



Islington college
(इरिलिङ्टन कलेज)

Choose a Module

Choose Coursework Percentage **Individual Coursework**

2023-24 Choose an item.

Student Name: ASHISH BUDHA

London Met ID: 22067313

College ID: NP01AI4A220031

Assignment Due Date: 13 May 2024

Assignment Submission Date: Monday, May 13, 2024

Word Count: 1083

I confirm that I understand my coursework needs to be submitted online via MySecondTeacher under the relevant module page before the deadline in order for my assignment to be accepted and marked. I am fully aware that late submissions will be treated as non-submission and a marks of zero will be awarded.

Contents

Data Understanding	4
What is the dataset about?	4
Data Preparation	5
Python Program to load data into pandas Framework.....	5
Python Program to remove unnecessary columns i.e., salary and salary currency.	5
Write a python program to remove NaN missing values from updated dataframe..	6
Python Program to check duplicate values in the dataframe	7
Python program to see the unique values from all the columns in the dataframe...	7
Renaming experience level column	10
Data Analysis	13
Python Program to show summary statistics of sum, mean, standard deviation, skewness, and kurtosis of any chosen variable	13
Write a python program to calculate and show the correlation of all variables	14
Data Exploration.....	16
Write a python program to find out the top 15 jobs. Make a bar graph of sales as well.	16
Which job has the highest salary? Illustrate with bar graph	18
Python program to find out salaries based on experience level.....	19
Python program to show histogram and boxplot of any chosen different variables	21
Table 1 : Data Table.....	4

Figure 1: Load Data.....	5
Figure 2: Drop Data.....	6
Figure 3: Remove NaN.....	6
Figure 4: Check Duplicate.....	7
Figure 5: Unique Work Year	7
Figure 6: Unique Experience Level	7
Figure 7: Unique Employment Type	8
Figure 8: Unique Job Title	8
Figure 9: Unique Salary.....	9
Figure 10: Unique Employee Residence	9
Figure 11: Unique Remote Ratio	10
Figure 12: Unique Company Location	10
Figure 13: Unique Company Size	10
Figure 14: Replacing SE	11
Figure 15: Replacing MI	11
Figure 16: Replacing EN	12
Figure 17: Replacing EX	12
Figure 18: Summary Statistics	13
Figure 19: Skewness.....	13
Figure 20: Kurtosis	14
Figure 21: Correlation.....	15
Figure 22: Correlation Non-Numeric.....	15
Figure 23: Top 15 Jobs.....	16
Figure 24: Highest Salary	18
Figure 25: Experience Level Salary.....	19
Figure 26: Histogram.....	21
Figure 27: Box Plot.....	22

Data Understanding

What is the dataset about?

The provided dataset contains different data of employees that plays a part in determining the salary of the employee. These variables are : 'work_year' , 'experience_level', 'employment_type', 'job_title', 'salary', 'salary_currency', 'salary_in_USD', 'employee_residence', 'remote_ratio', 'company_location', 'company_size'. We are to identify correlations between the employee's salary and other variables.

