

# Providing data-driven suggestions for HR

## Description and Deliverables

This capstone project is 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, we will have two artifacts that we would be able to present to future employers. One is a brief one-page summary of this project that we would present to external stakeholders as the data professional in Salifort Motors. We 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 we only need to do one.

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

## PACE stages

### Pace: Plan

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

In this stage, consider the following:

### 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 me as a data analytics professional and ask me to provide data-driven suggestions based on my understanding of the data. They have the following question: what's likely to make the employee leave the company?

My 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 we 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), 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 we complete the plan stage.

- Who are the stakeholders for this project?
- What are we trying to solve or accomplish?
- What are my initial observations when I explore the data?
- What resources do I find yourself using as I complete this stage?
- Do I have any ethical considerations in this stage?

## Step 1. Imports

- Import packages
- Load dataset

## Import packages

```
[1]: # Import packages

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

# For data visualization
import matplotlib.pyplot as plt
import seaborn as sns

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

# For data modeling
from xgboost import XGBClassifier
from xgboost import XGBRegressor
from xgboost import plot_importance

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

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

# For saving models
import pickle
```

## Load dataset

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

# Load dataset into a dataframe
df0 = pd.read_csv("HR_capstone_dataset.csv")

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

Out[2]:

	satisfaction_level	last_evaluation	number_project	average_monthly_hours	time_spend_company
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 my variables
- Clean the dataset (missing data, redundant data, outliers)

### Gather basic information about the data

```
In [3]: # Gather basic information about the data
df0.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 14999 entries, 0 to 14998
Data columns (total 10 columns):
#   Column                      Non-Null Count  Dtype
---  -
0   satisfaction_level           14999 non-null  float64
1   last_evaluation              14999 non-null  float64
2   number_project               14999 non-null  int64
3   average_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
```

In

```
9 salary          14999 non-null object dtypes: float64(2),
int64(6), object(2) memory usage: 1.1+ MB
```

## Gather descriptive statistics about the data

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

Out[4]:

	satisfaction_level	last_evaluation	number_project	average_monthly_hours	time_spend_com
count	14999.000000	14999.000000	14999.000000	14999.000000	14999.00
mean	0.612834	0.716102	3.803054	201.050337	3.49
std	0.248631	0.171169	1.232592	49.943099	1.46
min	0.090000	0.360000	2.000000	96.000000	2.00
25%	0.440000	0.560000	3.000000	156.000000	3.00
50%	0.640000	0.720000	4.000000	200.000000	3.00
75%	0.820000	0.870000	5.000000	245.000000	4.00
max	1.000000	1.000000	7.000000	310.000000	10.00

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

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

```
In [6]: # 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
```

```
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]: # Check for missing values  
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]: # Check for duplicates  
df0.duplicated().sum()
```

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

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

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

```
Out[9]:
```

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

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?

We 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. we can proceed by dropping them.

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

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

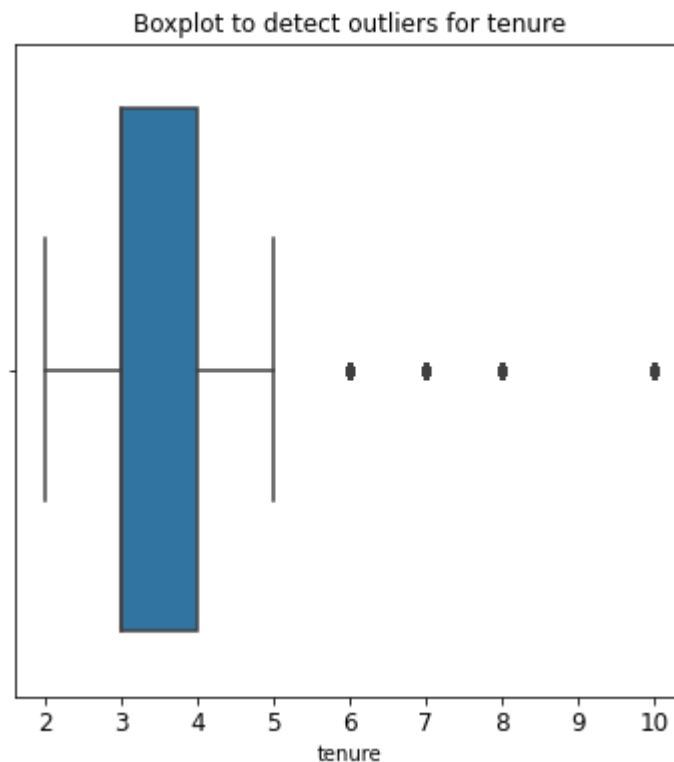
```
Out[10]:
```

	satisfaction_level	last_evaluation	number_project	average_monthly_hours	tenure	work_accide
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	

## Check outliers

Check for outliers in the data.

```
In [11]: # Create a boxplot to visualize distribution of `tenure` and detect any outli
plt.figure(figsize=(6,6))
plt.title('Boxplot to detect outliers for tenure', fontsize=12)
plt.xticks(fontsize=12)
plt.yticks(fontsize=12)
sns.boxplot(x=df1['tenure'])
plt.show()
```



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

# 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(outlie

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 we get to the stage of building the model, we consider whether to remove these outliers based on the type of model we decide to use.

