

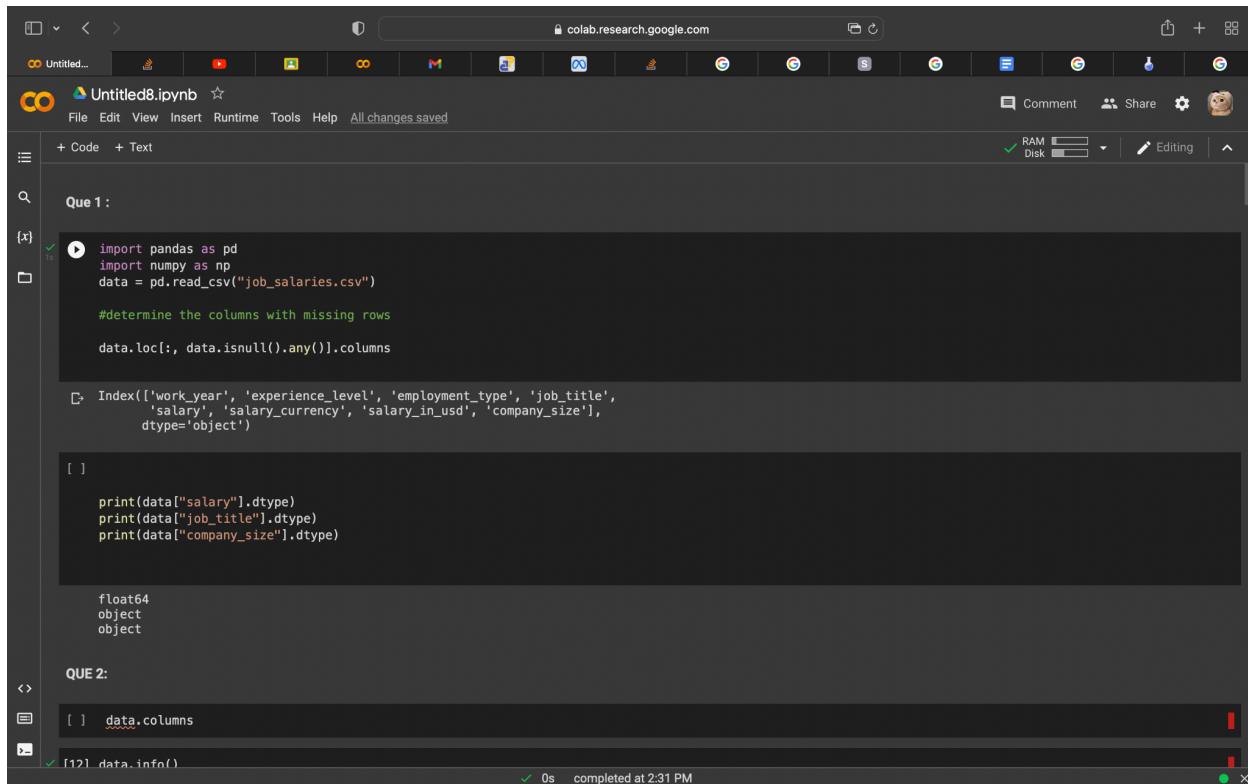
Smruti Dawale
B21BB007
LAB 2

Que 1)

By observation we can tell the data type of the given categories

- a) job_title = object
- b) Salary = float
- c) Company_size =object

We can verify this by using the code line filename[name_of_column].dtype



The screenshot shows a Google Colab notebook titled "Untitled8.ipynb". The code cell contains Python code to read a CSV file and print the data types of specific columns. The output shows that 'salary' is float64, 'job_title' is object, and 'company_size' is object. Below this, a new cell labeled "QUE 2:" is shown with the command "data.columns".

```
import pandas as pd
import numpy as np
data = pd.read_csv("job_salaries.csv")

#determine the columns with missing rows
data.loc[:, data.isnull().any()].columns

Index(['work_year', 'experience_level', 'employment_type', 'job_title',
       'salary', 'salary_currency', 'salary_in_usd', 'company_size'],
      dtype='object')

[ ]:

print(data["salary"].dtype)
print(data["job_title"].dtype)
print(data["company_size"].dtype)

float64
object
object

QUE 2:
[ ] data.columns
```

Que 2)

Almost all the algorithms in scikit learn are not capable of handling missing data , so before training your model you must handle your missing data.