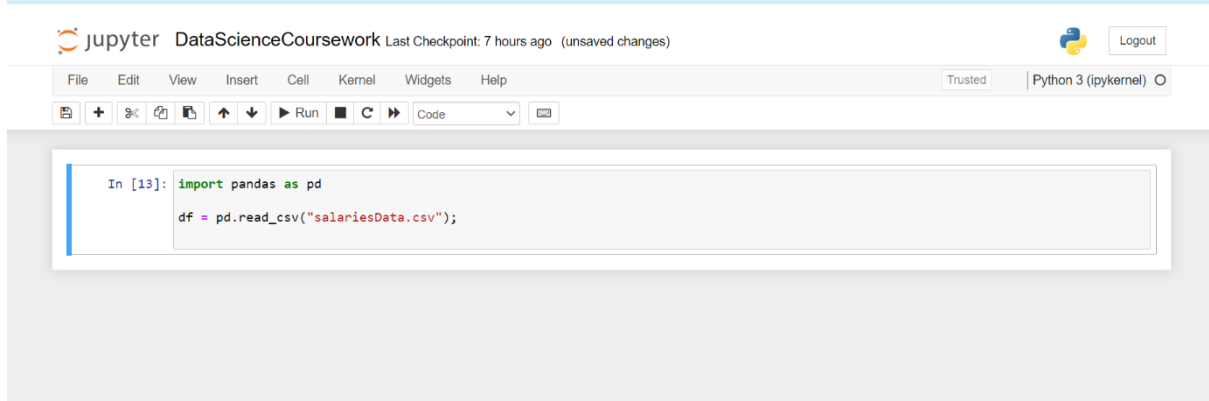
Table 1 : Data Table

S.No	Column Name	Description	Data Type
1	Work_year	This is a column that shows the year the employee worked in. We may be able to observe patterns in work year and salary, overtime a job may become more in demand or saturated, affecting the salary of the employee.	Integer
2	Experience_level	This is a column that shows the experience of the employee. The different experiences are : SE or Senior / Expert, MI or Medium Level / Intermediate, EN or Entry Level, EX or Executive Level.	String
3	Employment_type	This is a column showing the type of employment the employee is working as. Example it could be full type employment (FT) or a contract employment (CT)	
4	Job_title	This column shows us the job of the employee.	String
5	Salary	The Salary of the employee. This can be in any currency.	Integer
6	Salary Currency	The columns will tell us what the currency of the salary in the 5 th column is.	String
7	Salary_in_usd	This column will show the employees salary in USD, if in previous columns	Integer

		it was shown in another currency, it will be converted to USD.	
8	Employee_residence	This column shows the place the employee is currently residing in. People living in low cost of living could be paid less and vice versa.	String
9	Remote_ratio	This column shows the extent to which remote work is embraced by a company. This could affect the salary structure of the employees.	Integer
10	Company_location	This column shows the location of the country. Certain company locations could mean a higher / lower salary.	String
11	Company_size	The size of the company such as : S, M, and L.	String

Data Preparation

Python Program to load data into pandas Framework

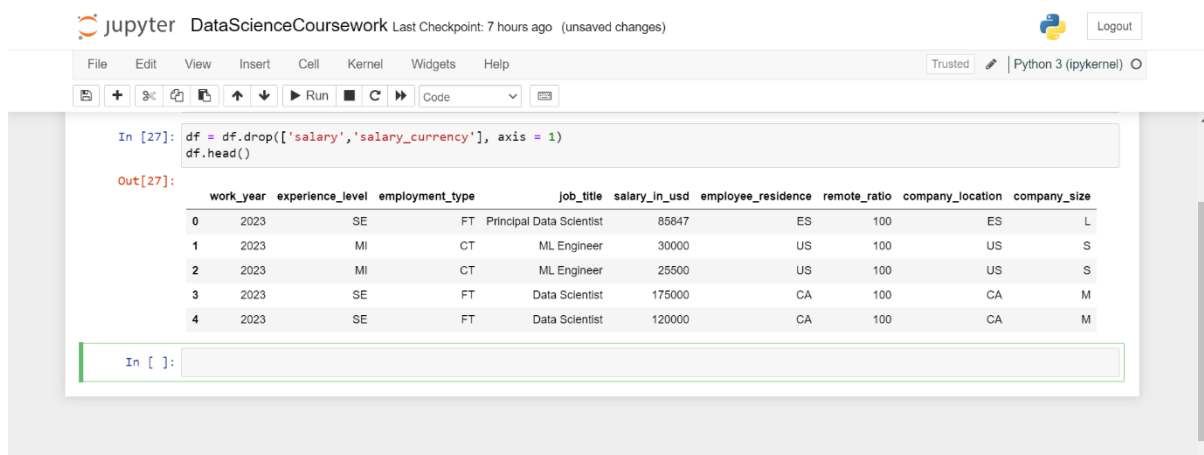


The screenshot shows a Jupyter Notebook interface. At the top, it says "jupyter DataScienceCoursework Last Checkpoint: 7 hours ago (unsaved changes)". There is a "Logout" button. Below the header is a menu bar with "File", "Edit", "View", "Insert", "Cell", "Kernel", "Widgets", and "Help". To the right of the menu bar are "Trusted" and "Python 3 (ipykernel)" labels. Below the menu bar is a toolbar with icons for file operations, running, and other notebook functions. The main area contains a code cell with the following Python code:

```
In [13]: import pandas as pd
df = pd.read_csv("salariesData.csv");
```

Figure 1: Load Data

Python Program to remove unnecessary columns i.e., salary and salary currency.



