

BA780: Introduction to Data Analytics

Team Project

Team 8:

- Jacinto Lemarroy
- Yongxian (Caroline) Lun
- Yipeng (Caroline) Guo
- Chris Chang
- Aash Gohil

Dataset: HR Analytics: Job Change in Data Scientists

(Data source from Analytics Vidhya: <https://datahack.analyticsvidhya.com/contest/janatahack-hr-analytics/True/#DiscussTab>)

The dataset is from JanataHack (Machine Learning Hackathon), a knowledge competition on Machine Learning & Data Science. It's powered by Analytics Vidhya Community, a data science community.

Background: A company which is active in Big Data and Data Science wants to hire data scientists among people who received training. "Company" wants to analyze the factors affecting candidates' decision on staying or looking for a new job after training.

Objective: The project goal is to predict whether a data scientist candidate will look for a new employment or wants to work for the company after training, which helps optimize HR costs and increase efficiencies. By using both descriptive and predictive analysis on a company's HR dataset, we seek to interpret affecting factors on employee decisions.

▼ Summary of the Dataset

To explore the factors that influence the employee decision, we chose to narrow the features of our dataset to the following 13 variables:

12,477 rows and 13 features

- enrollee_id : Unique ID for candidate
- city_development_index : Development index of the city (scaled)
- gender: Gender of candidate
- relevant_experience: Relevant experience of candidate

- `enrolled_university`: Type of University course enrolled if any
- `education_level`: Education level of candidate
- `major_discipline`: Education major discipline of candidate
- `experience`: Candidate total experience in years
- `company_size`: No of employees in current employer's company
- `company_type`: Type of current employer
- `lastnewjob`: Difference in years between previous job and current job
- `training_hours`: training hours completed
- `target`: 0 – Stay in the "Company" after Training, 1 – Looking for a New job after Training

Data preparationg for ML

To get dummies:

- `gender`->0,1,
- `experience`->0~21,
- `relevent_experience`: 0,1 binary
- `enrolled_university`: get dummy, 3 columns
- `education_level`: get dummy, 3 columns
- `major_discipline`: 0,1 STEM, NOT STEM
- `company_size`: small, large, medium, unknown
- `company_type`: Unknown, startupGroup, Private, public, NGO, other
- `lastnewjob`-> turn into integer: 0,1,2,3,4,5;

▼ 1. Data Cleaning and Processing

▼ 1.1 Null Value Reasoning

```
1 pip install squarify
```

```
Collecting squarify
  Downloading squarify-0.4.3-py3-none-any.whl (4.3 kB)
Installing collected packages: squarify
Successfully installed squarify-0.4.3
```

```
1 pip install circlify
```

```
Collecting circlify
  Downloading circlify-0.13-py2.py3-none-any.whl (10 kB)
```

```
Installing collected packages: circlify
Successfully installed circlify-0.13
```

```
1 # Importing libraries
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 import numpy as np
6 %matplotlib inline
7 import plotly.express as px
8 from plotly.subplots import make_subplots
9 import plotly.graph_objs as go
10 import squarify

1 # Importing Data
2 df = pd.read_csv('https://raw.githubusercontent.com/aashgohil/HR_Data_Science/main
3 df.head()

1 # Bucketing the company_size values into categorizes
2 df['size'] = np.where(df['company_size'].isin(['50-99','10-49','<10']), 'Small Com
3 df['size'] = np.where(df['company_size'].isin(['100-500','500-999']), 'Medium Comp
4 df['size'] = np.where(df['company_size'].isin(['10000+','1000-4999','5000-9999']),
5 sorted_counts=df['size'].value_counts()

1 df.drop(['company_size'], axis=1, inplace=True)
2 df.rename(columns={'size': 'company_size'}, inplace=True)

1 df.isnull().sum()

enrollee_id          0
city                 0
city_development_index  0
gender              4508
relevent_experience   0
enrolled_university  386
education_level      460
major_discipline     2813
experience           65
company_type         6140
last_new_job         423
training_hours       0
target              0
company_size         0
dtype: int64
```

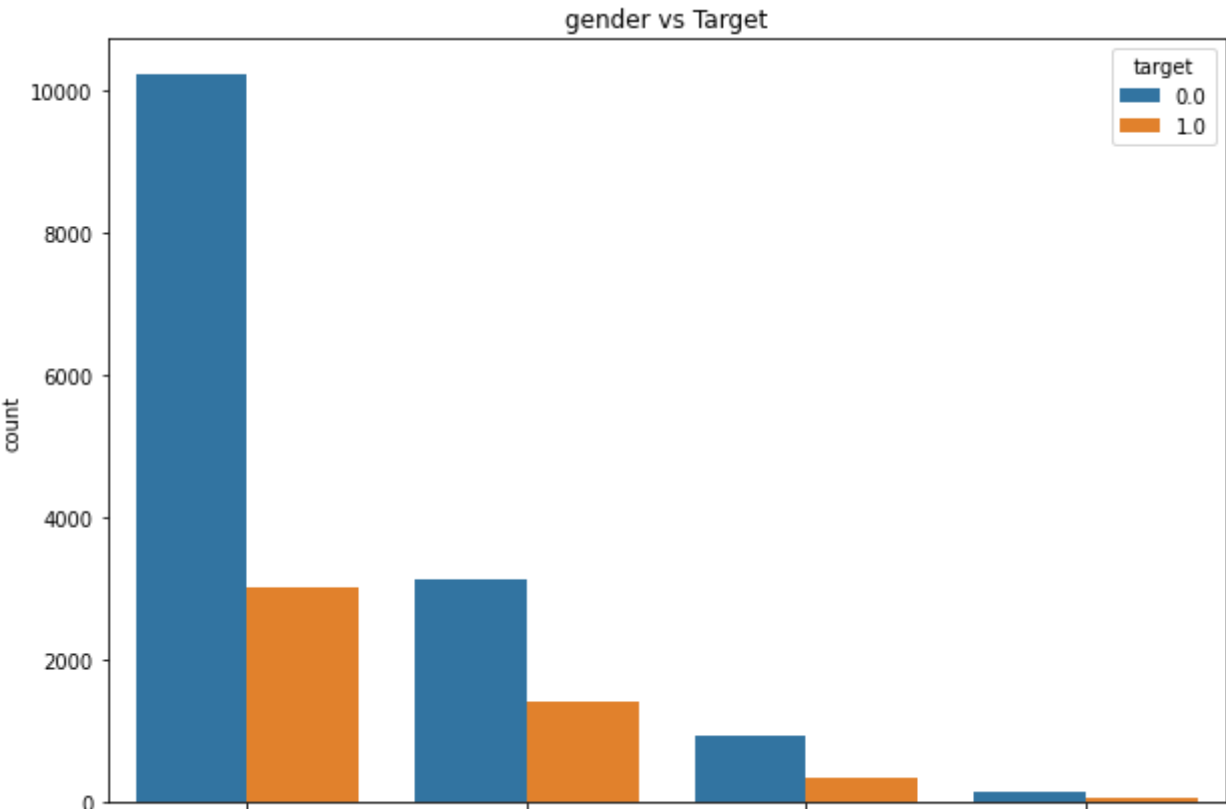
There are quite a few features with null values. We will investigate if null values have a significant impact on attrition, by encoding null with "Unknown". We will drop the null values for features where the ratio of attrition for null values is similar to that of non null values. On the contrary, we will keep

null values as "Unknown" for features where null values seems to have a significant impact on attrition .

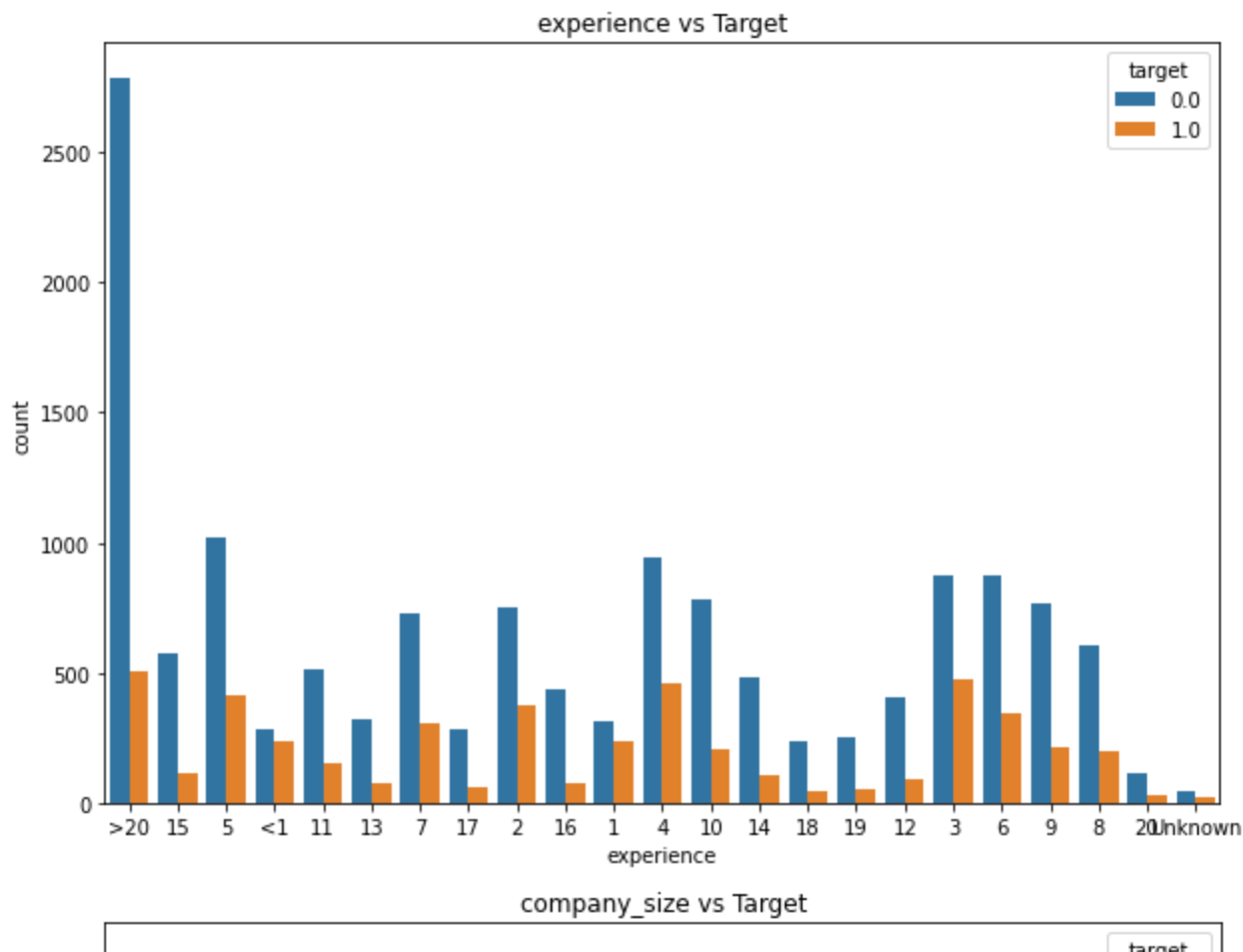
```
1 df_temp = df.fillna('Unknown')

1 ## This function takes in a list of features and plots a count plot against the ta
2 def plotsal(category):
3
4     for cat in category:
5         plt.figure(figsize=(10,7))
6         ax = sns.countplot(x=cat, hue='target', data=df_temp)
7         plt.title(cat + ' vs Target')
8
9         plt.show()

1 plotsal(['gender','enrolled_university','education_level','major_discipline'])
```



```
1 plotsal(['experience','company_size','company_type','last_new_job'])
```



From the graphs above we can conclude, only for the variables company_size and company_type does the unknown value play an impact.

▼ 1.2 Data Cleaning

Drop the null observations in specific columns:

- gender
- relevent_experience
- enrolled_university
- education_level
- major_discipline
- experience
- lastnewjob

```
1 # dropping null observations in specified columns and filling null values with "Un
2 df_Clean = df.dropna(subset = ['gender', 'enrolled_university', 'education_level', 'm
3 df_Clean = df_Clean.fillna('Unknown')
4 df_Clean['company_size'].replace({'10/49': '10-49'}, inplace=True)
```

```

1
2 df_Clean.info()

<class 'pandas.core.frame.DataFrame'>
Int64Index: 12477 entries, 0 to 19155
Data columns (total 14 columns):
 #   Column                                Non-Null Count  Dtype
---  -
 0   enrollee_id                          12477 non-null  int64
 1   city                                 12477 non-null  object
 2   city_development_index               12477 non-null  float64
 3   gender                               12477 non-null  object
 4   relevent_experience                  12477 non-null  object
 5   enrolled_university                 12477 non-null  object
 6   education_level                      12477 non-null  object
 7   major_discipline                    12477 non-null  object
 8   experience                           12477 non-null  object
 9   company_type                         12477 non-null  object
10  last_new_job                         12477 non-null  object
11  training_hours                       12477 non-null  int64
12  target                               12477 non-null  float64
13  company_size                         12477 non-null  object
dtypes: float64(2), int64(2), object(10)
memory usage: 1.4+ MB

```

```

1 df_Clean.isnull().sum()

```

```

enrollee_id      0
city              0
city_development_index  0
gender            0
relevent_experience  0
enrolled_university  0
education_level   0
major_discipline  0
experience         0
company_type      0
last_new_job      0
training_hours    0
target            0
company_size      0
dtype: int64

```

```

1 # Changing name of column to make it more interpretable to reader of notebook; Gra
2 df_Clean['education_level'].replace({'Graduate':'Undergraduate'}, inplace=True)

```

▼ 2. Exploratory Data Analysis

▼ 1. Demographics:

1. Gender
2. City (city, CDI)
3. Education (major_discipline, enrolled_university, education_level)
4. Job history (experience, relevent_experience, last_new_job)

▼ 1.1 Gender

1.1.1 Is the hiring of data scientists gender biased? What is the impact of gender on attrition?

```
1 ## Function to create pie plot, takes input as list of features and dataframe
2 def pie_plt(category, dataframe):
3
4     for cat in category:
5         values_m = dataframe[cat].value_counts()
6         labels_m = values_m.index
7         plt.subplots(figsize = (9,7))
8
9         plt.pie(values_m, labels=labels_m, wedgeprops = { 'linewidth' : 3, 'edgecolor
10             ,explode=(0.1, 0.1, 0.1), autopct='%0.2f%%' )
11         plt.title('Hiring of Data Scientist by Gender at Company')
12         plt.show()

1 ## PLOtting a pie chart for hiring of DS by gender
2 pie_plt(['gender'], df_Clean)
```


Hiring of Data Scientist by Gender at Company

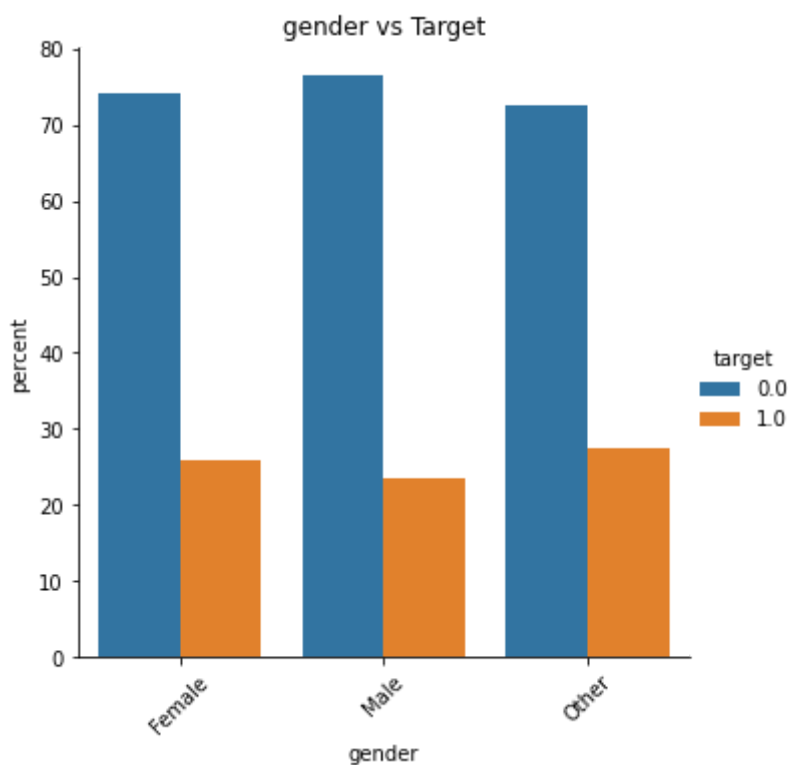
Approximately 90% of the hires are male, demonstrating a bias in the hiring of DS at this company. The industry average of male DS in USA is 65% according to the below source.