## Pace: Analyze Stage

- Perform EDA (analyze relationships between variables)

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

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

## Step 2. Data Exploration (Continue EDA)

Begin by understanding how many employees left and what percentage of all employees this fit

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

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

0    10000
1     1991
Name: left, dtype: int64

0    0.833959
1    0.166041
Name: left, dtype: float64
```

## Data visualizations

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

We 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, we could also plot a stacked histogram to visualize the distribution of `number_project` for those who stayed and those who left.

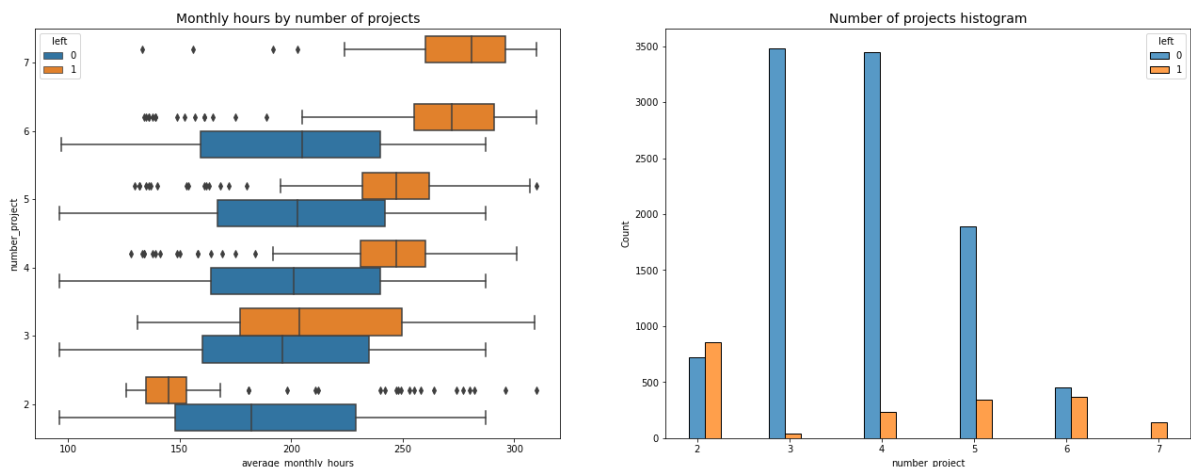
```
In [14]: # Create a plot as needed

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

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

# Create histogram showing distribution of `number_project`, comparing employee
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
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 we 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 =  $50 \text{ weeks} * 40 \text{ hours per week} / 12 \text{ months} = 166.67 \text{ hours per month}$ . 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

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

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

This confirms that all employees with 7 projects did leave.

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

In [16]: `# Create a plot as needed`

```
# Create scatterplot of `average_monthly_hours` versus `satisfaction_level`, co
plt.figure(figsize=(16, 9))
sns.scatterplot(data=df1, x='average_monthly_hours', y='satisfaction_level', hu
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 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.

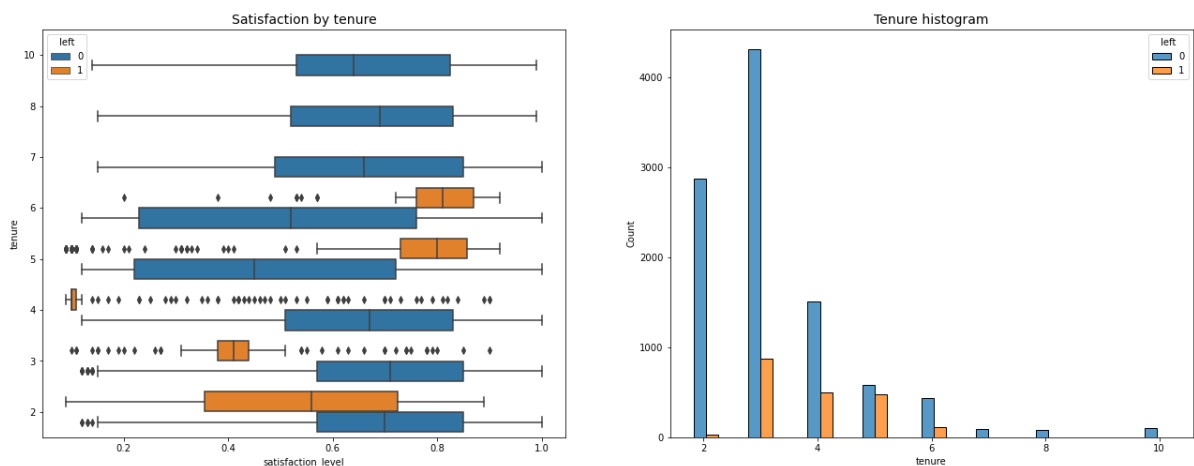
```
In [17]: # Create a plot as needed

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

# Create boxplot showing distributions of `satisfaction_level` by tenure, comparing employees who left and those who stayed
sns.boxplot(data=df1, x='satisfaction_level', y='tenure', hue='left', orient="horizontal", ax=ax[0].invert_yaxis())
ax[0].set_title('Satisfaction by tenure', fontsize='14')

# Create histogram showing distribution of `tenure`, comparing employees who left and those who stayed
tenure_stay = df1[df1['left']==0]['tenure']
tenure_left = df1[df1['left']==1]['tenure']
sns.histplot(data=df1, x='tenure', hue='left', multiple='dodge', shrink=5, ax=ax[1])
ax[1].set_title('Tenure histogram', fontsize='14')

plt.show();
```



There are many observations we 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, we could calculate the mean and median satisfaction scores of employees who left and those who didn't.

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

