

# This jupyter notebook is prepared by Matthias Rathbun.

import libraries: pandas, numpy, matplotlib (set %matplotlib inline), matplotlib's pyplot, seaborn, missingno, scipy's stats, sklearn (1 pt)

```
In [1]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import missingno as msn
from scipy import stats
%matplotlib inline
```

import the data to a dataframe and show how many rows and columns does it have (1 pt)

```
In [2]: df = pd.read_csv("hcddata.csv", index_col = "rec_num")
df.shape

Out[2]: (21287, 17)
```

call the describe method of dataframe to see some summary statistics of the numerical columns. (1 pt)

```
In [3]: df.describe()
```

	Unnamed: 0	enrollee_id	city	city_development_index	training_hours	target	city_development_matrices
count	21287.000000	21287.000000		21287.000000	21287.000000	19158.000000	21287.000000
mean	10643.000000	16873.983652		0.828462	65.328510	0.249348	8.284615
std	6145.171926	9612.131237		0.123537	60.075201	0.432647	1.235365
min	0.000000	1.000000		0.448000	1.000000	0.000000	4.480000
25%	5321.500000	8554.500000		0.739000	23.000000	0.000000	7.390000
50%	10643.000000	16967.000000		0.903000	47.000000	0.000000	9.030000
75%	15964.500000	25161.500000		0.920000	88.000000	0.000000	9.200000
max	21286.000000	33380.000000		0.949000	336.000000	1.000000	9.490000

Target is a categorical given by 0 and 1, but the number of values in target is less than the number of values in the rest of the column implying there are missing values. Probably need to drop those. The unnamed 0 is an index column which will be dropped as there is two index columns.

```
In [4]: df = df.iloc[:, 1:]
```

Show the top 5 rows and last 5 rows of the data frame (1 pt)

```
In [5]: df.head()
```

rec_num	enrollee_id	city	city_development_index	gender	relevant_experience	enrolled_university	education_level	major_discipline	exper
1	8949	city_103		0.920	Male	Has relevant experience	no_enrollment	Graduate	STEM
2	29725	city_40		0.776	Male	No relevant experience	no_enrollment	Graduate	STEM
3	11561	city_21		0.624	NaN	No relevant experience	Full time course	Graduate	STEM
4	33241	city_115		0.789	NaN	No relevant experience	NaN	Graduate	Business Degree
5	666	city_162		0.767	Male	Has relevant experience	no_enrollment	Masters	STEM

```
In [6]: df.tail()
```

```
Out[6]:
```

rec_num	enrollee_id	city	city_development_index	gender	relevant_experience	enrolled_university	education_level	major_discipline	exper
21283	1289	city_103		0.920	Male	No relevant experience	no_enrollment	Graduate	Humanities
21284	195	city_136		0.897	Male	Has relevant experience	no_enrollment	Masters	STEM
21285	31762	city_100		0.887	Male	No relevant experience	no_enrollment	Primary School	NaN
21286	7873	city_102		0.804	Male	Has relevant experience	Full time course	High School	NaN
21287	12215	city_102		0.804	Male	Has relevant experience	no_enrollment	Masters	STEM

List all the numerical columns (1 pt)

List all the categorical columns (1 pt)

```
In [7]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 21287 entries, 1 to 21287
Data columns (total 16 columns):
#   Column                Non-Null Count  Dtype
---  --
0   enrollee_id            21287 non-null   int64
1   city                   21287 non-null   object
2   city_development_index 21287 non-null   float64
3   gender                 16271 non-null   object
4   relevant_experience     21287 non-null   object
5   enrolled_university    20870 non-null   object
6   education_level        20775 non-null   object
7   major_discipline       18162 non-null   object
8   experience             14721 non-null   object
9   company_size           14721 non-null   object
10  company_type           14513 non-null   object
11  last_new_job           20824 non-null   object
12  training_hours         21287 non-null   int64
13  target                 19158 non-null   float64
14  state                  21287 non-null   object
15  city_development_matrices 21287 non-null  float64
dtypes: float64(3), int64(2), object(11)
memory usage: 2.8+ MB
```

Numeric Columns: city\_development\_index, training\_hours, training\_hours, training\_hours, city, gender, relevant\_experience, enrolled\_university, education\_level, major\_discipline, experience, company size, company type, last new job, target, its just 1's and 0's so its categorical), state

Examine missing values: (2 + 2 + 2 + 5 = 11 pt)

Show a list with column wise count of missing values and display the list in count wise descending order

```
In [8]: df.isna().sum().sort_values(ascending = False)
```

```
Out[8]: company_type        6774
company_size        6560
gender              5016
major_discipline    3125
target             2129
education_level     512
last_new_job        463
enrolled_university 417
experience           70
enrollee_id         0
city                0
city_development_index 0
relevant_experience  0
training_hours      0
state              0
city_development_matrices 0
dtype: int64
```

Show a list with column wise percentage of missing values and display the list in percentage wise descending order

```
In [9]: df.isna().mean().sort_values(ascending=False) * 100
```

```
Out[9]: company_type        31.822239
company_size        30.816931
gender              23.563677
major_discipline    14.609231
target             10.001409
education_level     2.405224
last_new_job        2.175036
enrolled_university 1.958842
experience           0.328839
enrollee_id         0.000000
city                0.000000
city_development_index 0.000000
relevant_experience  0.000000
training_hours      0.000000
state              0.000000
city_development_matrices 0.000000
dtype: float64
```