<https://www.zippia.com/data-scientist-jobs/demographics/>

```
1 ## This function takes in a list of features and plots a normalized count plot aga
2 def norm_cnt_plt(category, dataframe):
3     for cat in category:
4         plt.figure(figsize=(10,8))
5         x,y = cat, 'target'
6         df1 = dataframe.groupby(x)[y].value_counts(normalize=True) # Grouping by featu
7         df1 = df1.mul(100)
8         df1 = df1.rename('percent').reset_index()
9
10        ax = sns.catplot(x=x,y='percent',hue=y, kind='bar',data=df1) # Plotting the
11        plt.xticks(rotation= 45)
12        plt.title(cat + ' vs Target')
13        plt.show()
```

```
1 ## Checking the impact of attrition on gender using a normalized bar plot.
2 norm_cnt_plt(['gender'], df_Clean)
3
```

<Figure size 720x576 with 0 Axes>



From the above graph gender seems to have no impact on attrition.

▼ 1.2 City

CDI is ranked on a scale from 0 to 1.0, with 1.0 being the highest human development. HDI is broken down into four tiers: very high human development (0.8-1.0), high human development (0.7-0.79), medium human development (0.55-.70), and low human development (below 0.55).

```
1 ## Creating categorical values from CDI
2 cdi_bins = [0, 0.55, 0.70, 0.79, 1.0]
3 cdi_labels = ["low_human_development", "medium_human_development", "high_human_devel
4 df_Clean['cdi_bucket'] = pd.cut(df_Clean['city_development_index'], bins = cdi_bin
```

1.2.1 Which are the top 10 cities the company hires from? and their corresponding CDI.
(the higher the CDI, the more urban the city is)

```
1 df_Clean.groupby(['city', 'city_development_index', 'cdi_bucket'], as_index= False)[
2     .sort_values(by = 'enrollee_id', ascending = False).head(10).rese
```

	city	city_development_index	cdi_bucket	count
0	city_103	0.920	very_high_human_development	3262.0
1	city_21	0.624	medium_human_development	1480.0
2	city_16	0.910	very_high_human_development	1093.0
3	city_114	0.926	very_high_human_development	801.0
4	city_160	0.920	very_high_human_development	619.0
5	city_136	0.897	very_high_human_development	405.0
6	city_67	0.855	very_high_human_development	277.0
7	city_75	0.939	very_high_human_development	218.0
8	city_104	0.924	very_high_human_development	190.0
9	city_102	0.804	very_high_human_development	190.0

Company hires mainly from very high human development cities, with majority of the candidates coming from city_103. The only exception is city_21 which is the 2nd highest in terms of hiring but has medium CDI, it can be possible this is a University Town.

▼ 1.2.2 How is CDI correlated with an individual's education level?

```

1 ## Since CDI is a numeric variable and Education level is categorical, we use ANOV
2 df_cdi_edu = df_Clean[['education_level','city_development_index']]

1 # f_oneway() function takes the df_cdi_edu as input and returns F-statistic and P-
2 from scipy.stats import f_oneway
3
4 # Running the one-way anova test between City Development Index and Education Leve
5 # Assumption(H0) is that City Development Index and Education level are NOT correl
6
7 # Finds out the CDI data for each education level as a list
8 CategoryGroupLists = df_cdi_edu.groupby('education_level')['city_development_index
9
10 # Performing the ANOVA test
11 # We accept the Assumption(H0) only when P-Value > 0.05
12 AnovaResults = f_oneway(*CategoryGroupLists)
13 print('P-Value for Anova is: ', AnovaResults[1])

P-Value for Anova is:  2.8604593860651103e-18

```

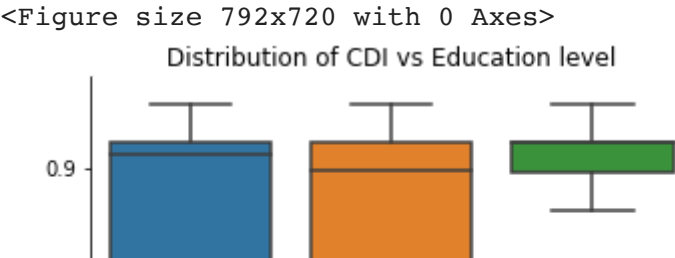
Since the P-value is less than 0.05, we can reject the null H0, and conclude there is a correlation between CDI and education level

Citation : <https://thinkingneuron.com/how-to-measure-the-correlation-between-a-numeric-and-a-categorical-variable-in-python/>

```

1 ## Plotting a Boxplot to check the distribution
2 plt.figure(figsize=(11,10))
3 sns.catplot(x="education_level", y="city_development_index", kind="box", data=df_c

```



In general, the candidates with higher qualifications belong to cities with higher CDI.

▼ 1.2.3 Relationship between CDI and city code

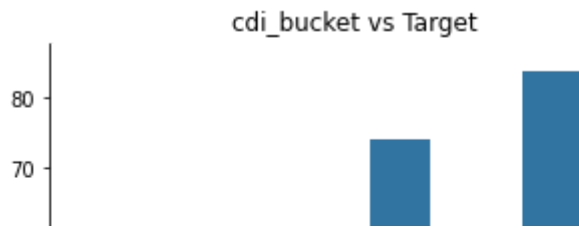
```
1 df_Clean.groupby('cdi_bucket', as_index=False)['city'].nunique().rename(columns={'
2                                     .sort_values(by = 'Number of cities', ascending = False).reset_ind
```

	cdi_bucket	Number of cities
0	very_high_human_development	47
1	medium_human_development	36
2	high_human_development	28
3	low_human_development	9

```
1 ### Impact of CDI on attrition using a normalized bar plot.
2 plt.figure(figsize=(13,10))
3 norm_cnt_plt(['cdi_bucket'], df_Clean)
4
```

<Figure size 936x720 with 0 Axes>

<Figure size 720x576 with 0 Axes>



Candidates from cities with lower and medium development index are more likely to look for a change.



▼ 1.3 Education



▼ 1.3.1 What are the top 5 education backgrounds for the data scientist (based on major discipline)

```
1 df_Clean.education_level.unique()

array(['Undergraduate', 'Masters', 'Phd'], dtype=object)

1 # import the circlify library
2 import circlify
3 values_e = df_Clean["education_level"].value_counts().sort_values(ascending=False)
4 labels_e = df_Clean["education_level"].value_counts().sort_values(ascending=True).
5 # compute circle positions:
6 circles = circlify.circlify(
7     values_e.tolist(),
8     show_enclosure=False,
9     target_enclosure=circlify.Circle(x=0, y=0, r=1)
10 )

1 # circular packing
2 # Create just a figure and only one subplot
3 fig, ax = plt.subplots(figsize=(10,10))
4
5 # Title
6 ax.set_title('Top Education Background')
7
8 # Remove axes
9 ax.axis('off')
10
11 # Find axis boundaries
12 lim = max(
13     max(
14         abs(circle.x) + circle.r,
```

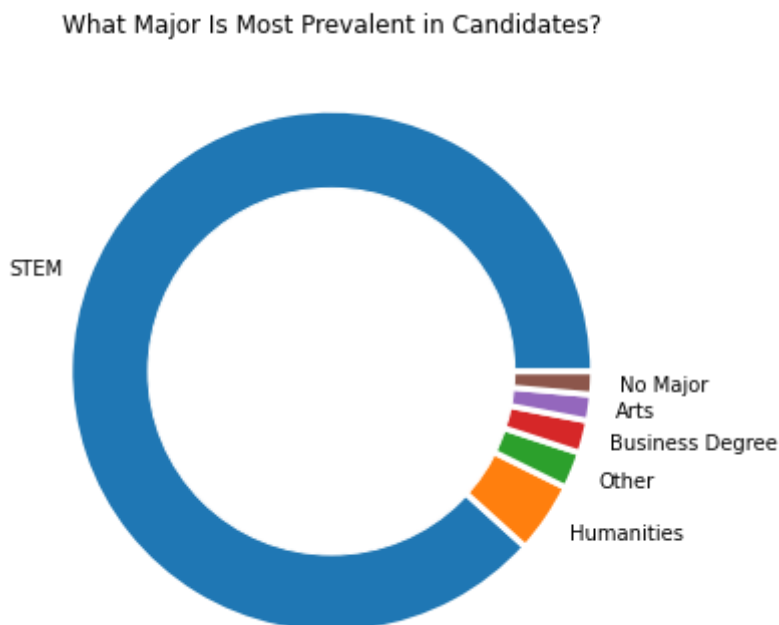
```
15         abs(circle.y) + circle.r,
16     )
17     for circle in circles
18 )
19 plt.xlim(-lim, lim)
20 plt.ylim(-lim, lim)
21
22 # list of labels
23 labels = labels_e
24
25 # print circles
26 for circle, label in zip(circles, labels):
27     x, y, r = circle
28     ax.add_patch(plt.Circle((x, y), r, alpha=0.2, color="orange", linewidth=2))
29     plt.annotate(
30         label,
31         (x,y) ,
32         va='center',
33         ha='center'
34     )
35
```

Top Education Background

The candidates mostly have undergraduate education background, and small portion of the candidates has Phd.

- Citation: <https://www.python-graph-gallery.com/circular-packing-several-levels-of-hierarchy>

```
1 # donut chart
2 # create data
3 plt.figure(figsize=(8,6))
4 values_m = df_Clean["major_discipline"].value_counts()
5 labels_m = values_m.index
6
7 # Create a circle at the center of the plot
8 my_circle = plt.Circle( (0,0), 0.7, color='white')
9
10 # Custom wedges
11 plt.pie(values_m, labels=labels_m, wedgeprops = { 'linewidth' : 3, 'edgecolor' : 'white' })
12 p = plt.gcf()
13 p.gca().add_artist(my_circle)
14 plt.title("What Major Is Most Prevalent in Candidates?")
15 plt.show()
```

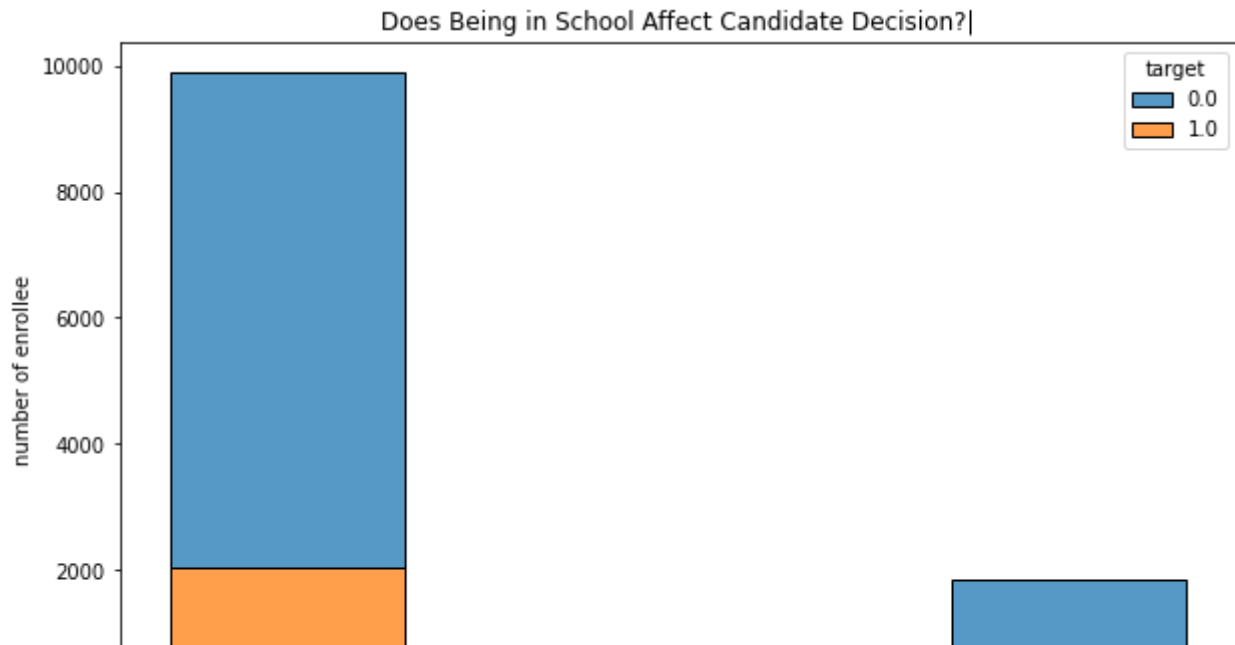


Most of the enrollees major in STEM, which is reasonable because the company is hiring data scientists.