One of the ways is “list wise deletion of cases” : discarding observations where value in any of the variables is missing .Here observations are your rows and variables are columns , you basically discard the row if any of the values in your column is missing.

I used `dataDropNa = data.dropna(axis=0)` , thereby you can see 497 rows x 12 columns , instead of the original data where we had 607 rows x 12 columns.

We can also replace all the NA / NaN values by specific value (I used 0) `data_fill = data.fillna(0)`, this basically replaces all of your missing values by 0.

The screenshot shows a Google Colab interface. The top bar includes tabs for 'Untitled...', 'Untitled8.ipynb' (which is the active tab), and various Google services like Sheets, Slides, and Drive. The toolbar has icons for file operations, search, and sharing. The code cell contains the following Python code:

```
data_fill = data.fillna(0)
```

The output cell displays a preview of a pandas DataFrame. The columns are labeled: work_year, experience_level, employment_type, job_title, salary, salary_currency, salary_in_usd, employee_residence, remote_ratio, company_location, and compa. The data shows rows from index 0 to 606, with some rows being ellipses (...). The preview includes a scroll bar on the right.

	work_year	experience_level	employment_type	job_title	salary	salary_currency	salary_in_usd	employee_residence	remote_ratio	company_location	compa
0	2020.0	MI	FT	Data Scientist	70000.0	EUR	79833.0	DE	0	DE	
1	2020.0	SE	FT	Machine Learning Scientist	260000.0	USD	260000.0	JP	0	JP	
2	2020.0	SE	FT	Big Data Engineer	85000.0	GBP	109024.0	GB	50	GB	
3	2020.0	MI	FT	Product Data Analyst	20000.0	USD	20000.0	HN	0	HN	
4	2020.0	SE	FT	Machine Learning Engineer	150000.0	USD	150000.0	US	50	US	
...	
602	2022.0	SE	FT	Data Engineer	154000.0	USD	154000.0	US	100	US	
603	2022.0	SE	FT	Data Engineer	126000.0	USD	126000.0	US	100	US	
604	2022.0	SE	FT	Data Analyst	129000.0	USD	129000.0	US	0	US	
605	2022.0	SE	FT	Data Analyst	150000.0	USD	150000.0	US	100	US	
606	2022.0	MI	FT	AI Scientist	200000.0	USD	200000.0	IN	100	US	