Display a bar plot to visualize only the columns with missing values and their count. The plot should display from less missing value columns in the left and then more missing value columns to the right side of the plot

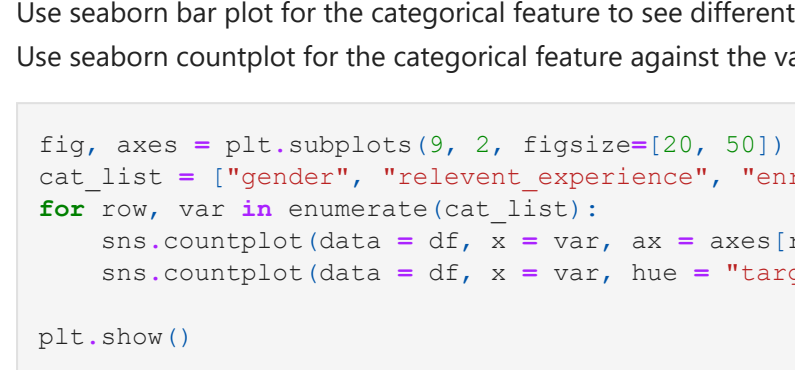
```
In [10]: missing_values = df.isna().sum().sort_values().to_frame()
missing_values = missing_values.loc[~missing_values.empty].all(axis=1)
```

```
In [11]: missing_values.columns = ["count"]
```

```
In [12]: missing_values.index.names = ["Name"]
```

```
In [13]: missing_values['Name'] = missing_values.index
```

```
In [13]: sns.set(style="whitegrid", color_codes=True)
sns.barplot(x = 'Name', y = 'count', data=missing_values)
plt.xticks(rotation = 30)
plt.show()
```

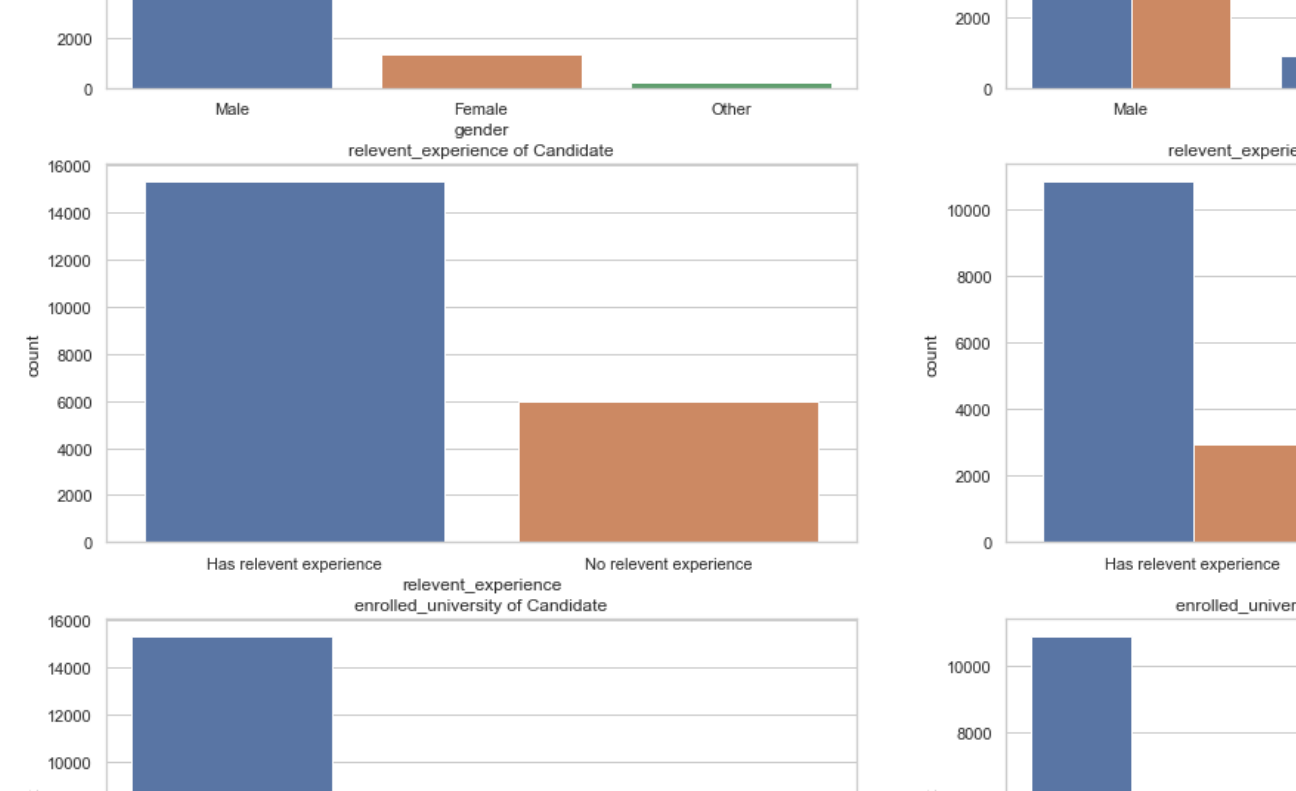


Use missingno's bar plot, matrix plot with 200 sample, and heatmap.

