

Activity_ Course 7 Salifort Motors project lab

July 14, 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)

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?

[Double-click to enter your responses here.]

2.2 Step 1. Imports

- Import packages

- Load dataset

2.2.1 Import packages

```
[50]: import numpy as np
import pandas as pd

import matplotlib.pyplot as plt
import seaborn as sns

pd.set_option('display.max_columns', None)

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.

```
[51]: # 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()
```

```
[51]: 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

```
[52]: # 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_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
```

2.3.2 Gather descriptive statistics about the data

```
[53]: # Gather descriptive statistics about the data
df0.describe()
```

```
[53]:
```

	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

2.3.3 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.

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

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

```
[55]: # Rename columns as needed  
df0 = df0.rename(columns={'Work_accident': 'work_accident',  
                          'average_monthly_hours': 'average_monthly_hours',  
                          'time_spend_company': 'tenure',  
                          'Department': 'department'})  
  
# Display all column names after the update  
df0.columns
```

```
[55]: 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.

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

```
[56]: 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.

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

[57]: 3008

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

```
[58]:
```

	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

```
[59]: # Drop duplicates and save resulting dataframe in a new variable as needed
df2 = df0.drop_duplicates(keep = 'first')

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

```
[59]:
```

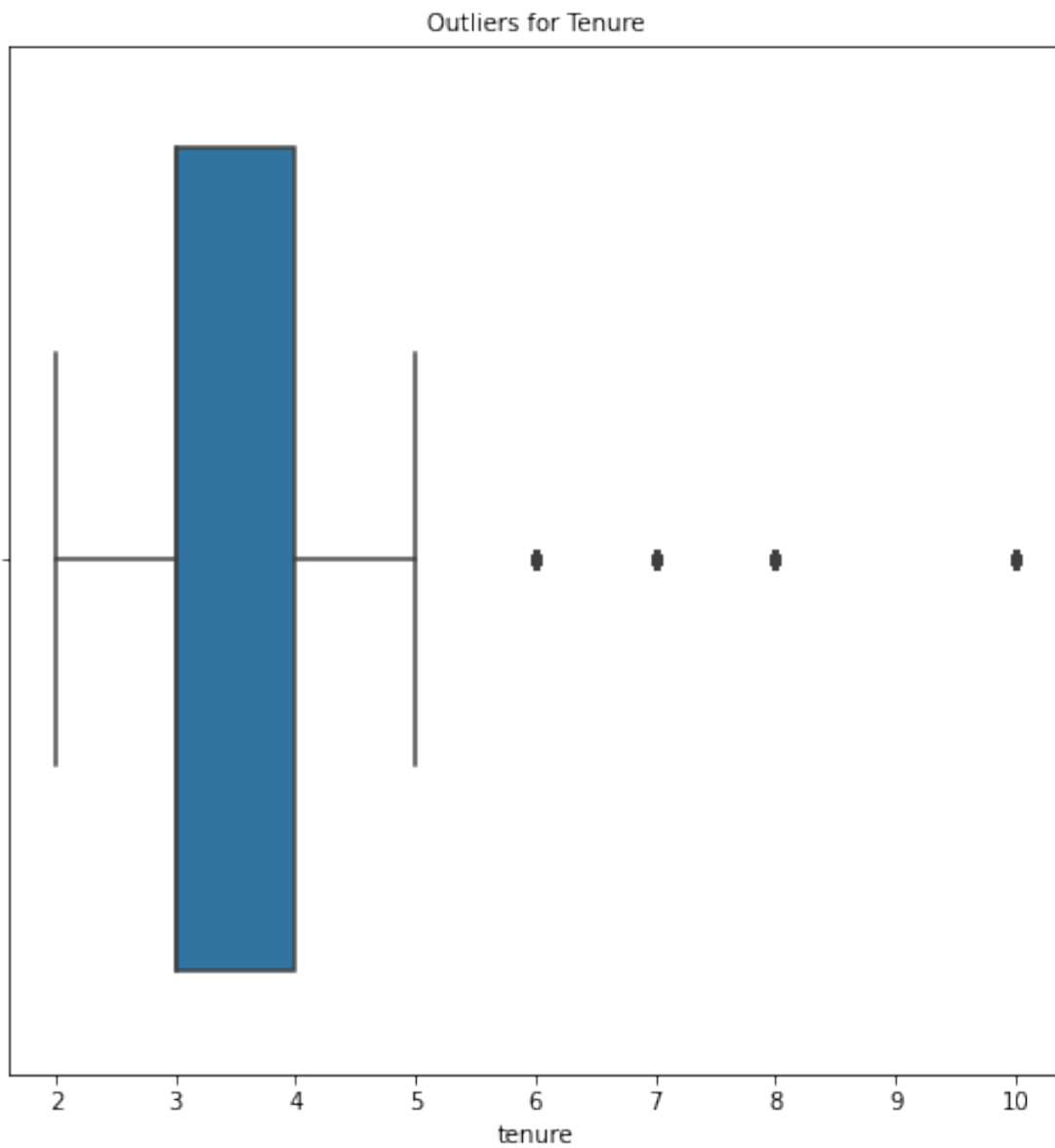
	satisfaction_level	last_evaluation	number_project	average_monthly_hours	\
0	0.38	0.53	2	157	
1	0.80	0.86	5	262	
2	0.11	0.88	7	272	
3	0.72	0.87	5	223	
4	0.37	0.52	2	159	

	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.