▼ 1.3.2 What are some education characteristics for those candidates who are staying?

```
1 df_education = df_Clean[["target", "enrolled_university", "major_discipline"]]

1 for i in df_education.drop("target", axis=1):
2     plt.figure(figsize=[10, 6])
3     hue_colors = {0: "red", 1: "black"}
4     plots = sns.histplot(data = df_Clean, x=i, hue="target", multiple="stack", shrink
5     # Iterrating over the bars one-by-one
6     # for bar in plots.patches:
7     #     # Using Matplotlib's annotate function and
8     #     # passing the coordinates where the annotation shall be done
9     #     # x-coordinate: bar.get_x() + bar.get_width() / 2
10    #     # y-coordinate: bar.get_height()
11    #     # free space to be left to make graph pleasing: (0, 8)
12    #     # ha and va stand for the horizontal and vertical alignment
13    #     plots.annotate(format(bar.get_height(), '.2f'),
14    #                     (bar.get_x() + bar.get_width() / 2,
15    #                      bar.get_height()), ha='center', va='center',
16    #                      size=10, xytext=(0, 8),
17    #                      textcoords='offset points')
18
19    plt.xlabel(i)
20    plt.ylabel("number of enrollee")
21    plt.title("Does Being in School Affect Candidate Decision?|")
22    plt.show()
```

```
1 test1 = df_Clean.groupby(["education_level", "target"])["target"].count()
```

```
1 test1
```

education_level	target	count
Masters	0.0	2643
	1.0	626
Phd	0.0	282
	1.0	43
Undergraduate	0.0	6586
	1.0	2297

Name: target, dtype: int64

```
1 test1
```

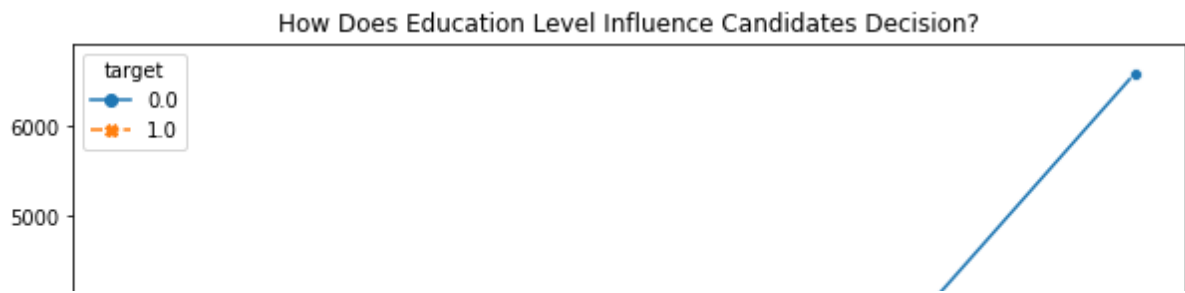
```
1 test1 = test1.reset_index(name='count')
```

```
1 test1
```

```
1 plt.figure(figsize=[10,6])
```

```
2 sns.lineplot(data=test1,x="education_level", y = "count", hue="target", style="tar
```

```
3 plt.title("How Does Education Level Influence Candidates Decision?");
```



- Master candidates who decide to stay have larger portion.
- Candidates who have no enrollment or enroll in part-time course tend to stay in the company.
- STEM majored candidates have higher chance to stay.

▼ 1.4 Job History

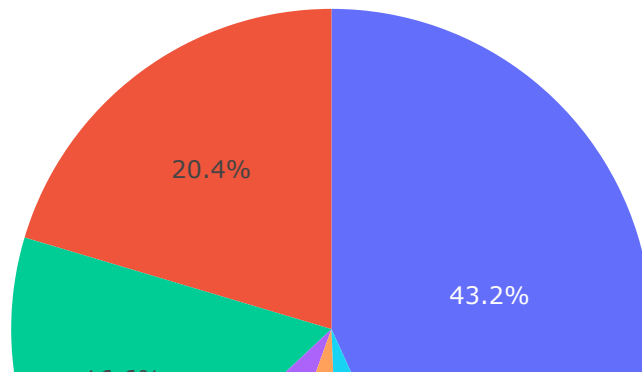
▼ 1.4.1 No Gap/Yes Gap: will the gap years of jobs affect whether a candidate is staying or leaving?

```

1 # First see the portion of the duration of job gap
2 ep = df_Clean['last_new_job'].value_counts().reset_index()
3 ep.columns = [
4     'last_new_job',
5     'percent'
6 ]
7 ep['percent'] /= len(df)
8 fig = px.pie(
9     ep,
10    names='last_new_job',
11    values='percent',
12    title='Difference in Years Between Previous Job and Current Job',
13    width=800,
14    height=500
15 )
16 fig.show()

```

Difference in Years Between Previous Job and Current Job



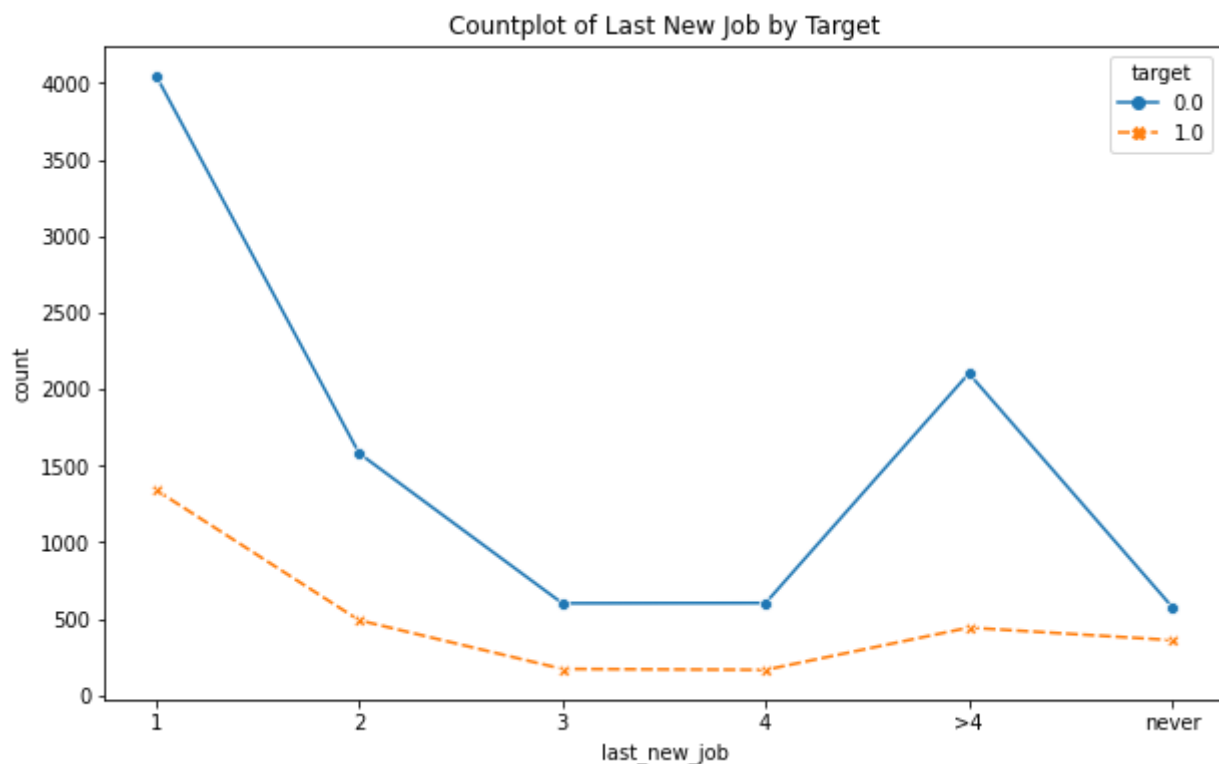
```
1 last_graph = df_Clean.groupby(["last_new_job", "target"])["target"].count()
```



```
1 last_graph = last_graph.reset_index(name='count')
```



```
1 plt.figure(figsize=[10,6])
2 sns.lineplot(data=last_graph,x="last_new_job", y = "count", hue="target", style="t
3 plt.title("Countplot of Last New Job by Target");
```



According to the above plot, we can easily tell that when the gap year(s) between last and new is(are) 1 or 2, the ratio between leave and stay is 1:3. When the gap duration increases to 3, 4 or

even more than 4, the ratio of leave and stay decreases to 1:4. This means if the last job and new job gap turns longer, more candidates tend to stay in this Big Data company correspondingly.

However, when we come to candidates never change their job -- the 'never' group, we can see that the ratio between leave and stay are 1:2, which means least candidates in this group will stay in this company rather than their original ones.

This will give HR a hint: When recruiting candidates to join this program, we can turn to candidates with longer gap between last job and new one. However, not the ones who never change their job,

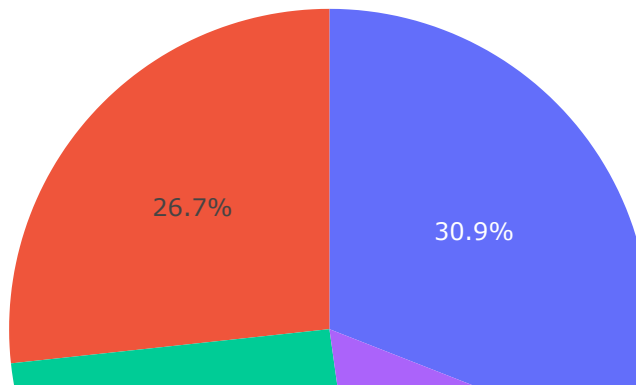
▼ 1.4.2 Is the hiring of Data Scientists impacted by their previous experience?

```

1 # Set bins for experience to make the graph clearer
2 df_Clean['cat_experience'] = np.where(df_Clean['experience'].isin(['<1','1','2','3
3 df_Clean['cat_experience'] = np.where(df_Clean['experience'].isin(['6','7','8','9'
4 df_Clean['cat_experience'] = np.where(df_Clean['experience'].isin(['11','12','13',
5 df_Clean['cat_experience'] = np.where(df_Clean['experience'].isin(['16','17','18',

1 # First see the portion of experience
2 ep = df_Clean['cat_experience'].value_counts().reset_index()
3 ep.columns = [
4     'cat_experience',
5     'percent'
6 ]
7 ep['percent'] /= len(df_Clean)
8
9 fig = px.pie(
10     ep,
11     names='cat_experience',
12     values='percent',
13     title='What Level of Experience do Most Candidates Have?',
14     width=800,
15     height=500
16 )
17
18 fig.show()
```

What Level of Experience do Most Candidates Have?



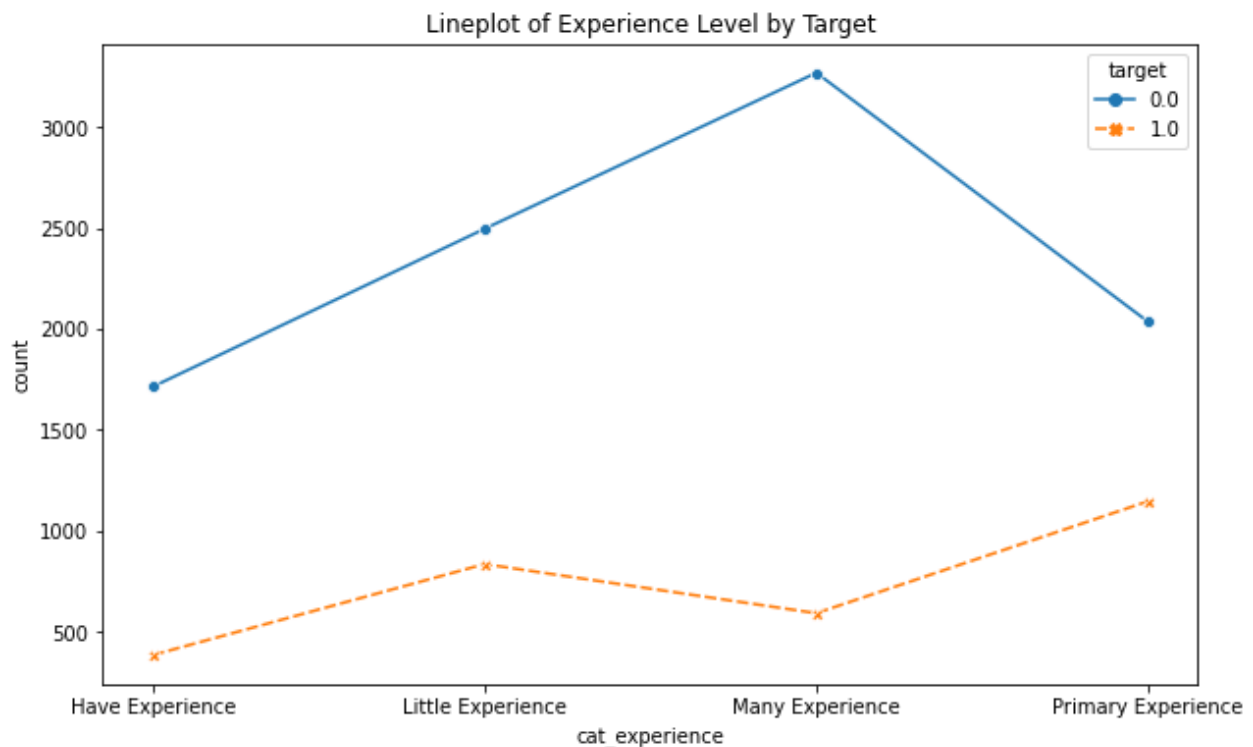
```
1 exp_graph = df_Clean.groupby(["cat_experience", "target"])["target"].count()
```



```
1 exp_graph = exp_graph.reset_index(name='count')
```



```
1 plt.figure(figsize=[10,6])
2 sns.lineplot(data=exp_graph,x="cat_experience", y = "count", hue="target", style="
3 plt.title("Lineplot of Experience Level by Target");
```



From the graph, we can easily tell that in the group with more experience candidates tend to stay in Big Data company after training. And this trend is monotonical.

This gives HR the hint -- Candidates with more experience joining this program will tend to stay.