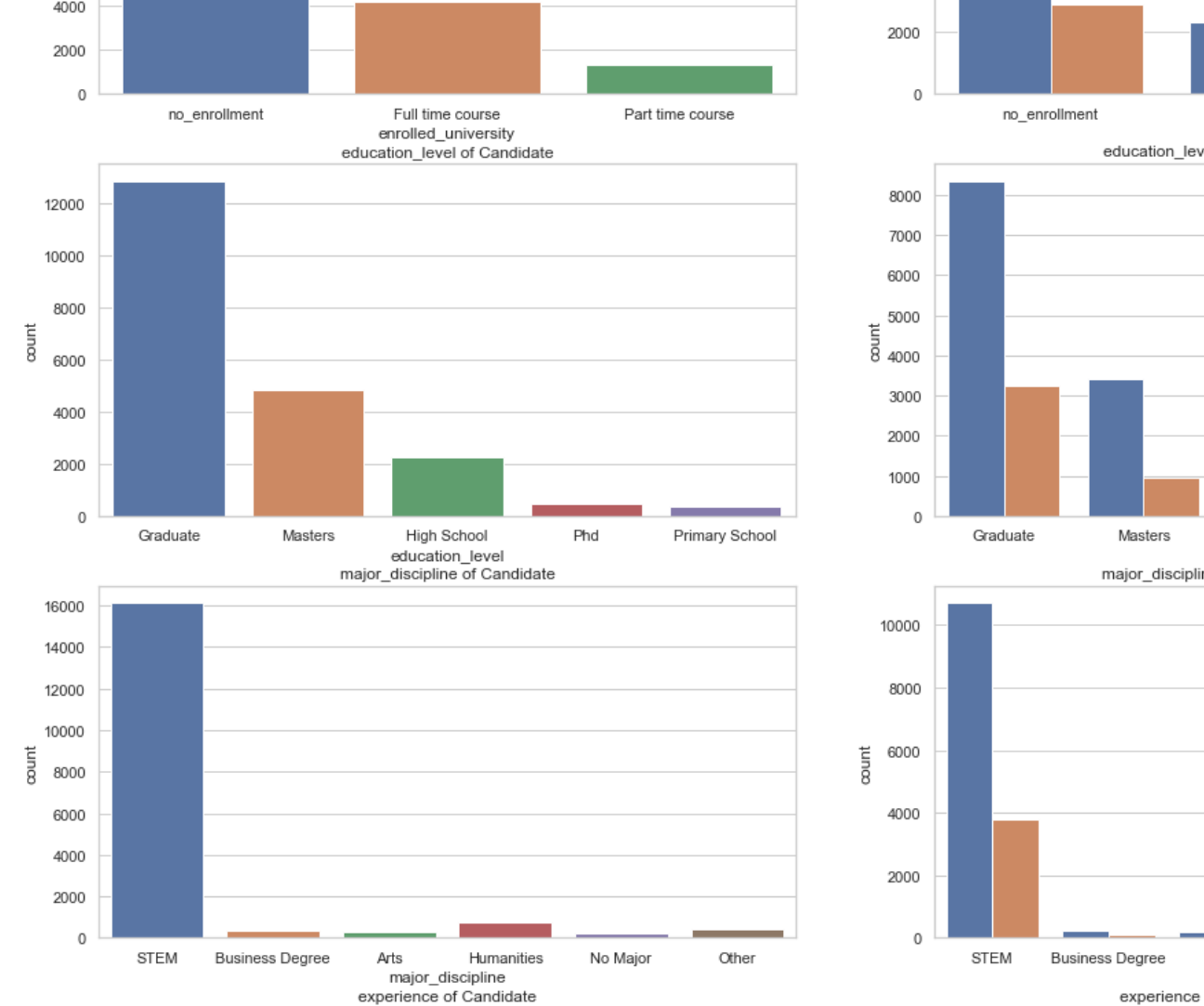
Interpret any interesting information you found in the heatmap and any one plot

```
In [14]: msn.bar(df, figsize=(10,5), fontsize=12)
```

