

Capstone project: Providing data-driven suggestions for HR

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](https://docs.google.com/presentation/d/1Pps5GKxi1V31y2oRHRzU-xhJubkEYzCgElFnjIEY3Og/templatepreview#slide=id.g1c07b9724f3_0_983) (https://docs.google.com/presentation/d/1Pps5GKxi1V31y2oRHRzU-xhJubkEYzCgElFnjIEY3Og/templatepreview#slide=id.g1c07b9724f3_0_983) 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.

PACE stages

- [Plan](#)
- [Analyze](#)
- [Construct](#)
- [Execute](#)



Pace: Plan Stage

- Understand your data in the problem context
- Consider how your data will best address the business need
- Contextualize & understand the data and the problem



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.

Familiarize yourself with the HR dataset

In this [dataset \(\[https://www.kaggle.com/datasets/mfaisalqureshi/hr-analytics-and-job-prediction?select=HR_comma_sep.csv\]\(https://www.kaggle.com/datasets/mfaisalqureshi/hr-analytics-and-job-prediction?select=HR_comma_sep.csv\)\)](https://www.kaggle.com/datasets/mfaisalqureshi/hr-analytics-and-job-prediction?select=HR_comma_sep.csv), there are 14,999 rows, 10 columns, and these variables:

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?

Step 1. Imports

- Import packages
- Load dataset

Import packages

```
In [1]: # Import packages
# ## YOUR CODE HERE ##

# For data manipulation
import numpy as np
import pandas as pd

# For data visualization
%matplotlib inline
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
```

Load dataset

```
In [2]: # Load dataset into a dataframe
# ## YOUR CODE HERE ##
df0 = pd.read_csv("HR_capstone_dataset.csv")
# Display first few rows of the dataframe
## YOUR CODE HERE ##
df0.head()
```

Out[2]:

	satisfaction_level	last_evaluation	number_project	average_montly_hours	time_spend_company	Work_accident
0	0.38	0.53	2	157		3
1	0.80	0.86	5	262		6
2	0.11	0.88	7	272		4
3	0.72	0.87	5	223		5
4	0.37	0.52	2	159		3

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

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

Gather basic information about the data

```
In [3]: # Gather basic information about the data
## YOUR CODE HERE ##
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
```

Gather descriptive statistics about the data

In [4]:

```
1 # Gather descriptive statistics about the data
2 ### YOUR CODE HERE ####
3 df0.describe()
```

Out[4]:

	satisfaction_level	last_evaluation	number_project	average_montly_hours	time_spend_company	Work_accident	left
count	14999.000000	14999.000000	14999.000000	14999.000000	14999.000000	14999.000000	14999.000000
mean	0.612834	0.716102	3.803054	201.050337	3.498233	0.000000	0.000000
std	0.248631	0.171169	1.232592	49.943099	1.460136	0.000000	0.000000
min	0.090000	0.360000	2.000000	96.000000	2.000000	0.000000	0.000000
25%	0.440000	0.560000	3.000000	156.000000	3.000000	0.000000	0.000000
50%	0.640000	0.720000	4.000000	200.000000	3.000000	0.000000	0.000000
75%	0.820000	0.870000	5.000000	245.000000	4.000000	0.000000	0.000000
max	1.000000	1.000000	7.000000	310.000000	10.000000	0.000000	0.000000



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.

In [5]:

```
1 # Display all column names
2 ### YOUR CODE HERE ####
3 df0.columns
```

Out[5]:

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

In [6]:

```
1 # Rename columns as needed
2 ### YOUR CODE HERE ####
3 df0 = df0.rename(columns={'Work_accident': 'work_accident',
                           'average_montly_hours': 'average_monthly_hours',
                           'time_spend_company': 'tenure',
                           'Department': 'department'})
4
5
6
7
8 # Display all column names after the update
9 ### YOUR CODE HERE ####
10 df0.columns
```

Out[6]:

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

Check missing values

Check for any missing values in the data.

```
In [7]: 1 # Check for missing values
2 ### YOUR CODE HERE ####
3 df0.isna().sum()
```

```
Out[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
```

There are no missing values in the data.

Check duplicates

Check for any duplicate entries in the data.

```
In [8]: 1 # Check for duplicates
2 ### YOUR CODE HERE ####
3 df0.duplicated().sum()
```

```
Out[8]: 3008
```

3,008 rows contain duplicates. That is 20% of the data.

```
In [9]: 1 # Inspect some rows containing duplicates as needed
2 ### YOUR CODE HERE ####
3 df0[df0.duplicated()].head()
```

```
Out[9]:
```

	satisfaction_level	last_evaluation	number_project	average_monthly_hours	tenure	work_accident	left
396	0.46	0.57	2		139	3	0
866	0.41	0.46	2		128	3	0
1317	0.37	0.51	2		127	3	0
1368	0.41	0.52	2		132	3	0
1461	0.42	0.53	2		142	3	0

The above output shows the first five occurrences of rows that are duplicated farther down in the dataframe. How likely is it that these are legitimate entries? In other words, how plausible is it that two employees self-reported the exact same response for every column?

You could perform a likelihood analysis by essentially applying Bayes' 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. You can proceed by dropping them.

In [10]:

```
1 # Drop duplicates and save resulting dataframe in a new variable as needed
2 ### YOUR CODE HERE ###
3 df1 = df0.drop_duplicates(keep='first')
4
5 # Display first few rows of new dataframe as needed
6 df1.head()
```

Out[10]:

	satisfaction_level	last_evaluation	number_project	average_monthly_hours	tenure	work_accident	left
0	0.38	0.53	2	157	3	0	1
1	0.80	0.86	5	262	6	0	1
2	0.11	0.88	7	272	4	0	1
3	0.72	0.87	5	223	5	0	1
4	0.37	0.52	2	159	3	0	1

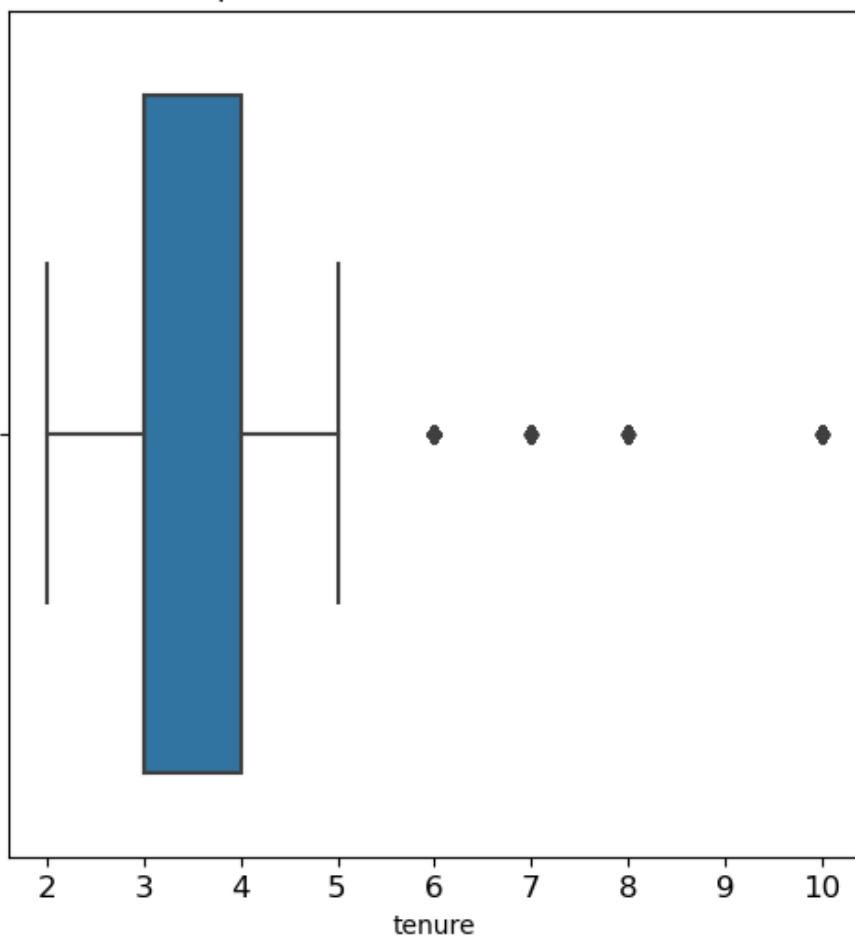
Check outliers

Check for outliers in the data.