The image shows a Jupyter Notebook interface with the title "DataScienceCoursework". The top bar indicates "Last Checkpoint: 7 hours ago (unsaved changes)" and a "Logout" button. The menu bar includes File, Edit, View, Insert, Cell, Kernel, Widgets, and Help. The toolbar shows various icons for file operations and execution. The code cell contains the following Python code:

```
In [27]: df = df.drop(['salary', 'salary_currency'], axis = 1)
df.head()
```

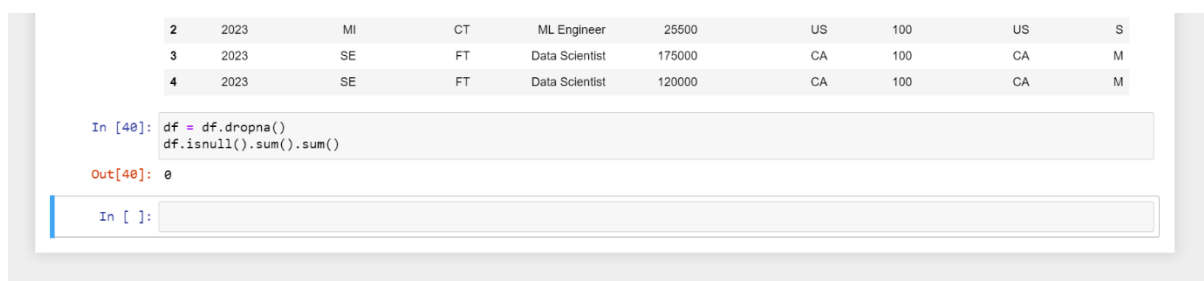
The output shows the first five rows of the DataFrame:

	work_year	experience_level	employment_type	job_title	salary_in_usd	employee_residence	remote_ratio	company_location	company_size
0	2023	SE	FT	Principal Data Scientist	85847	ES	100	ES	L
1	2023	MI	CT	ML Engineer	30000	US	100	US	S
2	2023	MI	CT	ML Engineer	25500	US	100	US	S
3	2023	SE	FT	Data Scientist	175000	CA	100	CA	M
4	2023	SE	FT	Data Scientist	120000	CA	100	CA	M

Figure 2: Drop Data

The columns “salary” and “salary_currency” are no longer seen on the DataFrame df.

Write a python program to remove NaN missing values from updated dataframe.



The image shows a Jupyter Notebook interface with the following Python code:

```
In [40]: df = df.dropna()
df.isnull().sum().sum()
```

