6]: # Rename columns as needed
     df0 = df0.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
df0.columns
```

[6]: 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.

[7]: ```
# Check for missing values
df0.isna().sum()
```

[7]: 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.

[8]: ```
# Check for duplicates
df0.duplicated().sum()
```

[8]: 3008

There are 3008 duplicated rows, that is 20 of the data.

[9]: ```
# Inspect some rows containing duplicates as needed
df0[df0.duplicated()].head(10)
```

```
[9]:          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
     1516                   0.40             0.50               2
     1616                   0.37             0.46               2
     1696                   0.39             0.56               2
     1833                   0.10             0.85               6
     12000                  0.38             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
     1516                     127       3              0     1
     1616                     156       3              0     1
     1696                     160       3              0     1
     1833                     266       4              0     1
     12000                    157       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
     1516                       0          IT     low
     1616                       0       sales     low
     1696                       0       sales     low
     1833                       0       sales     low
     12000                      0       sales     low
```

The above output shows the first five occurances of rows that are duplicated farther down in the dataframe. How plausiable is it that 2 employees self-reported the exact same response for every column?

We could perform a likelihood analysis by essentially applying Bayes's Theorem and multiplying the probabilities of finding each value in each column, but this does not seem necessary. With several continuous variables across 10 columns, it seems very unlikely that these observations are legitimate. We can drop these rows.

```
[10]:  # Drop duplicates and save resulting dataframe in a new variable as needed
       df1= df0.drop_duplicates(keep='first')
```

```
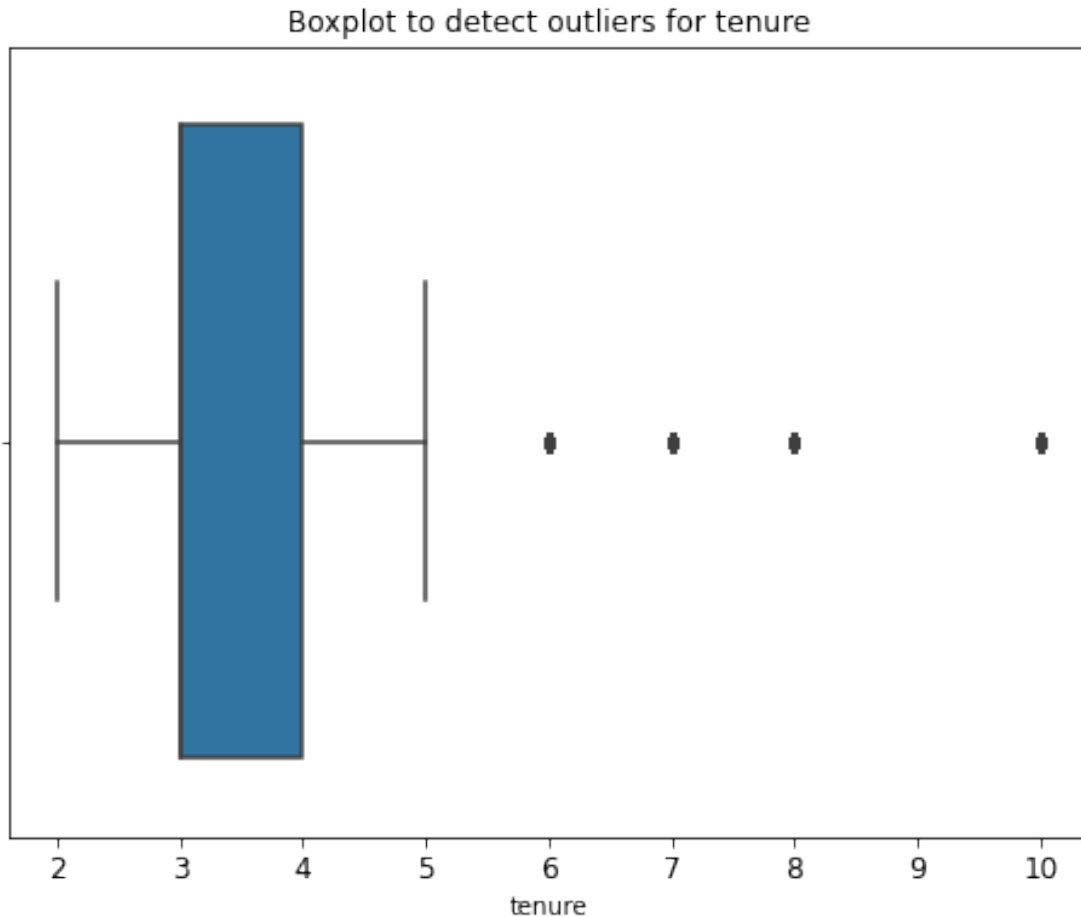# Display first few rows of new dataframe as needed
df1.head()
```

[10]:

| | 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.

[11]:
```
# Create a boxplot to visualize distribution of `tenure` and detect any outliers
plt.figure(figsize=(8,6))
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();
```

9

## Boxplot to detect outliers for tenure



The above boxplot shows that there are outliers in the **tenure** variable. It would be helpful to investigate how many rows in the data contain outliers in the **tenure** column.

```
[12]:  # Determine the number of rows containing outliers
       #Compute the 25th percentile value in 'tenure'
       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)

### Reflect on these questions as you complete the analyze stage.

- What did you observe about the relationships between variables?
- What do you observe about the distributions in the data?
- What transformations did you make with your data? Why did you chose to make those decisions?
- What are some purposes of EDA before constructing a predictive model?
- What resources do you find yourself using as you complete this stage? (Make sure to include the links.)
- Do you have any ethical considerations in this stage?

[Double-click to enter your responses here.]

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

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

```
[13]: # Get numbers of people who left vs. stayed
      print(df1['left'].value_counts())

      # Get percentages of people who left vs. stayed
      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.

```
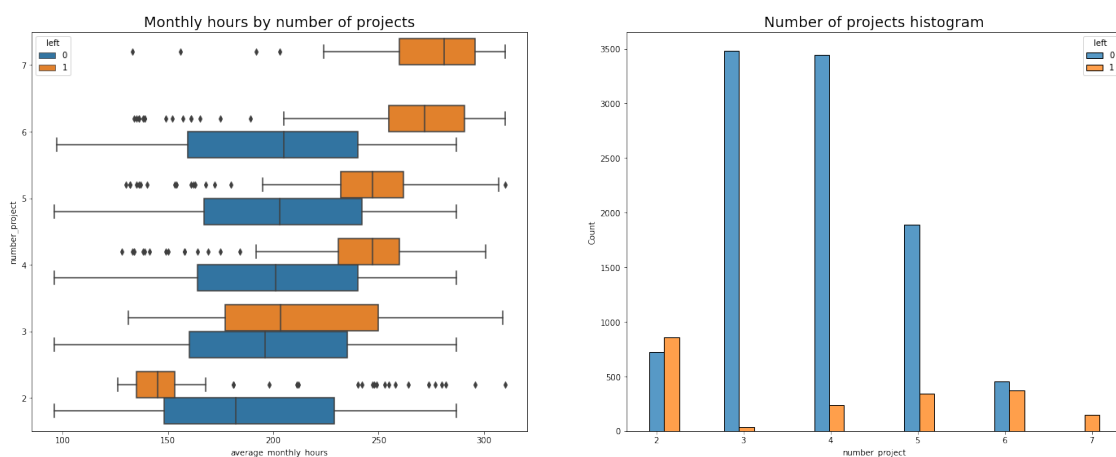[14]: # Create a plot as needed
      #set figure and axes
      fig, ax = plt.subplots(1, 2, figsize = (24, 9))

      #Create boxplot showing 'average_monthly_hours' distributions for␣
       ↪`number_project`, comparing employees who stayed versus those who left
      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 by number of projects', fontsize='18')

      #Create histogram showing distribution of 'number_project', comparing employees␣
       ↪who stayed versus those who left
      tenure_stay = df1[df1['left']==0]['number_project']
      tenure_left = df1[df1['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 projects histogram', fontsize='18')

      #Display the plots
      plt.show();
```



Boxplot: This shows how average monthly hours vary across different numbers of projects, and compares employees who left(orange)vs those who stayed(blue) * Employees who worked on 2-4

projects mostly stayed * Employees who worked on 6-7 projects had a much higher chance of leaving. * Those with extreme workloads were more likely to leave. * Moderate workloads with 3-5 projects seem more stable and associated with retention.