```
[60]: # Create a boxplot to visualize distribution of `tenure` and detect any outliers
plt.figure(figsize=(8,8))
plt.title('Outliers for Tenure', fontsize =10)
plt.xticks(fontsize = 10)
plt.yticks(fontsize = 10)
sns.boxplot(x = df2['tenure'])
plt.show()
```




```
[61]: # Determine the number of rows containing outliers

percentile25 = df2['tenure'].quantile(0.25)

percentile75 = df2['tenure'].quantile(0.75)

iqr = percentile75 - percentile25

upper_limit = percentile75 + 1.5 * iqr
lower_limit = percentile25 - 1.5 * iqr
print("Lower limit:", lower_limit)
print("Upper limit:", upper_limit)

outliers = (df2['tenure'] > upper_limit) | (df2['tenure'] < lower_limit)

print("Number of rows in the data containing outliers in `tenure`:", outliers.
      ↪sum())
```

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.

```
[62]: # Get numbers of people who left vs. stayed
print(df2['left'].value_counts())
print()

# Get percentages of people who left vs. stayed
print(df2['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.

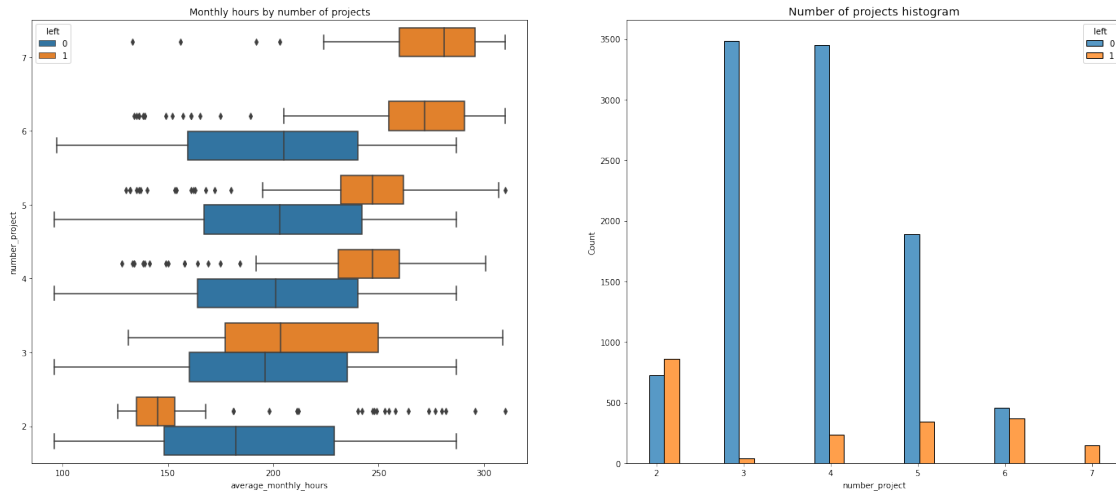
```
[63]: # Create a plot as needed
# First I am going to create a box plot showing the average_monthly_hours
# compared to the number_project
# variable to see the distribution between employees who stayed and who left.

fig, ax = plt.subplots(1, 2, figsize = (24, 10))

sns.boxplot(data = df2, 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 = '12')

tenure_stay = df2[df2['left']==0]['number_project']
tenure_left = df2[df2['left']==1]['number_project']
sns.histplot(data=df2, x='number_project', hue='left', multiple='dodge',
    shrink=2, ax=ax[1])
ax[1].set_title('Number of projects histogram', fontsize='14')

plt.show()
```