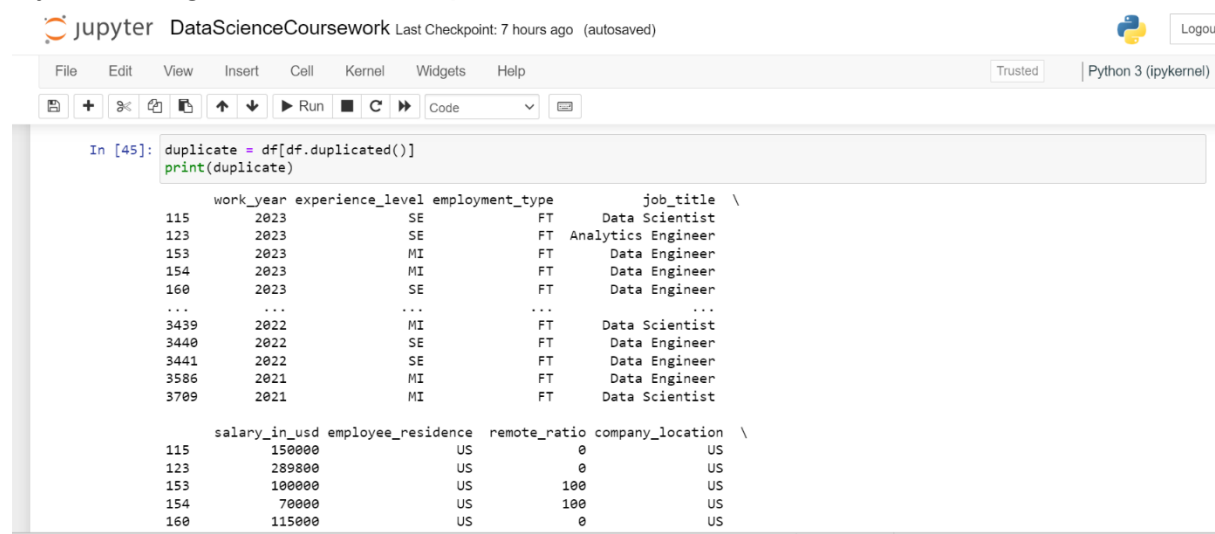
The output shows the result of the code:

```
Out[40]: 0
```

Figure 3: Remove NaN

Method called to drop NaN numbers. Checked to see if there are any NaN present in the new updated dataframe, there is no NaN present.

Python Program to check duplicate values in the dataframe



```
In [45]: duplicate = df[df.duplicated()]
print(duplicate)
```

	work_year	experience_level	employment_type	job_title
115	2023	SE	FT	Data Scientist
123	2023	SE	FT	Analytics Engineer
153	2023	MI	FT	Data Engineer
154	2023	MI	FT	Data Engineer
160	2023	SE	FT	Data Engineer
...
3439	2022	MI	FT	Data Scientist
3440	2022	SE	FT	Data Engineer
3441	2022	SE	FT	Data Engineer
3586	2021	MI	FT	Data Engineer
3709	2021	MI	FT	Data Scientist

	salary_in_usd	employee_residence	remote_ratio	company_location
115	150000	US	0	US
123	289000	US	0	US
153	100000	US	100	US
154	70000	US	100	US
160	115000	US	0	US

Figure 4: Check Duplicate

The `.duplicated()` method of the dataframe is used to check for duplicate values. It returns a boolean series, indicating a value of true if the row is a duplicate of a previous row.

This boolean Series is then used as condition for printing the dataframe rows.

Python program to see the unique values from all the columns in the dataframe.

```
[1171 rows x 9 columns]

In [47]: #Unique Values from column 'work_year'
print(df['work_year'].unique())

[2023 2022 2020 2021]
```

Figure 5: Unique Work Year

```
In [48]: #Unique values from column 'experience_level'
print(df['experience_level'].unique())

['SE' 'MI' 'EN' 'EX']

In [ ]:
```

Figure 6: Unique Experience Level

```
In [50]: #Unique values form column 'employment_type'
print(df['employment_type'].unique())

['FT' 'CT' 'FL' 'PT']
```

```
In [ ]:
```

Figure 7: Unique Employment Type

```
In [52]: #Unique values form column 'job_title'
print(df['job_title'].unique())

['Principal Data Scientist' 'ML Engineer' 'Data Scientist'
 'Applied Scientist' 'Data Analyst' 'Data Modeler' 'Research Engineer'
 'Analytics Engineer' 'Business Intelligence Engineer'
 'Machine Learning Engineer' 'Data Strategist' 'Data Engineer'
 'Computer Vision Engineer' 'Data Quality Analyst'
 'Compliance Data Analyst' 'Data Architect'
 'Applied Machine Learning Engineer' 'AI Developer' 'Research Scientist'
 'Data Analytics Manager' 'Business Data Analyst' 'Applied Data Scientist'
 'Staff Data Analyst' 'ETL Engineer' 'Data DevOps Engineer' 'Head of Data'
 'Data Science Manager' 'Data Manager' 'Machine Learning Researcher'
 'Big Data Engineer' 'Data Specialist' 'Lead Data Analyst'
 'BI Data Engineer' 'Director of Data Science'
 'Machine Learning Scientist' 'MLOps Engineer' 'AI Scientist'
 'Autonomous Vehicle Technician' 'Applied Machine Learning Scientist'
 'Lead Data Scientist' 'Cloud Database Engineer' 'Financial Data Analyst'
 'Data Infrastructure Engineer' 'Software Data Engineer' 'AI Programmer'
 'Data Operations Engineer' 'BI Developer' 'Data Science Lead']
```