Histogram: This shows how many employees worked on each number of projects, comparing those who left vs stayed. * Most employees who stayed had 3, 4, or 5 projects. * Very few employees with 2 0r 6+ projects stayed. * A large portion of leavers had 2 or 6-7 projects, supporting the previous point.

If you assume a work week of 40 hours and 2 weeks of vaction 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 2 projects, every group- even those who didn't leave the company-worked considerably more hours than this. It seems that employees here are overworked.

```python
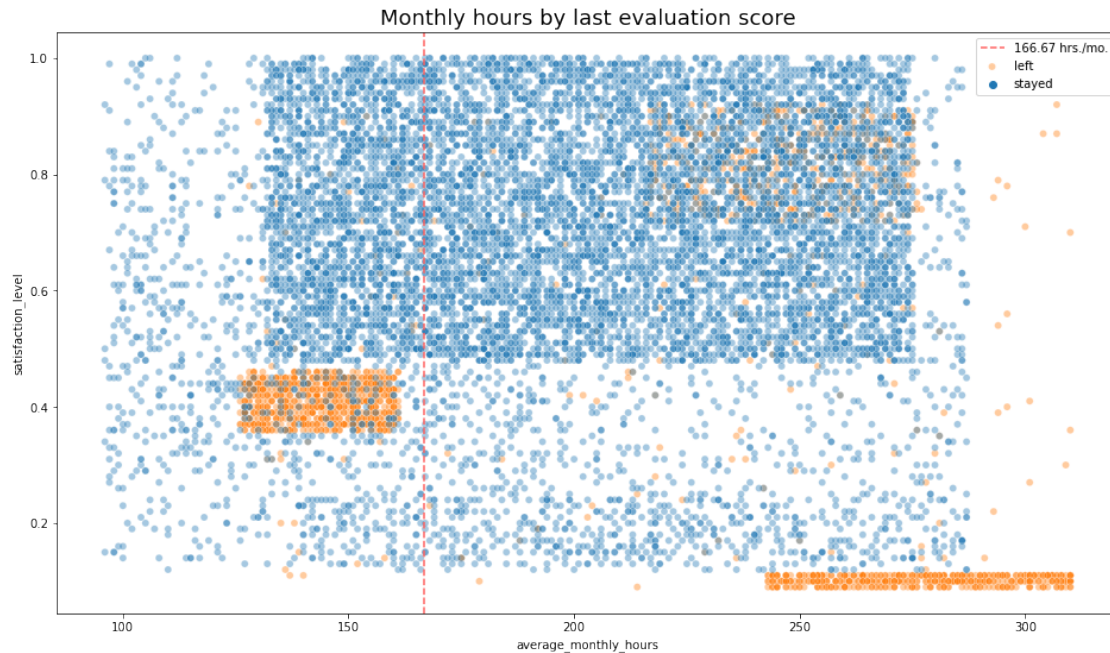[15]: #Get value counts of stayed/left for employees with 7 projects
      df1[df1['number_project']==7]['left'].value_counts()
```

```
[15]: 1    145
      Name: left, dtype: int64
```

This confirms that all employess with 7 projects did leave. Next, you could examine the average monthly hours vs the satisfaction levels.

```python
[16]: # Create a plot as needed
      #Create scatterplot of 'average_monthly_hours' vs 'satisfaction_level',␣
       ↪comparing employees who stayed versus those who left
      plt.figure(figsize=(16, 9))
      sns.scatterplot(data=df1, x='average_monthly_hours', y='satisfaction_level',␣
       ↪hue='left', alpha=0.4)
      plt.axvline(x=166.67, color='#ff6361', label='166.67 hrs./mo.', ls='--')
      plt.legend(labels=['166.67 hrs./mo.', 'left', 'stayed'])
      plt.title('Monthly hours by last evaluation score', fontsize='18');
```

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.

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

```
[17]: # Create a plot as needed - Visualize satisfaction levels by tenure
      #Set figure and axes
      fig, ax = plt.subplots(1, 2, figsize = (22, 9))

      #Create boxplot showing distributions of 'satisfaction_level' by tenure,␣
       ↪comparing employees who stayed versus those who left
      sns.boxplot(data=df1, x='satisfaction_level', y='tenure', hue='left',␣
       ↪orient='h', ax=ax[0])
      ax[0].invert_yaxis()
      ax[0].set_title('Satisfaction by tenure', fontsize='18')

      #Create histogram showing distribution of 'tenure', comparing employees who␣
       ↪stayed vs those who left
      tenure_stay = df1[df1['left']==0]['tenure']
      tenure_left = df1[df1['left']==0]['tenure']
```

```
sns.histplot(data=df1, x='tenure', hue='left', multiple='dodge', shrink=5,␣
 ↪ax=ax[1])
ax[1].set_title('Tenure histogram', fontsize='18')

plt.show();
```



There are many observations you could make from this plot. * Employess who left fall into two general categories, dissatified employees with shorter tenures and very satisfied employees with medium-length tenures. * Four-year employees who left seem to have an unusually low satisfaction level. It's worth investigating changes to company policy that might have affected people specifically at the four-year mark, if possible. * The longest-tenured employees didn't leave. Their satisfaction level aligned with those of newer employees who stayed. * The histogram shows that there are relatively few longer-tenured employees. It's possible that they're the higher-ranking, higher-paid employees.

```
[18]: # Calculate mean and median satisfaction scores of employees who left and those␣
 ↪who stayed
df1.groupby(['left'])['satisfaction_level'].agg([np.mean,np.median])
```

```
[18]:          mean   median
    left
    0      0.667365    0.69
    1      0.440271    0.41
```

The mean, and median satisfaction scores of employees who left are lower than those of employees who stayed. Among employees who stayed, the mean satisfaction score appears to be slightly below the median score. This indicates that satisfaction level amoung whose who stayed might be skewed to the left.

```
[19]: # Create a plot as needed – Examine salary levels for different tenures.
#Set figure and axes
```

```
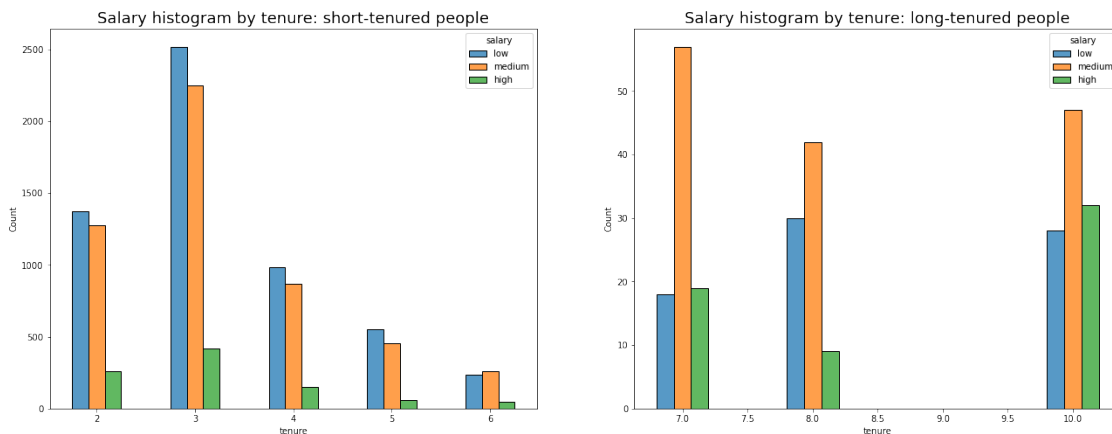fig, ax = plt.subplots(1, 2, figsize = (22, 8))

#Define short-tenured emplyoees
tenure_short = df1[df1['tenure'] < 7]

#Define long-tenured employees
tenure_long = df1[df1['tenure'] > 6]

#Plot short_tenured histogram
sns.histplot(data=tenure_short, x='tenure', hue='salary', discrete=1,
             hue_order=['low', 'medium', 'high'], multiple='dodge', shrink=.5,␣
 ↪ax=ax[0])
ax[0].set_title('Salary histogram by tenure: short-tenured people',␣
 ↪fontsize='18')

#Plot long_tenured histogram
sns.histplot(data=tenure_long, x='tenure', hue='salary', discrete=1,
             hue_order=['low', 'medium', 'high'], multiple='dodge', shrink=.4,␣
 ↪ax=ax[1])
ax[1].set_title('Salary histogram by tenure: long-tenured people',␣
 ↪fontsize='18');
```