In [11]:

```
1 # Create a boxplot to visualize distribution of `tenure` and detect any outliers
2 plt.figure(figsize=(6,6))
3 plt.title('Boxplot to detect outliers for tenure', fontsize=12)
4 plt.xticks(fontsize=12)
5 plt.yticks(fontsize=12)
6 sns.boxplot(x=df1['tenure'])
7 plt.show()
```

Boxplot to detect outliers for tenure



The boxplot above 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.

```
In [12]: # Determine the number of rows containing outliers
# YOUR CODE HERE ##

# 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 these outliers based on the type of model you decide to use.



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?

Step 2. Data Exploration (Continue EDA)

Data visualizations

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

You could start by creating a stacked boxplot showing `average_monthly_hours` distributions for `number_project`, comparing the distributions of employees who stayed versus those who left.

Box plots are very useful in visualizing distributions within data, but they can be deceiving without the context of how big the sample sizes that they represent are. So, you could also plot a stacked histogram to visualize the distribution of `number_project` for those who stayed and those who left.

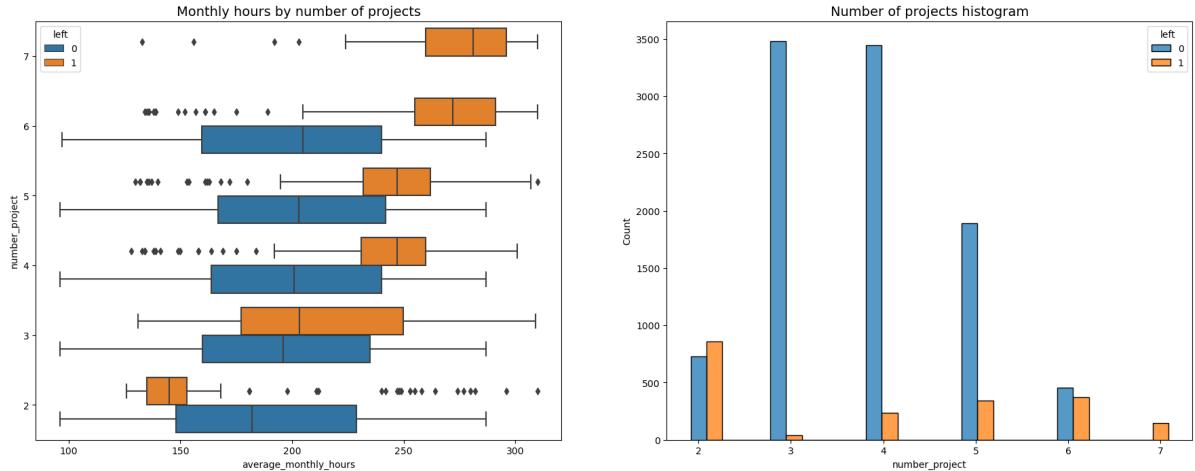
```
In [13]: # Create a plot as needed
# YOUR CODE HERE ##

# Set figure and axes
fig, ax = plt.subplots(1, 2, figsize = (22,8))

# Create boxplot showing `average_monthly_hours` distributions for `number_project`
sns.boxplot(data=df1, x='average_monthly_hours', y='number_project', hue='left', order=[2, 3, 4, 5, 6, 7])
ax[0].invert_yaxis()
ax[0].set_title('Monthly hours by number of projects', fontsize='14')

# Create histogram showing distribution of `number_project`, comparing employees who stayed and 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[1].set_title('Number of projects histogram', fontsize='14')

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



It might be natural that people who work on more projects would also work longer hours. This appears to be the case here, with the mean hours of each group (stayed and left) increasing with number of projects worked. However, a few things stand out from this plot.

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

As the next step, you could confirm that all employees with seven projects left.

```
In [14]: ┶ 1 # Get value counts of stayed/Left for employees with 7 projects
          2 df1[df1['number_project']==7]['left'].value_counts()
```

```
Out[14]: left
          1    145
          Name: count, dtype: int64
```

This confirms that all employees with 7 projects did leave.

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

In [15]:

```

1 # Create a plot as needed
2 ### YOUR CODE HERE ####
3
4 # Create scatterplot of `average_monthly_hours` versus `satisfaction_level`, compo-
5 plt.figure(figsize=(16, 9))
6 sns.scatterplot(data=df1, x='average_monthly_hours', y='satisfaction_level', hue=
7 plt.axvline(x=166.67, color='#ff6361', label='166.67 hrs./mo.', ls='--')
8 plt.legend(labels=['166.67 hrs./mo.', 'left', 'stayed'])
9 plt.title('Monthly hours by last evaluation score', fontsize=14);

```



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

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

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

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

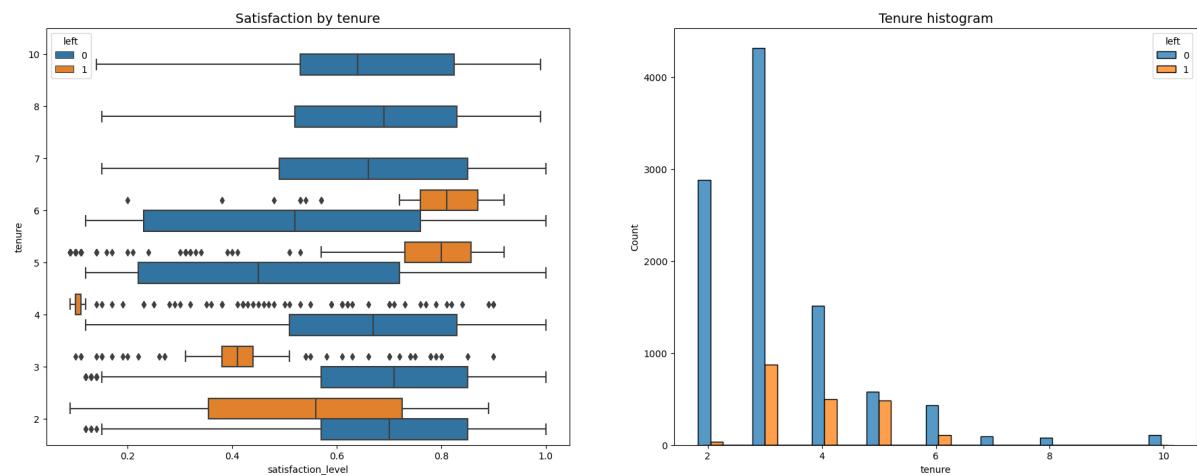
For the next visualization, it might be interesting to visualize satisfaction levels by tenure.

In [16]:

```

1 # Create a plot as needed
2 ### YOUR CODE HERE ####
3
4 # Set figure and axes
5 fig, ax = plt.subplots(1, 2, figsize = (22,8))
6
7 # Create boxplot showing distributions of `satisfaction_Level` by tenure, comparing employees who stayed vs left
8 sns.boxplot(data=df1, x='satisfaction_level', y='tenure', hue='left', orient="h",
9 ax[0].invert_yaxis()
10 ax[0].set_title('Satisfaction by tenure', fontsize='14')
11
12 # Create histogram showing distribution of `tenure`, comparing employees who stayed vs left
13 tenure_stay = df1[df1['left']==0]['tenure']
14 tenure_left = df1[df1['left']==1]['tenure']
15 sns.histplot(data=df1, x='tenure', hue='left', multiple='dodge', shrink=5, ax=ax[1])
16 ax[1].set_title('Tenure histogram', fontsize='14')
17
18 plt.show();

```



There are many observations you could make from this plot.

- Employees who left fall into two general categories: dissatisfied 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 levels 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.

As the next step in analyzing the data, you could calculate the mean and median satisfaction scores of employees who left and those who didn't.

In [17]:

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

Out[17]:

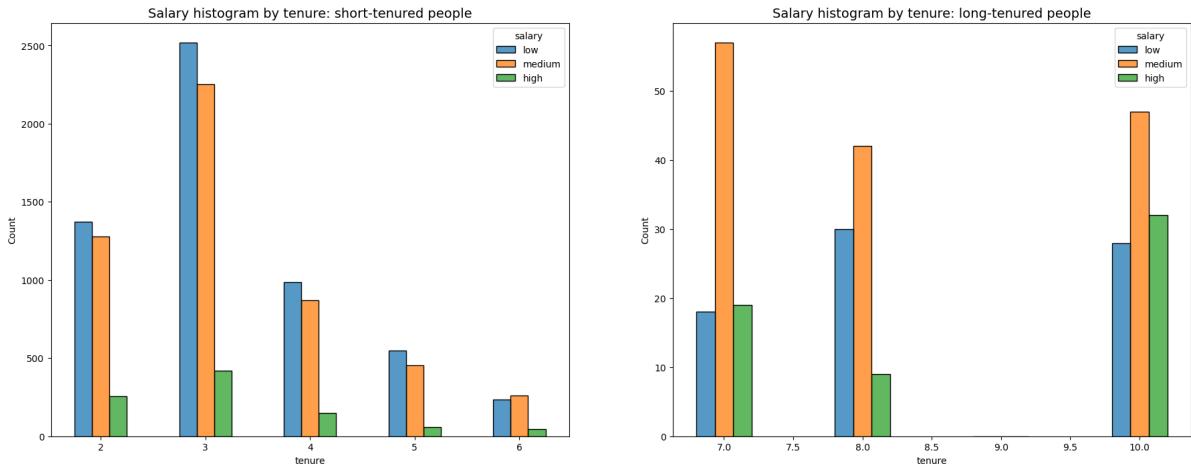
	mean	median
left		
0	0.667365	0.69
1	0.440271	0.41

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

Next, you could examine salary levels for different tenures.