```
Out[18]:
```

	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, we could examine salary levels for different tenures.

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

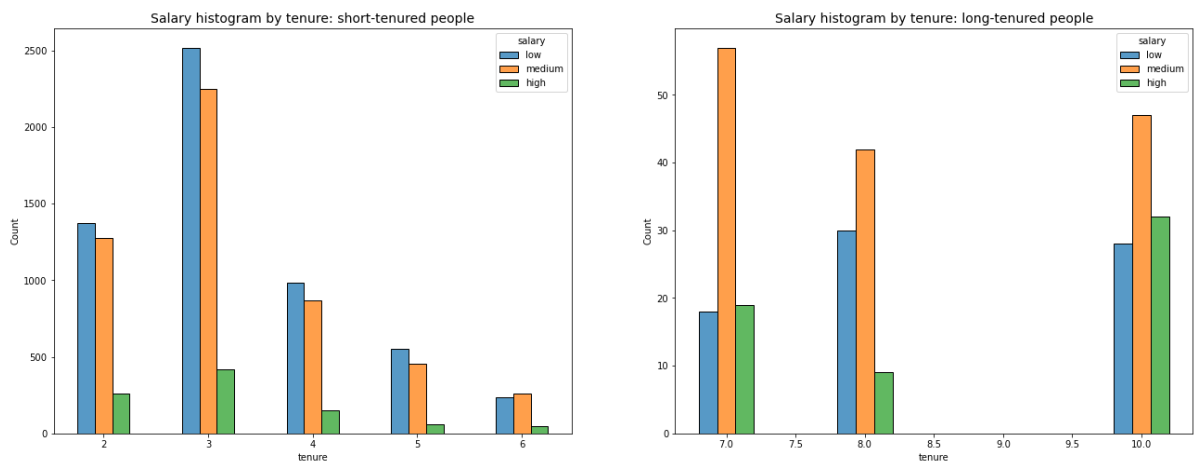
# Set figure and axes
fig, ax = plt.subplots(1, 2, figsize = (22,8))

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

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

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

# Plot Long-tenured histogram
sns.histplot(data=tenure_long, x='tenure', hue='salary', discrete=1,
             hue_order=['low', 'medium', 'high'], multiple='dodge', shrink=.4,
             ax[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, we could explore whether there's a correlation between working long hours and receiving high evaluation scores. We could create a scatterplot of average\_monthly\_hours versus



last\_evaluation .

In [20]: # Create a plot as needed

```
# 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=
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');
```



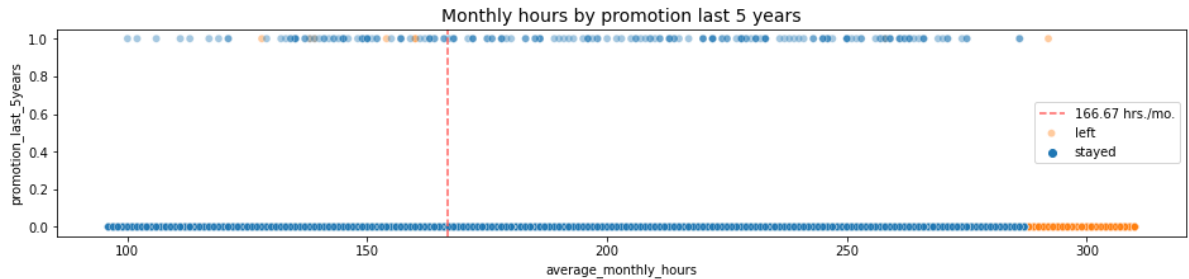
The following observations can be made from the scatterplot above:

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

```
In [21]: # Create a plot as needed

# Create plot to examine relationship between `average_monthly_hours` and `prom
plt.figure(figsize=(16, 3))
sns.scatterplot(data=df1, x='average_monthly_hours', y='promotion_last_5years',
plt.axvline(x=166.67, color='#ff6361', ls='--')
plt.legend(labels=['166.67 hrs./mo.', 'left', 'stayed'])
plt.title('Monthly hours by promotion last 5 years', fontsize='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, we could inspect how the employees who left are distributed across departments.

```
In [22]: # Display counts for each department
df1["department"].value_counts()
```