The plots above show that long_tenured employees were not disproportionately comprised of higher-paid employees.

```
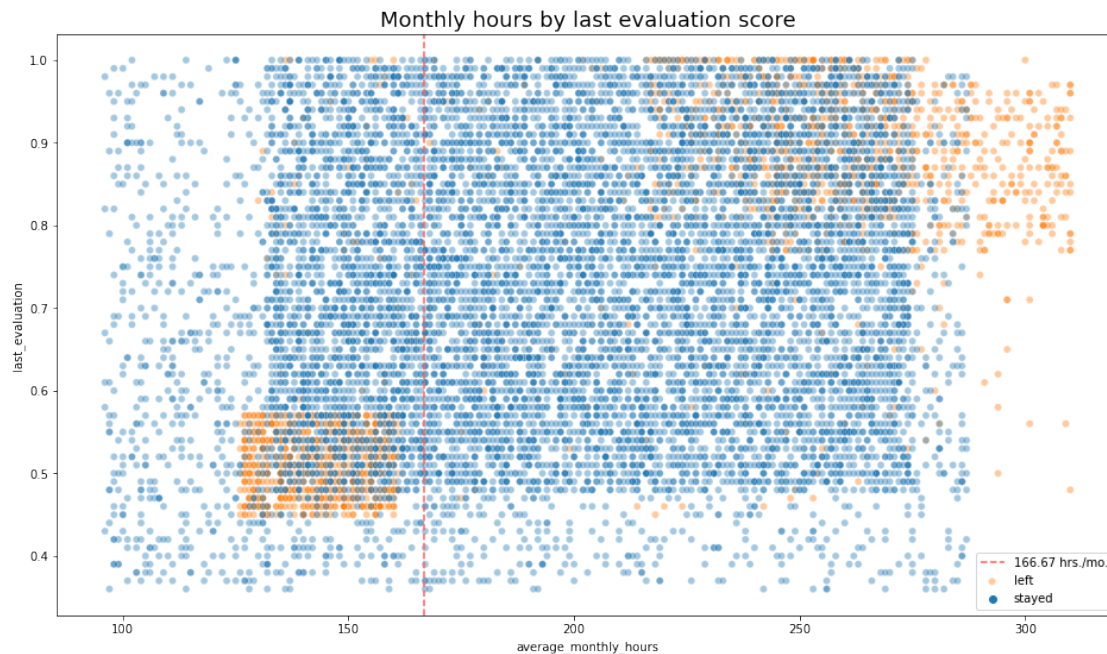[20]: # Create a plot as needed - create scatterplot of 'average_monthly_scores' vs␣
 ↪'last_evaluation'
plt.figure(figsize=(16, 9))
sns.scatterplot(data=df1, x='average_monthly_hours', y='last_evaluation',␣
 ↪hue='left', alpha=0.4)
plt.axvline(x=166.67, color='#ff6361', label='166.67 hrs./mo.', ls='--')
plt.legend(labels=['166.67 hrs./mo.', 'left', 'stayed'])
```

```
plt.title('Monthly hours by last evaluation score', fontsize='18');
```



The following observations can be made from the scatterplot: * The scatterplot indicates two groups of employees who left: overworked employees who performed very well and employees who worked slightly under the nominal monthly average of 166.67 hours with lower evaluation scores. * There seems to be a correlation between hours worked and evaluation score. * These isn't a high percentage of employees in the upper left quadrant of this plot, but working long hours doesn't guarantee a good evaluation score. * Most of the employees in this company work well over 167 hours per month.

[21]:
```
# Create a plot as needed- examine whether employees who worked very long hours
 →were promoted in the last five years.
#Create plot to examine relationship between 'average_monthly_hours', and
 →'promotion_last_5years'
plt.figure(figsize=(16, 3))
sns.scatterplot(data=df1, x='average_monthly_hours', y='promotion_last_5years',
 →hue='left', alpha=0.4)
plt.axvline(x=166.67, color='#ff6361', ls='--')
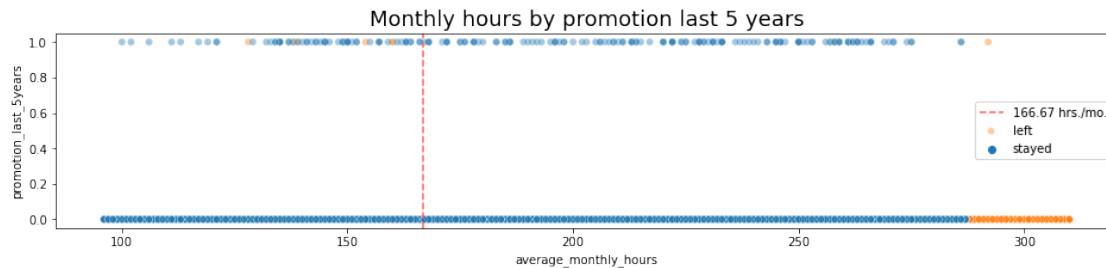plt.legend(labels=['166.67 hrs./mo.', 'left', 'stayed'])
plt.title('Monthly hours by promotion last 5 years', fontsize='18');
```

Monthly hours by promotion last 5 years

The above plot shows the following: * Very few employees who were promoted in the last five years left * Very few employees who worked the most hours were promoted * All of the employees who left were working the longest hours

```
[22]: #Inspect how employees who left are distributed across departments
      #Display counts foe each department.
      df1['department'].value_counts()
```

```
[22]: 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
```

```
[23]: # Create a plot as needed
      # Create stacked histogram to compare department distribution of employees who␣
       ↪left to that of employees who didn't
      plt.figure(figsize=(11,8))
      sns.histplot(data=df1, x='department', hue='left', discrete=1,
                   hue_order=[0, 1], multiple='dodge', shrink=.5)
      plt.xticks(rotation='45')
      plt.title('Counts of stayed/left by department', fontsize='18');
```

18

## 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.

```
[24]: # Create a plot as needed- check for strong correlation between variables in
      ⮑the data
      plt.figure(figsize=(16, 9))
      heatmap = sns.heatmap(df0.corr(), vmin=-1, vmax=1, annot=True, cmap=sns.
      ⮑color_palette("vlag", as_cmap=True))
      heatmap.set_title('Correlation Heatmap', fontdict={'fontsize':'18'}, pad=12);
```

Correlation Heatmap

The correlation heatmap confirms that the number of projects, monthly hours, and evaluation scores all have some positive correlation with each other, and whether an employee leaves is negatively correlated with their satisfaction level.

### 3.1.2 Insights

[What insights can you gather from the plots you created to visualize the data?

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 level. 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 6 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

### Reflect on these questions as you complete the constructing stage.

- Do you notice anything odd?
- Which independent variables did you choose for the model and why?
- Are each of the assumptions met?
- How well does your model fit the data?
- Can you improve it? Is there anything you would change about the model?
- What resources do you find yourself using as you complete this stage? (Make sure to include the links.)
- Do you have any ethical considerations in this stage?

[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.

The 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(employee left) or 0(employee didn't leave)

### 4.1.2 Identify the types of models most appropriate for this task.

The variables you want to predict is categorical, you could either build a LogisticRegression model, or a Tree-based Machine learning model. We could proceed with one of the two following approaches, or if we could implement both and determine how they compare.