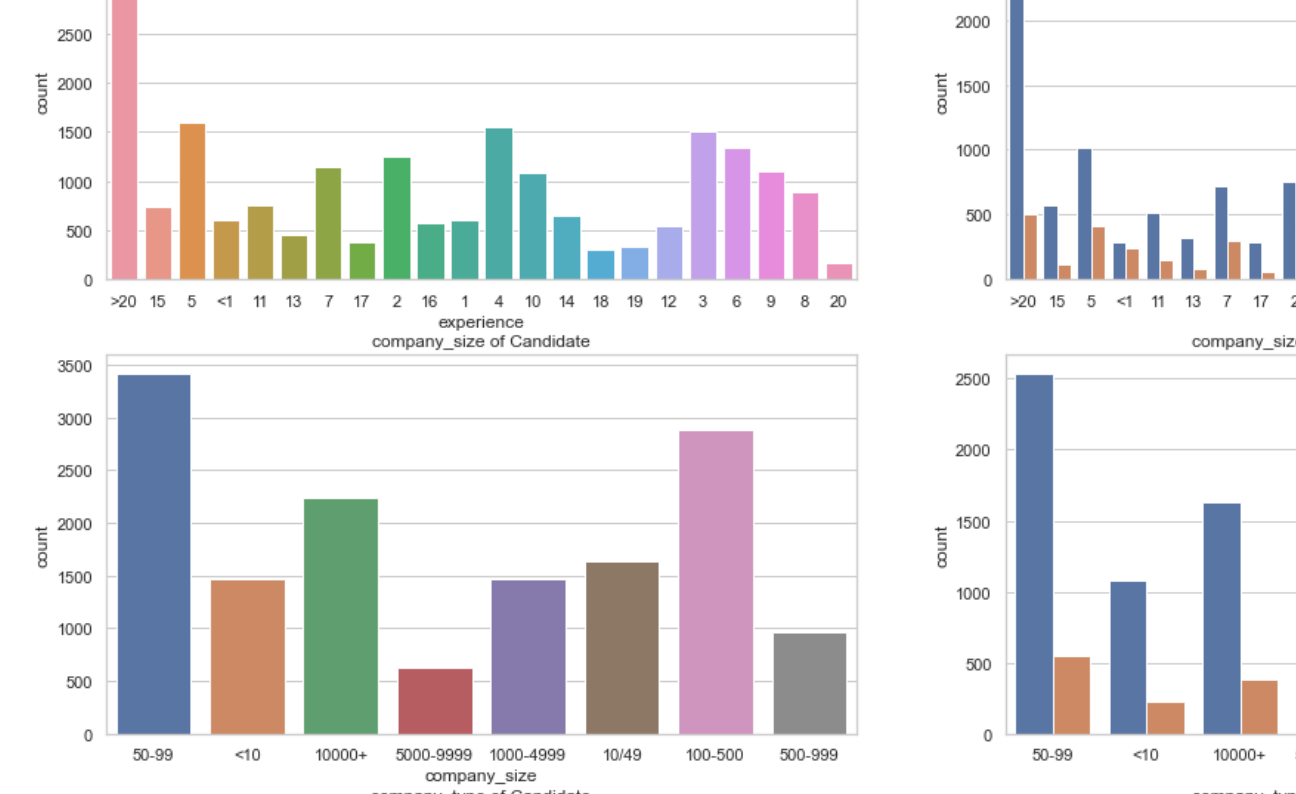
```
Out[14]: <AxesSubplot>
```



```
In [15]: msn.matrix(df, figsize=(10,5), fontsize=12);
```



```
In [16]: msn.heatmap(df, figsize=(10,5), fontsize=12);
```



The values we want to predict have some missing values, potential to separate. Also cols with massive missing values should be dropped as a whole as they are categorical and cannot be interpolated.

## Note to grader, I am using countplot for both as I can then pair each of the plots by variable they represent

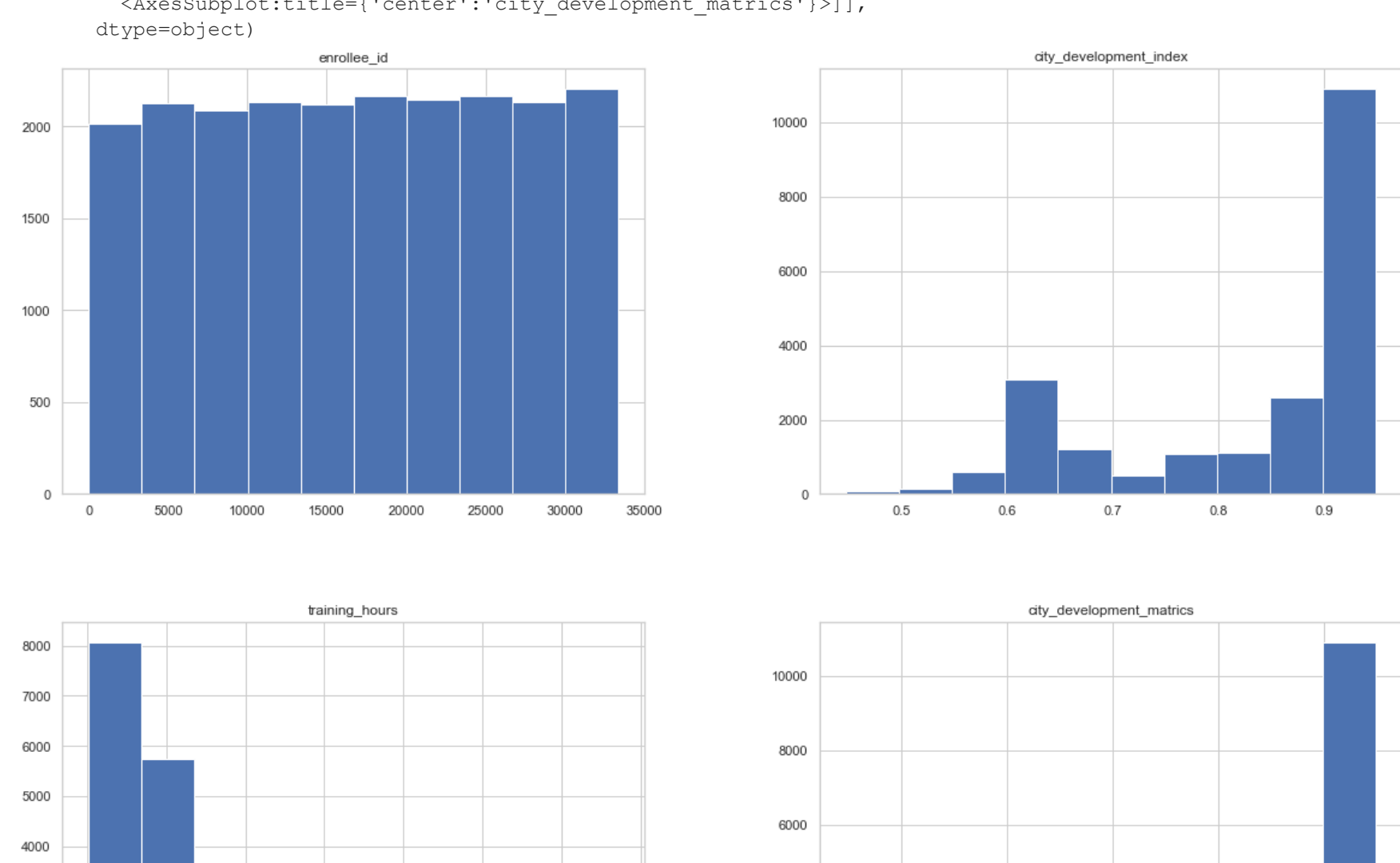
Understanding Categorical attributes (this part may require you to make 20+ plots ) [26 pts]