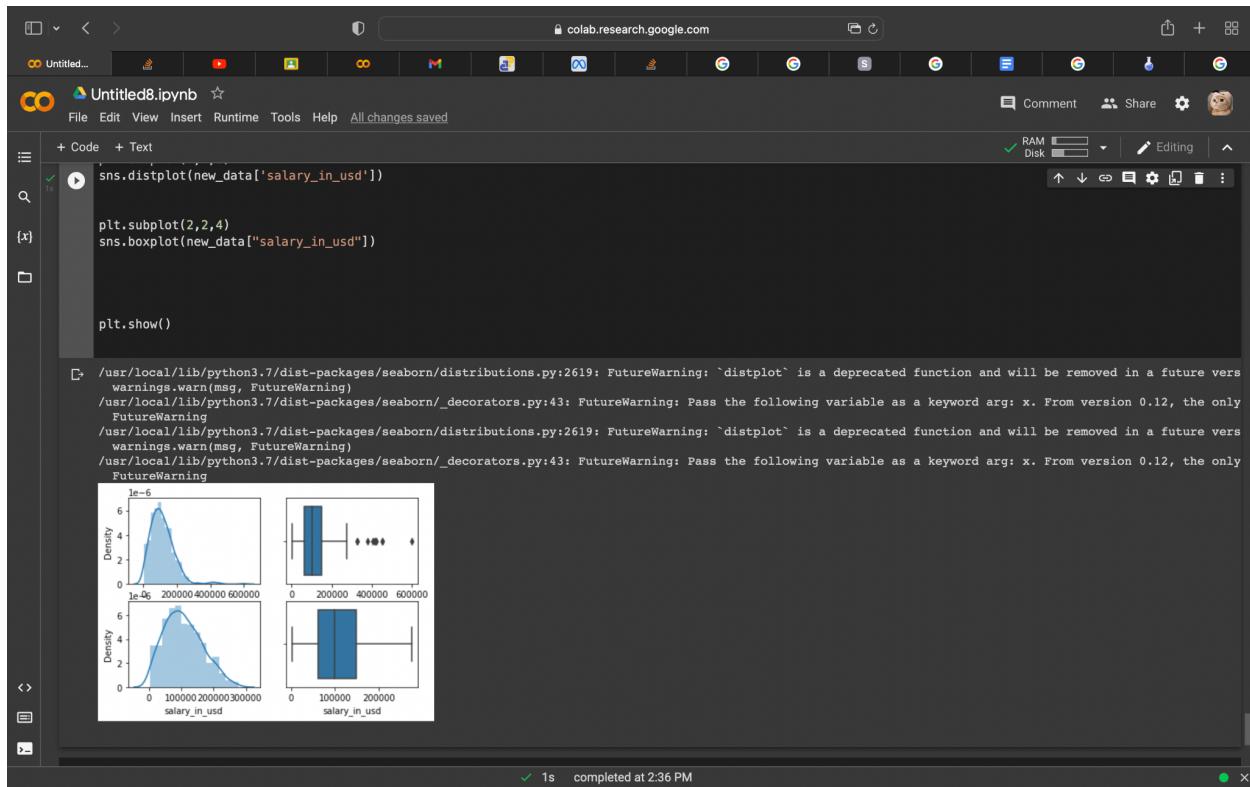
Below the preview, it says '607 rows x 11 columns'. At the bottom of the screen, there's a status bar with '0s completed at 2:31 PM'.