### 4.1.3 Modeling Approach A: Logistic Regression Model

This approach covers implementation of Logisitic Regression.

**Logistic Regression**   The binomial logistic regression suits the task because it involves binary classification. Before splitting the data, we encode the non-numeric variables. There are 2: **department** and **salary**. **department** is a categorical variable, which means we can dummy it for modeling. **salary** is categorical too, but its ordinal. There's hierarchy to the categories, so its better not to dummy this column, but rather to convert the levels to numbers, 0-2.

```
[25]:  #copy the dataframe
       df_enc = df1.copy()

       #Encode the 'salary' column as an ordinal numeric category
       df_enc['salary'] = (
           df_enc['salary'].astype('category')
           .cat.set_categories(['low', 'medium', 'high'])
       )

       #Dummy encode the 'department' column
       df_enc = pd.get_dummies(df_enc, drop_first=False)

       #Display the new dataframe
       df_enc.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  department_IT  \
       0       3              0     1                      0              0
       1       6              0     1                      0              0
       2       4              0     1                      0              0
       3       5              0     1                      0              0
       4       3              0     1                      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  salary_low  \
       0                 1                   0                     0           1
       1                 1                   0                     0           0
       2                 1                   0                     0           0

| | | | | | |
|---|---|---|---|---|---|
| 3 | 1 | 0 | 0 | 1 | |
| 4 | 1 | 0 | 0 | 1 | |

| | salary_medium | salary_high |
|---|---|---|
| 0 | 0 | 0 |
| 1 | 1 | 0 |
| 2 | 1 | 0 |
| 3 | 0 | 0 |
| 4 | 0 | 0 |

```
[26]: #Create a heatmap to visualize how correlated variables are
      plt.figure(figsize=(8, 6))
      sns.heatmap(df_enc[['satisfaction_level', 'last_evaluation', 'number_project',␣
       ↪'number_project', 'average_monthly_hours', 'tenure']]
                   .corr(), annot=True, cmap='crest')
      plt.title('Heatmap of the dataset')
      plt.show();
```



Heatmap of the dataset

```
[27]: # Create a stacked bart plot to visualize number of employees across␣
      ↪department, comparing those who left with those who didn't
      # In the legend, 0 (purple color) represents employees who did not leave, 1␣
      ↪(red color) represents employees who left
      pd.crosstab(df1['department'], df1['left']).plot(kind ='bar',color='mr')
      plt.title('Counts of employees who left versus stayed across department')
      plt.ylabel('Employee count')
      plt.xlabel('Department')
      plt.show()
```



```
[28]: # Select rows without outliers in `tenure` and save resulting dataframe in a␣
      ↪new variable
      df_logreg = df_enc[(df_enc['tenure'] >= lower_limit) & (df_enc['tenure'] <=␣
      ↪upper_limit)]

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

```
[28]:     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  department_IT  \
     0       3              0     1                      0              0
     2       4              0     1                      0              0
     3       5              0     1                      0              0
     4       3              0     1                      0              0
     5       3              0     1                      0              0

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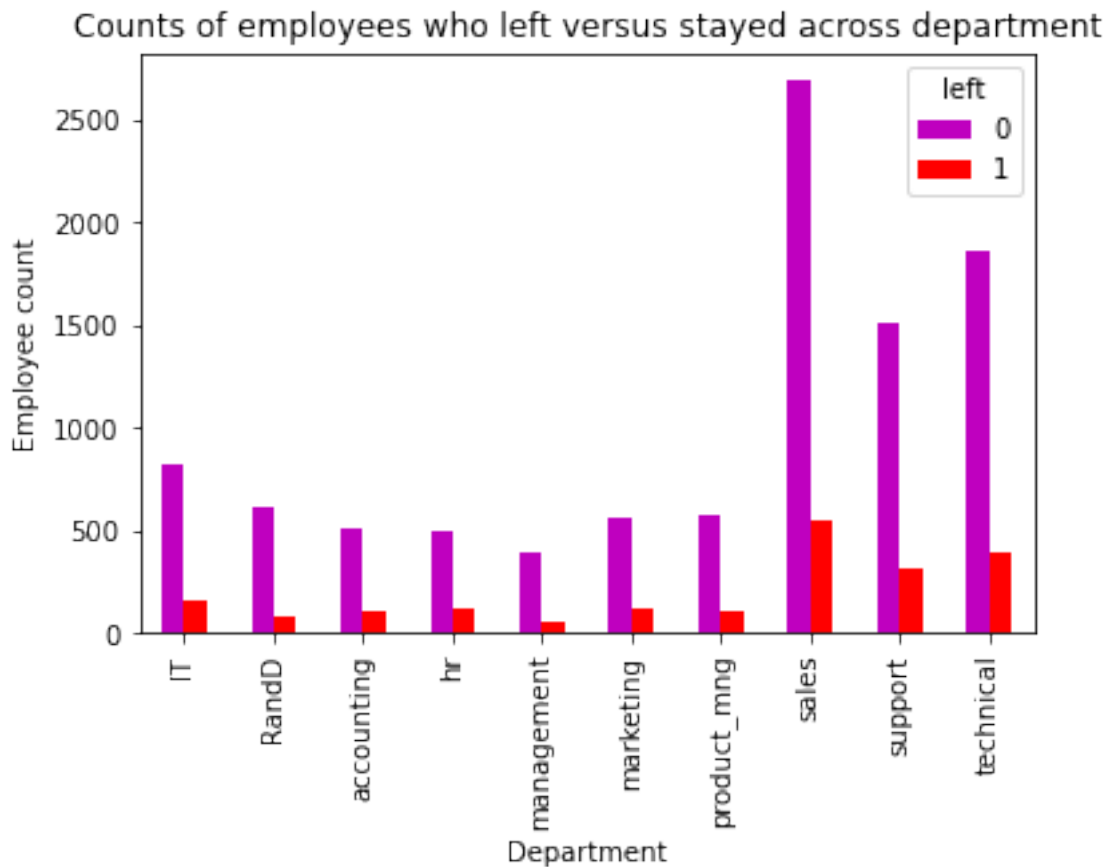     4                 0                      0              0
     5                 0                      0              0

        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  salary_low  \
     0                 1                   0                     0           1
     2                 1                   0                     0           0
     3                 1                   0                     0           1
     4                 1                   0                     0           1
     5                 1                   0                     0           1

        salary_medium  salary_high
     0              0            0
     2              1            0
     3              0            0
     4              0            0
     5              0            0
```

```python
[29]: #Isolate the outcome variable
      y = df_logreg['left']

      #Display first few rows of the outcome variables
      y.head()
```

```
[29]: 0    1
      2    1
      3    1
      4    1
      5    1
      Name: left, dtype: int64
```

```
[30]: #select the features we want to use in the model
      X = df_logreg.drop('left', axis=1)

      #Display the first few rows of the selected features
      X.head()
```

```
[30]:    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  promotion_last_5years  department_IT  \
      0       3              0                      0              0
      2       4              0                      0              0
      3       5              0                      0              0
      4       3              0                      0              0
      5       3              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  salary_low  \
      0                 1                   0                     0           1
      2                 1                   0                     0           0
      3                 1                   0                     0           1
      4                 1                   0                     0           1
      5                 1                   0                     0           1
```

```
      salary_medium  salary_high
0                0            0
2                1            0
3                0            0
4                0            0
5                0            0
```

[31]: *#split the data into training set and testing set*
```
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.25,␣
 ↪stratify=y, random_state=42)
```

[32]: *#Construct a logistiv regression model and fit it to the training dataset*
```
log_clf = LogisticRegression(random_state=42, max_iter=500).fit(X_train,␣
 ↪y_train)
```

[33]: *#Use logistic regression model to get predictions on the test set*
```
y_pred = log_clf.predict(X_test)
```

[34]: *#Compute values for confusion matrix*
```
log_cm = confusion_matrix(y_test, y_pred, labels=log_clf.classes_)
```

*#Create display of confusion matrix*
```
log_disp = ConfusionMatrixDisplay(confusion_matrix=log_cm,
                                  display_labels=log_clf.classes_)
```

*#plot confusion matrix*
```
log_disp.plot(values_format='')
```

*#Display plot*
```
plt.show();
```

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.