In general when people work on more projects this means they will have a greater number of hours worked. In the case above this seems to hold true.

```
[64]: # Create a plot as needed
plt.figure(figsize=(10,6))
sns.histplot(df2['satisfaction_level'], bins=20, kde=True, color='skyblue')
plt.title('Distribution of Job Satisfaction Levels', fontsize=14)
plt.xlabel('Satisfaction Level', fontsize=12)
plt.ylabel('Frequency', fontsize=12)
plt.show()
```



This figure above shows the satisfaction level of the employees on a scale from 0-1. From this plot you can tell where the majority of workers are concentrated and how satisfied they are with the company. From this it shows a positive sign with a higher concentration on the right side. The kde curve allows us to see the peaks of the graph and the overall distribution of satisfaction levels

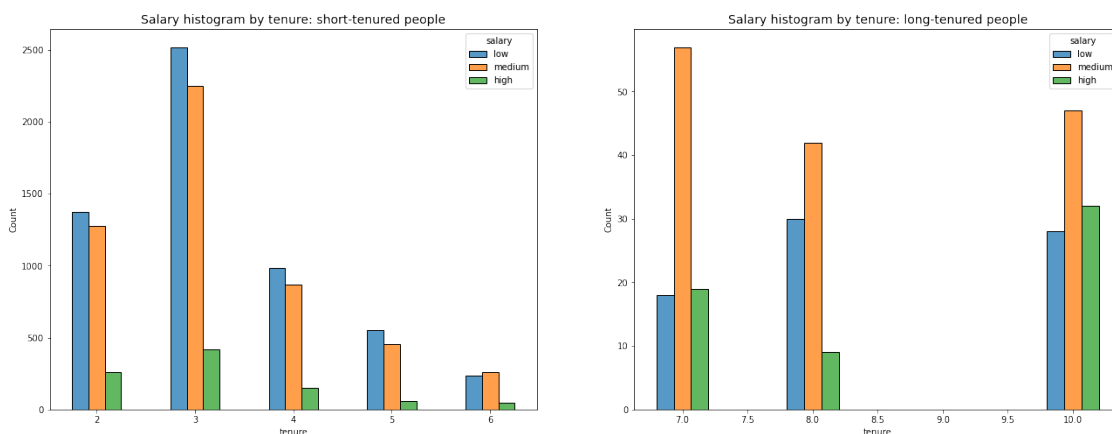
```
[65]: fig, ax = plt.subplots(1, 2, figsize = (22,8))

tenure_short = df2[df2['tenure'] < 7]

tenure_long = df2[df2['tenure'] > 6]

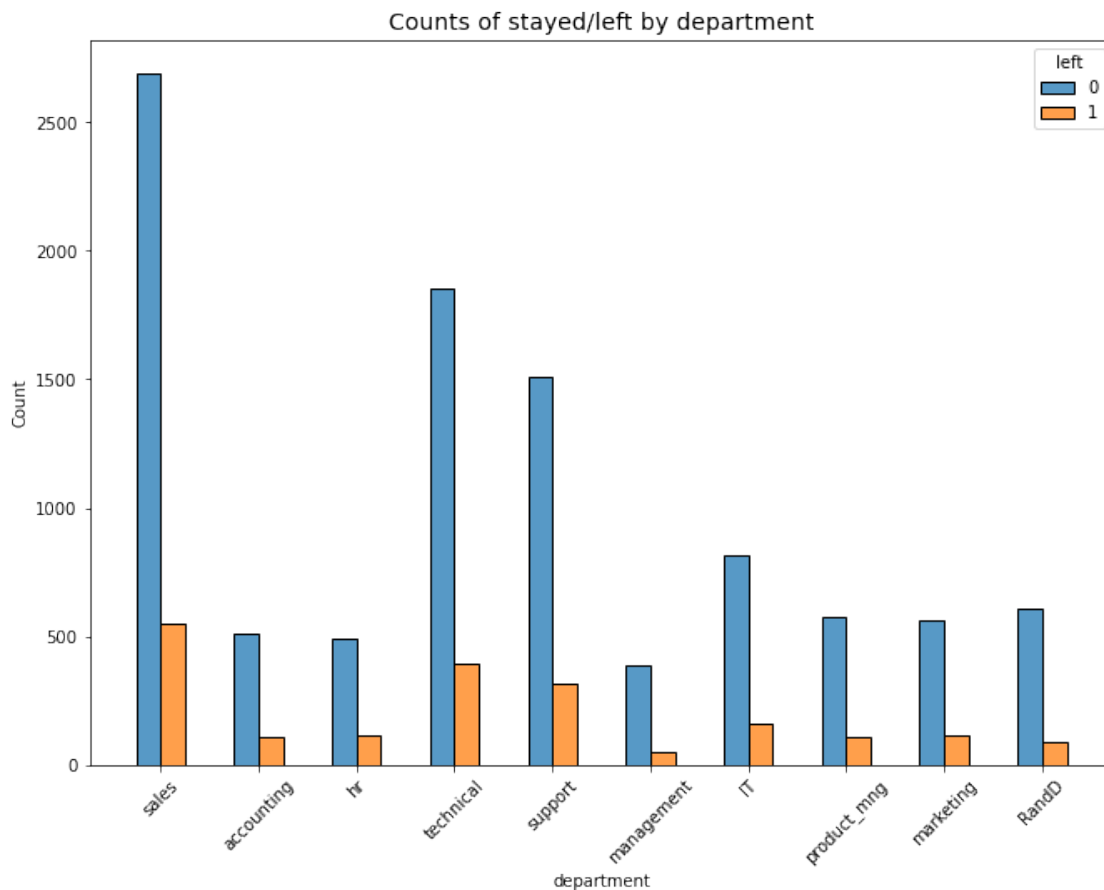
sns.histplot(data=tenure_short, x='tenure', hue='salary', discrete=1,
             hue_order=['low', 'medium', 'high'], multiple='dodge', shrink=.5,
             →ax=ax[0])
ax[0].set_title('Salary histogram by tenure: short-tenured people',
             →fontsize='14')

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='14');
```