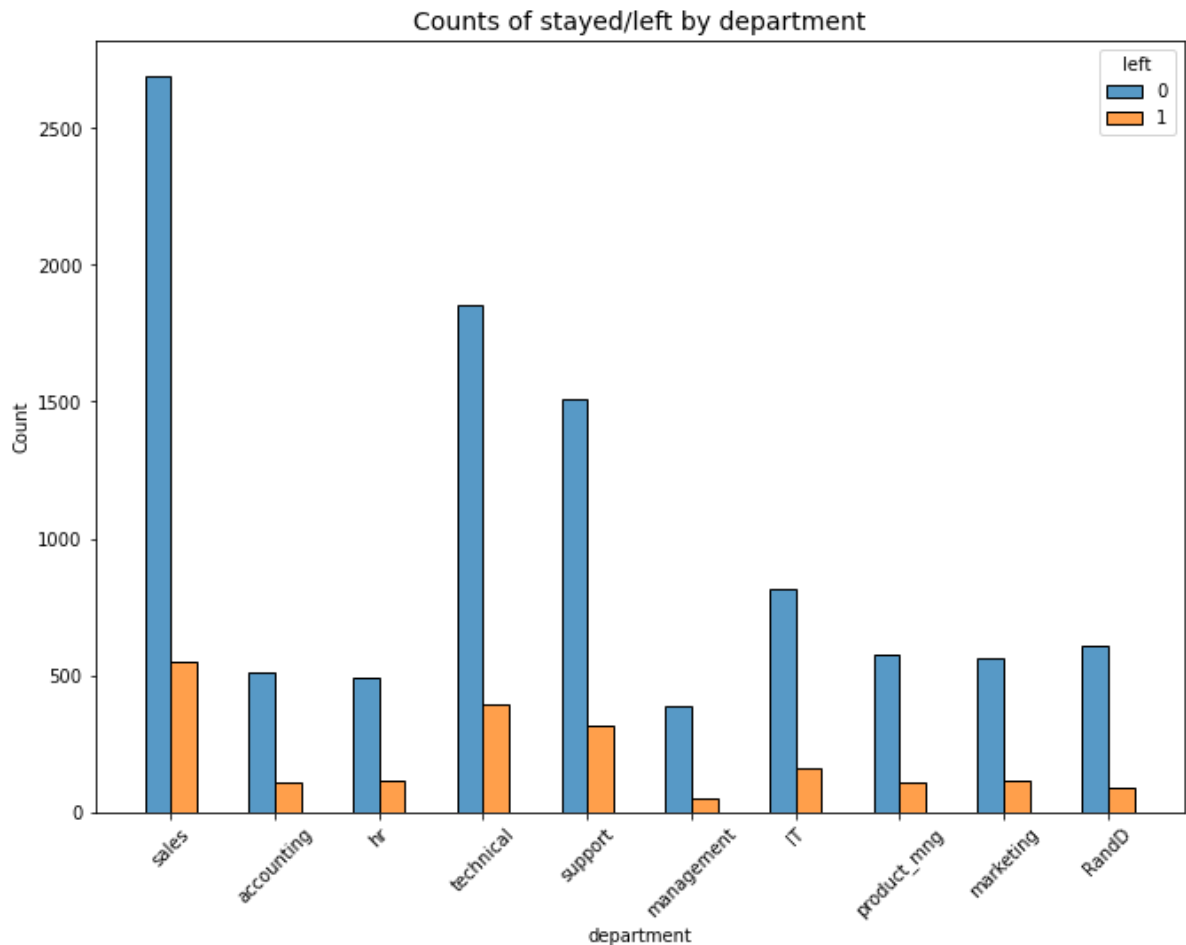
```
Out[22]: sales          3239
technical    2244
support      1821
IT           976
RandD        694
product_mng  686
marketing     673
accounting    621
hr            601
management   436
Name: department, dtype: int64
```



In

```
[23]: # Create a plot as needed

# Create stacked histogram to compare department distribution of employees who
plt.figure(figsize=(11,8))
sns.histplot(data=df1, x='department', hue='left', discrete=1,
             hue_order=[0, 1], multiple='dodge', shrink=.5)
plt.xticks(rotation='45')
plt.title('Counts of stayed/left by department', fontsize=14);
```



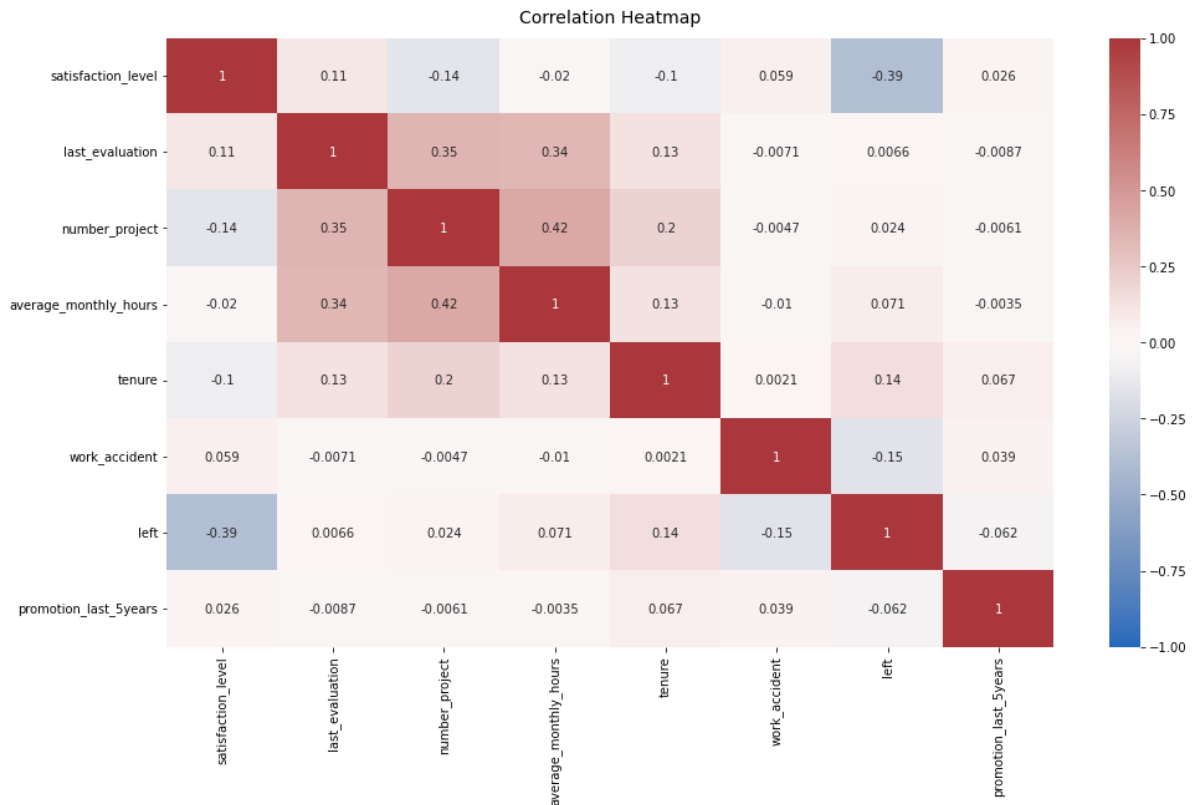
There doesn't seem to be any department that differs significantly in its proportion of employees who left to those who stayed.