Create a classification report that includes precision, recall, f1-score, and accuracy metrics to evaluate the performance of the logistic regression model.

Check the class balance in the data. In other words, check the value counts in the left column. Since this is a binary classification task, the class balance informs the way you interpret accuracy metrics.

```
[35]: df_logreg['left'].value_counts(normalize=True)
```

```
[35]: 0    0.831468
      1    0.168532
      Name: left, dtype: float64
```

There is an approximately 83%-17% split. So the data is not prefectly balanced, but it is not too imbalanced. If it was more severely imbalanced, we might want to resample the data to make it more balanced. We can use this data without modifying the class balance and continue evaluating the model.

```
[36]: #Create classification report for logistic regression model
      target_names = ['Predicted would not leave', 'Predicted would leave']
      print(classification_report(y_test, y_pred, target_names=target_names))
```

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| Predicted would not leave | 0.86 | 0.94 | 0.90 | 2321 |
| Predicted would leave | 0.45 | 0.26 | 0.33 | 471 |
| accuracy |  |  | 0.82 | 2792 |
| macro avg | 0.66 | 0.60 | 0.62 | 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 or 80%, and accuracy of 82%. However, if it'S most important to predict employees who leave, then the scores are signigicantly lower.

**Modelling Approach B: Tree-based Model**  This approach covers implementation of Decision Tree and Random Forest.

```
[37]: #Isolate the outcome variable
      y = df_enc['left']

      #Display the first few rows of 'y'
      y.head()
```

```
[37]: 0    1
      1    1
      2    1
      3    1
      4    1
      Name: left, dtype: int64
```

```
[38]: #select the features
      X = df_enc.drop('left', axis=1)

      #Display the first few rows of 'X'
      X.head()
```

```
[38]:    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  department_IT  \
0        3              0                      0              0
1        6              0                      0              0
2        4              0                      0              0
3        5              0                      0              0
4        3              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  salary_low  \
0                  1                   0                     0           1
1                  1                   0                     0           0
2                  1                   0                     0           0
3                  1                   0                     0           1
4                  1                   0                     0           1

    salary_medium  salary_high
0              0            0
1              1            0
2              1            0
3              0            0
4              0            0
```

[39]:
```python
#Split the data
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.25,␣
 ↪stratify=y, random_state=0)
```

[40]:
```python
#Decision tree - Round 1
# Instantiate model
tree = DecisionTreeClassifier(random_state=0)
```

```
# Assign a dictionary of hyperparameters to search over
cv_params = {'max_depth':[4, 6, 8, None],
             'min_samples_leaf': [2, 5, 1],
             'min_samples_split': [2, 4, 6]
             }

# Assign a dictionary of scoring metrics to capture
scoring = {'accuracy', 'precision', 'recall', 'f1', 'roc_auc'}

# Instantiate GridSearch
tree1 = GridSearchCV(tree, cv_params, scoring=scoring, cv=4, refit='roc_auc')
```

Fit the decision tree model to the training data

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

```
CPU times: user 3.5 s, sys: 227 ms, total: 3.73 s
Wall time: 3.73 s
```

```
[41]: 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={'roc_auc', 'precision', 'f1', 'recall', 'accuracy'},
                   verbose=0)
```

```
[42]: #Check best parameters
      tree1.best_params_
```

```
[42]: {'max_depth': 4, 'min_samples_leaf': 5, 'min_samples_split': 2}
```

```
[43]: #Check best AUC score on CV
      tree1.best_score_
```

```
[43]: 0.969819392792457
```

This is a strong AUC score, which shows that this model can predict employees who will leave very well.

```python
[44]: #extract all the scores from the grid search
      def make_results(model_name:str, model_object, metric:str):

          # 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
```

```python
[45]: #Get all CV scores
      tree1_cv_results = make_results('decision tree cv', tree1, 'auc')
      tree1_cv_results
```

```
[45]:              model  precision    recall        F1  accuracy       auc
      0  decision tree cv   0.914552  0.916949  0.915707  0.971978  0.969819
```

All of these scores from the decision tree model are strong indicators of good model performance.

```python
[46]: # Random forest - Round 1
      # Instantiate model
      rf = RandomForestClassifier(random_state=0)

      # Assign a dictionary of hyperparameters to search over
      cv_params = {'max_depth': [3,5, None],
                   'max_features': [1.0],
                   'max_samples': [0.7, 1.0],
                   'min_samples_leaf': [1,2,3],
                   'min_samples_split': [2,3,4],
                   'n_estimators': [300, 500],
                   }

      # Assign a dictionary of scoring metrics to capture
      scoring = {'accuracy', 'precision', 'recall', 'f1', 'roc_auc'}

      # Instantiate GridSearch
      rf1 = GridSearchCV(rf, cv_params, scoring=scoring, cv=4, refit='roc_auc')
```

```python
[47]: %%time
      rf1.fit(X_train, y_train)
```

```
CPU times: user 10min 33s, sys: 9.47 s, total: 10min 43s
Wall time: 10min 43s
```

```
[47]: GridSearchCV(cv=4, error_score=nan,
                   estimator=RandomForestClassifier(bootstrap=True, ccp_alpha=0.0,
                                                    class_weight=None,
                                                    criterion='gini', max_depth=None,
                                                    max_features='auto',
                                                    max_leaf_nodes=None,
                                                    max_samples=None,
                                                    min_impurity_decrease=0.0,
                                                    min_impurity_split=None,
                                                    min_samples_leaf=1,
                                                    min_samples_split=2,
                                                    min_weight_fraction_leaf=0.0,
                                                    n_estimators=100, n_jobs=None,…
                                                    verbose=0, warm_start=False),
                   iid='deprecated', n_jobs=None,
                   param_grid={'max_depth': [3, 5, None], 'max_features': [1.0],
                               'max_samples': [0.7, 1.0],
                               'min_samples_leaf': [1, 2, 3],
```

```
                    'min_samples_split': [2, 3, 4],
                    'n_estimators': [300, 500]},
            pre_dispatch='2*n_jobs', refit='roc_auc', return_train_score=False,
            scoring={'roc_auc', 'precision', 'f1', 'recall', 'accuracy'},
            verbose=0)
```

[58]:
```python
# Define a path to the folder where you want to save the model
path = '/home/jovyan/work/'
```

[59]:
```python
def write_pickle(path, model_object, save_as:str):
    '''
    In:
        path:         path of folder where you want to save the pickle
        model_object: a model you want to pickle
        save_as:      filename for how you want to save the model

    Out: A call to pickle the model in the folder indicated
    '''

    with open(path + save_as + '.pickle', 'wb') as to_write:
        pickle.dump(model_object, to_write)
```

[60]:
```python
def read_pickle(path, saved_model_name:str):
    '''
    In:
        path:             path to folder where you want to read from
        saved_model_name: filename of pickled model you want to read in

    Out:
        model: the pickled model
    '''
    with open(path + saved_model_name + '.pickle', 'rb') as to_read:
        model = pickle.load(to_read)

    return model
```

[61]:
```python
# Write pickle
write_pickle(path, rf1, 'hr_rf1')
```

[62]:
```python
# Read pickle
rf1 = read_pickle(path, 'hr_rf1')
```

[63]:
```python
# Check best AUC score on CV
rf1.best_score_
```

[63]: 0.980416028921927

```
[64]: # Check best params
      rf1.best_params_
```

```
[64]: {'max_depth': 5,
       'max_features': 1.0,
       'max_samples': 0.7,
       'min_samples_leaf': 2,
       'min_samples_split': 2,
       'n_estimators': 500}
```

```
[65]: # Get all CV scores
      rf1_cv_results = make_results('random forest cv', rf1, 'auc')
      print(tree1_cv_results)
      print(rf1_cv_results)
```

```
                 model  precision    recall        F1  accuracy       auc
0   decision tree cv   0.914552  0.916949  0.915707  0.971978  0.969819
                 model  precision    recall        F1  accuracy       auc
0   random forest cv   0.949886  0.912933  0.931017  0.977538  0.980416
```