▼ 1.4.3 Does the type of the candidate's current company affect his decision?

```
1 com_graph = df_Clean.groupby(["company_type", "target"])["target"].count()
```

```
1 com_graph
```

company_type	target	
Early Stage Startup	0.0	307
	1.0	79
Funded Startup	0.0	677
	1.0	108
NGO	0.0	313
	1.0	59
Other	0.0	60
	1.0	17
Public Sector	0.0	518
	1.0	128
Pvt Ltd	0.0	5789
	1.0	1183
Unknown	0.0	1847
	1.0	1392

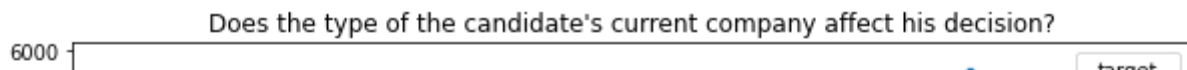
Name: target, dtype: int64

```
1 com_graph = com_graph.reset_index(name='count')
```

```
1 plt.figure(figsize=[10,6])
```

```
2 sns.lineplot(data=com_graph,x="company_type", y = "count", hue="target", style="ta
```

```
3 plt.title("Does the type of the candidate's current company affect his decision?")
```



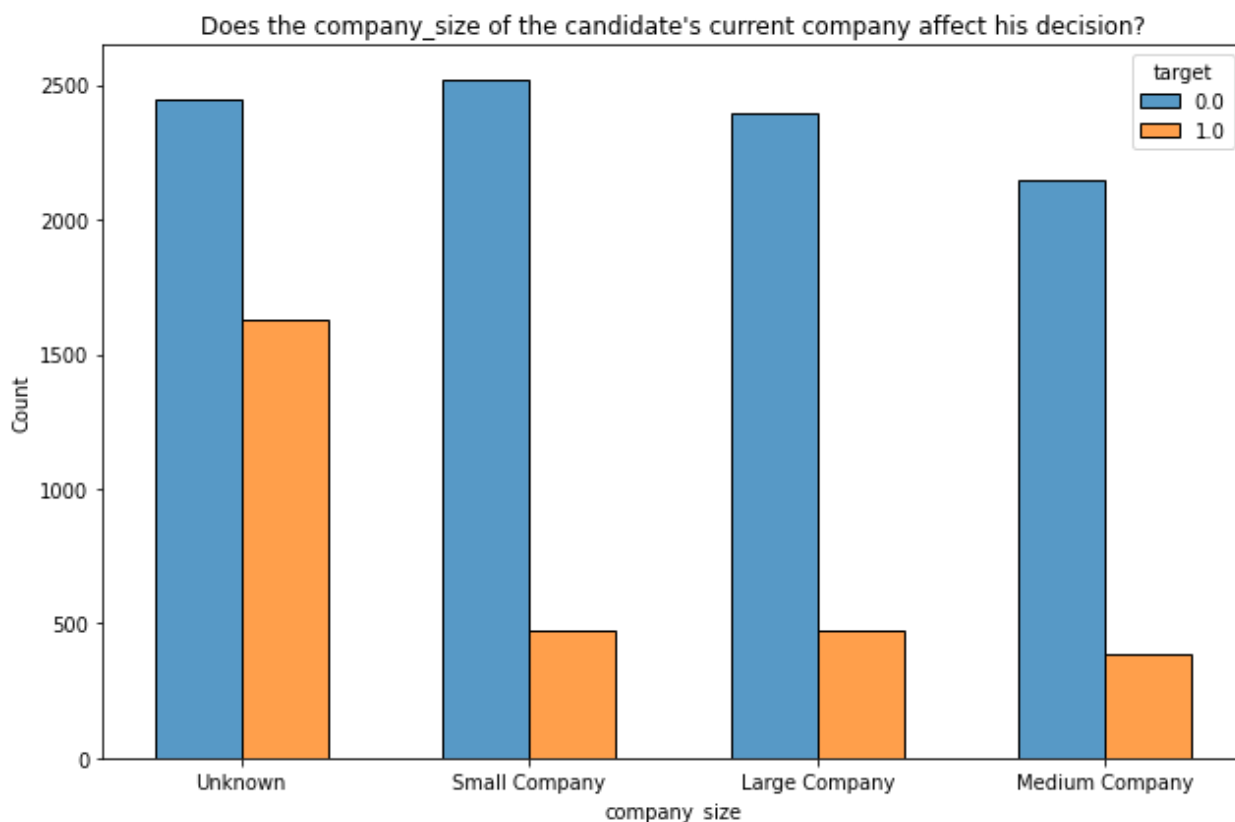
It seems like candidates from private companies tend to stay with the "company" after training. However, 75% of the candidates come from private companies. We need to further investigate into why most data scientist candidates come from a private company. Is there a reason? Perhaps "company" should focus its efforts into targeting only that company type. We could isolate those who come from pvt ltd.



▼ 1.4.4 Does the size of the candidate's current company affect his decision?



```
1 plt.figure(figsize=[9,6])
2 sns.histplot(data = df_Clean,x="company_size",hue="target", multiple="dodge", shri
3 plt.title("Does the {} of the candidate's current company affect his decision?".fo
4
5 plt.tight_layout()
6 plt.show()
```

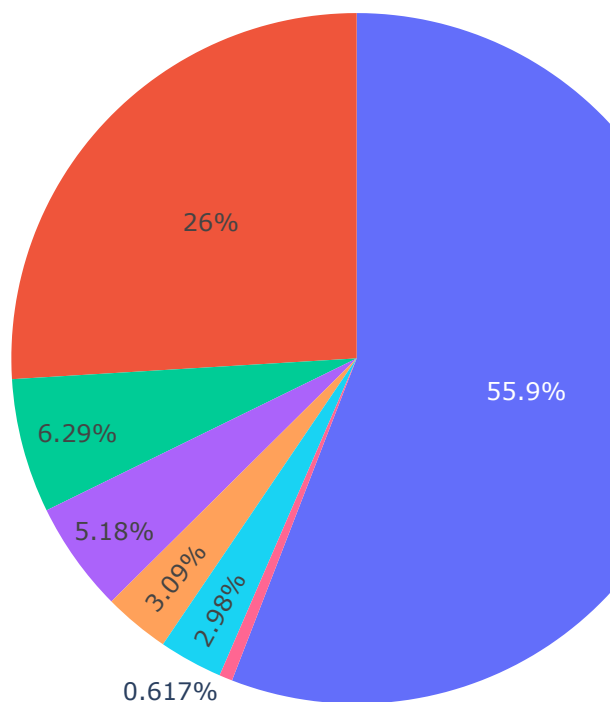


When comparing the sizes of the candidates' current companies, the candidates coming from small companies tend to stay with the "company" the most. However, there is not a clear relationship between the company size and the candidate's decision to stay after training.

▼ 1.4.5 What type of company do most company hires come from?

```
1 fig = px.pie(df_Clean['company_type'].value_counts(), values='company_type',
2             names = df_Clean['company_type'].value_counts().index,title = 'What i
3             )
4 fig.show()
```

What is the most prevalent company type in the dataset?

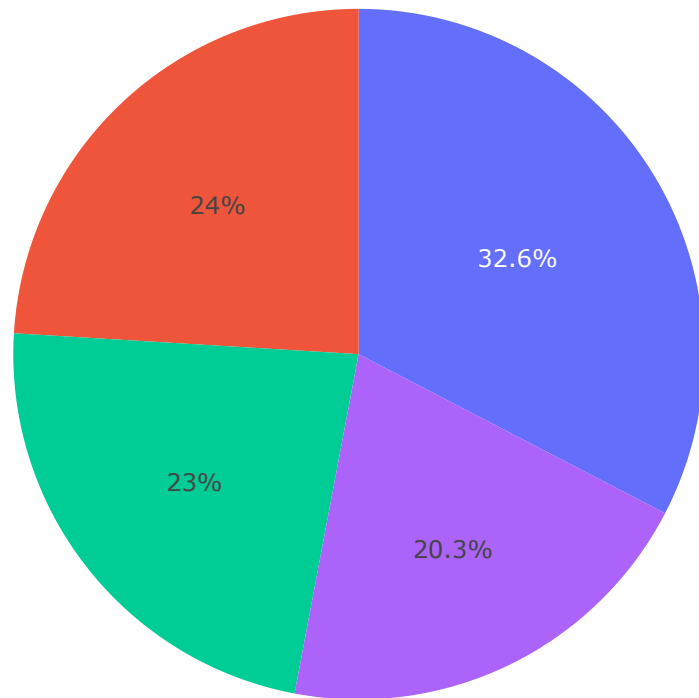


Most candidates come from private limited companies.

▼ 1.4.6 What company size do most candidates come from?

```
1 fig = px.pie(df_Clean['company_size'].value_counts(), values='company_size',
2             names = df_Clean['company_size'].value_counts().index,title = 'What i
3             )
4 fig.show()
```


What is the most prevalent company size in the dataset?

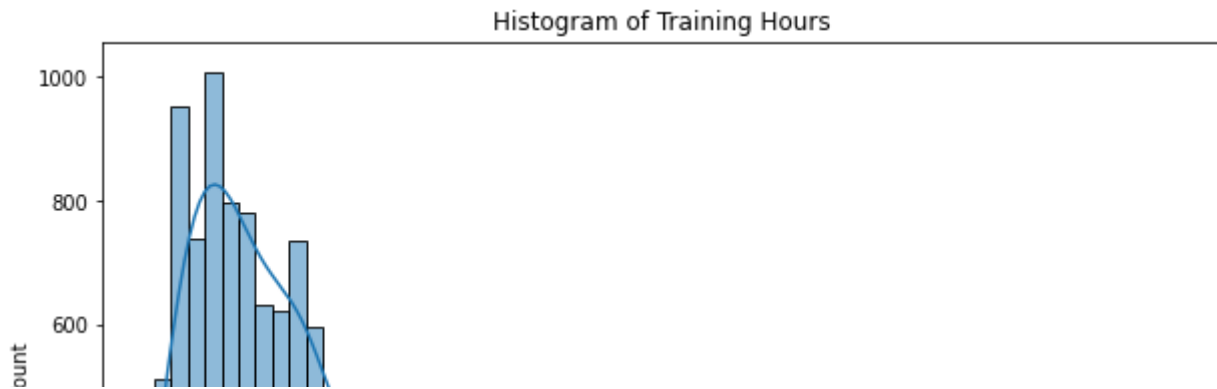


Not a clear winner; hires come from different size companies.

▼ 2. Engagement and Retention

▼ 2.1 How much training hours does the company invest in its future employees?

```
1 # We can plot a histogram of training hours
2 plt.figure(figsize=[10,6])
3 sns.histplot(data=df_Clean, x="training_hours", kde=True)
4 plt.title("Histogram of Training Hours");
```



```
1 df_Clean['training_hours'].median()
```

47.0



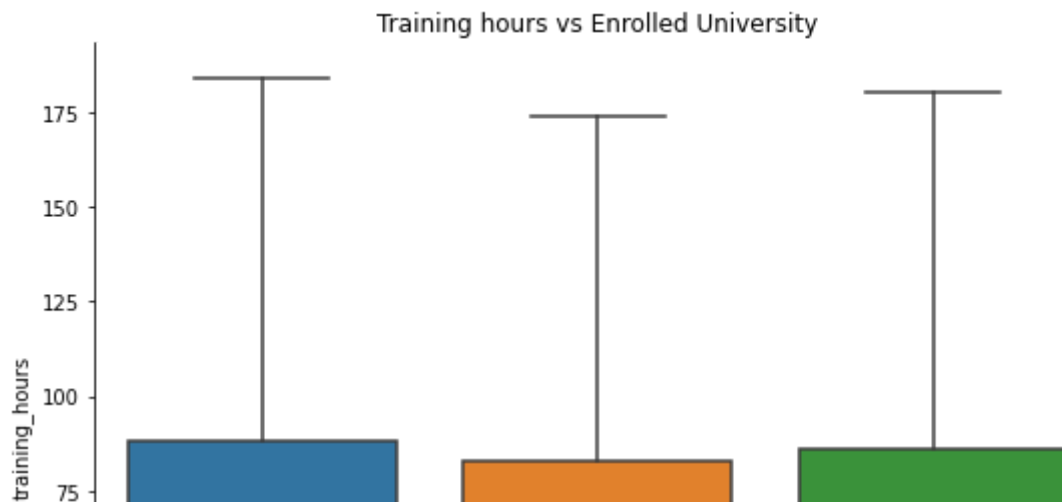
```
1 df_Clean['training_hours'].describe()
```

```
count      12477.000000
mean         64.927306
std          59.732622
min           1.000000
25%          23.000000
50%          47.000000
75%          88.000000
max         336.000000
Name: training_hours, dtype: float64
```

Distribution is right skewed. On average, the company spends approximately 47 hours training its future employees.

▼ 2.1.1 Will candidates enrolled in university experience longer or shorter training hours?

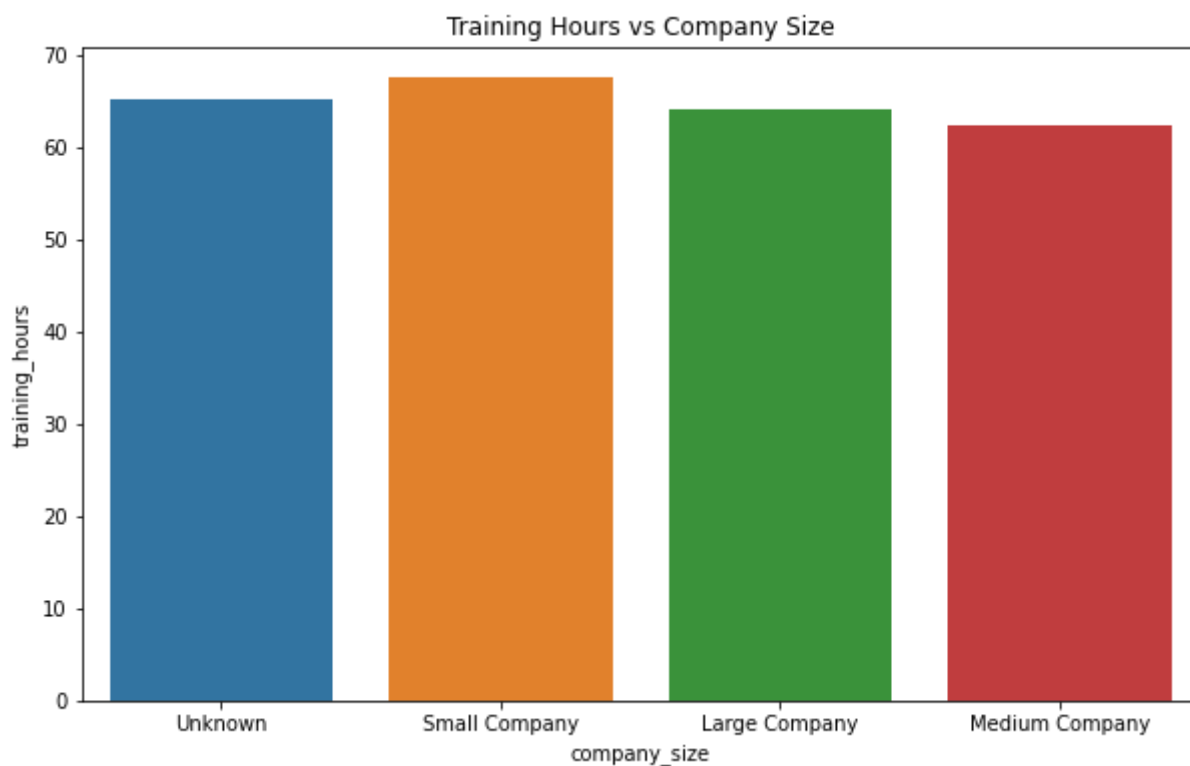
```
1 plt.figure(figsize=[9,7])
2 sns.boxplot(x="enrolled_university", y="training_hours",
3             data= df_Clean, showfliers = False)
4 sns.despine()
5 plt.title("Training hours vs Enrolled University");
```



The distribution is quite similar for all three enrolled university categories. However, the no enrollment category seems to have slightly higher training hours compared to those of the other two. This makes sense as people who are not enrolled in university may have more time for training.