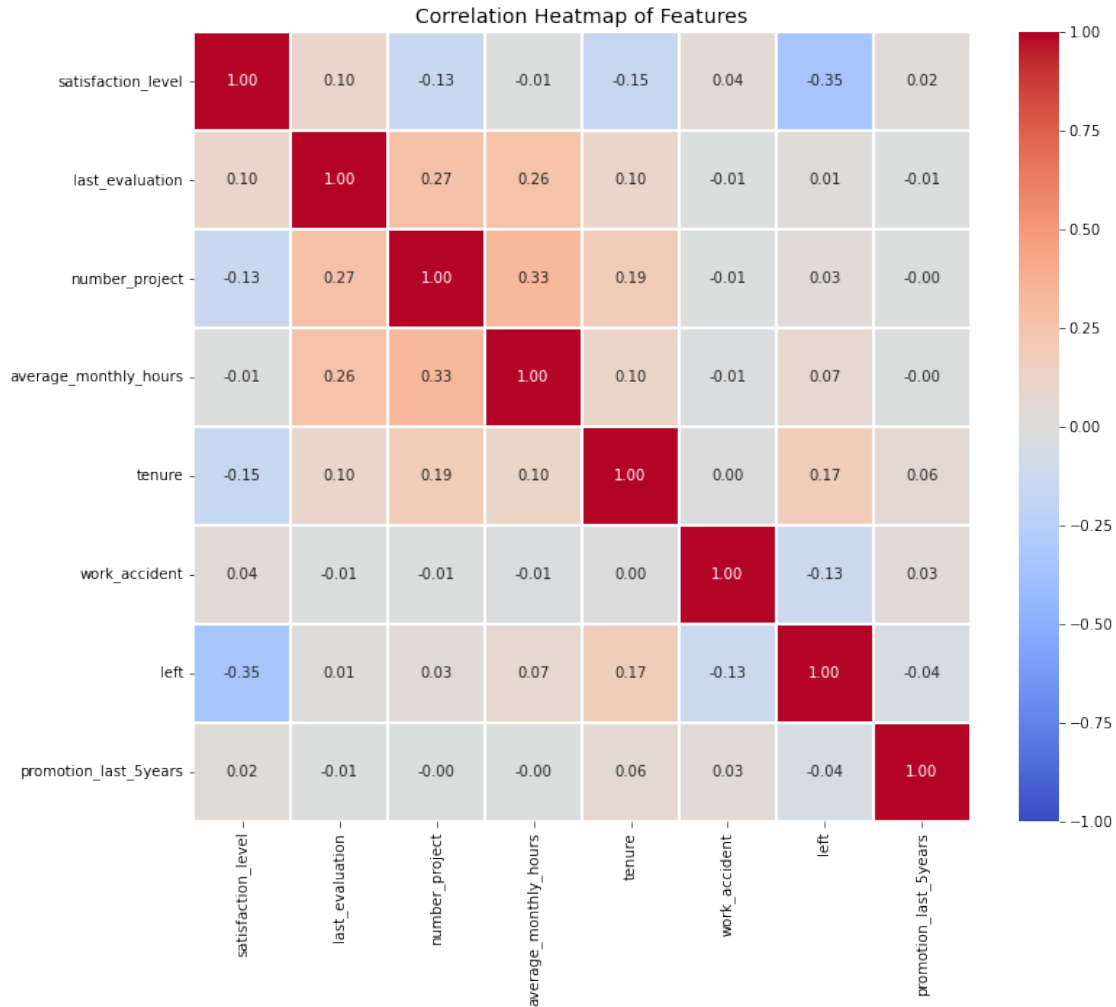
```
[66]: # An important measurement you would want to see is the number of employees
      →that stayes vs left based on their
      # position at the company.
plt.figure(figsize=(11,8))
sns.histplot(data=df2, x='department', hue='left', discrete=1,
             hue_order=[0, 1], multiple='dodge', shrink=.5)
```

```
plt.xticks(rotation='45')
plt.title('Counts of stayed/left by department', fontsize=14);
```