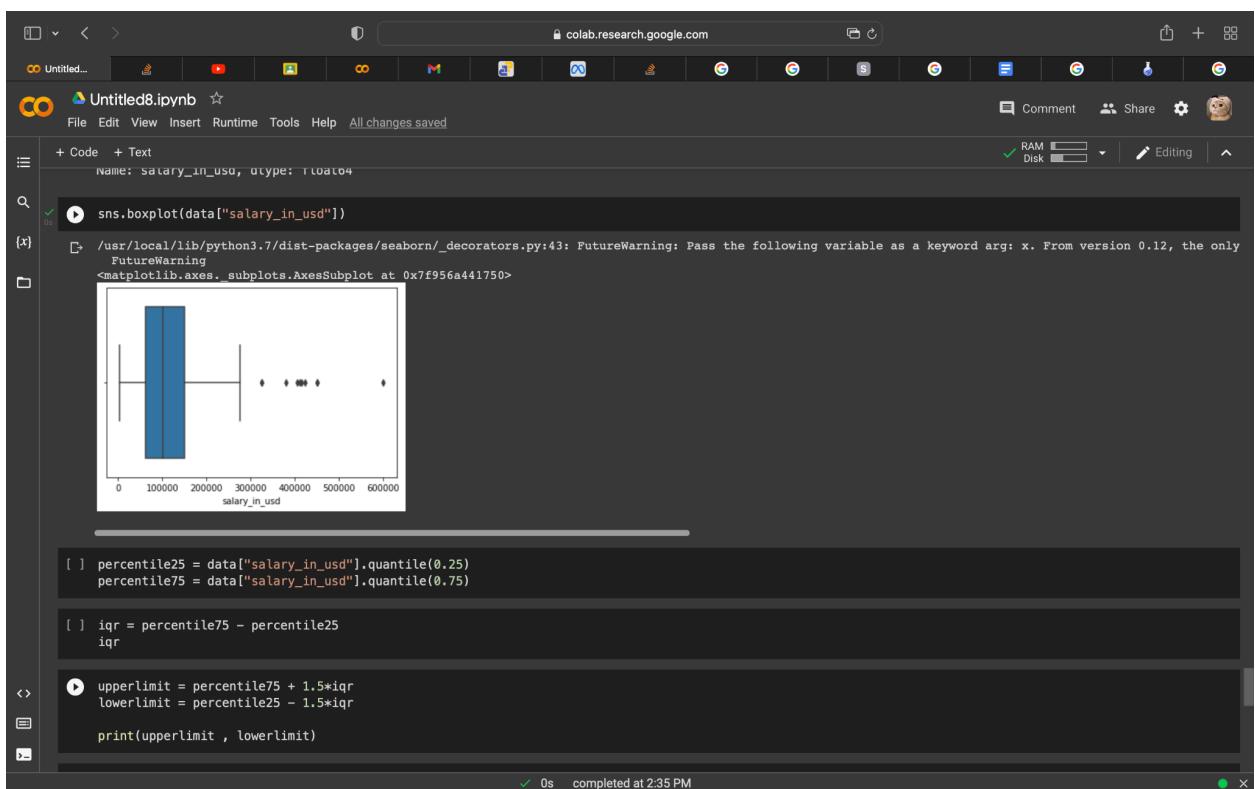
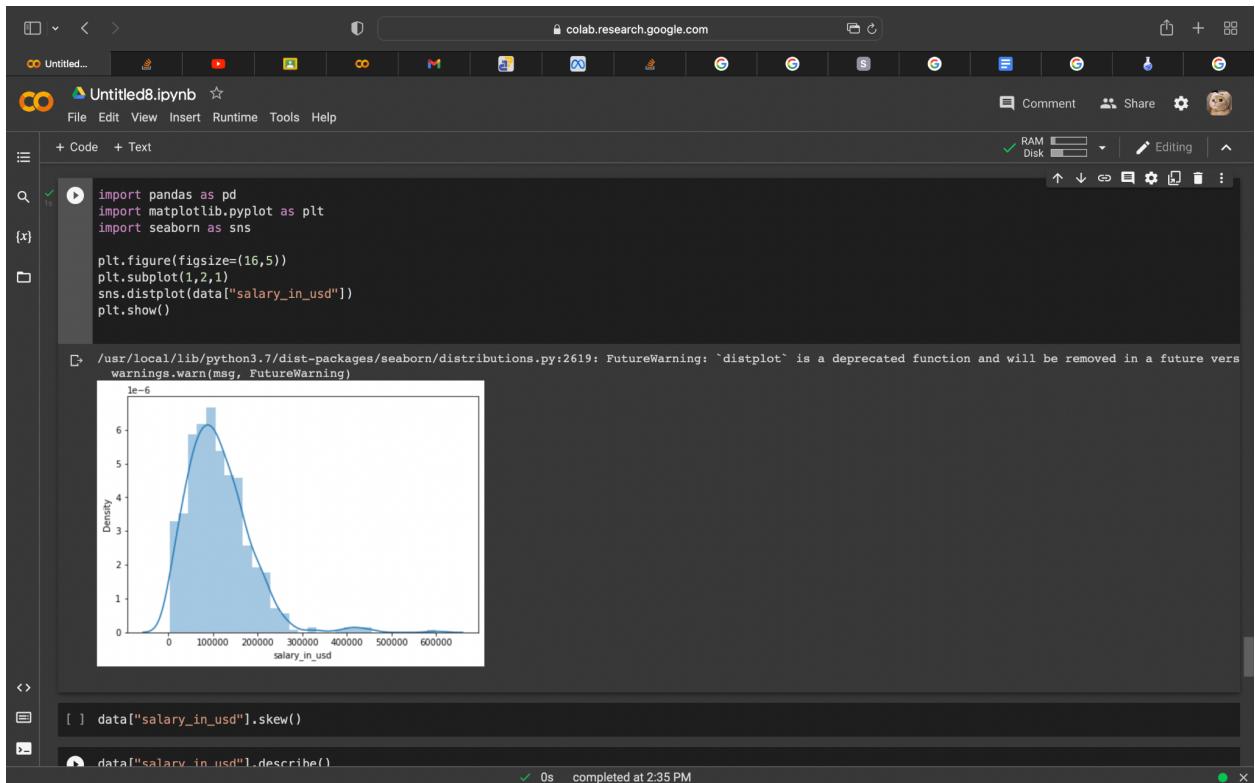
Que 3)

You can reduce noise in various ways , one of them is removing the outliers. Outliers are basically those **data points that are significantly different from the rest of the dataset**. They are often abnormal observations that skew the data distribution, and arise due to inconsistent data entry, or erroneous observations.