For each categorical attribute perform the following:

Use seaborn bar plot for the categorical feature to see different values and count

Use seaborn countplot for the categorical feature against the values of the target

```
In [17]: fig, axes = plt.subplots(9, 2, figsize=(20, 50))
cat_list = ["gender", "relevant_experience", "enrolled_university", "education_level", "major_discipline", "experience", "company_size", "company_type", "last_new_job"]
for row, var in enumerate(cat_list):
    sns.countplot(data = df, x = var, ax = axes[row, 0]).set(title = var + " of Candidate")
    sns.countplot(data = df, x = var, hue = "target", ax = axes[row, 1]).set(title = var + " of Candidate separated by target")
plt.show()
```



Understanding Numerical attributes (16 pts) For each numerical features, perform the following: Interpret their distributions using histogram (removed the group by word) Plot the distribution using seaborn distplot Interpret any interesting information(Candidate has way too many variables that should have no effect on the value of target. The only state was California so that should be dropped as well. There are a lot more college graduates than any other class, so the dataset may be skewed towards them. Also those with more than 20 years of experience represent a lot of the data. It seems that those with less experience are more willing to change jobs. Almost all are stem degree holders, so it may be necessary to drop all non stem degrees as there isn't enough representation of them. Those in full time enrollment are also more likely to want to change jobs. Same with those with no relevant experience relative to those with relevant experience)

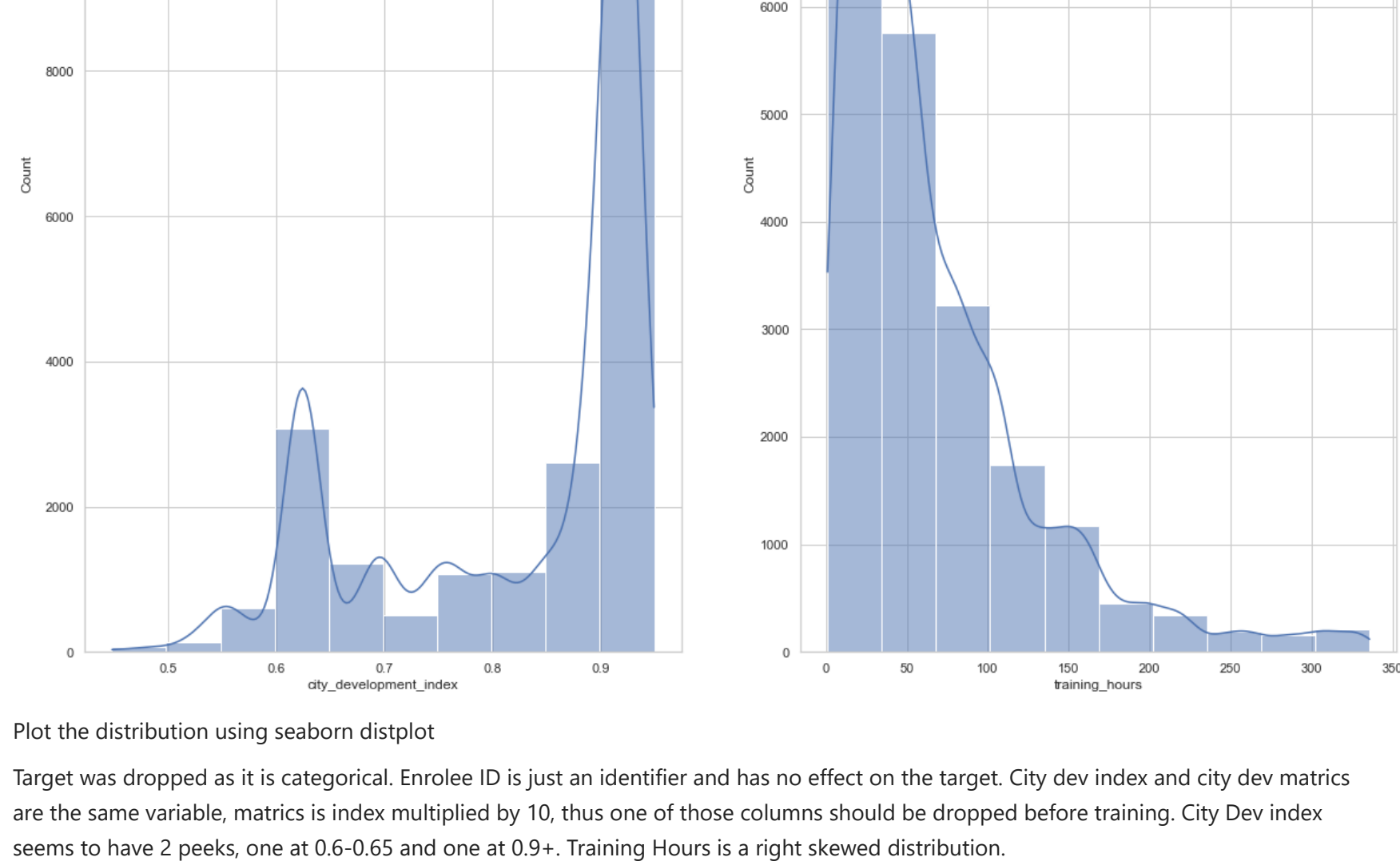
Understanding Numerical attributes (16 pts)