Based on the graph above there is a peak of employees leaving that are in the sales department and the smallest amount in the management department. Being able to look at this information can allow the company to make better decision moving forward to increase the satisfaction of the sales department. Seeing which sections have the most employees leaving allows them to make better decisions moving forward.

```
[67]: plt.figure(figsize=(12,10))
corr_matrix = df2.corr()
sns.heatmap(corr_matrix, annot=True, cmap='coolwarm', vmin=-1, vmax=1,
            linewidths=0.5, fmt='.2f')
plt.title('Correlation Heatmap of Features', fontsize=14)
plt.show()
```



3.1.2 Insights

Based on the various graphs created above people that are leaving the company are generally tied with working long hours, being unsatisfied, and having too many projects. It seems the employees become burned out after a certain amount of years and just become unsatisfied with the company. Too much work in a short time is causing them 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.

Whether an employee leaves the company.

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

Logistic Regression model or Tree-based Machine Learning model

4.1.3 Modeling

Add as many cells as you need to conduct the modeling process.

```
[68]: df_enc = df2.copy()

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

df_enc = pd.get_dummies(df_enc, drop_first=False)
```

```
df_enc.head()
```

```
[68]:
```

	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	

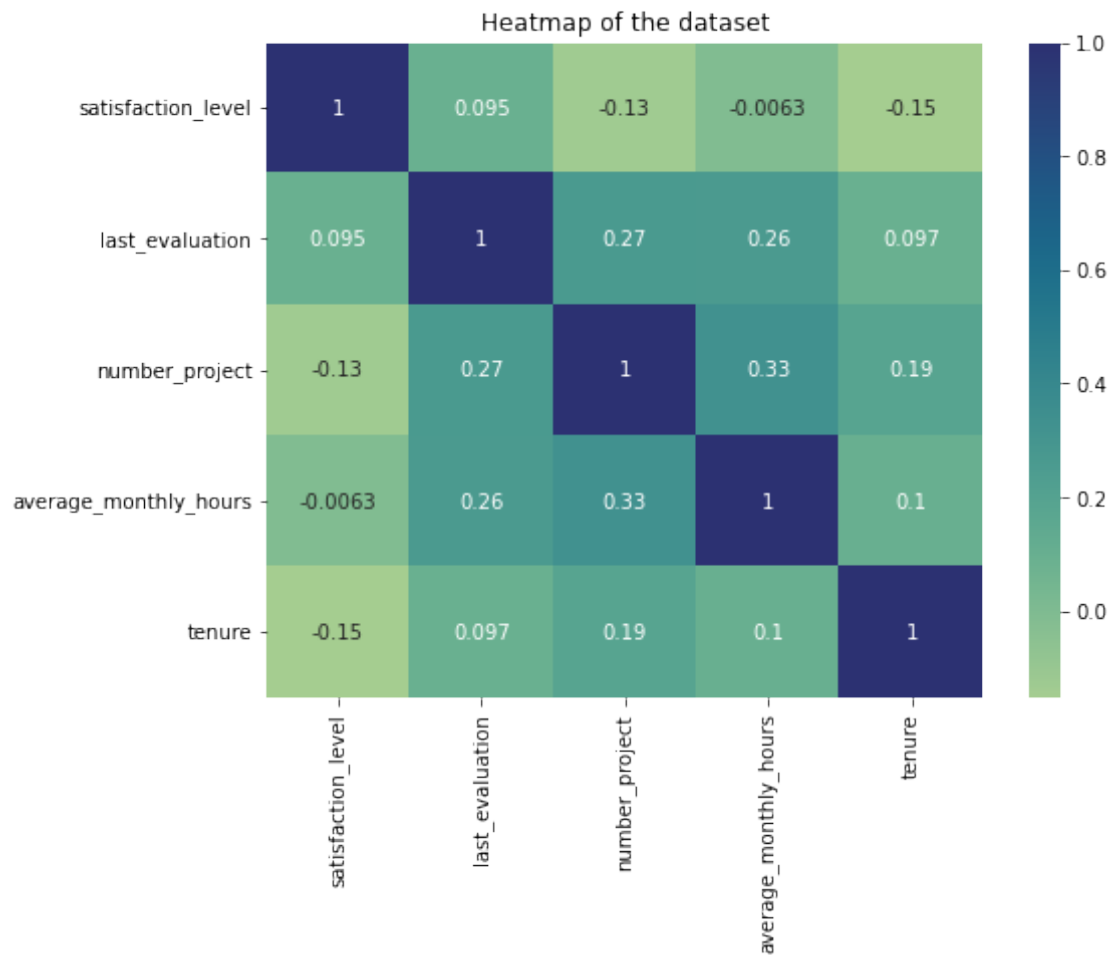
	tenure	work_accident	left	promotion_last_5years	salary	department_IT	\
0	3	0	1	0	0	0	
1	6	0	1	0	1	0	
2	4	0	1	0	1	0	
3	5	0	1	0	0	0	
4	3	0	1	0	0	0	

	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

```
[69]: plt.figure(figsize=(8, 6))
sns.heatmap(df_enc[['satisfaction_level', 'last_evaluation', 'number_project',
→ 'average_monthly_hours', 'tenure']]
            .corr(), annot=True, cmap="crest")
plt.title('Heatmap of the dataset')
plt.show()
```

```
[70]: pd.crosstab(df2['department'], df2['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()
```



```
[71]: df_logreg = df_enc[(df_enc['tenure'] >= lower_limit) & (df_enc['tenure'] <=
    ↳ upper_limit)]

df_logreg.head()
```

```
[71]:
```

	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

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

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

```
[72]: 0    1
      2    1
      3    1
      4    1
      5    1
      Name: left, dtype: int64
```

```
[73]: # Select the features you want to use in your model
X = df_logreg.drop('left', axis=1)

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

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