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

In

```
[24]: # Create a plot as needed

# Plot a correlation heatmap
plt.figure(figsize=(16, 9))
heatmap = sns.heatmap(df0.corr(), vmin=-1, vmax=1, annot=True, cmap=sns.color_p
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

In

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 we notice anything odd?
- Which independent variables did we choose for the model and why?
- Are each of the assumptions met?
- How well does our model fit the data?
- Can we improve it? Is there anything we would change about the model?
- What resources do I find myself using as I complete this stage?
- Do I 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.

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

So I could proceed with one of the two following approaches. Or, if we'd like, we could implement both and determine how they compare.

## 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.

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

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

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

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

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

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

# Display the new dataframe
df_enc.head()
```

```
Out[25]:
```

	satisfaction_level	last_evaluation	number_project	average_monthly_hours	tenure	work_accide
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	

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

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



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

```
[27]: # Create a stacked bar plot to visualize number of employees across department
# In the legend, 0 (purple color) represents employees who did not leave, 1 (red color) represents employees who left
pd.crosstab(df1['department'], df1['left']).plot(kind='bar', color=['purple', 'red'])
plt.title('Counts of employees who left versus stayed across department')
plt.ylabel('Employee count')
plt.xlabel('Department')
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]: # Select rows without outliers in `tenure` and save resulting dataframe in a new dataframe
df_logreg = df_enc[(df_enc['tenure'] >= lower_limit) & (df_enc['tenure'] <= upper_limit)]
# Display first few rows of new dataframe
df_logreg.head()
```

```
Out[28]:
```

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

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

In

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

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

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

Selecting the features I want to use in the model. We consider which variables will help us predict the outcome variable, left .

```
In [30]: # 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()
```

```
Out[30]:
```

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

Split the data into training set and testing set. Don't forget to stratify based on the values in y , since the classes are unbalanced.

```
In [31]: # Split the data into training set and testing set
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.25, strat
```

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

```
In [32]: # 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)
```

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

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

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

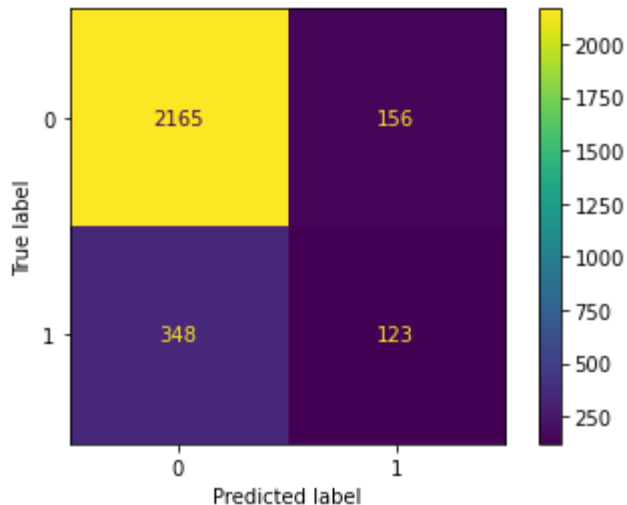


```
[34]: # Compute values for confusion matrix
log_cm = confusion_matrix(y_test, y_pred, labels=log_clf.classes_)

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

# Plot confusion matrix
log_disp.plot(values_format='')

# Display plot
plt.show()
```



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

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

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

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

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

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

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

In

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

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

```
Out[35]: 0    0.831468 1    0.168532
```

```
Name: left, dtype: float64
```

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

```
In [36]: # Create classification report for logistic regression model
target_names = ['Predicted would not leave', 'Predicted would leave']
print(classification_report(y_test, y_pred, target_names=target_names))
precision    recall  f1-score   support
```

Predicted would not leave	0.86	0.93	0.90	2321
Predicted would leave	0.44	0.26	0.33	471
accuracy				
macro avg	0.65	0.60	0.61	2792
avg	0.79	0.82	0.80	2792
				weighted
				0.82
				2792

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

## Modeling Approach B: Tree-based Model

This approach covers implementation of Decision Tree and Random Forest.

Isolate the outcome variable.

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

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

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

Select the features.

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

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

```
Out[38]:
```

	satisfaction_level	last_evaluation	number_project	average_monthly_hours	tenure	work_accide
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	

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

```
In [39]: # Split the data
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.25, strat
```

## 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 [40]: # Instantiate model tree =
DecisionTreeClassifier(random_state=0)

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

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

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

Fit the decision tree model to the training data.