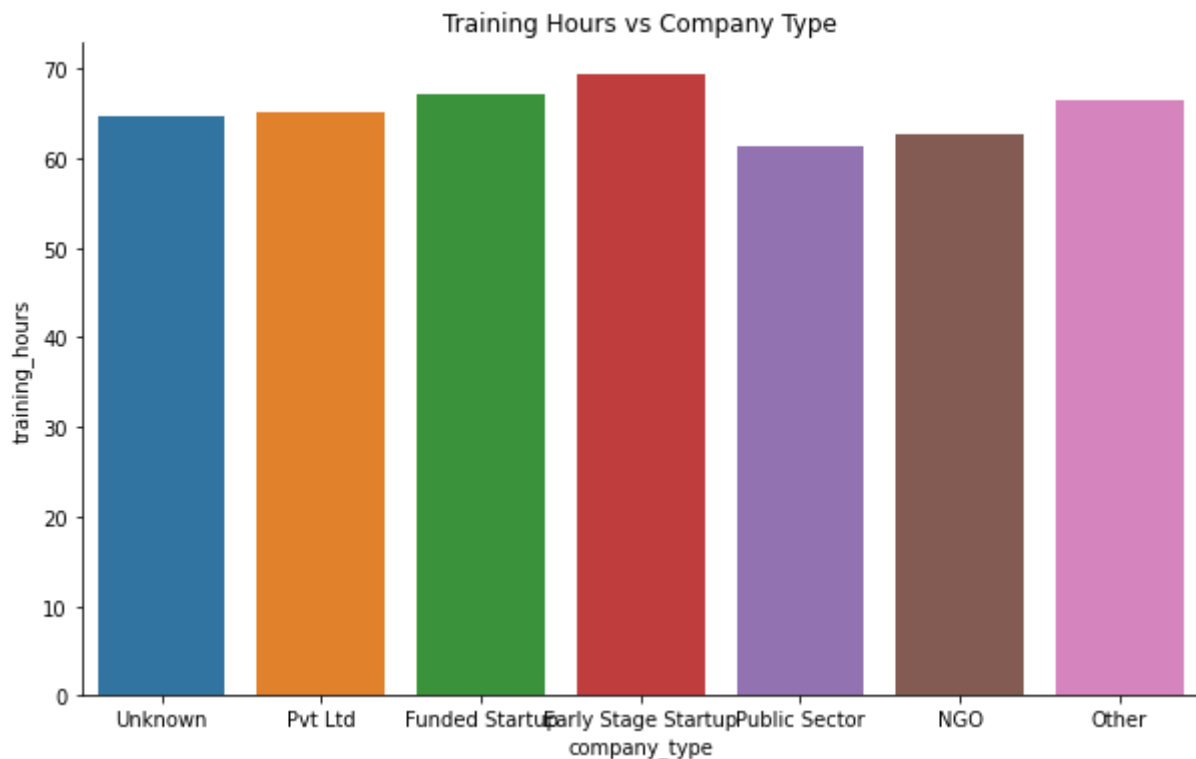
▼ 2.1.2 Does company size/type affect the number of training hours?

```
1 plt.figure(figsize=[10,6])
2
3 sns.barplot(x="company_size", y="training_hours", data=df_Clean, ci=None)
4 plt.title("Training Hours vs Company Size");
```



Small company tends to have higher training hours compared to the others, which makes sense considering the fact that people from smaller companies may not have as much experience compared to people from larger and more established companies, and thus would required more training hours.

```
1 plt.figure(figsize=[10,6])
2
3 sns.barplot(x="company_type", y="training_hours",
4             data= df_Clean, ci = None)
5 sns.despine()
6 plt.title("Training Hours vs Company Type");
```



Startups have highest training hours, followed by Pvt Ltd and NGO, and public sector has lowest training hours. People from startups may need more training as they may not have been exposed to the conduct of established companies.

2.1.3 How long should the candidates get trained if the STEM candidates have relevant experience?

```
1 # Set new conditions and values to match & save to a new column
2 condition = [
3     (df_Clean['major_discipline'] == 'STEM') & (df_Clean['relevent_experi
4     (df_Clean['major_discipline'] == 'STEM') & (df_Clean['relevent_experi
```

```

5         (df_Clean['major_discipline'] != 'STEM') & (df_Clean['relevent_experi
6         (df_Clean['major_discipline'] != 'STEM') & (df_Clean['relevent_experi
7 ]

```

```
1 values = ['STEM_rel', 'STEM_nonrel', "Non_STEM_rel", "Non_STEM_nonrel"]

```

```

1 df_Clean['educational_condition'] = np.select(condition, values)
2 df_Clean.head()

```

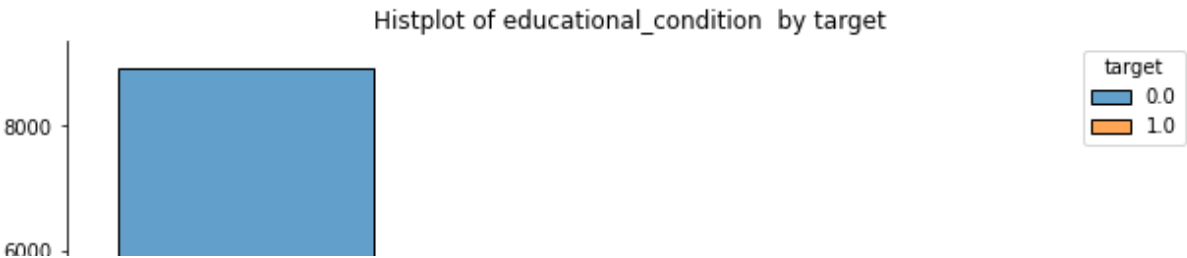
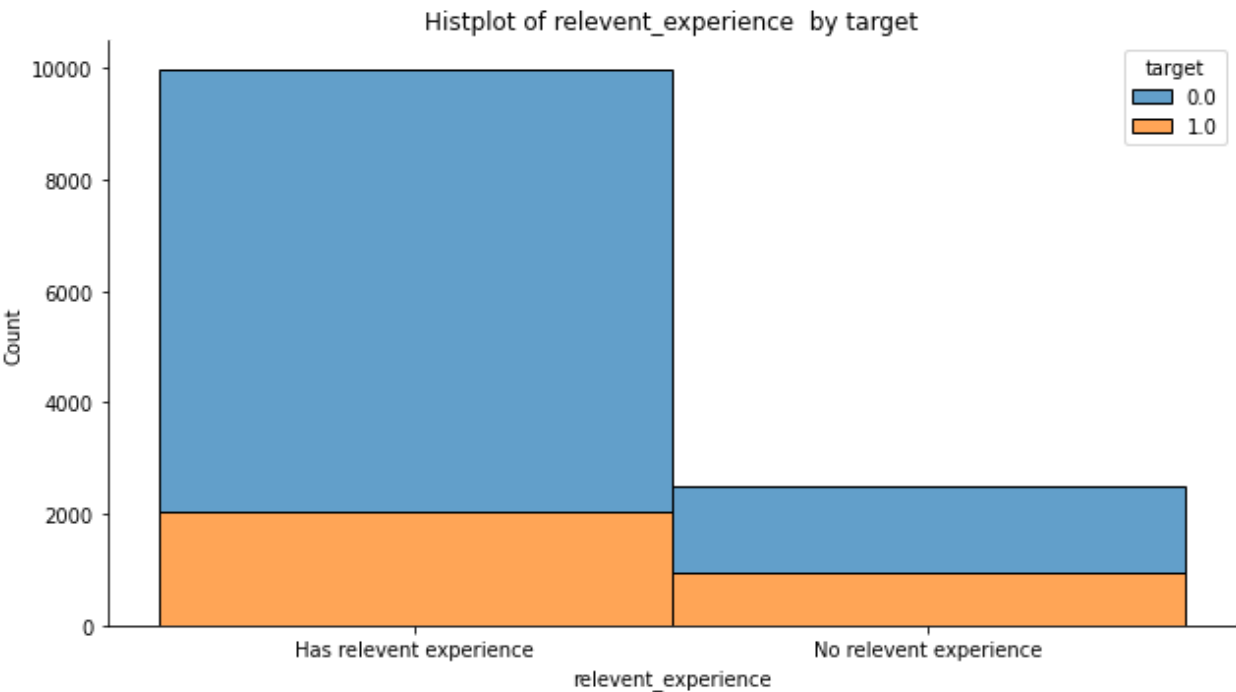
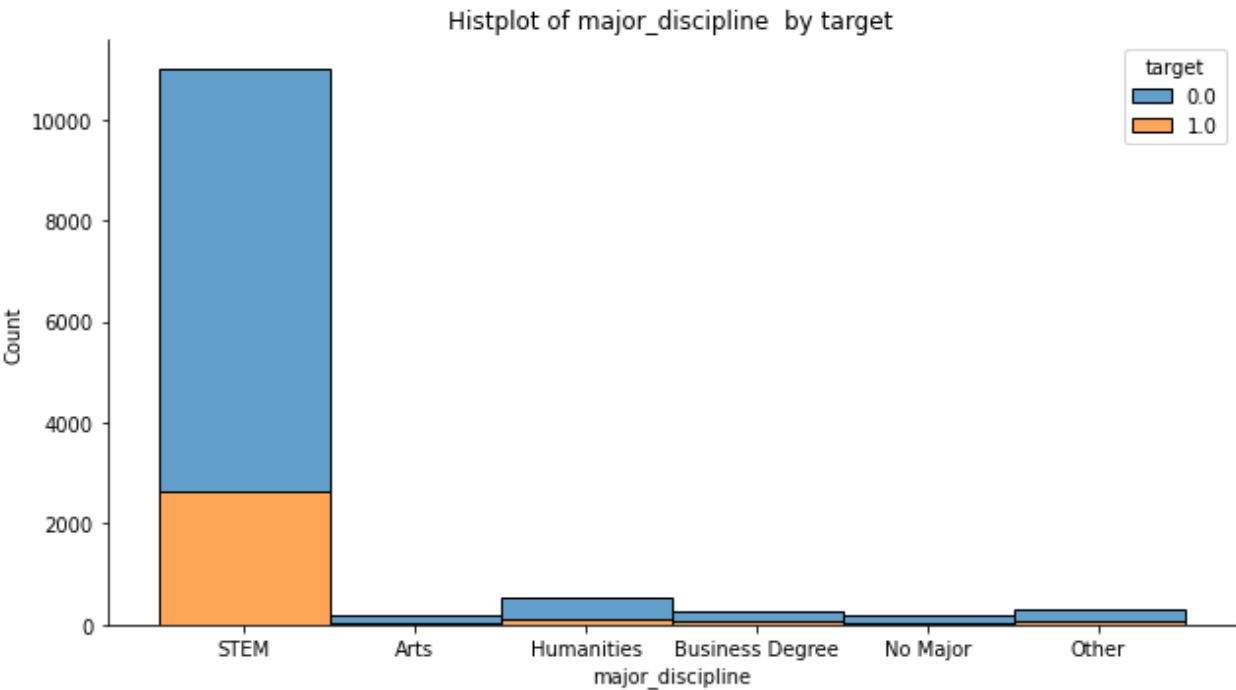
	enrollee_id	city	city_development_index	gender	relevent_experience
0	8949	city_103	0.920	Male	Has relevent experience
1	29725	city_40	0.776	Male	No relevent experience
4	666	city_162	0.767	Male	Has relevent experience
7	402	city_46	0.762	Male	Has relevent experience
8	27107	city_103	0.920	Male	Has relevent experience

▼ 2.1.3.1 Educational Condition by target

```

1 # Plot the educational condition by target to see the different
2 plt.figure(figsize=[9,15])
3 plot=["major_discipline", "relevent_experience", "educational_condition"]
4 n=1
5 for f in plot:
6     plt.subplot(3,1,n)
7     sns.histplot(x=f, hue='target', edgecolor="black", multiple="stack", alpha=0.7
8     sns.despine()
9     plt.title("Histplot of {} by target".format(f))
10    n=n+1
11 plt.tight_layout()
12 plt.show()

```



Through this plot, it is undeniable that candidates with relevant experiecnce tend to stay in this company, especially when we compare 'relevent' with 'non_relevent' groups.

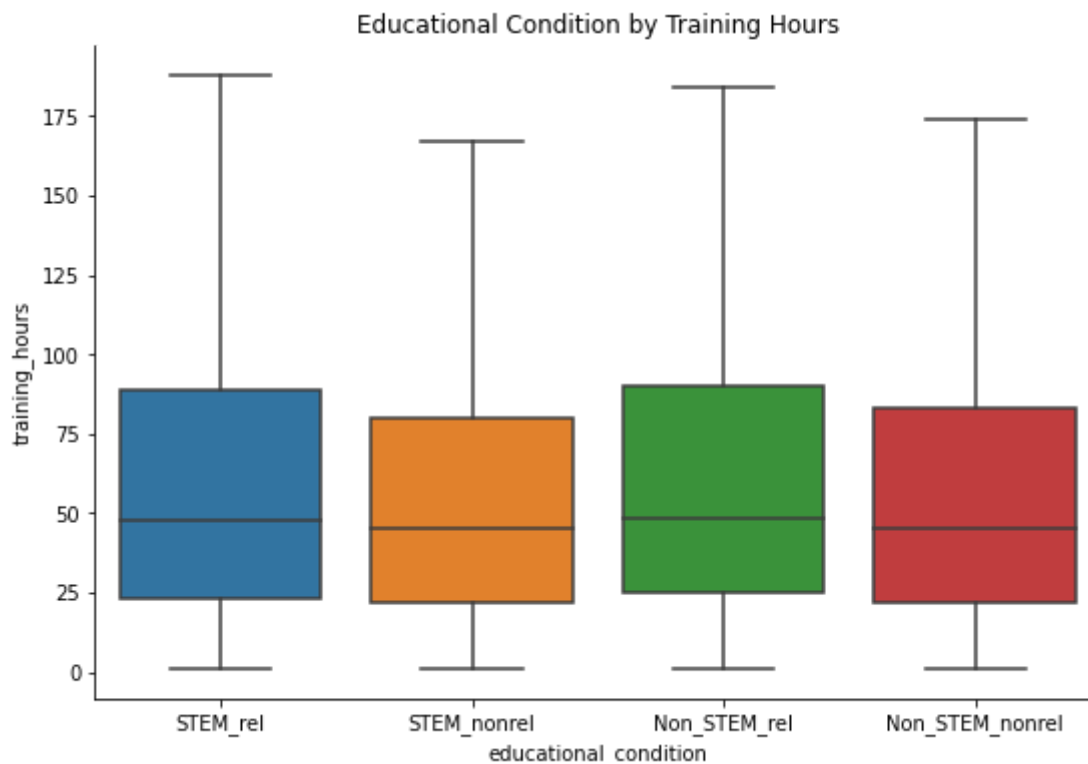
▼ 2.1.3.1 Educational Condition by training hours



```

1 plt.figure(figsize=[9,6])
2 sns.boxplot(x="educational_condition", y="training_hours",
3             data=df_Clean, showfliers = False)
4 #plt.legend(loc='upper right', title='Target')
5 sns.despine()
6 plt.title("Educational Condition by Training Hours");

```



We are guessing that maybe candidates who are majored in STEM and have relevant experience need less training hours compared with other groups.

However, apparently the result is not like our guessing – Training hours does not rely on the candidates' educational condition.