The evaluation scores of the random forest model are better than those of the decision tree model, with the exception of recall (the recall score of the random forest model is approximately 0.001 lower, which is a negligible amount). This indicates that the random forest model mostly outperforms the decision tree model.

```
[66]: def get_scores(model_name:str, model, X_test_data, y_test_data):
          '''
          Generate a table of test scores.

          In:
              model_name (string):   How you want your model to be named in the output␣
      ↪table
              model:                 A fit GridSearchCV object
              X_test_data:           numpy array of X_test data
              y_test_data:           numpy array of y_test data

          Out: pandas df of precision, recall, f1, accuracy, and AUC scores for your␣
      ↪model
          '''

          preds = model.best_estimator_.predict(X_test_data)

          auc = roc_auc_score(y_test_data, preds)
          accuracy = accuracy_score(y_test_data, preds)
          precision = precision_score(y_test_data, preds)
          recall = recall_score(y_test_data, preds)
          f1 = f1_score(y_test_data, preds)
```

```
        table = pd.DataFrame({'model': [model_name],
                              'precision': [precision],
                              'recall': [recall],
                              'f1': [f1],
                              'accuracy': [accuracy],
                              'AUC': [auc]
                            })

        return table
```

[67]:
```
# Get predictions on test data
rf1_test_scores = get_scores('random forest1 test', rf1, X_test, y_test)
rf1_test_scores
```

[67]:
| | model | precision | recall | f1 | accuracy | AUC |
|---|---|---|---|---|---|---|
| 0 | random forest1 test | 0.964135 | 0.917671 | 0.940329 | 0.980654 | 0.955435 |

The test scores are very similar to the validation scores, which is good. This appears to be a strong model. Since this test set was only used for this model, we can be more confident that your model's performance on this data is representative of how it will perform on new, unseeen data.

**Feature Engineering** It might be skeptical of the high evaluation scores. There is a chance that there is some data leakage occurring. Data leakage is when we use data to train the model that should not be used during training, either because it appears in the test data or because it's not data that you'd expect to have when the model is actually deployed. Training a model with leaked data can give an unrealistic score that is not replicated in production.

In this case, it's likely that the company won't have satisfaction levels reported for all of its employees. It's also possible that the average_monthly_hours column is a source of some data leakage. If employees have already decided upon quitting, or have already been identified by management as people to be fired, they may be working fewer hours.

The first round of decision tree and random forest models included all variables as features. This next round will incorporate feature engineering to build improved models.

We can proceed by dropping satisfaction_level and creating a new feature that roughly captures whether an employee is overworked. we could call this new feature overworked. It will be a binary variable.

[68]:
```
#Drop 'satisfaction_level' and save resulting dataframe in new variable
df2 = df_enc.drop('satisfaction_level', axis=1)

#Display first few rows of new dataframe
df2.head()
```

[68]:
| | last_evaluation | number_project | average_monthly_hours | tenure \ |
|---|---|---|---|---|
| 0 | 0.53 | 2 | 157 | 3 |
| 1 | 0.86 | 5 | 262 | 6 |

|   |       | 7 | 272 | 4 |
|---|-------|---|-----|---|
| 2 | 0.88  | 7 | 272 | 4 |
| 3 | 0.87  | 5 | 223 | 5 |
| 4 | 0.52  | 2 | 159 | 3 |

|   | work_accident | left | promotion_last_5years | department_IT \ |
|---|---------------|------|-----------------------|-----------------|
| 0 | 0 | 1 | 0 | 0 |
| 1 | 0 | 1 | 0 | 0 |
| 2 | 0 | 1 | 0 | 0 |
| 3 | 0 | 1 | 0 | 0 |
| 4 | 0 | 1 | 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 | salary_low \ |
|---|------------------|--------------------|----------------------|--------------|
| 0 | 1 | 0 | 0 | 1 |
| 1 | 1 | 0 | 0 | 0 |
| 2 | 1 | 0 | 0 | 0 |
| 3 | 1 | 0 | 0 | 1 |
| 4 | 1 | 0 | 0 | 1 |

|   | salary_medium | salary_high |
|---|---------------|-------------|
| 0 | 0 | 0 |
| 1 | 1 | 0 |
| 2 | 1 | 0 |
| 3 | 0 | 0 |
| 4 | 0 | 0 |

```
[69]:  # Create `overworked` column. For now, it's identical to average monthly hours.
       df2['overworked'] = df2['average_monthly_hours']

       # Inspect max and min average monthly hours values
       print('Max hours:', df2['overworked'].max())
       print('Min hours:', df2['overworked'].min())
```

```
Max hours: 310
```

```
Min hours: 96
```

166.67 is approximately the average number of monthly hours for someone who works 50 weeks per year, 5 days per week, 8 hours per day.

We could define being overworked as working more than 175 hours per month on average. To make the **overworked** column binary, you could reassign the column using a boolean mask. * df3['overworked'] > 175 creates a series of booleans, consisting of True for every value>175 and **False** for every values   175 * .astype(int) converts all **True** to 1 and all **False** to 0.

```python
[70]: #Define 'overworked' as working > 175 hrs/week
      df2['overworked'] = (df2['overworked'] > 175).astype(int)

      #Display first few rows of new column
      df2['overworked'].head()
```

```
[70]: 0    0
      1    1
      2    1
      3    1
      4    0
      Name: overworked, dtype: int64
```

```python
[71]: #Drop the 'average_monthly_hours' column
      df2 = df2.drop('average_monthly_hours', axis=1)

      #Display first few rows of resulting dataframe
      df2.head()
```

```
[71]:    last_evaluation  number_project  tenure  work_accident  left  \
      0             0.53               2       3              0     1
      1             0.86               5       6              0     1
      2             0.88               7       4              0     1
      3             0.87               5       5              0     1
      4             0.52               2       3              0     1

         promotion_last_5years  department_IT  department_RandD  \
      0                      0              0                 0
      1                      0              0                 0
      2                      0              0                 0
      3                      0              0                 0
      4                      0              0                 0

         department_accounting  department_hr  department_management  \
      0                      0              0                      0
      1                      0              0                      0
      2                      0              0                      0
      3                      0              0                      0
```

```
4                        0                0                          0

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

    department_support  department_technical  salary_low  salary_medium  \
0                    0                     0           1              0
1                    0                     0           0              1
2                    0                     0           0              1
3                    0                     0           1              0
4                    0                     0           1              0

    salary_high  overworked
0             0           0
1             0           1
2             0           1
3             0           1
4             0           0
```

```python
[72]: #Isolate the outcome variable
      y = df2['left']

      #select the features
      X = df2.drop('left', axis=1)
```

```python
[73]: #Create test data
      X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.25,
      ↪stratify=y, random_state=0)
```

```python
[74]: #Decision Tree- Round 2
      # Instantiate model
      tree = DecisionTreeClassifier(random_state=0)

      # Assign a dictionary of hyperparameters to search over
      cv_params = {'max_depth':[4, 6, 8, None],
                   'min_samples_leaf': [2, 5, 1],
                   'min_samples_split': [2, 4, 6]
                   }

      # Assign a dictionary of scoring metrics to capture
      scoring = {'accuracy', 'precision', 'recall', 'f1', 'roc_auc'}

      # Instantiate GridSearch
```

```
tree2 = GridSearchCV(tree, cv_params, scoring=scoring, cv=4, refit='roc_auc')
```

[75]:
```
%%time
tree2.fit(X_train, y_train)
```

CPU times: user 2.69 s, sys: 10.4 ms, total: 2.7 s
Wall time: 2.7 s

[75]: 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={'roc_auc', 'precision', 'f1', 'recall', 'accuracy'},
                verbose=0)
```

[76]:
```
#Check best params
tree2.best_params_
```

[76]: {'max_depth': 6, 'min_samples_leaf': 2, 'min_samples_split': 6}

[77]:
```
#Check best AUC score on CV
tree2.best_score_
```

[77]: 0.9577889277573082

[78]:
```
# Get all CV scores
tree2_cv_results = make_results('decision tree2 cv', tree2, 'auc')
print(tree1_cv_results)
print(tree2_cv_results)
```

|   | model | precision | recall | F1 | accuracy | auc |
|---|-------|-----------|--------|-----|----------|-----|
| 0 | decision tree cv | 0.914552 | 0.916949 | 0.915707 | 0.971978 | 0.969819 |

|   | model | precision | recall | F1 | accuracy | auc |
|---|-------|-----------|--------|-----|----------|-----|
| 0 | decision tree2 cv | 0.854952 | 0.903553 | 0.878003 | 0.958189 | 0.957789 |

Some of the other scores fell. That's to be expected given fewer features were taken into account in this round of the model. The scores are very good.