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

%%time

CPU times: user 2.7 s, sys: 91 ms, total: 2.79 s
Wall time: 2.8 s
```

```
Out[41]: GridSearchCV(cv=4, error_score=nan,
    estimator=DecisionTreeClassifier(ccp_alpha=0.0, class_weight=None,
    criterion='gini', max_depth=None,
    max_features=None,
    max_leaf_nodes=None,
    min_impurity_decrease=0.0,
    min_impurity_split=None,
    min_samples_leaf=1,
    min_samples_split=2,
    min_weight_fraction_leaf=0.0,
    presort='deprecated',
    random_state=0, splitter='best'),
    iid='deprecated', n_jobs=None,
    param_grid={'max_depth': [4, 6, 8, None],
    'min_samples_leaf': [2, 5, 1],
    'min_samples_split': [2, 4, 6]},
    pre_dispatch='2*n_jobs', refit='roc_auc', return_train_score=False,
    scoring={'f1', 'precision', 'accuracy', 'roc_auc', 'recall'},
    verbose=0)
```

Identify the optimal values for the decision tree parameters.

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

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

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

```
In [43]: # Check best AUC score on CV
```

```
tree1.best_score_
```

```
Out[43]: 0.969819392792457
```

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

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

```
[44]: 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    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 score
    for the model with the best mean 'metric' score across all validation folds
    ...

    # Create dictionary that maps input metric to actual metric name in GridSea
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]].id
    # 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

```

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

```

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

```

```

Out[45]:

```

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

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

Recall that decision trees can be vulnerable to overfitting, and random forests avoid overfitting by incorporating multiple trees to make predictions. We could construct a random forest model next.

## 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 [46]: # 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')
```

Fit the random forest model to the training data.

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

CPU times: user 9min 7s, sys: 3.01 s, total: 9min 10s  
Wall time: 9min 10s

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

```

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

```

Specify path to where we want to save your model.

```

In [48]: # Define a path to the folder where you want to save the model
path = '/home/jovyan/work/'

```

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

```

In [49]: 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)

In [50]: 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

```

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

```

In [51]: # Write pickle
write_pickle(path, rf1, 'hr_rf1')

```

```

In [52]: # Read pickle
rf1 = read_pickle(path, 'hr_rf1')

```

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

```

In [53]: # Check best AUC score on CV

```

```
rf1.best_score_
```

```
Out[53]: 0.9804250949807172
```

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

In

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

```
Out[54]: {'max_depth': 5,
          'max_features': 1.0,
          'max_samples': 0.7,
          'min_samples_leaf': 1,
          'min_samples_split': 4,
          'n_estimators': 500}
```

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

```
In [55]: # Get all CV scores
rf1_cv_results = make_results('random forest cv', rf1, 'auc')
print(tree1_cv_results)
print(rf1_cv_results)
```

	model	precision	recall	F1	accuracy	auc
0	decision tree cv	0.914552	0.916949	0.915707	0.971978	0.969819
	model	precision	recall	F1	accuracy	auc
	0	random				
	forest cv	0.950023	0.915614	0.932467	0.977983	0.980425



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

Next, we can evaluate the final model on the test set.

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

```
In [56]: 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          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
        """

        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
```

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

```
In [57]: # Get predictions on test data
rf1_test_scores = get_scores('random forest1 test', rf1, X_test, y_test)
rf1_test_scores
```

```
Out[57]:
```

	model	precision	recall	f1	accuracy	AUC
0	random forest1 test	0.964211	0.919679	0.941418	0.980987	0.956439

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

## Feature Engineering

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

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

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

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

```
In [58]: # 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()
```

```
Out[58]:
```

	last_evaluation	number_project	average_monthly_hours	tenure	work_accident	left	promotio
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 [59]: # Create `overworked` column. For now, it's identical to average monthly hours.
df2['overworked'] = df2['average_monthly_hours']

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

Max hours: 310

Min hours: 96

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

We could define being overworked as working more than 175 hours per month on average.

To make the `overworked` column binary, we 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 [60]: # Define `overworked` as working > 175 hrs/week
df2['overworked'] = (df2['overworked'] > 175).astype(int)

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

```
Out[60]: 0    0
         1    1
         2    1
         3    1
         4    0
Name: overworked, dtype: int64
```

Drop the `average_monthly_hours` column.

```
In [61]: # Drop the `average_monthly_hours` column
df2 = df2.drop('average_monthly_hours', axis=1)

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

```
Out[61]:
```

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

Again, isolate the features and target variables

```
In [62]: # Isolate the outcome variable
y = df2['left']

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

Split the data into training and testing sets.

[63]:

#

Create test data

```
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.25, strat
```

## Decision tree - Round 2

In

In [64]:

```
# 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')
```

```
tree2.fit(X_train, y_train)
```

[65]: %%time

```
CPU times: user 2.49 s, sys: 1.36 ms, total: 2.49 s
Wall time: 2.49 s
```

```
Out[65]: GridSearchCV(cv=4, error_score=nan,
                      estimator=DecisionTreeClassifier(ccp_alpha=0.0, class_weight=Non
e,
e,
                      criterion='gini', max_depth=Non
e,
                      max_features=None,
                      max_leaf_nodes=None,
                      min_impurity_decrease=0.0,
                      min_impurity_split=None,
                      min_samples_leaf=1,
                      min_samples_split=2,
                      min_weight_fraction_leaf=0.0,
                      presort='deprecated',
                      random_state=0, splitter='bes
t'),
                      iid='deprecated', n_jobs=None,
                      param_grid={'max_depth': [4, 6, 8, None],
                      'min_samples_leaf': [2, 5, 1],
                      'min_samples_split': [2, 4, 6]},
                      pre_dispatch='2*n_jobs', refit='roc_auc', return_train_score=Fal
```

```
se,
        scoring={'f1', 'precision', 'accuracy', 'roc_auc', 'recall'},
        verbose=0)
```