2.2 Is it a fact that: the longer the candidate is being trained, the higher chance he/she will stay? How should we adjust the training hours as a HR?

```

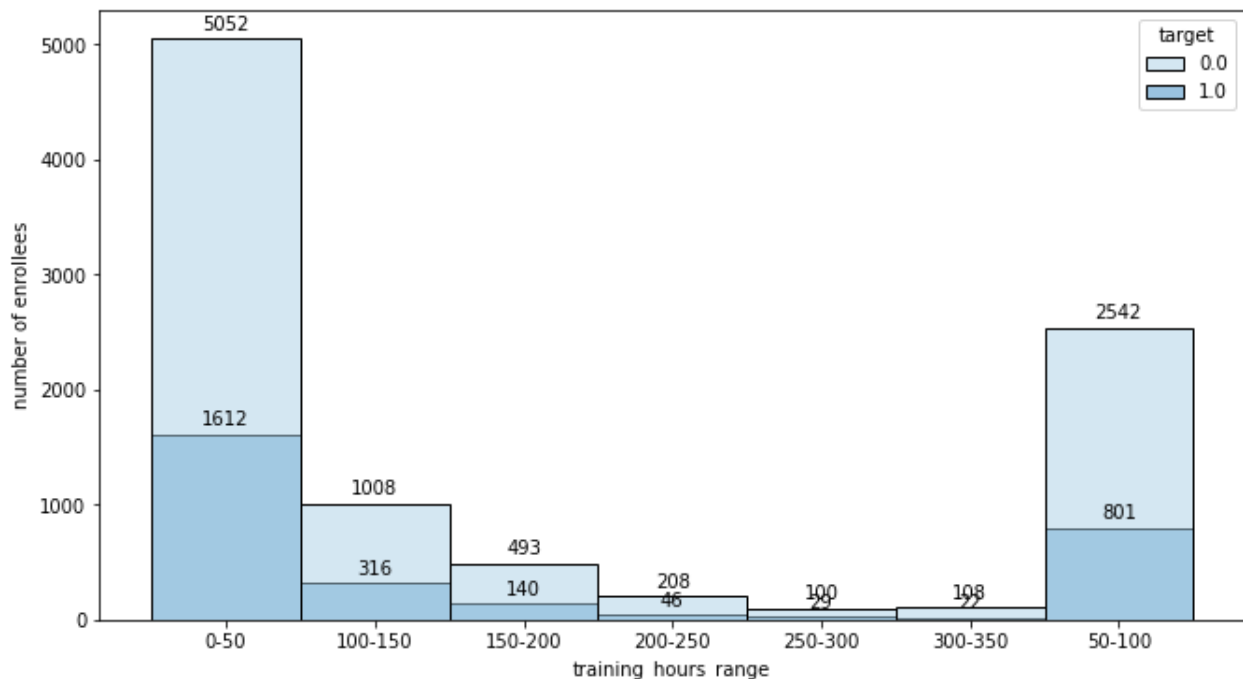
1 df_Clean['training_hours_range'] = np.where(df_Clean['training_hours'].isin(range(
2 df_Clean['training_hours_range'] = np.where(df_Clean['training_hours'].isin(range(
3 df_Clean['training_hours_range'] = np.where(df_Clean['training_hours'].isin(range(
4 df_Clean['training_hours_range'] = np.where(df_Clean['training_hours'].isin(range(
5 df_Clean['training_hours_range'] = np.where(df_Clean['training_hours'].isin(range(
6 df_Clean['training_hours_range'] = np.where(df_Clean['training_hours'].isin(range(
7 df_Clean['training_hours_range'] = np.where(df_Clean['training_hours'].isin(range(
8 df_Clean=df_Clean.sort_values('training_hours_range', ascending=True)

```

```

1 # How much training hours does the company invest in its future employees? we can
2 plt.figure(figsize=[11,6])
3
4 plots_thrs = sns.histplot(data=df_Clean, x='training_hours_range', hue="target", p
5
6 # Iterrating over the bars one-by-one
7 for bar in plots_thrs.patches:
8     # Using Matplotlib's annotate function and
9     # passing the coordinates where the annotation shall be done
10    # x-coordinate: bar.get_x() + bar.get_width() / 2
11    # y-coordinate: bar.get_height()
12    # free space to be left to make graph pleasing: (0, 8)
13    # ha and va stand for the horizontal and vertical alignment
14    plots_thrs.annotate(format(bar.get_height(), 'd'),
15                        (bar.get_x() + bar.get_width() / 2,
16                         bar.get_height()), ha='center', va='center',
17                        size=10, xytext=(0, 8),
18                        textcoords='offset points')
19
20 plt.ylabel("number of enrollees")
21 plt.show()

```



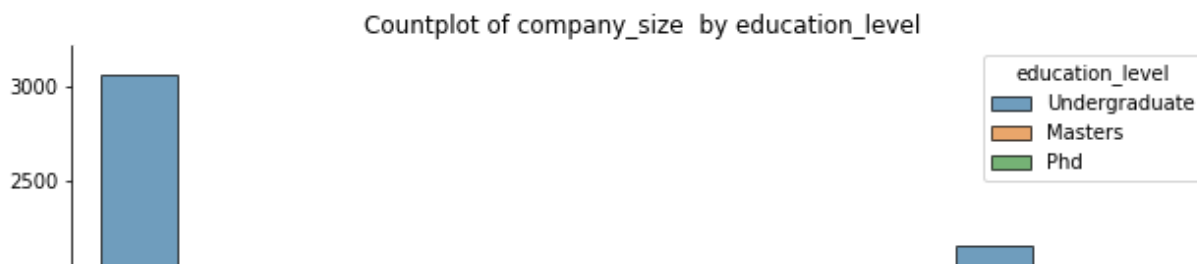
Based on the hisplot for training hours, we can observe that it's not necessary that the more time the candidate is being trained, the higher chance he/she will stay. From a HR perspective, 50-150hrs training hour range has higher percentage of retention about 76%. It's the ideal range for the employees.

- 100-150:76.1%

- 50-100: 76.04%

2.3 Do people with higher education levels affect the company size and type they work in?

```
1 plt.figure(figsize=[9,15])
2 categories=["company_size", "company_type"]
3 n=1
4 for f in categories:
5     plt.subplot(3,1,n)
6     sns.countplot(x=f, hue='education_level', edgecolor="black", alpha=0.7, data=d)
7     sns.despine()
8     plt.title("Countplot of {} by education_level".format(f))
9     n=n+1
10 plt.tight_layout()
11 plt.show()
```



- People with different education levels have no difference in choosing company size.
- Undergraduates candidates favor private company the most.



3. Data Preparation for Machine Learning



```
1 from sklearn.neural_network import MLPClassifier
2 from sklearn.datasets import make_classification
3 from sklearn.model_selection import train_test_split
4 import pandas as pd
5 import sklearn.metrics as metrics
6 from sklearn.ensemble import RandomForestClassifier
7 from xgboost import XGBClassifier
8 from sklearn.metrics import confusion_matrix
9 from sklearn.metrics import classification_report
10 from sklearn.datasets import make_classification
```

```
3000 {
```

```
1 df_Clean.head()
```

	enrollee_id	city	city_development_index	gender	relevent_experience
0	8949	city_103	0.920	Male	Has relevent experience
1	29725	city_40	0.776	Male	No relevent experience
4	666	city_162	0.767	Male	Has relevent experience
7	402	city_46	0.762	Male	Has relevent experience
8	27107	city_103	0.920	Male	Has relevent experience

```
1 # Bucketing the company_size values into categorizes
2 df_Clean['size'] = np.where(df_Clean['company_size'].isin(['50-99','10-49','<10']))
3 df_Clean['size'] = np.where(df_Clean['company_size'].isin(['100-500','500-999']),
4 df_Clean['size'] = np.where(df_Clean['company_size'].isin(['10000+','1000-4999'],'5
5 sorted_counts=df_Clean['size'].value_counts())
```

```
1 df_Clean.drop(['company_size'], axis=1, inplace=True)
2 df_Clean.rename(columns={'size': 'company_size'}, inplace=True)
```

```
1 df_Clean.head()
```

	enrollee_id	city	city_development_index	gender	relevent_experience
0	8949	city_103	0.920	Male	Has relevent experience
1	29725	city_40	0.776	Male	No relevent experience
4	666	city_162	0.767	Male	Has relevent experience
7	402	city_46	0.762	Male	Has relevent experience
8	27107	city_103	0.920	Male	Has relevent experience

```
1 df["company_type"].unique()
```

```
array([nan, 'Pvt Ltd', 'Funded Startup', 'Early Stage Startup', 'Other',
       'Public Sector', 'NGO'], dtype=object)
```

```
1 df_Clean.drop(['company_size'], axis=1, inplace=True)
```

```
2 df_Clean.rename(columns={'size': 'company_size'}, inplace=True)
```

```
1 df_Clean['company_type'].isin(['Funded Startup','Early Stage Startup'])
```

```
0      False
```

```
1      False
```

```
4       True
```

```
7      False
```

```
8      False
```

```
...
```

```
19150  False
```

```
19152   True
```

```
19153  False
```

```
19154  False
```

```
19155  False
```

```
Name: company_type, Length: 12477, dtype: bool
```

```
1 df_Clean['type'] = np.where(df_Clean['company_type'].isin(['Funded Startup','Early
```

```
1 df_Clean.info()
```

```
<class 'pandas.core.frame.DataFrame'>
```

```
Int64Index: 12477 entries, 0 to 19155
```

```
Data columns (total 14 columns):
```

#	Column	Non-Null Count	Dtype
0	enrollee_id	12477 non-null	int64
1	city	12477 non-null	object
2	city_development_index	12477 non-null	float64

```

3   gender                12477 non-null object
4   relevent_experience    12477 non-null object
5   enrolled_university    12477 non-null object
6   education_level        12477 non-null object
7   major_discipline       12477 non-null object
8   experience             12477 non-null object
9   company_type           12477 non-null object
10  last_new_job           12477 non-null object
11  training_hours         12477 non-null int64
12  target                 12477 non-null float64
13  type                   12477 non-null object
dtypes: float64(2), int64(2), object(10)
memory usage: 1.4+ MB

```

```

1 df_Clean.drop(['company_type'], axis=1, inplace=True)
2 df_Clean.rename(columns={'type': 'company_type'}, inplace=True)

```

```
1 df_Clean.head()
```

	enrollee_id	city	city_development_index	gender	relevent_experience
0	8949	city_103	0.920	Male	Has relevent experience
1	29725	city_40	0.776	Male	No relevent experience
4	666	city_162	0.767	Male	Has relevent experience
7	402	city_46	0.762	Male	Has relevent experience
8	27107	city_103	0.920	Male	Has relevent experience

```
1 df_Clean["last_new_job"].unique()
```

```
array(['1', '>4', '4', '3', '2', 'never'], dtype=object)
```

```

1 df_Clean["last_new_job"] = df_Clean["last_new_job"].replace(to_replace = "never",
2                       value = 0)
3 df_Clean["last_new_job"] = df_Clean["last_new_job"].replace(to_replace = ">4",
4                       value = 5)
5 df_Clean["last_new_job"] = df_Clean["last_new_job"].astype("int64")

```

```
1 df_Clean["last_new_job"].unique()
```

```
array([1, 5, 4, 3, 2, 0])
```

```
1 df_Clean.head()
```

	enrollee_id	city	city_development_index	gender	relevent_experience
0	8949	city_103	0.920	Male	Has relevent experi
10803	24658	city_103	0.920	Male	Has relevent experi
10797	17847	city_21	0.624	Male	Has relevent experi
10796	1122	city_46	0.762	Male	Has relevent experi
10794	24653	city_160	0.920	Male	Has relevent experi

```
1 df_Clean["experience"].unique()
```

```
array(['>20', '15', '13', '7', '5', '16', '4', '11', '<1', '18', '19',  
      '12', '10', '9', '2', '6', '14', '3', '8', '20', '17', '1'],  
      dtype=object)
```

```
1 df_Clean["experience"].head()
```

```
0      >20  
10803    >20  
10797      4  
10796     17  
10794      7  
Name: experience, dtype: object
```

```
1 df_Clean.info()
```

```
<class 'pandas.core.frame.DataFrame'>  
Int64Index: 12477 entries, 0 to 12878  
Data columns (total 17 columns):  
#   Column                                Non-Null Count  Dtype  
---  -  
0   enrollee_id                          12477 non-null  int64  
1   city                                  12477 non-null  object  
2   city_development_index               12477 non-null  float64  
3   gender                               12477 non-null  object  
4   relevent_experience                  12477 non-null  object  
5   enrolled_university                 12477 non-null  object  
6   education_level                     12477 non-null  object  
7   major_discipline                    12477 non-null  object  
8   experience                           12477 non-null  object  
9   last_new_job                         12477 non-null  int64  
10  training_hours                       12477 non-null  int64  
11  target                               12477 non-null  float64  
12  cdi_bucket                           12477 non-null  category  
13  cat_experience                        12477 non-null  object  
14  educational_condition                12477 non-null  object  
15  training_hours_range                 12477 non-null  object  
16  company_type                         12477 non-null  object  
dtypes: category(1), float64(2), int64(3), object(11)  
memory usage: 1.9+ MB
```

```
1 df_Clean.head()
```

	enrollee_id	city	city_development_index	gender	relevent_experience
0	8949	city_103	0.920	Male	Has relevent experi
10803	24658	city_103	0.920	Male	Has relevent experi
10797	17847	city_21	0.624	Male	Has relevent experi
10796	1122	city_46	0.762	Male	Has relevent experi
10794	24653	city_160	0.920	Male	Has relevent experi

```
1
```

```
1
```

```
1 df_Clean["experience"] =df_Clean["experience"].replace(to_replace = ">20",
2                 value =21)
3 df_Clean["experience"] =df_Clean["experience"].replace(to_replace = "<1",
4                 value =0)
5 df_Clean["experience"] = df_Clean["experience"].astype("int64")
```

```
1 df_Clean["experience"].unique()
```

```
array([21,  4, 17,  7,  8,  6, 15, 18,  3, 14, 11, 16,  1,  5, 10, 13,  9,
        20, 12, 19,  2,  0])
```

```
1 df_Clean["major_discipline"].unique()
```

```
array(['STEM', 'No Major', 'Humanities', 'Business Degree', 'Other',
       'Arts'], dtype=object)
```

```
1 df_Clean['major'] = np.where(df_Clean['major_discipline'].isin(['Arts', 'Humanitie
2       'Other']), 'Non-STEM', df_Clean['major_discipline'])
```

```
1 df_Clean['major'].unique()
```

```
array(['STEM', 'Non-STEM'], dtype=object)
```

```
1 df_Clean['major_discipline'].unique()
```

```
array(['STEM', 'No Major', 'Humanities', 'Business Degree', 'Other',
       'Arts'], dtype=object)
```

```
1 df_Clean.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 12477 entries, 0 to 12878
Data columns (total 18 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   enrollee_id                          12477 non-null  int64
1   city                                 12477 non-null  object
2   city_development_index              12477 non-null  float64
3   gender                              12477 non-null  object
4   relevent_experience                 12477 non-null  object
5   enrolled_university                12477 non-null  object
6   education_level                    12477 non-null  object
7   major_discipline                   12477 non-null  object
8   experience                          12477 non-null  int64
9   last_new_job                       12477 non-null  int64
10  training_hours                     12477 non-null  int64
11  target                             12477 non-null  float64
12  cdi_bucket                         12477 non-null  category
13  cat_experience                      12477 non-null  object
14  educational_condition               12477 non-null  object
15  training_hours_range                12477 non-null  object
16  company_type                       12477 non-null  object
17  major                              12477 non-null  object
dtypes: category(1), float64(2), int64(4), object(11)
memory usage: 2.0+ MB
```

```
1 df_Clean.drop(['major_discipline'], axis=1, inplace=True)
2
```

```
1 df_Clean.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 12477 entries, 0 to 12878
Data columns (total 17 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   enrollee_id                          12477 non-null  int64
1   city                                 12477 non-null  object
2   city_development_index              12477 non-null  float64
3   gender                              12477 non-null  object
4   relevent_experience                 12477 non-null  object
5   enrolled_university                12477 non-null  object
6   education_level                    12477 non-null  object
7   experience                          12477 non-null  int64
8   last_new_job                       12477 non-null  int64
9   training_hours                     12477 non-null  int64
10  target                             12477 non-null  float64
11  cdi_bucket                         12477 non-null  category
12  cat_experience                      12477 non-null  object
13  educational_condition               12477 non-null  object
14  training_hours_range                12477 non-null  object
```

```
15 company_type      12477 non-null object
16 major             12477 non-null object
dtypes: category(1), float64(2), int64(4), object(10)
memory usage: 1.9+ MB
```

```
1 df_Clean.head()
```