I have used the IQR method , since the data(i.e the salary_in_usd) is skewed in nature . We define a Interquartile Range which is (75 - 25)percentile and whiskers using the formula.all the points lying outside the whiskers are our outliers.

Here in the histogram plot you can observe our data is right skewed (skewness can be determined using `(data["salary_in_usd"].skew())`) In the box plot you can see the outliers on the right side (`1.664144558187621`).now we define a new data set whose value should be less than the upper limit . And basically after plotting the new data you can see that the outliers got removed .





Que 4)

Categorical data can be divided into 2 types Nominal data (here you cannot define a relationship between the data) and Ordinal data (here you can define a relationship between the data).Categorical data is present mostly in the form of strings and machine learning algorithms expect numbers , so it's your duty to convert these categories to numbers .there are various type of encoding , one is Ordinal encoding which is used for the Ordinal type of data and OneHotEncoding is used for Nominal data .

Firstly i splitted the data into test and train proving the job_title as x and y as salary_in_usd
The OneHotEncoder converts the data into arrays by using the hstack i stacked the x_train and the previous data column.

Que 4:

```
[36] dum_df = pd.get_dummies( data, drop_first=True)
dum_df

[7] #Nominal Encoding
from sklearn.model_selection import train_test_split

x_train,x_test , y_train , y_test  = train_test_split( data.iloc[:,5] , data.iloc[:,7] , test_size = 0.3 , random_state = 2021)

[8] x_train

[9] from sklearn.preprocessing import OneHotEncoder
ohe = OneHotEncoder(drop = "first" , sparse = False)

x_train_new = ohe.fit_transform(x_train[['job_title']])
x_test_new = ohe.fit_transform(x_test[['job_title']] )

[10] np.hstack((x_train[['job_title']].values , x_train_new))

array([['Director of Data Science', 0.0, 0.0, ..., 0.0, 0.0, 0.0],
       ['Data Engineer', 0.0, 0.0, ..., 0.0, 0.0, 0.0],
       ['ML Engineer', 0.0, 0.0, ..., 0.0, 0.0, 0.0],
       ...,
       ['Data Scientist', 0.0, 0.0, ..., 0.0, 0.0, 0.0],
       ['Data Engineer', 0.0, 0.0, ..., 0.0, 0.0, 0.0],
       ['Data Scientist', 0.0, 0.0, ..., 0.0, 0.0, 0.0], dtype=object)

[94] #ordinal Encoding
data1 = dataDropNa.iloc[:, -1 ]
```

```
[94] #ordinal Encoding
data1 = dataDropNa.iloc[:, -1 ]

data1.head
<bound method NDFrame.head of 0      L
1      S
2      M
3      S
5      L
.
602     M
603     M
604     M
605     M
606     L
Name: company_size, Length: 497, dtype: object>

[84] from sklearn.preprocessing import OrdinalEncoder
oe = OrdinalEncoder(categories=[['L' , 'M' , 'S']])

[90] oe.fit(np.array(data1).reshape(-1 ,1))
OrdinalEncoder(categories=[['L' , 'M' , 'S']]

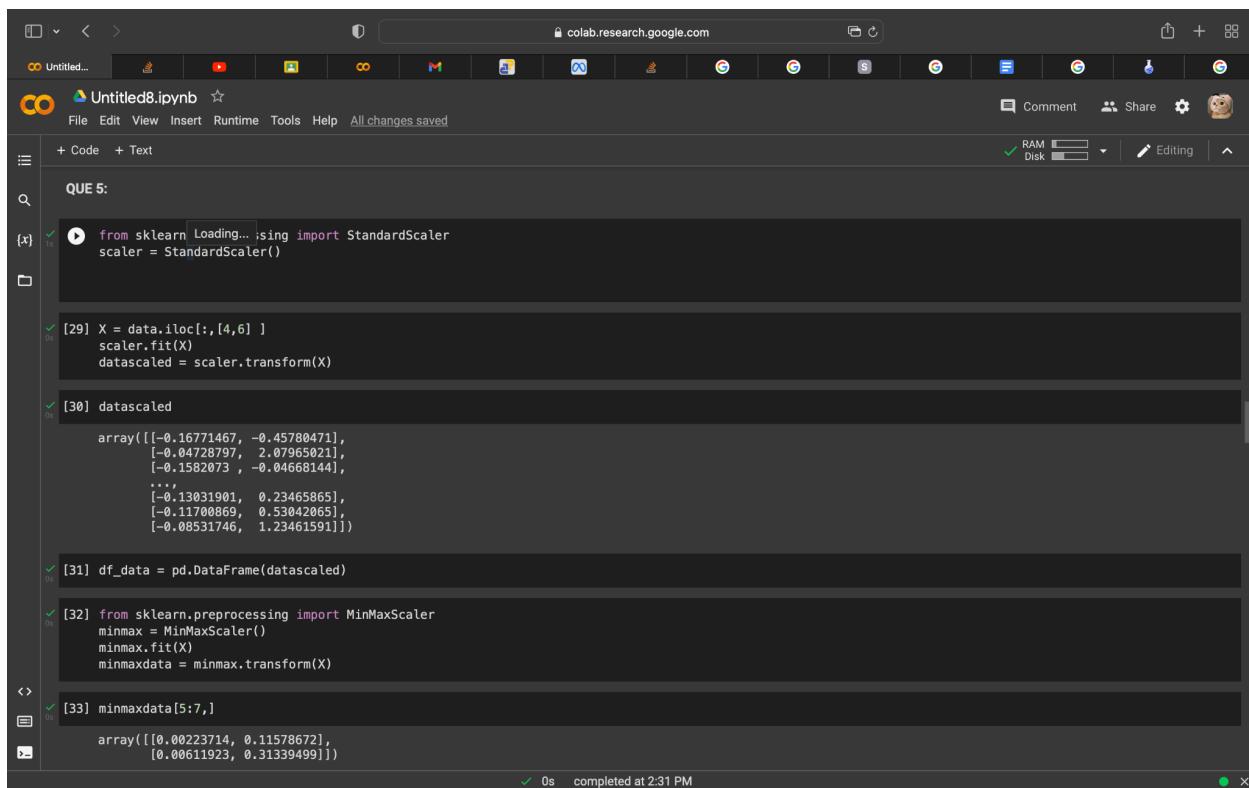
[91] x = oe.transform(np.array(data1).reshape(-1,1))
x

array([[0.1,
       [2.1,
        [1.1,
        [2.1,
        [0.1,
        [2.1,
```