In [18]:

```
1 # Create a plot as needed
2 ### YOUR CODE HERE ####
3
4 # Set figure and axes
5 fig, ax = plt.subplots(1, 2, figsize = (22,8))
6
7 # Define short-tenured employees
8 tenure_short = df1[df1['tenure'] < 7]
9
10 # Define long-tenured employees
11 tenure_long = df1[df1['tenure'] > 6]
12
13 # Plot short-tenured histogram
14 sns.histplot(data=tenure_short, x='tenure', hue='salary', discrete=1,
15               hue_order=['low', 'medium', 'high'], multiple='dodge', shrink=.5, ax=
16 ax[0].set_title('Salary histogram by tenure: short-tenured people', fontsize='14')
17
18 # Plot Long-tenured histogram
19 sns.histplot(data=tenure_long, x='tenure', hue='salary', discrete=1,
20               hue_order=['low', 'medium', 'high'], multiple='dodge', shrink=.4, ax=
21 ax[1].set_title('Salary histogram by tenure: long-tenured people', fontsize='14')
```



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

Next, you could explore whether there's a correlation between working long hours and receiving high evaluation scores. You could create a scatterplot of `average_monthly_hours` versus `last_evaluation`.

```
In [19]: # Create a plot as needed
### YOUR CODE HERE ###

# Create scatterplot of `average_monthly_hours` versus `last_evaluation`
plt.figure(figsize=(16, 9))
sns.scatterplot(data=df1, x='average_monthly_hours', y='last_evaluation', hue='left')
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='14');
```



The following observations can be made from the scatterplot above:

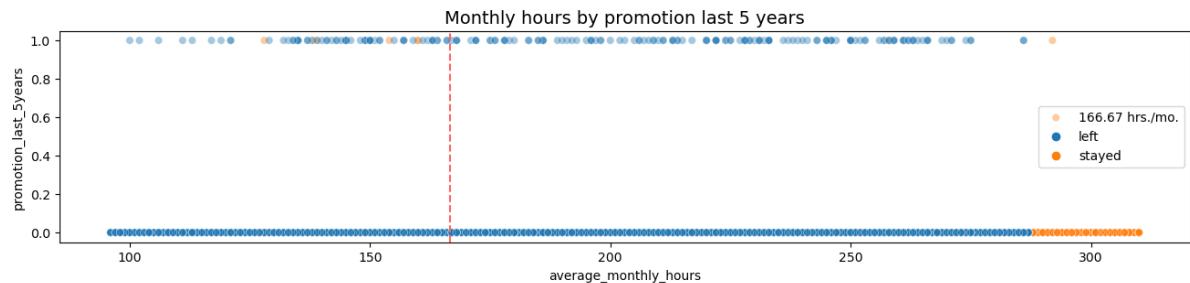
- The scatterplot indicates two groups of employees who left: overworked employees who performed very well and employees who worked slightly under the nominal monthly average of 166.67 hours with lower evaluation scores.
- There seems to be a correlation between hours worked and evaluation score.
- There 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.

Next, you could examine whether employees who worked very long hours were promoted in the last five years.

In [20]:

```

1 # Create a plot as needed
2 ### YOUR CODE HERE ####
3
4 # Create plot to examine relationship between `average_monthly_hours` and `promotion_
5 plt.figure(figsize=(16, 3))
6 sns.scatterplot(data=df1, x='average_monthly_hours', y='promotion_last_5years', h
7 plt.axvline(x=166.67, color='#ff6361', ls='--')
8 plt.legend(labels=['166.67 hrs./mo.', 'left', 'stayed'])
9 plt.title('Monthly hours by promotion last 5 years', fontsize='14');
```



The plot above 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

Next, you could inspect how the employees who left are distributed across departments.

In [21]:

```

1 # Display counts for each department
2 df1["department"].value_counts()
```

Out[21]:

sales	3239
technical	2244
support	1821
IT	976
RandD	694
product_mng	686
marketing	673
accounting	621
hr	601
management	436

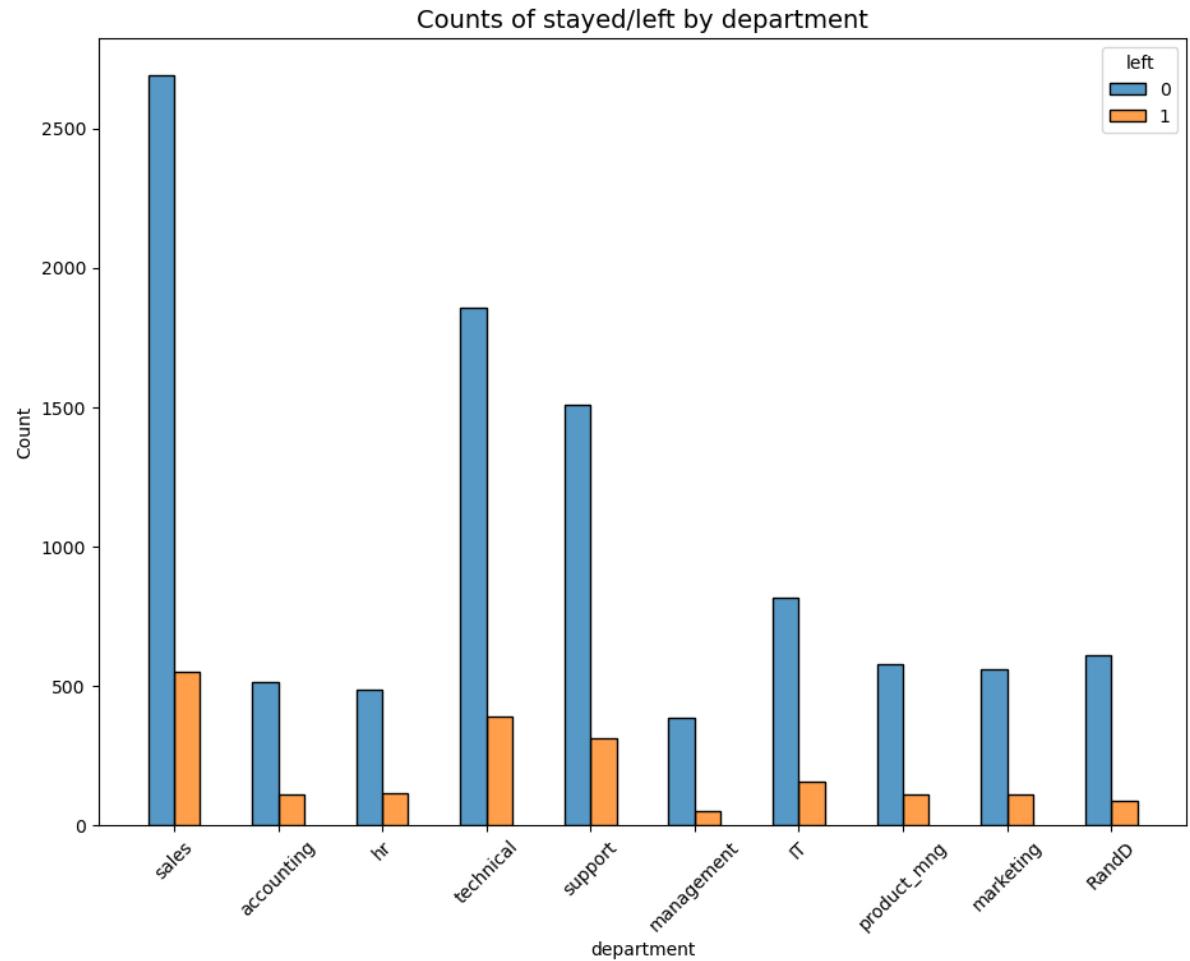
Name: count, dtype: int64

In [22]:

```

1 # Create a plot as needed
2 ### YOUR CODE HERE ####
3
4 # Create stacked histogram to compare department distribution of employees who Le
5 plt.figure(figsize=(11,8))
6 sns.histplot(data=df1, x='department', hue='left', discrete=1,
7               hue_order=[0, 1], multiple='dodge', shrink=.5)
8 plt.xticks(rotation='45')
9 plt.title('Counts of stayed/left by department', fontsize=14);
10

```



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

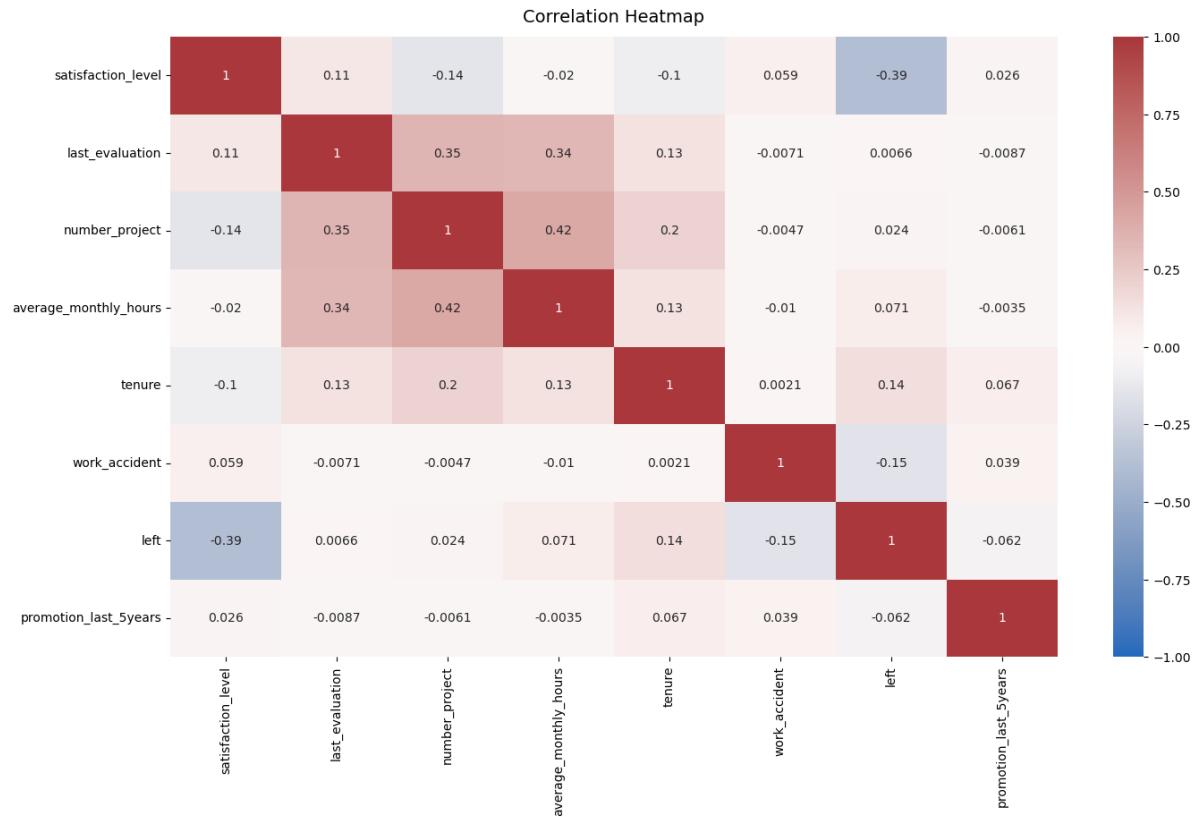
Lastly, you could check for strong correlations between variables in the data.