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

```
Out[66]: {'max_depth': 6, 'min_samples_leaf': 2, 'min_samples_split': 6}
```

```
In [67]: # Check best AUC score on CV
```

```
tree2.best_score_
```

```
Out[67]: 0.9586752505340426
```

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

Next, check the other scores.

```
In [68]: # Get all CV scores
      tree2_cv_results = make_results('decision tree2 cv', tree2, 'auc')
      print(tree1_cv_results)
      print(tree2_cv_results)
```

	model	precision	recall	F1	accuracy	auc
0	decision tree cv	0.914552	0.916949	0.915707	0.971978	0.969819
	model precision	recall	F1	accuracy	auc	0
	tree2 cv	0.856693	0.903553	0.878882	0.958523	0.958675

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 [69]: # 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')
```

```
[70]: %%time rf2.fit(X_train, y_train) # --> Wall time:
```

```
7min 5s
```

```
CPU times: user 7min 10s, sys: 1.17 s, total: 7min 11s
```

```
Wall time: 7min 11s
```

```
Out[70]: GridSearchCV(cv=4, error_score=nan,
 estimator=RandomForestClassifier(bootstrap=True, ccp_alpha=0.0,
 class_weight=None, criterion='gini',
 max_depth=None,
 e,
 max_features='auto',
 max_leaf_nodes=None,
 max_samples=None,
 min_impurity_decrease=0.0,
 min_impurity_split=None,
 min_samples_leaf=1,
 min_samples_split=2,
 min_weight_fraction_leaf=0.0,
 n_estimators=100, n_jobs=None,
 e,... verbose=0,
 warm_start=False), iid='deprecated', n_jobs=None,
 param_grid={'max_depth': [3, 5, None], 'max_features': [1.0],
 'max_samples': [0.7, 1.0],
 'min_samples_leaf': [1, 2, 3],
 'min_samples_split': [2, 3, 4],
 'n_estimators': [300, 500]},
 pre_dispatch='2*n_jobs', refit='roc_auc', return_train_score=False,
 scoring={'f1', 'precision', 'accuracy', 'roc_auc', 'recall'},
 verbose=0)
```

In

```
In [71]: # Write pickle
write_pickle(path, rf2, 'hr_rf2')
```

```
In [72]: # Read in pickle
rf2 = read_pickle(path, 'hr_rf2')
```

```
rf2.best_params_
```

```
Out[73]: {'max_depth': 5,
 'max_features': 1.0,
 [73]: # Check best params
 'max_samples': 0.7,
 'min_samples_leaf': 2,
 'min_samples_split': 2,
```

```
'n_estimators': 300}
```

```
In [74]: # Check best AUC score on CV
```

```
rf2.best_score_
```

```
Out[74]: 0.9648100662833985
```

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

	model	precision	recall	F1	accuracy	auc
0	decision tree2 cv	0.856693	0.903553	0.878882	0.958523	0.958675
	model precision	recall	F1	accuracy	auc	0 random
	forest2 cv	0.866758	0.878754	0.872407	0.957411	0.96481

Again, the scores dropped slightly, but the random forest performs better than the decision tree if using AUC as the deciding metric.

Score the champion model on the test set now.

```
In [76]: # Get predictions on test data
rf2_test_scores = get_scores('random forest2 test', rf2, X_test, y_test)
rf2_test_scores
```

```
Out[76]:
```

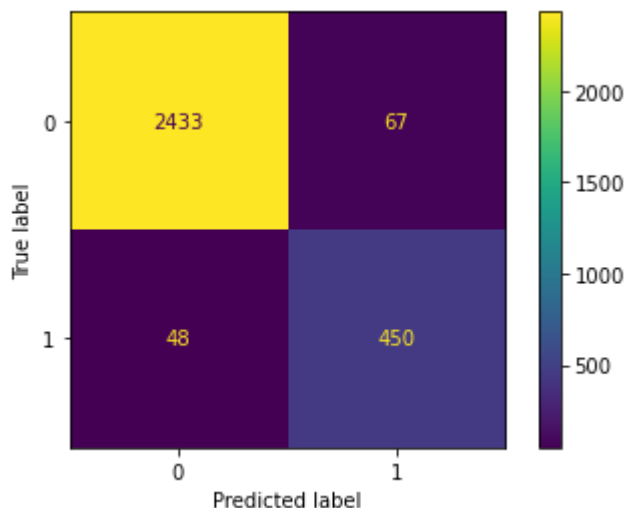
	model	precision	recall	f1	accuracy	AUC
0	random forest2 test	0.870406	0.903614	0.8867	0.961641	0.938407

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

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



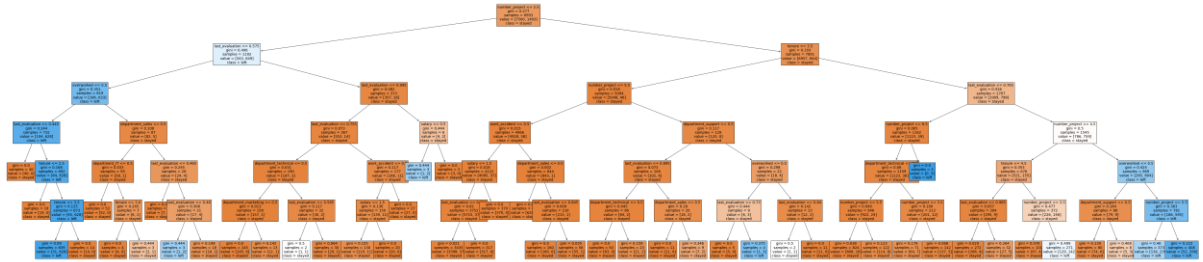
The model predicts more false positives than false negatives, which means that some employees may be identified as at risk of quitting or getting fired, when that's actually not the case. But this is still a strong model.



For exploratory purpose, We 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 [78]: # 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()
```



Note that we can double-click on the tree image to zoom in on it and inspect the splits.