Que 5)

You do scaling of the data to minimize the errors , making the numerical data relatable so the computer can minimize the errors.

Minmax basically converts the data into the range 0 to 1 .



The screenshot shows a Google Colab notebook titled "Untitled8.ipynb". The code cell at index [x] contains the following imports:

```
from sklearn.preprocessing import StandardScaler
scaler = StandardScaler()
```

The code cell at index [29] contains the following code to scale the data:

```
X = data.iloc[:,[4,6]]
scaler.fit(X)
datascaled = scaler.transform(X)
```

The code cell at index [30] displays the scaled data:

```
[30] datascaled
array([[-0.16771467, -0.45780471],
       [-0.04728797,  2.07965021],
       [-0.1582073 , -0.04668144],
       ...,
       [-0.13031901,  0.23465865],
       [-0.11700869,  0.53042065],
       [-0.08531746,  1.23461591]])
```

The code cell at index [31] creates a DataFrame from the scaled data:

```
[31] df_data = pd.DataFrame(datascaled)
```

The code cell at index [32] uses MinMaxScaler to scale the data:

```
from sklearn.preprocessing import MinMaxScaler
minmax = MinMaxScaler()
minmax.fit(X)
minmaxdata = minmax.transform(X)
```

The code cell at index [33] displays the scaled data from step 32:

```
[33] minmaxdata[5:7,:]
array([[0.00223714, 0.11578672],
       [0.00611923, 0.31339499]])
```

At the bottom of the interface, it shows "0s completed at 2:31 PM".

Que 6)
b)

Train test split , here you basically split the data into 2 sets , where you use the training set for the training purpose rather than training the whole data.

For dividing the data into 3 parts , I first put the train and validation set together (main)and then the splitted the main into train and validation.