In [24]:

```

1 # Create a plot as needed
2 ### YOUR CODE HERE ####
3
4 # Plot a correlation heatmap
5 plt.figure(figsize=(16, 9))
6 heatmap = sns.heatmap(df0.corr(numeric_only=True), vmin=-1, vmax=1, annot=True, cbar_kws={"shrink": 0.5})
7 heatmap.set_title('Correlation Heatmap', fontdict={'fontsize':14}, pad=12);

```



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.

Insights

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



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?

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

Identify the type of prediction task.

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

Identify the types of models most appropriate for this task.

Since the variable you want to predict (whether an employee leaves the company) is categorical, you could either build a Logistic Regression model, or a Tree-based Machine Learning model.

Modeling Approach A: Logistic Regression Model

This approach covers implementation of Logistic Regression.

Logistic regression

Note that binomial logistic regression suits the task because it involves binary classification.

Start by one-hot encoding the categorical variables as needed.

```
In [25]: # One-hot encode the categorical variables as needed and save resulting dataframe
          df_enc = pd.get_dummies(df1, prefix=['salary', 'dept'], columns = ['salary', 'dept'])
          # Display the new dataframe
          df_enc.head()
```

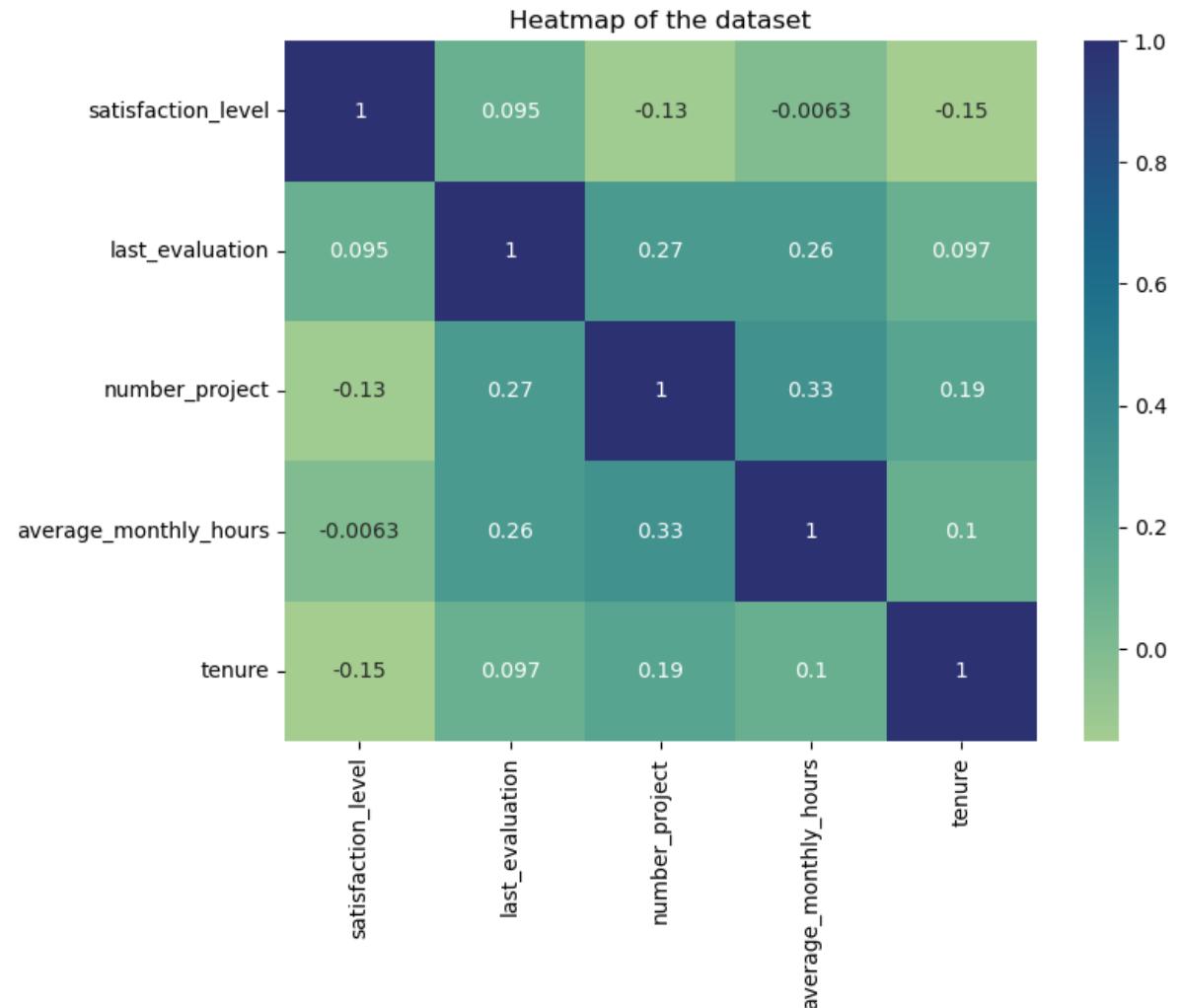
Out[25]:

	satisfaction_level	last_evaluation	number_project	average_monthly_hours	tenure	work_accident	left
0	0.38	0.53	2	157	3	0	1
1	0.80	0.86	5	262	6	0	1
2	0.11	0.88	7	272	4	0	1
3	0.72	0.87	5	223	5	0	1
4	0.37	0.52	2	159	3	0	1

Create a heatmap to visualize how correlated variables are. Consider which variables you're interested in examining correlations between.

In [26]:

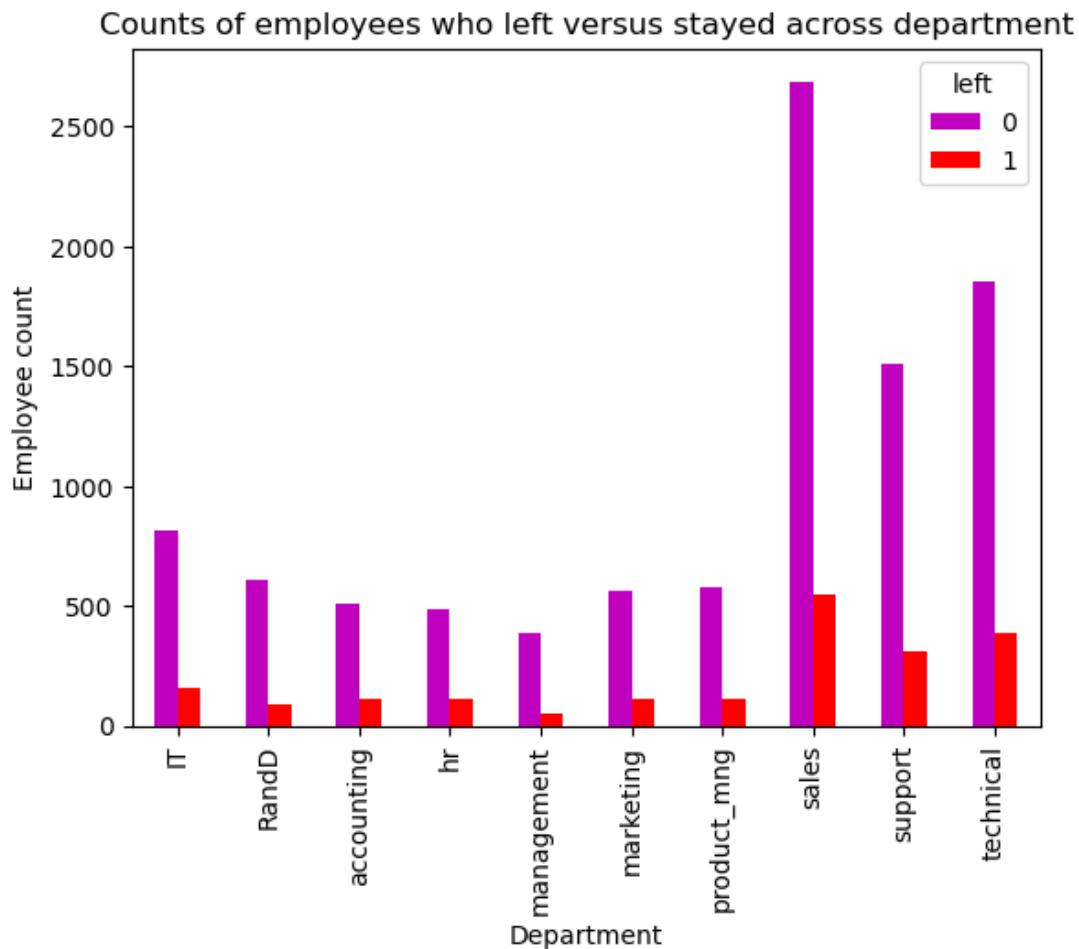
```
1 # Create a heatmap to visualize how correlated variables are
2 plt.figure(figsize=(8, 6))
3 sns.heatmap(df_enc[['satisfaction_level', 'last_evaluation', 'number_project', 'a
4 plt.title('Heatmap of the dataset')
5 plt.show()
```



Create a stacked bar plot to visualize number of employees across department, comparing those who left with those who didn't.