For each numerical features, perform the following:

Plot their distributions using histogram (removed the group by word)

```
In [18]: df.drop(columns = ["target"]).hist(figsize = (20,15))
```

```
Out[18]: array([[<AxesSubplot:title='center': 'enrollee_id'>],
[<AxesSubplot:title='center': 'city_development_index'>],
[<AxesSubplot:title='center': 'training_hours'>],
[<AxesSubplot:title='center': 'city_development_matrices'>]],
dtype=object)
```



Note to grader: Histplot was used over dist plot so I could put both graphs in the same figure.

Plot the distribution using seaborn distplot

```
In [19]: fig, axes = plt.subplots(1, 2, figsize=(20, 15))
dist_cols = ["city_development_index", "training_hours"]
for col, var in enumerate(dist_cols):
    sns.histplot(x = df[var], ax = axes[col], kde = True, bins = 10).set(title = "Histogram of " + var)
```



Plot the distribution using seaborn distplot

Target was dropped as it is categorical. Enrollee ID is just an identifier and has no effect on the target. City dev index and city dev matrices are the same variable, matrices is index multiplied by 10, thus one of those columns should be dropped before training. City Dev index seems to have 2 peaks, one at 0.6-0.65 and one at 0.9+. Training Hours is a right skewed distribution.

Correlation: (15 pts)

For the numerical attributes, use heatmap to show the correlation

If you find any interesting short list of columns, create another heatmap with them and show the correlations inside the heatmap as well Show scatter plots between columns to show the relationships with the target Interpret and explain any finding and next course of action from there

```
In [20]: corr = df.drop(columns = ["target", "enrollee_id", "city_development_matrices"]).corr()
plt.figure(figsize = (10,10))
sns.heatmap(corr).set(title = "Correlation between City Development Index and Training Hours (Relevant Numerical Columns)")
```

```
Out[20]: [Text(0.5, 1.0, 'Correlation between City Development Index and Training Hours (Relevant Numerical Columns)')]
```



```
In [21]: plt.figure(figsize = (10,10))
sns.scatterplot(data = df, x = "city_development_index", y = "training_hours", hue = "target", palette = "deep")
plt.show()
```

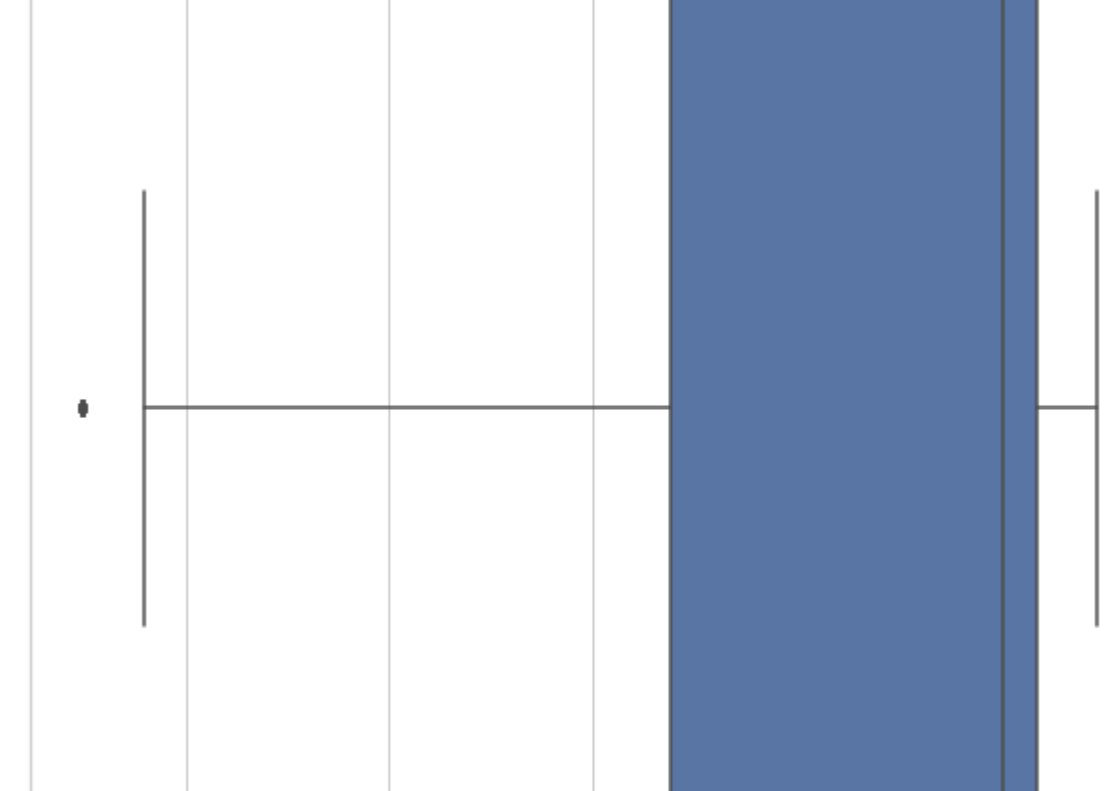


Both the Histogram and the scatterplot show no correlation between the city development index and training hours. Both distributions could be normalized in the future, but the features should not be combined.