```
[74]: # 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)
```

```
[75]: # Construct a logistic regression model and fit it to the training dataset
log_clf = LogisticRegression(random_state=42, max_iter=500).fit(X_train,
↳y_train)
```

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

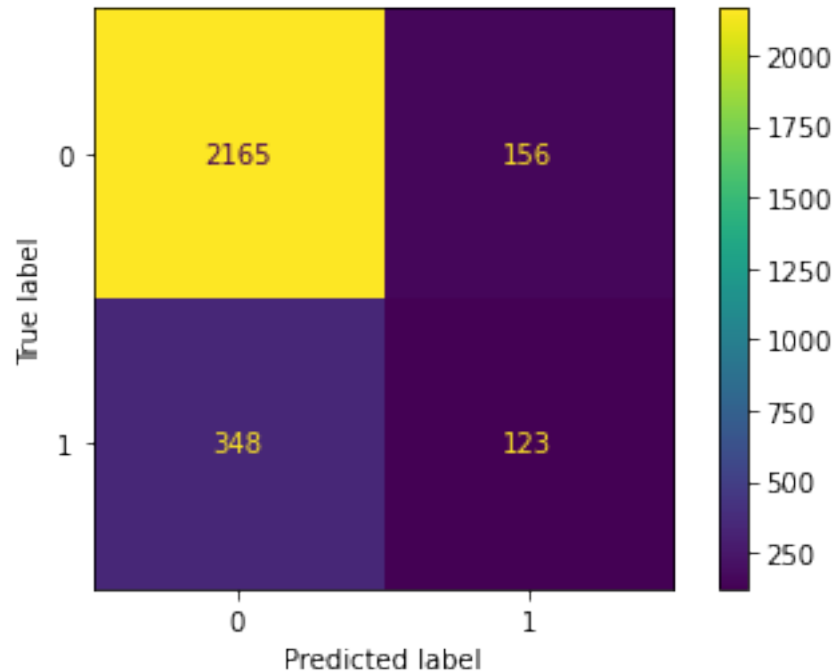
```
[77]: # Compute values for confusion matrix
log_cm = confusion_matrix(y_test, y_pred, labels=log_clf.classes_)

# 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()
```



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

```
[78]: 0    0.831468
      1    0.168532
      Name: left, dtype: float64
```

```
[ ]: # 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))
```

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

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

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

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

```
[ ]: # 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)
```

```
[ ]: # 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')
```

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

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

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

```
[ ]: 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

```

```

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

```

```

[ ]: # 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')
```

```
[ ]: %%time
rf1.fit(X_train, y_train) # --> Wall time: ~10min
```

```
[ ]: path = '/home/jovyan/work/'
```

```
[ ]: 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)
```

```
[ ]: 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
```

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

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

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



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

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

```
[ ]: 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
```

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

```
[ ]: # 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()
```

```
[ ]: # 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())
```

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

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

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

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

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

# Select the features
X = df2.drop('left', axis=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
tree2 = GridSearchCV(tree, cv_params, scoring=scoring, cv=4, refit='roc_auc')
```

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

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

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

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

[ ]: # 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')

[ ]: %%time
rf2.fit(X_train, y_train) # --> Wall time: 7min 5s

[ ]: # Write pickle
write_pickle(path, rf2, 'hr_rf2')

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

[ ]: # Check best params
rf2.best_params_

[ ]: # Check best AUC score on CV
rf2.best_score_

[ ]: # Get all CV scores
rf2_cv_results = make_results('random forest2 cv', rf2, 'auc')
print(tree2_cv_results)
```

```
print(rf2_cv_results)
```

```
[ ]: # Get predictions on test data
rf2_test_scores = get_scores('random forest2 test', rf2, X_test, y_test)
rf2_test_scores
```

```
[ ]: # 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='');
```

```
[ ]: # 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()
```

```
[ ]: #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
```

```
[ ]: 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()
```

```
[ ]: # Get feature importances
feat_impt = rf2.best_estimator_.feature_importances_
```

```

# Get indices of top 10 features
ind = np.argsort(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()

```

[]:

[]:

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

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.

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 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 eval-

uations 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.

For another project, you could try building a K-means model on this data and analyzing the clusters. This may yield valuable insight.

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.