In [27]:

```
1 # Create a stacked bar plot to visualize number of employees across department, 
2 # In the Legend, 0 (purple color) represents employees who did not Leave, 1 (red color) represents employees who left
3 pd.crosstab(df1["department"], df1["left"]).plot(kind ='bar',color='mr')
4 plt.title('Counts of employees who left versus stayed across department')
5 plt.ylabel('Employee count')
6 plt.xlabel('Department')
7 plt.show()
```



Since logistic regression is quite sensitive to outliers, it would be a good idea at this stage to remove the outliers in the tenure column that were identified earlier.

In [28]:

```

1 # Select rows without outliers in `tenure` and save resulting dataframe in a new variable
2 df_logreg = df_enc[(df_enc['tenure'] >= lower_limit) & (df_enc['tenure'] <= upper_limit)]
3
4 # Display first few rows of new dataframe
5 df_logreg.head()

```

Out[28]:

	satisfaction_level	last_evaluation	number_project	average_monthly_hours	tenure	work_accident	left
0	0.38	0.53	2		157	3	0 1
2	0.11	0.88	7		272	4	0 1
3	0.72	0.87	5		223	5	0 1
4	0.37	0.52	2		159	3	0 1
5	0.41	0.50	2		153	3	0 1

Isolate the outcome variable, which is the variable you want your model to predict.

In [29]:

```

1 # Isolate the outcome variable
2 y = df_logreg['left']
3
4 # Display first few rows of the outcome variable
5 y.head()

```

Out[29]:

```

0    1
1    1
2    1
3    1
4    1
5    1
Name: left, dtype: int64

```

Select the features you want to use in your model. Consider which variables will help you predict the outcome variable, `left`.

In [30]:

```

1 # Select the features you want to use in your model
2 X = df_logreg[['satisfaction_level', 'last_evaluation', 'number_project', 'average_monthly_hours', 'tenure', 'work_accident']]
3
4 # Display the first few rows of the selected features
5 X.head()

```

Out[30]:

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

Split the data into training set and testing set.

```
In [31]: 1 # Split the data into training set and testing set
2 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_s...
```

Construct a logistic regression model and fit it to the training dataset.

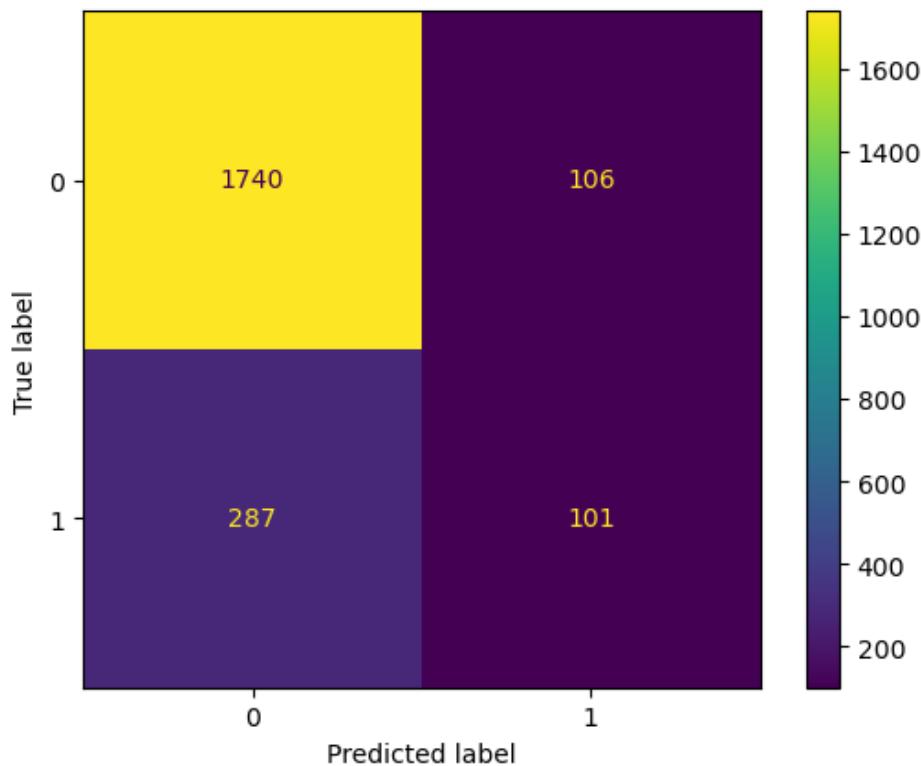
```
In [32]: 1 # Construct a Logistic regression model and fit it to the training dataset
2 log_clf = LogisticRegression(random_state=42, max_iter=500).fit(X_train, y_train)
```

Test the logistic regression model: use the model to make predictions on the test set.

```
In [33]: 1 # Use the logistic regression model to get predictions on the test set
2 y_pred = log_clf.predict(X_test)
```

Create a confusion matrix to visualize the results of the logistic regression model.

```
In [34]: 1 # Compute values for confusion matrix
2 log_cm = confusion_matrix(y_test, y_pred, labels=log_clf.classes_)
3
4 # Create display of confusion matrix
5 log_disp = ConfusionMatrixDisplay(confusion_matrix=log_cm, display_labels=log_clf
6
7 # Plot confusion matrix
8 log_disp.plot()
9
10 # Display plot
11 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.

```
In [35]: 1 df_logreg['left'].value_counts(normalize=True)
```

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

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

```
In [36]: 1 # Create classification report for Logistic regression model
2 target_names = ['Predicted would not leave', 'Predicted would leave']
3 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	1846
Predicted would leave	0.49	0.26	0.34	388
accuracy			0.82	2234
macro avg	0.67	0.60	0.62	2234
weighted avg	0.79	0.82	0.80	2234

The classification report above shows that the logistic regression model achieved a precision of 80%, recall of 83%, f1-score of 80% (all weighted averages), and accuracy of 83%.

Modeling Approach B: Tree-based Model

This approach covers implementation of Decision Tree and Random Forest.

Encode the categorical variables.

```
In [37]: 1 # Encode categorical variables
2 df2 = pd.get_dummies(df1)
```

Isolate the outcome variable.

```
In [38]: 1 # Isolate the outcome variable
2 y = df2['left']
3
4 # Display the first few rows of `y`
5 y.head()
```

```
Out[38]: 0    1
1    1
2    1
3    1
4    1
Name: left, dtype: int64
```

Select the features.

```
In [39]: 1 # Select the features
2 X = df2.drop('left', axis=1)
3
4 # Display the first few rows of `X`
5 X.head()
```

```
Out[39]:
   satisfaction_level  last_evaluation  number_project  average_monthly_hours  tenure  work_accident  prorr
0                  0.38          0.53             2                 157       3            0
1                  0.80          0.86             5                 262       6            0
2                  0.11          0.88             7                 272       4            0
3                  0.72          0.87             5                 223       5            0
4                  0.37          0.52             2                 159       3            0
```

Split the data into training, validating, and testing sets.

```
In [40]: 1 # Create test data
2 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, stratify=y)
3
4 # Create train & validate data
5 X_tr, X_val, y_tr, y_val = train_test_split(X_train, y_train, test_size=0.25, strati
```

Decision tree - Round 1

Construct a decision tree model and set up cross-validated grid-search to exhaustively search for the best model parameters.

```
In [41]: 1 # Instantiate model
          2 tree = DecisionTreeClassifier(random_state=0)
          3
          4 # Assign a dictionary of hyperparameters to search over
          5 cv_params = {'max_depth':[4, 6, 8, None],
          6             'min_samples_leaf': [2, 5, 1],
          7             'min_samples_split': [2, 4, 6]
          8             }
          9
         10 # Assign a dictionary of scoring metrics to capture
         11 scoring = {'accuracy', 'precision', 'recall', 'f1', 'roc_auc'}
         12
         13 # Instantiate GridSearch
         14 tree1 = GridSearchCV(tree, cv_params, scoring=scoring, cv=4, refit='roc_auc')
```

Fit the decision tree model to the training data.

```
In [42]: 1 %%time
          2 tree1.fit(X_tr, y_tr)
```

Wall time: 10.1 s

```
Out[42]: GridSearchCV(cv=4, estimator=DecisionTreeClassifier(random_state=0),
                      param_grid={'max_depth': [4, 6, 8, None],
                                  'min_samples_leaf': [2, 5, 1],
                                  'min_samples_split': [2, 4, 6]},
                      refit='roc_auc',
                      scoring={'f1', 'roc_auc', 'accuracy', 'precision', 'recall'})
```

Identify the optimal values for the decision tree parameters.

```
In [43]: 1 # Check best parameters
          2 tree1.best_params_
```

```
Out[43]: {'max_depth': 4, 'min_samples_leaf': 1, 'min_samples_split': 2}
```

Identify the best AUC score achieved by the decision tree model on the training set.

```
In [44]: 1 # Check best AUC score on CV
          2 tree1.best_score_
```

```
Out[44]: 0.9703634179699269
```

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

Next, you can write a function that will help you extract all the scores from the grid search.

```
In [57]: 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 = {'Model': model_name,
              'AUC': auc,
              'Precision': precision,
              'Recall': recall,
              'F1': f1,
              'Accuracy': accuracy,
              }
    table = pd.DataFrame(table, index=range(1))

    return table
```

Use the function just defined to get all the scores from grid search.

```
In [58]: # Get all CV scores
tree1_cv_results = make_results('decision tree cv', tree1, 'auc')
tree1_cv_results
```