```
[79]: #Random forest - Round 2
      # Instantiate model
      rf = RandomForestClassifier(random_state=0)

      # Assign a dictionary of hyperparameters to search over
      cv_params = {'max_depth': [3,5, None],
                   'max_features': [1.0],
                   'max_samples': [0.7, 1.0],
                   'min_samples_leaf': [1,2,3],
                   'min_samples_split': [2,3,4],
                   'n_estimators': [300, 500],
                   }

      # Assign a dictionary of scoring metrics to capture
      scoring = {'accuracy', 'precision', 'recall', 'f1', 'roc_auc'}

      # Instantiate GridSearch
      rf2 = GridSearchCV(rf, cv_params, scoring=scoring, cv=4, refit='roc_auc')
```

```
[80]: %%time
      rf2.fit(X_train, y_train)
```

```
CPU times: user 8min, sys: 1.87 s, total: 8min 2s
Wall time: 8min 2s
```

```
[80]: GridSearchCV(cv=4, error_score=nan,
                   estimator=RandomForestClassifier(bootstrap=True, ccp_alpha=0.0,
                                                     class_weight=None,
                                                     criterion='gini', max_depth=None,
                                                     max_features='auto',
                                                     max_leaf_nodes=None,
                                                     max_samples=None,
                                                     min_impurity_decrease=0.0,
                                                     min_impurity_split=None,
                                                     min_samples_leaf=1,
                                                     min_samples_split=2,
                                                     min_weight_fraction_leaf=0.0,
                                                     n_estimators=100, n_jobs=None,…
                                                     verbose=0, warm_start=False),
                   iid='deprecated', n_jobs=None,
                   param_grid={'max_depth': [3, 5, None], 'max_features': [1.0],
                               'max_samples': [0.7, 1.0],
                               'min_samples_leaf': [1, 2, 3],
                               'min_samples_split': [2, 3, 4],
                               'n_estimators': [300, 500]},
```

```
                pre_dispatch='2*n_jobs', refit='roc_auc', return_train_score=False,
                scoring={'roc_auc', 'precision', 'f1', 'recall', 'accuracy'},
                verbose=0)
```

[81]: 
```
#write pickle
write_pickle(path, rf2, 'hr_rf2')
```

[82]: 
```
#Read in pickle
rf2 = read_pickle(path, 'hr_rf2')
```

[83]: 
```
#Check best params
rf2.best_params_
```

[83]: 
```
{'max_depth': 5,
 'max_features': 1.0,
 'max_samples': 0.7,
 'min_samples_leaf': 1,
 'min_samples_split': 2,
 'n_estimators': 300}
```

[84]: 
```
#check best AUC score on CV
rf2.best_score_
```

[84]: 0.9648895083941449

[85]: 
```
#Get all CV scores
rf2_cv_results = make_results('random forests cv', rf2, 'auc')
print(tree2_cv_results)
print(rf2_cv_results)
```

```
                model  precision    recall        F1  accuracy       auc
0  decision tree2 cv   0.854952  0.903553  0.878003  0.958189  0.957789
                model  precision    recall        F1  accuracy       auc
0  random forests cv   0.866353  0.880093   0.87287  0.957523   0.96489
```

The scores dropped slightly, but the random forest performs better than the decision tree if using AUC as the deciding metric.
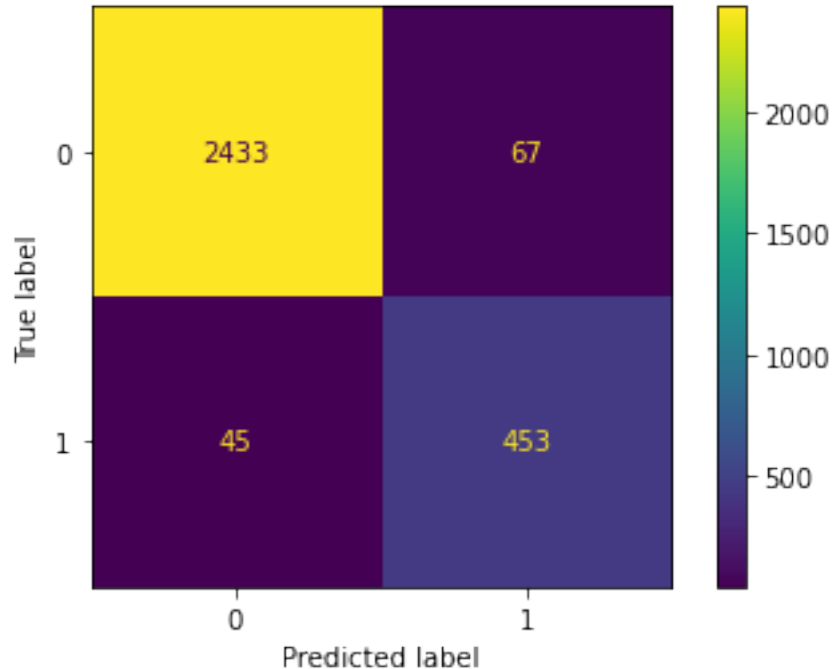
[86]: 
```
#Get predictions on test data
rf2_test_scores = get_scores('random forest2 test', rf2, X_test, y_test)
rf2_test_scores
```

[86]: 
```
                 model  precision    recall       f1  accuracy       AUC
0  random forest2 test   0.871154  0.909639  0.88998  0.962642  0.941419
```

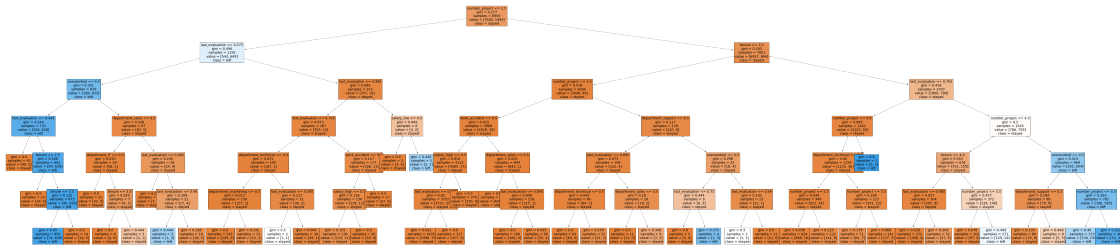This seems to be a stable, well-performing final model.

```
[88]: #Generate array of values for confusion matrix
      preds = rf2.best_estimator_.predict(X_test)
      cm = confusion_matrix(y_test, preds, labels=rf2.classes_)

      # Plot confusion matrix
      disp = ConfusionMatrixDisplay(confusion_matrix=cm,
                                    display_labels=rf2.classes_)
      disp.plot(values_format='');
```



The model predicts more false positives than false negatives, which means that some employees may be identified as at risk of quitting or getting fired, when that's actually not the case. But this is still a strong model.

```
[89]: #Decision tree splits
      #plot the tree
      plt.figure(figsize=(85,20))
      plot_tree(tree2.best_estimator_, max_depth=6, fontsize=14, feature_names=X.
       ↪columns,
                class_names={0:'stayed', 1:'left'}, filled=True);
      plt.show()
```

Double-click on the tree image to zoom in on it and inspect the splits.