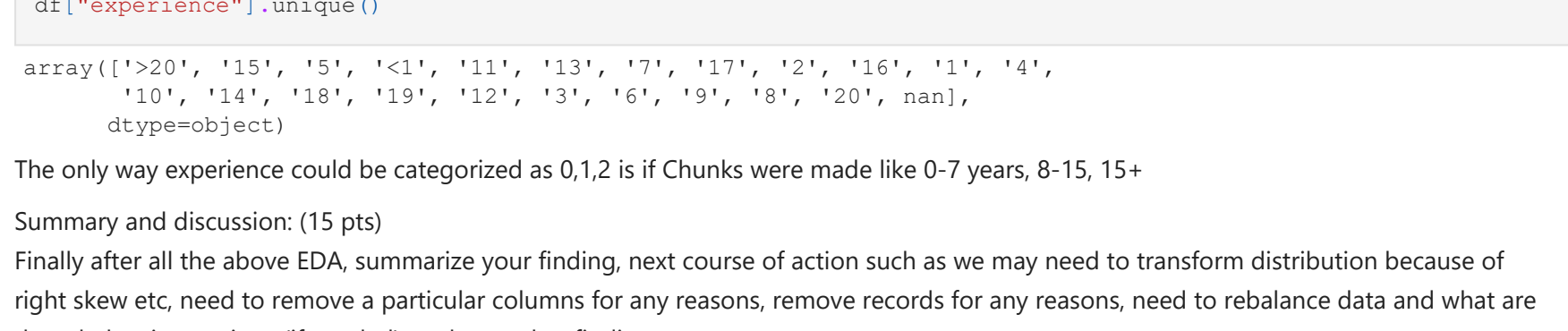
Outliers: (5) Use boxplot or any other strategies to find outliers

```
In [22]: plt.figure(figsize = (10,10))
sns.boxplot(data = df, x = "training_hours")
```

```
Out[22]: <AxesSubplot:xlabel='training hours'>
```



```
In [23]: plt.figure(figsize = (10,10))
sns.boxplot(data = df, x = "city_development_index")
```