Out[58]:

	Model	AUC	Precision	Recall	F1	Accuracy
0	decision tree cv	0.970363	0.922167	0.921337	0.921719	0.974007

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

Recall that decision trees can be vulnerable to overfitting and random forests avoid overfitting by

Random forest - Round 1

Construct a random forest model and set up cross-validated grid-search to exhaustively search for the best model parameters.

```
In [59]: # Instantiate model
1 rf = RandomForestClassifier(random_state=0)
2
3
4 # Assign a dictionary of hyperparameters to search over
5 cv_params = {'max_depth': [3,5, None],
6               'max_features': [1.0],
7               'max_samples': [0.7, 1.0],
8               'min_samples_leaf': [1,2,3],
9               'min_samples_split': [2,3,4],
10              'n_estimators': [300, 500],
11            }
12
13 # Assign a dictionary of scoring metrics to capture
14 scoring = {'accuracy', 'precision', 'recall', 'f1', 'roc_auc'}
15
16 # Instantiate GridSearch
17 rf1 = GridSearchCV(rf, cv_params, scoring=scoring, cv=4, refit='roc_auc')
```

Fit the random forest model to the training data.

```
In [60]: %time
1 rf1.fit(X_tr, y_tr) # --> Wall time: ~22min
```

Wall time: 26min 1s

```
Out[60]: GridSearchCV(cv=4, estimator=RandomForestClassifier(random_state=0),
1   param_grid={'max_depth': [3, 5, None], 'max_features': [1.0],
2               'max_samples': [0.7, 1.0],
3               'min_samples_leaf': [1, 2, 3],
4               'min_samples_split': [2, 3, 4],
5               'n_estimators': [300, 500]},
6   refit='roc_auc',
7   scoring={'f1', 'roc_auc', 'accuracy', 'precision', 'recall'})
```

Specify path to where you want to save your model.

```
In [62]: # Define a path to the folder where you want to save the model
1 path = 'C:\\Users\\USER\\Google Advanced Data Analytics Certificate\\Capstone Course'
```

Define functions to pickle the model and read in the model.

```
In [63]: 1 def write_pickle(path, model_object, save_as:str):
2     ...
3     In:
4         path:           path of folder where you want to save the pickle
5         model_object: a model you want to pickle
6         save_as:       filename for how you want to save the model
7
8     Out: A call to pickle the model in the folder indicated
9     ...
10
11    with open(path + save_as + '.pickle', 'wb') as to_write:
12        pickle.dump(model_object, to_write)
```

```
In [64]: 1 def read_pickle(path, saved_model_name:str):
2     ...
3     In:
4         path:           path to folder where you want to read from
5         saved_model_name: filename of pickled model you want to read in
6
7     Out:
8         model: the pickled model
9     ...
10
11    with open(path + saved_model_name + '.pickle', 'rb') as to_read:
12        model = pickle.load(to_read)
13
14    return model
```

Use the functions defined above to save the model in a pickle file and then read it in.

```
In [65]: 1 # Write pickle
2 write_pickle(path, rf1, 'hr_rf1')
```

```
In [66]: 1 # Read pickle
2 rf1 = read_pickle(path, 'hr_rf1')
```

Identify the best AUC score achieved by the random forest model on the training set.

```
In [67]: 1 # Check best AUC score on CV
2 rf1.best_score_
```

Out[67]: 0.9795790967836983

Identify the optimal values for the parameters of the random forest model.

```
In [68]: 1 # Check best params
2 rf1.best_params_
```

Out[68]: {'max_depth': 5,
 'max_features': 1.0,
 'max_samples': 0.7,
 'min_samples_leaf': 1,
 'min_samples_split': 3,
 'n_estimators': 500}

Collect the evaluation scores on the training set for the decision tree and random forest models.

```
In [69]: # Get all CV scores
1 rf1_cv_results = make_results('random forest cv', rf1, 'auc')
2 print(tree1_cv_results)
3 print(rf1_cv_results)
```

	Model	AUC	Precision	Recall	F1	Accuracy
0	decision tree cv	0.970363	0.922167	0.921337	0.921719	0.974007
	Model	AUC	Precision	Recall	F1	Accuracy
0	random forest cv	0.979579	0.943576	0.923021	0.933145	0.978037

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 0.008 lower, which is a negligible amount). This indicates that the random forest model mostly outperforms the decision tree model.

Next, you can evaluate these models on the validation set.

Define a function that gets all the scores from a model's predictions.

```
In [70]: def get_scores(model_name:str, model, X_test_data, y_test_data):
1     """
2         Generate a table of test scores.
3
4         In:
5             model_name (string): How you want your model to be named in the output table
6             model: A fit GridSearchCV object
7             X_test_data: numpy array of X_test data
8             y_test_data: numpy array of y_test data
9
10            Out: pandas df of precision, recall, f1, accuracy, and AUC scores for your model
11            """
12
13
14     preds = model.best_estimator_.predict(X_test_data)
15
16     auc = round(roc_auc_score(y_test_data, preds), 3)
17     accuracy = round(accuracy_score(y_test_data, preds), 3)
18     precision = round(precision_score(y_test_data, preds), 3)
19     recall = round(recall_score(y_test_data, preds), 3)
20     f1 = round(f1_score(y_test_data, preds), 3)
21
22     table = pd.DataFrame({'model': [model_name],
23                           'AUC': [auc],
24                           'precision': [precision],
25                           'recall': [recall],
26                           'f1': [f1],
27                           'accuracy': [accuracy]
28                           })
29
30
31     return table
```

Apply the function defined above to get scores for the decision tree model and the random forest model.

In [71]:

```

1 # Get the results on validation set for both models
2 tree1_val_results = get_scores('decision tree1 val', tree1, X_val, y_val)
3 rf1_val_results = get_scores('random forest1 val', rf1, X_val, y_val)
4
5 # Concatenate validation scores into table
6 all_val_results1 = [tree1_val_results, rf1_val_results]
7 all_val_results1 = pd.concat(all_val_results1).sort_values(by='AUC', ascending=False)
8 all_val_results1

```

Out[71]:

	model	AUC	precision	recall	f1	accuracy
0	random forest1 val	0.954	0.955	0.917	0.936	0.979
0	decision tree1 val	0.952	0.924	0.920	0.922	0.974

On the validation set, the random forest model outperforms the decision tree model across most metrics.

Now use the best performing model to predict on the test set.

In [72]:

```

1 # Get predictions on test data
2 rf1_test_scores = get_scores('random forest1 test', rf1, X_test, y_test)
3 rf1_test_scores

```

Out[72]:

	model	AUC	precision	recall	f1	accuracy
0	random forest1 test	0.955	0.961	0.917	0.938	0.98

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, you can be more confident that your model's performance on this data is representative of how it will perform on new, unseen data.

Feature Engineering

You might be skeptical of the high evaluation scores. There is a chance that there is some data leakage occurring. Data leakage is when you use data to train your 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 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.

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

In [73]:

```

1 # Drop `satisfaction_Level` and save resulting dataframe in new variable
2 df3 = df1.drop('satisfaction_level', axis=1)
3
4 # Display first few rows of new dataframe
5 df3.head()

```

Out[73]:

	last_evaluation	number_project	average_monthly_hours	tenure	work_accident	left	promotion_last_5y
0	0.53	2		157	3	0	1
1	0.86	5		262	6	0	1
2	0.88	7		272	4	0	1
3	0.87	5		223	5	0	1
4	0.52	2		159	3	0	1

In [74]:

```

1 # Create `overworked` column. For now, it's identical to average monthly hours.
2 df3['overworked'] = df3['average_monthly_hours']
3
4 # Inspect max and min average monthly hours values
5 print('Max hours:', df3['overworked'].max())
6 print('Min hours:', df3['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.

You 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

In [75]:

```

1 # Define `overworked` as working > 175 hrs/week
2 df3['overworked'] = (df3['overworked'] > 175).astype(int)
3
4 # Display first few rows of new column
5 df3['overworked'].head()

```

Out[75]:

0	0
1	1
2	1
3	1
4	0

Name: overworked, dtype: int32

Drop the `average_monthly_hours` column.

```
In [76]: 1 # Drop the `average_monthly_hours` column
2 df3 = df3.drop('average_monthly_hours', axis=1)
3
4 # Display first few rows of resulting dataframe
5 df3.head()
```

Out[76]:

	last_evaluation	number_project	tenure	work_accident	left	promotion_last_5years	department	salary	
0	0.53	2	3	0	1		0	sales	low
1	0.86	5	6	0	1		0	sales	medium
2	0.88	7	4	0	1		0	sales	medium
3	0.87	5	5	0	1		0	sales	low
4	0.52	2	3	0	1		0	sales	low

Start by one-hot encoding the categorical variables as needed.

```
In [77]: 1 # One-hot encode the categorical variables as needed and save resulting dataframe
2 df4 = pd.get_dummies(df3)
3
4 # Display the new dataframe
5 df4.head()
```

Out[77]:

	last_evaluation	number_project	tenure	work_accident	left	promotion_last_5years	overworked	depart
0	0.53	2	3	0	1		0	0
1	0.86	5	6	0	1		0	1
2	0.88	7	4	0	1		0	1
3	0.87	5	5	0	1		0	1
4	0.52	2	3	0	1		0	0

Isolate the outcome variable.