	enrollee_id	city	city_development_index	gender	relevent_experience
0	8949	city_103	0.920	Male	Has relevent experience
10803	24658	city_103	0.920	Male	Has relevent experience
10797	17847	city_21	0.624	Male	Has relevent experience
10796	1122	city_46	0.762	Male	Has relevent experience
10794	24653	city_160	0.920	Male	Has relevent experience

```
1 df_Clean["major"].unique()
```

```
array(['STEM', 'Non-STEM'], dtype=object)
```

```
1 df_Clean = pd.get_dummies(df_Clean,columns=["gender", "relevent_experience", "enro
```

```
1 df_ML = df_Clean.drop(["enrollee_id", "city","cat_experience","educational_conditi
```

```
-----
-
NameError                                Traceback (most recent call
last)
<ipython-input-11-b8190ea87665> in <module>()
----> 1 df_ML = df_Clean.drop(["enrollee_id",
"city","cat_experience","educational_condition","training_hours_range"],
axis=1)

NameError: name 'df_Clean' is not defined
```

```
1 df_Clean.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 12477 entries, 0 to 12878
Data columns (total 26 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   enrollee_id                          12477 non-null  int64
1   city                                  12477 non-null  object
2   city_development_index               12477 non-null  float64
3   experience                           12477 non-null  int64
```



```

4   last_new_job          12477 non-null   int64
5   training_hours        12477 non-null   int64
6   target                12477 non-null   float64
7   cat_experience         12477 non-null   object
8   educational_condition  12477 non-null   object
9   training_hours_range  12477 non-null   object
10  gender_Male            12477 non-null   uint8
11  gender_Other           12477 non-null   uint8
12  relevent_experience_No relevent experience  12477 non-null   uint8
13  enrolled_university_Part time course      12477 non-null   uint8
14  enrolled_university_no_enrollment         12477 non-null   uint8
15  education_level_PhD      12477 non-null   uint8
16  education_level_Undergraduate             12477 non-null   uint8
17  cdi_bucket_medium_human_development       12477 non-null   uint8
18  cdi_bucket_high_human_development         12477 non-null   uint8
19  cdi_bucket_very_high_human_development    12477 non-null   uint8
20  company_type_Other       12477 non-null   uint8
21  company_type_Public Sector                 12477 non-null   uint8
22  company_type_Pvt Ltd      12477 non-null   uint8
23  company_type_Startup Company               12477 non-null   uint8
24  company_type_Unknown     12477 non-null   uint8
25  major_STEM               12477 non-null   uint8
dtypes: float64(2), int64(4), object(4), uint8(16)
memory usage: 1.5+ MB

```

```
1 df_ML.head()
```

	city_development_index	experience	last_new_job	training_hours	target
0	0.920	21	1	36	
10803	0.920	21	1	26	
10797	0.624	4	1	28	
10796	0.762	17	2	27	
10794	0.920	7	1	43	

Train and Test Dataset Split

```

1 ## Creating a train test split
2 from sklearn.model_selection import train_test_split
3 X = df_ML.drop(['target'], axis = 1)
4 y = df_ML['target']
5
6 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.3, random_

1 X.head()

```

	city_development_index	experience	last_new_job	training_hours	gender
0	0.920	21	1	36	
10803	0.920	21	1	26	
10797	0.624	4	1	28	
10796	0.762	17	2	27	
10794	0.920	7	1	43	

▼ 3.1 Gaussian Naive Bayes

```
1 from sklearn.naive_bayes import GaussianNB # 1. choose model class
2 model = GaussianNB() # 2. instantiate model
3 model.fit(X_train, y_train) # 3. fit model to data
4 y_model = model.predict(X_test) # 4. predict on new data
```

```
1 from sklearn.metrics import accuracy_score
2 accuracy_score(y_test, y_model)
```

```
0.7475961538461539
```

```
1 confusion_matrix(y_test, y_model)
```

```
array([[2269, 606],
       [ 339, 530]])
```

```
1 print(2357/(2357 + 517))
2 print("Sensitivity")
```

```
0.8201113430758524
Sensitivity
```

```
1 print(501/(501+369))
2 print("Specificity")
```

```
0.5758620689655173
Specificity
```

```
1 test = X_test.join(y_test).reset_index()
2 test.join(pd.Series(y_model, name='predicted')).head(20)
```

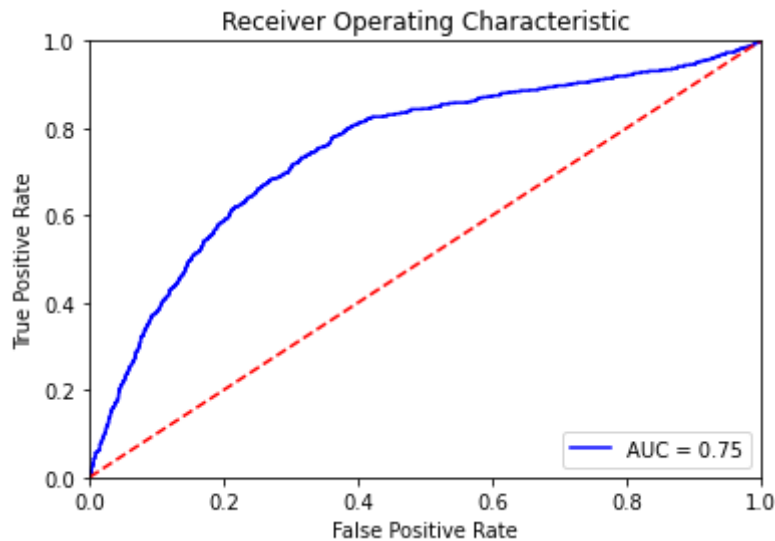
	index	city_development_index	experience	last_new_job	training_hours
0	9296	0.878	5	4	69
1	5358	0.512	9	1	68
2	8352	0.698	2	1	10
3	7679	0.920	21	5	153
4	7855	0.920	15	1	136
5	1533	0.740	21	0	31
6	16677	0.795	5	1	80
7	18739	0.920	4	3	87
8	11548	0.762	13	1	106
9	14143	0.926	12	1	322
10	1150	0.897	2	1	46
11	8176	0.897	3	2	28
12	8399	0.624	4	1	25
13	8183	0.920	4	2	84
14	3868	0.840	6	1	94
15	221	0.926	21	5	10
16	8952	0.920	6	2	15
17	9238	0.887	10	1	160
18	15628	0.920	5	1	200
19	18104	0.790	11	1	179

```

1 # creating AUC/ROC graph (code from stackoverflow)
2 import sklearn.metrics as metrics
3 probs = model.predict_proba(X_test)
4 preds = probs[:,1]
5 fpr, tpr, threshold = metrics.roc_curve(y_test, preds)
6 roc_auc = metrics.auc(fpr, tpr)
7
8 import matplotlib.pyplot as plt
9 plt.title('Receiver Operating Characteristic')
10 plt.plot(fpr, tpr, 'b', label = 'AUC = %0.2f' % roc_auc)
11 plt.legend(loc = 'lower right')
12 plt.plot([0, 1], [0, 1], 'r--')
13 plt.xlim([0, 1])
14 plt.ylim([0, 1])
15 plt.ylabel('True Positive Rate')

```

```
16 plt.xlabel('False Positive Rate')
17 plt.show()
```



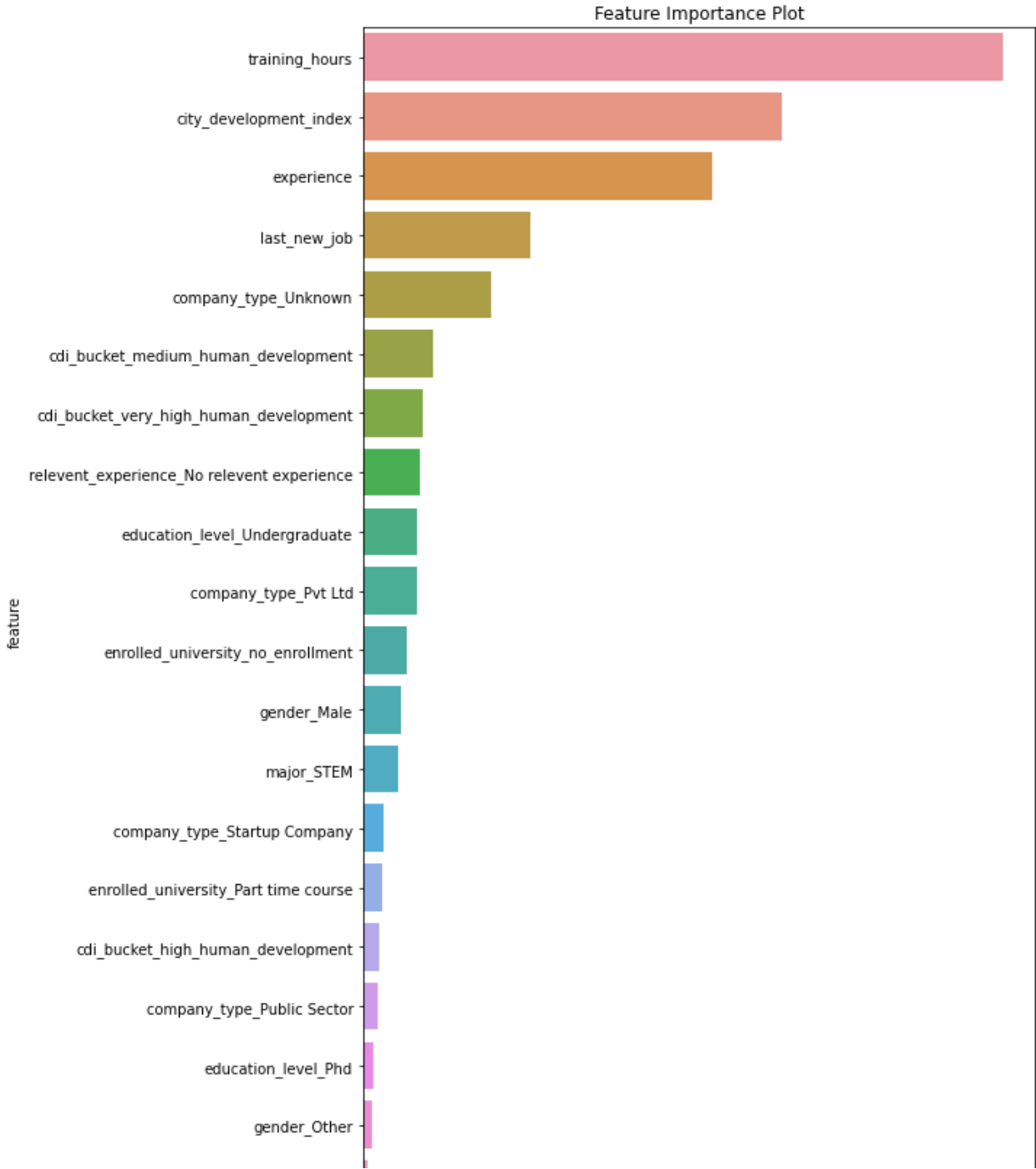
▼ 3.2 Random Forest

```
1 # initiate the model
2 rfc = RandomForestClassifier()
3 rfc.fit(X_train, y_train)
```

```
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, oob_score=False, random_state=None,
                        verbose=0, warm_start=False)
```

```
1 # explore feature importance
2 feature_importance = pd.DataFrame({'feature':X.columns, 'importance':rfc.feature_i
3 fig, ax = plt.subplots()
4 fig.set_size_inches(8.27,15)
5 plt.title('Feature Importance Plot')
6 sns.barplot(x='importance',y='feature',ax=ax,data=feature_importance[:50])
```

<matplotlib.axes._subplots.AxesSubplot at 0x7f3e806a3690>



```
1 y_pred = rfc.predict(X_test)
2 print(classification_report(y_test, y_pred))
3
```

	precision	recall	f1-score	support
0.0	0.83	0.89	0.86	2875
1.0	0.53	0.42	0.47	869
accuracy			0.78	3744
macro avg	0.68	0.65	0.66	3744
weighted avg	0.76	0.78	0.77	3744

```
1 from sklearn.metrics import roc_auc_score
2 y_pred_rfc = rfc.predict_proba(X_test)[: , 1]
3 print(roc_auc_score(y_test, y_pred_rfc))
4
```