### Decision tree feature importance

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

```
In [79]: #tree2_importances = pd.DataFrame(tree2.best_estimator_.feature_importances_, c
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
```

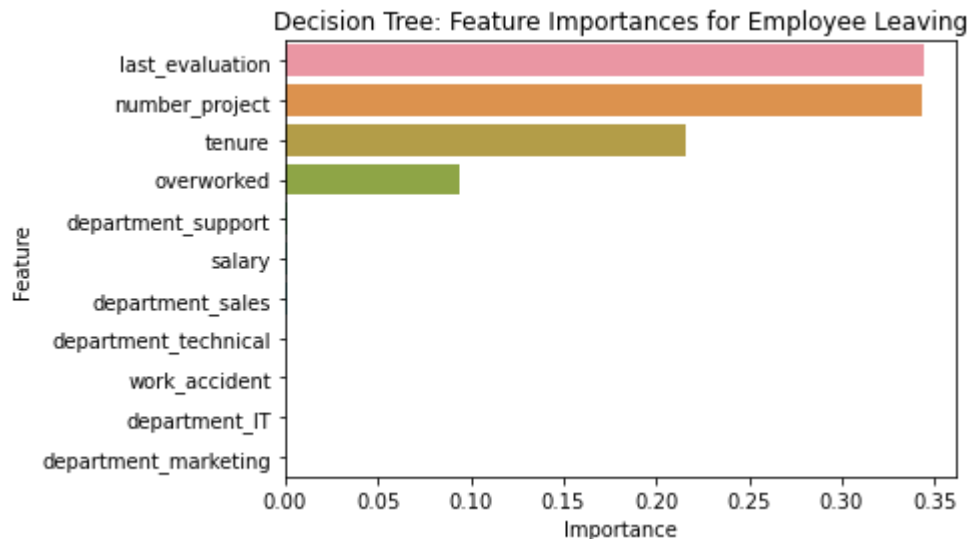
```
Out[79]:
```

	gini_importance
last_evaluation	0.343958
number_project	0.343385
tenure	0.215681
overworked	0.093498
department_support	0.001142
salary	0.000910
department_sales	0.000607
department_technical	0.000418
work_accident	0.000183

department_IT	0.000139
department_marketing	0.000078

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

```
In [80]: sns.barplot(data=tree2_importances, x="gini_importance", y=tree2_importances.in
plt.title("Decision Tree: Feature Importances for Employee Leaving", fontsize=1
plt.ylabel("Feature")
plt.xlabel("Importance")
plt.show()
```



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 [81]: # 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]

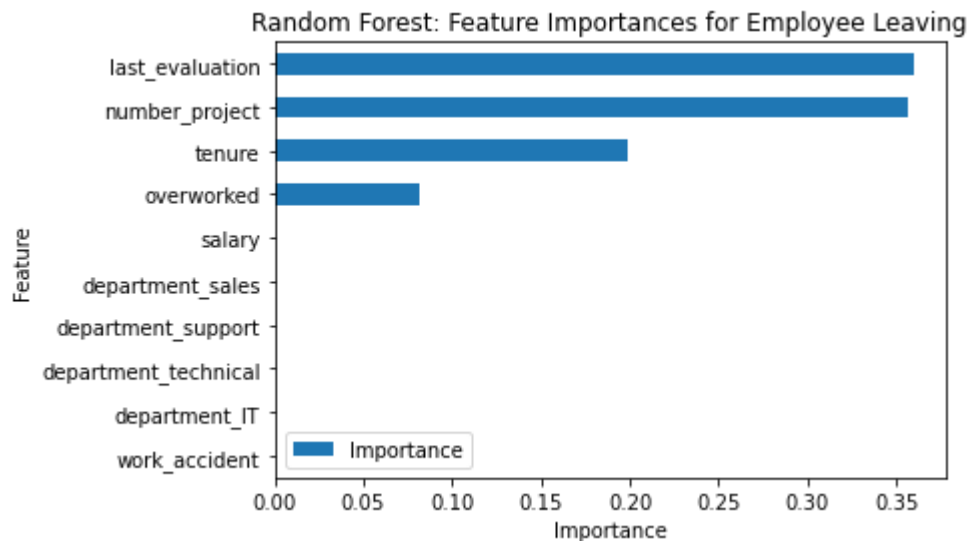
# Filter `feat_impt` to consist of top 10 feature importances
feat_impt = feat_impt[ind]

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", fontsi
ax1.set_ylabel("Feature")
ax1.set_xlabel("Importance")

plt.show()
```



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

## 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 the model(s)?
- What business recommendations do we propose based on the models built?
- What potential recommendations would we make to the manager/company?
- Do I think the model could be improved? Why or why not? How?
- Given what I know about the data and the models we were using, what other questions could we address for the team?
- What resources do I find yourself using as I complete this stage?
- Do I 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

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.

### 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, we could try building a K-means model on this data and analyzing the clusters. This may yield valuable insight.