```
[90]:  #Decision tree feature importance
       #tree2_importances = pd.DataFrame(tree2.best_estimator_.feature_importances_,
        ↪columns=X.columns)
       tree2_importances = pd.DataFrame(tree2.best_estimator_.feature_importances_,
                                        columns=['gini_importance'],
                                        index=X.columns
                                        )
       tree2_importances = tree2_importances.sort_values(by='gini_importance',
        ↪ascending=False)

       # Only extract the features with importances > 0
       tree2_importances = tree2_importances[tree2_importances['gini_importance'] != 0]
       tree2_importances
```
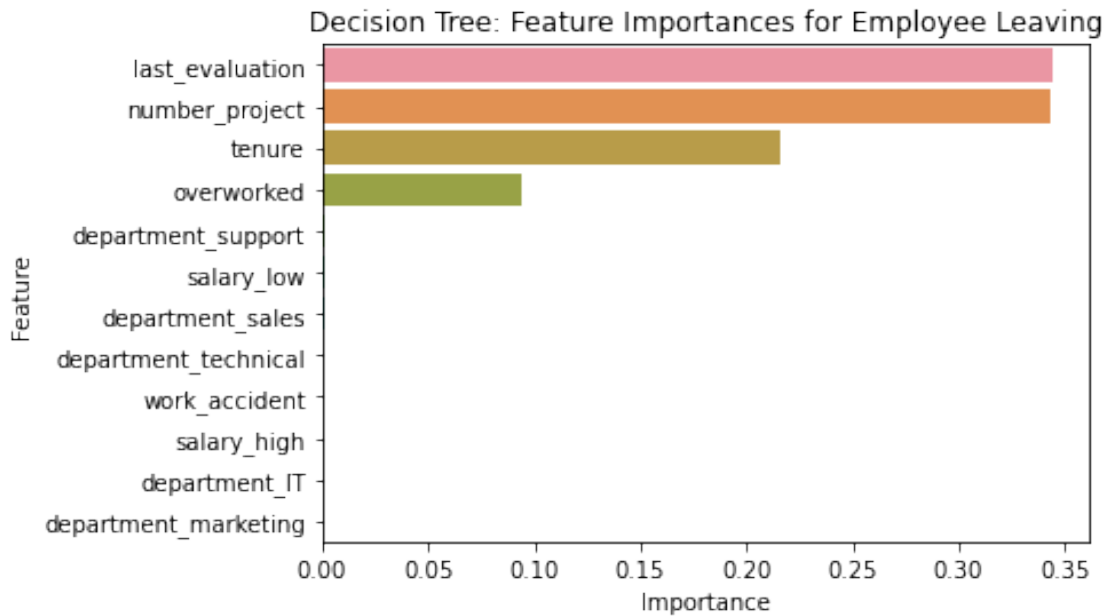
```
[90]:                        gini_importance
       last_evaluation              0.343958
       number_project               0.343385
       tenure                       0.215681
       overworked                   0.093498
       department_support           0.001142
       salary_low                   0.000737
       department_sales             0.000607
       department_technical         0.000418
       work_accident                0.000183
       salary_high                  0.000173
       department_IT                0.000139
       department_marketing         0.000078
```

```
[91]:  #create barplot to visualize the decision tree feature importances.
       sns.barplot(data=tree2_importances, x="gini_importance", y=tree2_importances.
        ↪index, orient='h')
       plt.title("Decision Tree: Feature Importances for Employee Leaving",
        ↪fontsize=12)
       plt.ylabel("Feature")
       plt.xlabel("Importance")
```

```
plt.show()
```



Decision Tree: Feature Importances for Employee Leaving

The **barplot** above shows that in this decision tree model, **last_evaluation**, **number_project**, **tenure**, and **overworked** have the highest importance, in that order. These variables are most helpful in prediciting the outcome variable, **left**.

```
[92]:  #Random forest feature importance
       # Get feature importances
       feat_impt = rf2.best_estimator_.feature_importances_

       # Get indices of top 10 features
       ind = np.argpartition(rf2.best_estimator_.feature_importances_, -10)[-10:]

       # Get column labels of top 10 features
       feat = X.columns[ind]

       # Filter `feat_impt` to consist of top 10 feature importances
       feat_impt = feat_impt[ind]

       y_df = pd.DataFrame({"Feature":feat,"Importance":feat_impt})
       y_sort_df = y_df.sort_values("Importance")
       fig = plt.figure()
       ax1 = fig.add_subplot(111)

       y_sort_df.plot(kind='barh',ax=ax1,x="Feature",y="Importance")
```
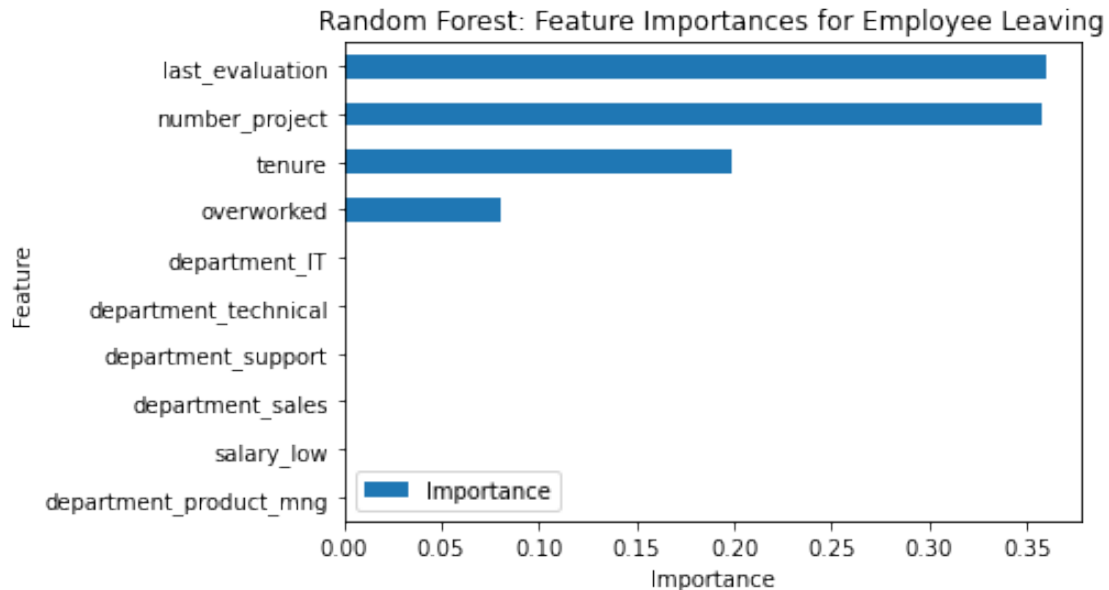
```
ax1.set_title("Random Forest: Feature Importances for Employee Leaving",␣
  ↪fontsize=12)
ax1.set_ylabel("Feature")
ax1.set_xlabel("Importance")

plt.show()
```



Random Forest: Feature Importances for Employee Leaving

The plot above shows that in this random forest model, last_evaluation, number_project, tenure, and overworked have the highest importance, in that order. These variables are most helpful in predicting the outcome variable, left, and they are the same as the ones used by the decision tree model.

[ ]:

[ ]:

# 5  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.

### Reflect on these questions as you complete the executing stage.

- What key insights emerged from your model(s)?
- What business recommendations do you propose based on the models built?
- What potential recommendations would you make to your manager/company?
- Do you think your model could be improved? Why or why not? How?
- Given what you know about the data and the models you were using, what other questions could you address for the team?
- What resources do you find yourself using as you complete this stage? (Make sure to include the links.)
- Do you have any ethical considerations in this stage?

Double-click to enter your responses here.

## 5.1 Step 4. Results and Evaluation

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

### 5.1.1 Summary of model results

[Double-click to enter your summary here.]

### 5.1.2 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 atleast 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.

**Congratulations!** You've completed this lab. However, you may not notice a green check mark next to this item on Coursera's platform. Please continue your progress regardless of the check mark. Just click on the "save" icon at the top of this notebook to ensure your work has been logged.