0.7681150747986191

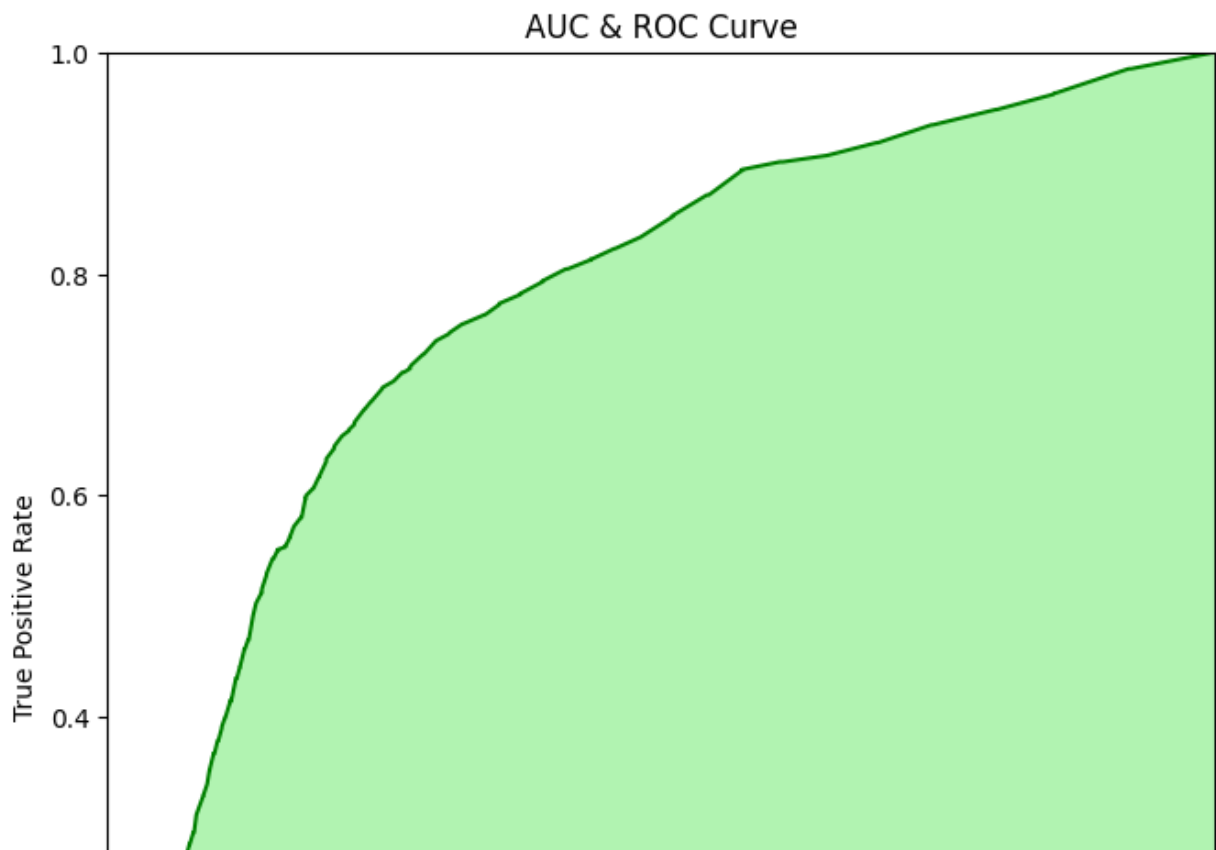
```
1 from sklearn.metrics import confusion_matrix
2 tn, fp, fn, tp = confusion_matrix(y_test, y_pred).ravel()
3 specificity = tn / (tn+fp)
4 sensitivity = tp / (tp+fn)
5 print("specificity", specificity, "sensitivity", sensitivity )
```

specificity 0.888 sensitivity 0.4154200230149597

1

```
1 from sklearn import metrics
2 auc = metrics.roc_auc_score(y_test, y_pred_rfc)
3 print(auc)
4 false_positive_rate, true_positive_rate, thresholds = metrics.roc_curve(y_test, y_p
5 #print(false_positive_rate,true_positive_rate)
6 plt.figure(figsize=(10, 8), dpi=100)
7 plt.axis('scaled')
8 plt.xlim([0, 1])
9 plt.ylim([0, 1])
10 plt.title("AUC & ROC Curve")
11 plt.plot(false_positive_rate, true_positive_rate, 'g')
12 plt.fill_between(false_positive_rate, true_positive_rate, facecolor='lightgreen',
13 plt.text(0.95, 0.05, 'AUC = %0.4f' % auc, ha='right', fontsize=12, weight='bold',
14 plt.xlabel("False Positive Rate")
15 plt.ylabel("True Positive Rate")
16 plt.show()
```

0.7681150747986191



▼ 3.3 XGBoost

```

1 def auc_roc_graph(ytest,preds):# creating AUC/ROC graph
2     ## Calculating metrics from ytest and predicted value
3     fpr, tpr, threshold = metrics.roc_curve(ytest, preds)
4     roc_auc = metrics.auc(fpr, tpr)
5     ## Plotting AUC-ROC
6     plt.title('Receiver Operating Characteristic')
7     plt.plot(fpr, tpr, 'b', label = 'AUC = %0.2f' % roc_auc)
8     plt.legend(loc = 'lower right')
9     plt.plot([0, 1], [0, 1],'r--')
10    plt.xlim([0, 1])
11    plt.ylim([0, 1])
12    plt.ylabel('True Positive Rate')
13    plt.xlabel('False Positive Rate')
14    plt.show()

1 def train_and_predict(X_train_model, X_test_model, y_train_model, y_test_model, cl
2     ## Training the Model
3     classifier.fit(X_train_model, y_train_model)
4
5     ## Predicting the model
6     y_model_tp = classifier.predict(X_test_model)
7

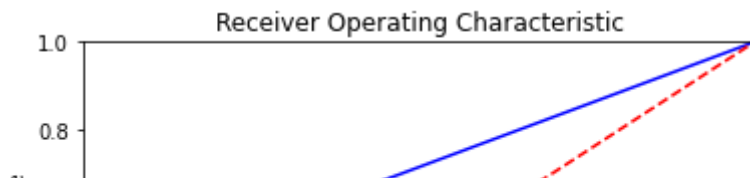
```

```
8  ## Accuracy score
9  a_score = metrics.accuracy_score(y_test_model, y_model_tp)
10 print("The Accuracy is {}".format(a_score))
11
12 ## Calculating sensitivity
13 TP = sum((y_test_model == 1) & (y_model_tp == 1))
14 p = sum((y_test_model == 1))
15 TPR = TP/p
16 print("The Sensitivity / True Positive Rate (TPR) is {}".format(TPR))
17
18 ## Calculating Specifictiy
19 N = sum(y_test_model == 0)
20 TN = sum((y_test_model == 0) & (y_model_tp == 0))
21 TNR = TN/N
22 print("The Specifictiy / TNR is {}".format(TNR))
23
24 # calculate recall
25 recall = metrics.recall_score(y_test_model, y_model_tp, average='binary')
26 print('Recall: %.3f' % recall)
27
28 # calculating precision
29 precision = metrics.precision_score(y_test_model, y_model_tp, average='binary')
30 print('Precision: {}'.format(precision))
31
32 ## Calculating F1 score
33 score = metrics.f1_score(y_test_model, y_model_tp, average='binary')
34 print('F1-score: {}'.format(score))
35
36 ##Printing the roc_auc
37 auc_roc_graph(y_test_model,y_model_tp)
38
39 return [a_score,TPR,TNR, recall, precision,score]

1 ## Instantiating the XGboost model
2
3 XgB = XGBClassifier(use_label_encoder=False, eval_metric='mlogloss')

1 xg_metrics = train_and_predict(X_train, X_test, y_train, y_test,XgB)
2
```


The Accuracy is 0.8036858974358975
 The Sensitivity / True Positive Rate (TPR) is 0.48331415420023016
 The Specificity / TNR is 0.9005217391304348
 Recall: 0.483
 Precision: 0.5949008498583569
 F1-score: 0.5333333333333333

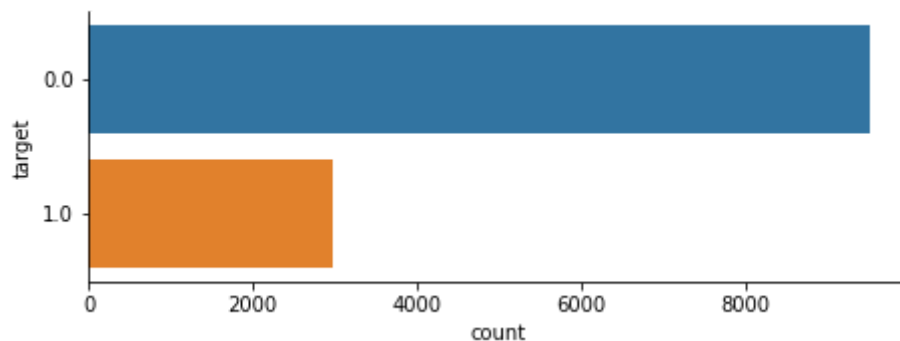


▼ 3.4 Logistic Model

```
df_ML.groupby('target').size()/df_ML.shape[0]
```

```
target
0.0    0.762283
1.0    0.237717
dtype: float64
```

```
1 sns.catplot(y="target", kind="count", data=df_ML, height=2.6, aspect=2.5);
```



We see that the target variable is very unbalanced

```
1 def train_and_predict(X_train_model, X_test_model, y_train_model, y_test_model, cl
2   ## Training the Model
3   classifier.fit(X_train_model, y_train_model)
4   ## Predicting the model
5   y_model_tp = classifier.predict(X_test_model)
6   ## Accuracy score
7   a_score = accuracy_score(y_test_model, y_model_tp)
8   print("The Accuracy is {}".format(a_score))
9   ## Calculating sensitivity
10  TP = sum((y_test_model == 1) & (y_model_tp == 1))
11  p = sum((y_test_model == 1))
12  TPR = TP/p
13  print("The Sensitivity / True Positive Rate (TPR) is {}".format(TPR))
```

```

14  ## Calculating Specifictiy
15  N = sum(y_test_model == 0)
16  TN = sum((y_test_model == 0) & (y_model_tp == 0))
17  TNR = TN/N
18  print("The Specifictiy / TNR is {}".format(TNR))
19  return [a_score,TPR,TNR]

```

```

1 from sklearn.linear_model import LogisticRegression
2 model = LogisticRegression(solver="liblinear")
3 model.fit(X_train, y_train)

```

```

    LogisticRegression(C=1.0, class_weight=None, dual=False, fit_intercept=True,
        intercept_scaling=1, l1_ratio=None, max_iter=100,
        multi_class='auto', n_jobs=None, penalty='l2',
        random_state=None, solver='liblinear', tol=0.0001, verbose=0,
        warm_start=False)

```

```

1 y_model = model.predict(X_test)

```

```

1 from sklearn.metrics import accuracy_score
2 print("Acurracy score:",accuracy_score(y_test, y_model))

```

```

    Acurracy score: 0.7769764957264957

```

```

1 #Sensitivity
2 P = sum(y_test == 1)
3 P

```

```

    869

```

```

1 TP = sum((y_test == 1) & (y_model == 1))
2 TP

```

```

    170

```

```

1 print("Sensitivity:",TP/P)

```

```

    Sensitivity: 0.1956271576524741

```

```

1 #Specificity
2 N = sum(y_test == 0)
3 N

```

```

    2875

```

```

1 TN = sum((y_test == 0) & (y_model == 0))
2 TN

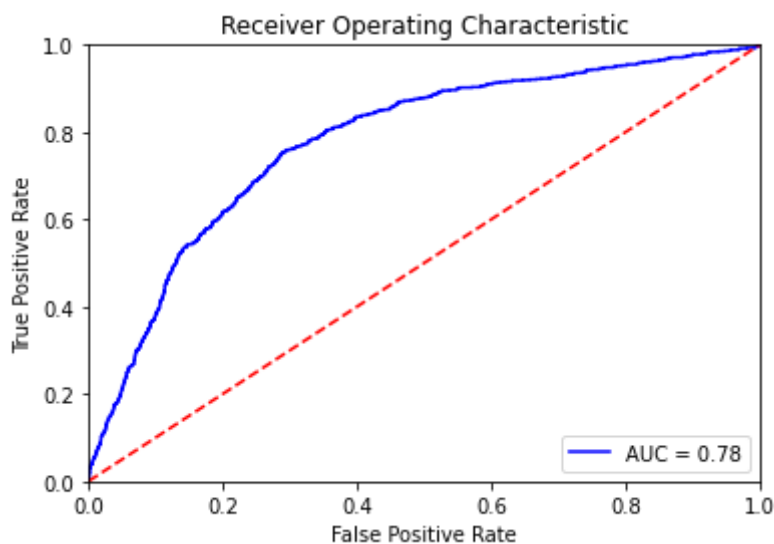
```

2739

```
1 print("Specificity:",TN/N)
```

```
Specificity: 0.952695652173913
```

```
1 # creating AUC/ROC graph
2 import sklearn.metrics as metrics
3 probs = model.predict_proba(X_test)
4 preds = probs[:,1]
5 fpr, tpr, threshold = metrics.roc_curve(y_test, preds)
6 roc_auc = metrics.auc(fpr, tpr)
7 import matplotlib.pyplot as plt
8 plt.title('Receiver Operating Characteristic')
9 plt.plot(fpr, tpr, 'b', label = 'AUC = %0.2f' % roc_auc)
10 plt.legend(loc = 'lower right')
11 plt.plot([0, 1], [0, 1], 'r--')
12 plt.xlim([0, 1])
13 plt.ylim([0, 1])
14 plt.ylabel('True Positive Rate')
15 plt.xlabel('False Positive Rate')
16 plt.show()
```



▼ 3.5 Multilayer Perceptron Classifier

- Rationale: The most typical MLP includes three layers: an input layer, a hidden layer and an output layer. The different layers of the MLP neural network are fully connected (full connection means: any neuron in the upper layer and all the neurons in the next layer). Neurons are connected)

- ReLU is a relatively popular activation function recently. When the input signal is less than 0, the output is 0; when the input signal is greater than 0, the output is equal to the input; the specific activation function used depends on the specific situation.

```
1 fit_MLP = MLPClassifier(random_state=1, max_iter=300).fit(X_train, y_train)
2 fit_MLP.predict_proba(X_test[:1])
```

```
array([[0.8507274, 0.1492726]])
```

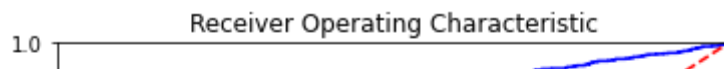
```
1 fit_MLP.predict(X_test)
```

```
array([0., 0., 1., ..., 0., 0., 1.])
```

```
1 fit_MLP.score(X_test, y_test)
```

```
0.7751068376068376
```

```
1 # creating AUC/ROC graph
2 import sklearn.metrics as metrics
3 probs = fit_MLP.predict_proba(X_test)
4 preds = probs[:,1]
5 fpr, tpr, threshold = metrics.roc_curve(y_test, preds)
6 roc_auc = metrics.auc(fpr, tpr)
7 import matplotlib.pyplot as plt
8 plt.title('Receiver Operating Characteristic')
9 plt.plot(fpr, tpr, 'b', label = 'AUC = %0.2f' % roc_auc)
10 plt.legend(loc = 'lower right')
11 plt.plot([0, 1], [0, 1], 'r--')
12 plt.xlim([0, 1])
13 plt.ylim([0, 1])
14 plt.ylabel('True Positive Rate')
15 plt.xlabel('False Positive Rate')
16 plt.show()
```



Conclusion

We are thinking from the perspective of HRs. Based on the EDA we conducted, we have several suggestions for HR regarding the optimal candidates selection.

- Hint 1: Gender seems to have no impact on candidates' decision. It is paramount to note that most candidates with relevant experience are male.
- Hint 2: In general, the candidates with higher qualifications belong to cities with higher CDI. However, candidates from cities with lower and medium development index are more likely to look for a change.
- Hint 3: Most of the enrollees major in STEM, which is reasonable because the company is hiring data scientists. However, the major discipline and having relevant experience or not do not affect the training hours -- which are related to the budget of the company. So HRs can lower down their "major bias". They all need training!
- Hint 4: STEM major, Masters and candidates with more experience joining this program will tend to stay.
- Hint 5: When recruiting candidates to join this program, HRs can turn to candidates with longer gap between last job and new one. However, not the ones who never change their job, those may experience nostalgia and don't want any change.
- Hint 6: Starting from the company size & type criteria, candidates from private limited company type and small or medium company size tend to stay.
- Hint 7: Considering the training hours which is highly related to the budget, candidates from small companies tend to have higher training hours compared to the others.
- Hint 8: 50-150hrs training hour range has higher percentage of retention, about 76%. This is the ideal training hour range for the employees.

Machine Learning

- After plotting 5 different machine models, we observed that Logistic Model and Multilayer Perceptron Classifier has the best AUC score of 0.78.
- The feature importance chart shows that training hours, city development and candidates' years of experience play a vital role in predicting whether the candidates stay or leave the company after training.
- Education background, gender, and company type/size acutally shows less significance in this problem.