Outliers are outside the whiskers of the box plots

What are the different values of experience, can you categorize them in to 0, 1, and 27 (5 pts)

```
In [24]: df["experience"].unique()
```

```
Out[24]: array(['120.0', '13.0', '15.0', '16.0', '17.0', '18.0', '19.0', '20.0', '21.0', '22.0', '23.0', '24.0', '25.0', '26.0', '27.0', '28.0', '29.0', '30.0', '31.0', '32.0', '33.0', '34.0', '35.0', '36.0', '37.0', '38.0', '39.0', '40.0', '41.0', '42.0', '43.0', '44.0', '45.0', '46.0', '47.0', '48.0', '49.0', '50.0', '51.0', '52.0', '53.0', '54.0', '55.0', '56.0', '57.0', '58.0', '59.0', '60.0', '61.0', '62.0', '63.0', '64.0', '65.0', '66.0', '67.0', '68.0', '69.0', '70.0', '71.0', '72.0', '73.0', '74.0', '75.0', '76.0', '77.0', '78.0', '79.0', '80.0', '81.0', '82.0', '83.0', '84.0', '85.0', '86.0', '87.0', '88.0', '89.0', '90.0', '91.0', '92.0', '93.0', '94.0', '95.0', '96.0', '97.0', '98.0', '99.0', '100.0', '101.0', '102.0', '103.0', '104.0', '105.0', '106.0', '107.0', '108.0', '109.0', '110.0', '111.0', '112.0', '113.0', '114.0', '115.0', '116.0', '117.0', '118.0', '119.0', '120.0', '121.0', '122.0', '123.0', '124.0', '125.0', '126.0', '127.0', '128.0', '129.0', '130.0', '131.0', '132.0', '133.0', '134.0', '135.0', '136.0', '137.0', '138.0', '139.0', '140.0', '141.0', '142.0', '143.0', '144.0', '145.0', '146.0', '147.0', '148.0', '149.0', '150.0', '151.0', '152.0', '153.0', '154.0', '155.0', '156.0', '157.0', '158.0', '159.0', '160.0', '161.0', '162.0', '163.0', '164.0', '165.0', '166.0', '167.0', '168.0', '169.0', '170.0', '171.0', '172.0', '173.0', '174.0', '175.0', '176.0', '177.0', '178.0', '179.0', '180.0', '181.0', '182.0', '183.0', '184.0', '185.0', '186.0', '187.0', '188.0', '189.0', '190.0', '191.0', '192.0', '193.0', '194.0', '195.0', '196.0', '197.0', '198.0', '199.0', '200.0', '201.0', '202.0', '203.0', '204.0', '205.0', '206.0', '207.0', '208.0', '209.0', '210.0', '211.0', '212.0', '213.0', '214.0', '215.0', '216.0', '217.0', '218.0', '219.0', '220.0', '221.0', '222.0', '223.0', '224.0', '225.0', '226.0', '227.0', '228.0', '229.0', '230.0', '231.0', '232.0', '233.0', '234.0', '235.0', '236.0', '237.0', '238.0', '239.0', '240.0', '241.0', '242.0', '243.0', '244.0', '245.0', '246.0', '247.0', '248.0', '249.0', '250.0', '251.0', '252.0', '253.0', '254.0', '255.0', '256.0', '257.0', '258.0', '259.0', '260.0', '261.0', '262.0', '263.0', '264.0', '265.0', '266.0', '267.0', '268.0', '269.0', '270.0', '271.0', '272.0', '273.0', '274.0', '275.0', '276.0', '277.0', '278.0', '279.0', '280.0', '281.0', '282.0', '283.0', '284.0', '285.0', '286.0', '287.0', '288.0', '289.0', '290.0', '291.0', '292.0', '293.0', '294.0', '295.0', '296.0', '297.0', '298.0', '299.0', '300.0', '301.0', '302.0', '303.0', '304.0', '305.0', '306.0', '307.0', '308.0', '309.0', '310.0', '311.0', '312.0', '313.0', '314.0', '315.0', '316.0', '317.0', '318.0', '319.0', '320.0', '321.0', '322.0', '323.0', '324.0', '325.0', '326.0', '327.0', '328.0', '329.0', '330.0', '331.0', '332.0', '333.0', '334.0', '335.0', '336.0', '337.0', '338.0', '339.0', '340.0', '341.0', '342.0', '343.0', '344.0', '345.0', '346.0', '347.0', '348.0', '349.0', '350.0', '351.0', '352.0', '353.0', '354.0', '355.0', '356.0', '357.0', '358.0', '359.0', '360.0', '361.0', '362.0', '363.0', '364.0', '365.0', '366.0', '367.0', '368.0', '369.0', '370.0', '371.0', '372.0', '373.0', '374.0', '375.0', '376.0', '377.0', '378.0', '379.0', '380.0', '381.0', '382.0', '383.0', '384.0', '385.0', '386.0', '387.0', '388.0', '389.0', '390.0', '391.0', '392.0', '393.0', '394.0', '395.0', '396.0', '397.0', '398.0', '399.0', '400.0', '401.0', '402.0', '403.0', '404.0', '405.0', '406.0', '407.0', '408.0', '409.0', '410.0', '411.0', '412.0', '413.0', '414.0', '415.0', '416.0', '417.0', '418.0', '419.0', '420.0', '421.0', '422.0', '423.0', '424.0', '425.0', '426.0', '427.0', '428.0', '429.0', '430.0', '431.0', '432.0', '433.0', '434.0', '435.0', '436.0', '437.0', '438.0', '439.0', '440.0', '441.0', '442.0', '443.0', '444.0', '445.0', '446.0', '447.0', '448.0', '449.0', '450.0', '451.0', '452.0', '453.0', '454.0', '455.0', '456.0', '457.0', '458.0', '459.0', '460.0', '461.0', '462.0', '463.0', '464.0', '465.0', '466.0', '467.0', '468.0', '469.0', '470.0', '471.0', '472.0', '473.0', '474.0', '475.0', '476.0', '477.0', '478.0', '479.0', '480.0', '481.0', '482.0', '483.0', '484.0', '485.0', '486.0', '487.0', '488.0', '489.0', '490.0', '491.0', '492.0', '493.0', '494.0', '495.0', '496.0', '497.0', '498.0', '499.0', '500.0', '501.0', '502.0', '503.0', '504.0', '505.0', '506.0', '507.0', '508.0', '509.0', '510.0', '511.0', '512.0', '513.0', '514.0', '515.0', '516.0', '517.0', '518.0', '519.0', '520.0', '521.0', '522.0', '523.0', '524.0', '525.0', '526.0', '527.0', '528.0', '529.0', '530.0', '531.0', '532.0', '533.0', '534.0', '535.0', '536.0', '537.0', '538.0', '539.0', '540.0', '541.0', '542.0', '543.0', '544.0', '545.0', '546.0', '547.0', '548.0', '549.0', '550.0', '551.0', '552.0', '553.0', '554.0', '555.0', '556.0', '557.0', '558.0', '559.0', '560.0', '561.0', '562.0', '563.0', '564.0', '565.0', '566.0', '567.0', '568.0', '569.0', '570.0', '571.0', '572.0', '573.0', '574.0', '575.0', '576.0', '577.0', '578.0', '579.0', '580.0', '581.0', '582.0', '583.0', '584.0', '585.0', '586.0', '587.0', '588.0', '589.0', '590.0', '591.0', '592.0', '593.0', '594.0', '595.0', '596.0', '597.0', '598.0', '599.0', '600.0', '601.0', '602.0', '603.0', '604.0', '605.0', '606.0', '607.0', '608.0', '609.0', '610.0', '611.0', '612.0', '613.0', '614.0', '615.0', '616.0', '617.0', '618.0', '619.0', '620.0', '621.0', '622.0', '623.0', '624.0', '625.0', '626.0', '627.0', '628.0', '629.0', '630.0', '631.0', '632.0', '633.0', '634.0', '635.0', '636.0', '637.0', '638.0', '639.0', '640.0', '641.0', '642.0', '643.0', '644.0', '645.0', '646.0', '647.0', '648.0', '649.0', '650.0', '651.0', '652.0', '653.0', '654.0', '655.0', '656.0', '657.0', '658.0', '659.0', '660.0', '661.0', '662.0', '663.0', '664.0', '665.0', '666.0', '667.0', '668.0', '669.0
```