```
In [78]: 1 # Isolate the outcome variable
2 y = df4['left']
3
4 # Display the first few rows of `y`
5 y.head()
```

```
Out[78]: 0    1
1    1
2    1
3    1
4    1
Name: left, dtype: int64
```

Select the features.

In [79]:

```

1 # Select the features
2 X = df4.drop('left', axis=1)
3
4 # Display the first few rows of `X`
5 X.head()

```

Out[79]:

	last_evaluation	number_project	tenure	work_accident	promotion_last_5years	overworked	department
0	0.53	2	3	0		0	0 Fa
1	0.86	5	6	0		0	1 Fa
2	0.88	7	4	0		0	1 Fa
3	0.87	5	5	0		0	1 Fa
4	0.52	2	3	0		0	0 Fa

Split the data into training, validating, and testing sets.

In [80]:

```

1 # Create test data
2 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, stratify=y)
3
4 # Create train & validate data
5 X_tr, X_val, y_tr, y_val = train_test_split(X_train, y_train, test_size=0.25, stratify=y)

```

Decision tree - Round 2

In [81]:

```

1 # Instantiate model
2 tree = DecisionTreeClassifier(random_state=0)
3
4 # Assign a dictionary of hyperparameters to search over
5 cv_params = {'max_depth': [4, 6, 8, None],
6             'min_samples_leaf': [2, 5, 1],
7             'min_samples_split': [2, 4, 6]
8           }
9
10 # Assign a dictionary of scoring metrics to capture
11 scoring = {'accuracy', 'precision', 'recall', 'f1', 'roc_auc'}
12
13 # Instantiate GridSearch
14 tree2 = GridSearchCV(tree, cv_params, scoring=scoring, cv=4, refit='roc_auc')

```

In [82]:

```

1 %time
2 tree2.fit(X_tr, y_tr)

```

Wall time: 3.51 s

Out[82]:

```
GridSearchCV(cv=4, estimator=DecisionTreeClassifier(random_state=0),
            param_grid={'max_depth': [4, 6, 8, None],
                        'min_samples_leaf': [2, 5, 1],
                        'min_samples_split': [2, 4, 6]},
            refit='roc_auc',
            scoring={'f1', 'roc_auc', 'accuracy', 'precision', 'recall'})
```

```
In [83]: 1 # Check best params
2 tree2.best_params_
```

Out[83]: {'max_depth': 6, 'min_samples_leaf': 1, 'min_samples_split': 4}

```
In [84]: 1 # Check best AUC score on CV
2 tree2.best_score_
```

Out[84]: 0.9534827790328492

This model performs very well, even without satisfaction levels and detailed hours worked data.

Next, check the other scores.

```
In [85]: 1 # Get all CV scores
2 tree2_cv_results = make_results('decision tree2 cv', tree2, 'auc')
3 tree2_cv_results
```

Out[85]:

	Model	AUC	Precision	Recall	F1	Accuracy
0	decision tree2 cv	0.953483	0.864826	0.902098	0.883024	0.960245

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

Random forest - Round 2

```
In [86]: 1 # Instantiate model
2 rf = RandomForestClassifier(random_state=0)
3
4 # Assign a dictionary of hyperparameters to search over
5 cv_params = {'max_depth': [3,5, None],
6               'max_features': [1.0],
7               'max_samples': [0.7, 1.0],
8               'min_samples_leaf': [1,2,3],
9               'min_samples_split': [2,3,4],
10              'n_estimators': [300, 500],
11              }
12
13 # Assign a dictionary of scoring metrics to capture
14 scoring = {'accuracy', 'precision', 'recall', 'f1', 'roc_auc'}
15
16 # Instantiate GridSearch
17 rf2 = GridSearchCV(rf, cv_params, scoring=scoring, cv=4, refit='roc_auc')
```

```
In [87]: 1 %%time
2 rf2.fit(X_tr, y_tr) # --> Wall time: 17min 5s
```

Wall time: 20min 10s

```
Out[87]: GridSearchCV(cv=4, estimator=RandomForestClassifier(random_state=0),
                      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]},
                      refit='roc_auc',
                      scoring={'f1', 'roc_auc', 'accuracy', 'precision', 'recall'})
```

```
In [88]: 1 # Write pickle
2 write_pickle(path, rf2, 'hr_rf2')
```

```
In [89]: 1 # Read in pickle
2 rf2 = read_pickle(path, 'hr_rf2')
```

```
In [90]: 1 # Check best params
2 rf2.best_params_
```

```
Out[90]: {'max_depth': None,
          'max_features': 1.0,
          'max_samples': 0.7,
          'min_samples_leaf': 3,
          'min_samples_split': 2,
          'n_estimators': 300}
```

```
In [91]: 1 # Check best AUC score on CV
2 rf2.best_score_
```

```
Out[91]: 0.9656664586139387
```

```
In [92]: 1 # Get all CV scores
2 rf2_cv_results = make_results('random forest2 cv', rf2, 'auc')
3 print(tree2_cv_results)
4 print(rf2_cv_results)
```

	Model	AUC	Precision	Recall	F1	Accuracy
0	decision tree2 cv	0.953483	0.864826	0.902098	0.883024	0.960245
	Model	AUC	Precision	Recall	F1	Accuracy
0	random forest2 cv	0.965666	0.909473	0.879506	0.894049	0.965388

Again, the scores dropped slightly, but the random forest performs better than the decision tree.

Test the models on the validation set now.

In [93]:

```

1 # Collect validation scores
2 tree2_val_results = get_scores('decision tree2 val', tree2, X_val, y_val)
3 rf2_val_results = get_scores('random forest2 val', rf2, X_val, y_val)
4
5 # Concatenate validation scores into table
6 all_val_results2 = [tree2_val_results, rf2_val_results]
7 all_val_results2 = pd.concat(all_val_results2).sort_values(by='AUC', ascending=False)
8 all_val_results2

```

Out[93]:

	model	AUC	precision	recall	f1	accuracy
0	decision tree2 val	0.942	0.883	0.907	0.895	0.965
0	random forest2 val	0.933	0.905	0.884	0.895	0.965

It appears that the random forest performs slightly better than the decision tree, across most of the evaluation metrics.

Use this random forest model to predict on the test set now.

In [94]:

```

1 # Get predictions on test data
2 rf2_test_scores = get_scores('random forest2 test', rf2, X_test, y_test)
3 rf2_test_scores

```

Out[94]:

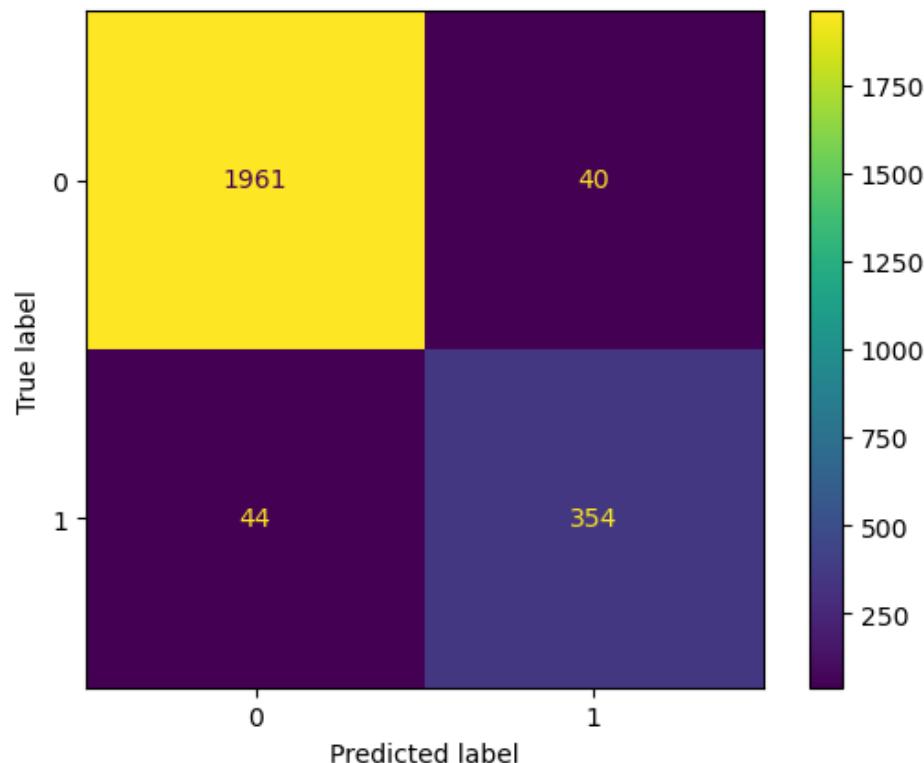
	model	AUC	precision	recall	f1	accuracy
0	random forest2 test	0.935	0.898	0.889	0.894	0.965

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

Plot a confusion matrix to visualize how well it predicts on the test set.

In [95]:

```
1 # Generate array of values for confusion matrix
2 preds = rf2.best_estimator_.predict(X_test)
3 cm = confusion_matrix(y_test, preds, labels=rf2.classes_)
4
5 # Plot confusion matrix
6 disp = ConfusionMatrixDisplay(confusion_matrix=cm,
7                                display_labels=rf2.classes_)
8 disp.plot();
```



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. who are not actually at risk of doing so. But this is still a strong model.

For exploratory purpose, you might want to inspect the splits of the decision tree model and the most important features in the random forest model.

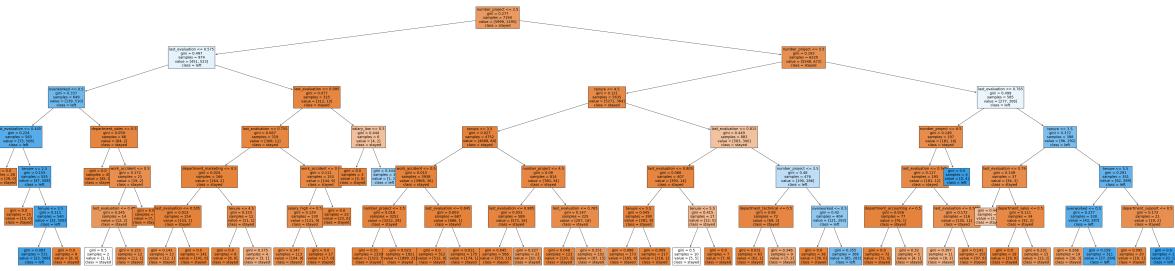
Decision tree splits

In [96]:

```

1 # Plot the tree
2 plt.figure(figsize=(85,20))
3 plot_tree(tree2.best_estimator_, max_depth=6, fontsize=14, feature_names=X.columns,
4           class_names={0:'stayed', 1:'left'}, filled=True);
5 plt.show()

```



Decision tree feature importance

You can also get feature importance from decision trees (see the [DecisionTreeClassifier scikit-learn documentation](https://scikit-learn.org/stable/modules/generated/sklearn.tree.DecisionTreeClassifier.html#sklearn.tree.DecisionTreeClassifier) (<https://scikit-learn.org/stable/modules/generated/sklearn.tree.DecisionTreeClassifier.html#sklearn.tree.DecisionTreeClassifier> for details).

In [97]:

```

1 #tree2_importances = pd.DataFrame(tree2.best_estimator_.feature_importances_, columns=['feature'])
2 tree2_importances = pd.DataFrame(tree2.best_estimator_.feature_importances_, columns=['feature'])
3 tree2_importances = tree2_importances.sort_values(by='gini_importance', ascending=False)
4
5 # Only extract the features with importances > 0
6 tree2_importances = tree2_importances[tree2_importances['gini_importance'] != 0]
7 tree2_importances

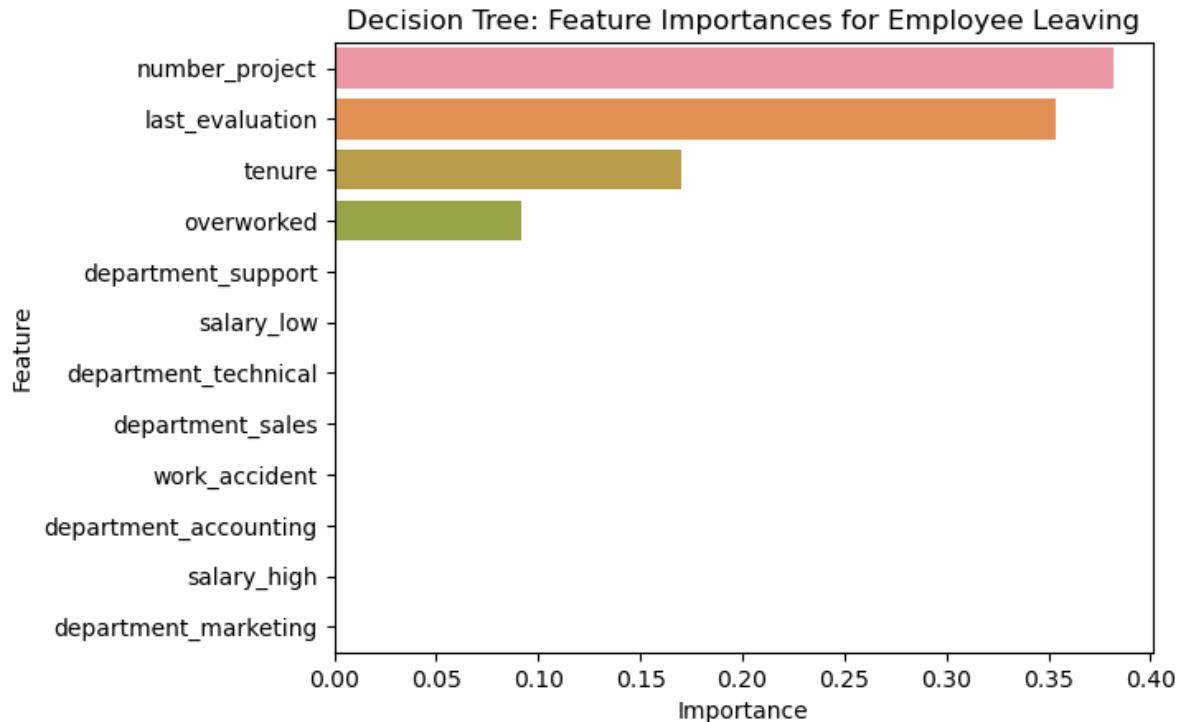
```

Out[97]:

	gini_importance
number_project	0.381767
last_evaluation	0.353482
tenure	0.169771
overworked	0.091517
department_support	0.001099
salary_low	0.000853
department_technical	0.000429
department_sales	0.000357
work_accident	0.000282
department_accounting	0.000239
salary_high	0.000120
department_marketing	0.000084

You can then create a barplot to visualize the decision tree feature importances.

```
In [98]: 1 sns.barplot(data=tree2_importances, x="gini_importance", y=tree2_importances.index)
2 plt.title("Decision Tree: Feature Importances for Employee Leaving", fontsize=12)
3 plt.ylabel("Feature")
4 plt.xlabel("Importance")
5 plt.show()
```



```
# This is formatted as code
```

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 predicting the outcome variable, `left`.

Random forest feature importance

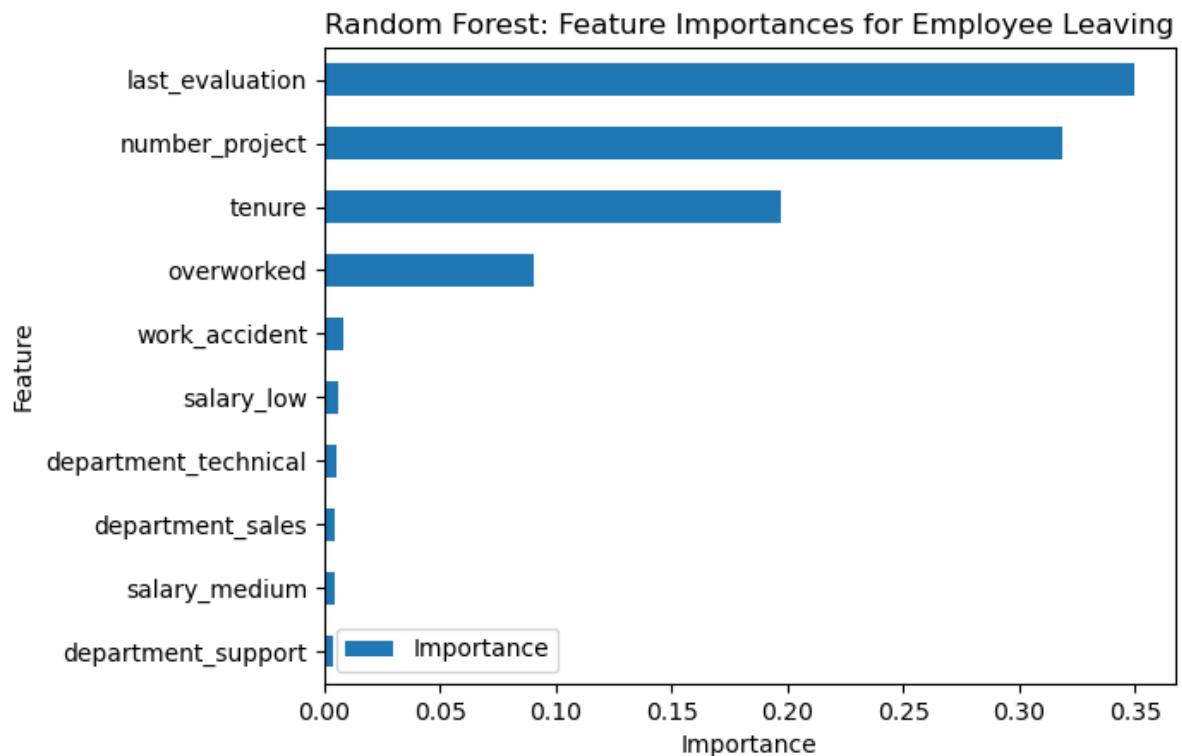
Now, plot the feature importances for the random forest model.

In [99]:

```

1 # Get feature importances
2 feat_impt = rf2.best_estimator_.feature_importances_
3
4 # Get indices of top 10 features
5 ind = np.argpartition(rf2.best_estimator_.feature_importances_, -10)[-10:]
6
7 # Get column labels of top 10 features
8 feat = X.columns[ind]
9
10 # Filter `feat_impt` to consist of top 10 feature importances
11 feat_impt = feat_impt[ind]
12
13 y_df = pd.DataFrame({"Feature":feat,"Importance":feat_impt})
14 y_sort_df = y_df.sort_values("Importance")
15 fig = plt.figure()
16 ax1 = fig.add_subplot(111)
17
18 y_sort_df.plot(kind='barh',ax=ax1,x="Feature",y="Importance")
19
20 ax1.set_title("Random Forest: Feature Importances for Employee Leaving", fontsize=14)
21 ax1.set_ylabel("Feature")
22 ax1.set_xlabel("Importance")
23
24 plt.show()

```



The plot above shows that in this random forest model, `last_evaluation`, `number_project`, `tenure`, `overworked`, `work_accident`, and `salary_low` have the highest importance, in that order. These variables are most helpful in predicting the outcome variable, `left`.



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?
- What business recommendations do you propose based on the models built?
- 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?

Step 4. Results and Evaluation

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

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

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.

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

In []: 1