Figure 8: Unique Job Title


```

'Data Infrastructure Engineer' 'Software Data Engineer' 'AI Programmer'
'Data Operations Engineer' 'BI Developer' 'Data Science Lead'
'Deep Learning Researcher' 'BI Analyst' 'Data Science Consultant'
'Data Analytics Specialist' 'Machine Learning Infrastructure Engineer'
'BI Data Analyst' 'Head of Data Science' 'Insight Analyst'
'Deep Learning Engineer' 'Machine Learning Software Engineer'
'Big Data Architect' 'Product Data Analyst'
'Computer Vision Software Engineer' 'Azure Data Engineer'
'Marketing Data Engineer' 'Data Analytics Lead' 'Data Lead'
'Data Science Engineer' 'Machine Learning Research Engineer'
'NLP Engineer' 'Manager Data Management' 'Machine Learning Developer'
'3D Computer Vision Researcher' 'Principal Machine Learning Engineer'
'Data Analytics Engineer' 'Data Analytics Consultant'
'Data Management Specialist' 'Data Science Tech Lead'
'Data Scientist Lead' 'Cloud Data Engineer' 'Data Operations Analyst'
'Marketing Data Analyst' 'Power BI Developer' 'Product Data Scientist'
'Principal Data Architect' 'Machine Learning Manager'
'Lead Machine Learning Engineer' 'ETL Developer' 'Cloud Data Architect'
'Lead Data Engineer' 'Head of Machine Learning' 'Principal Data Analyst'
'Principal Data Engineer' 'Staff Data Scientist' 'Finance Data Analyst']

```

Figure 8.1: Unique Job Title

```

In [55]: #Unique values from column 'salary_in_usd'
print(df['salary_in_usd'].unique())

[ 85847  30000  25500 ... 28369 412000  94665]

```

Figure 9: Unique Salary

```

In [57]: #Unique values from column 'employee_residence'
print(df['employee_residence'].unique())

['ES' 'US' 'CA' 'DE' 'GB' 'NG' 'IN' 'HK' 'PT' 'NL' 'CH' 'CF' 'FR' 'AU'
'FI' 'UA' 'IE' 'IL' 'GH' 'AT' 'CO' 'SG' 'SE' 'SI' 'MX' 'UZ' 'BR' 'TH'
'HR' 'PL' 'KW' 'VN' 'CY' 'AR' 'AM' 'BA' 'KE' 'GR' 'MK' 'LV' 'RO' 'PK'
'IT' 'MA' 'LT' 'BE' 'AS' 'IR' 'HU' 'SK' 'CN' 'CZ' 'CR' 'TR' 'CL' 'PR'
'DK' 'BO' 'PH' 'DO' 'EG' 'ID' 'AE' 'MY' 'JP' 'EE' 'HN' 'TN' 'RU' 'DZ'
'IQ' 'BG' 'JE' 'RS' 'NZ' 'MD' 'LU' 'MT']

```

Figure 10: Unique Employee Residence

```
In [58]: #Unique values from column 'remote_ratio'
print(df['remote_ratio'].unique())

[100    0   50]
```

Figure 11: Unique Remote Ratio

```
In [59]: #Unique values from column 'company_location'
print(df['company_location'].unique())

['ES' 'US' 'CA' 'DE' 'GB' 'NG' 'IN' 'HK' 'NL' 'CH' 'CF' 'FR' 'FI' 'UA'
 'IE' 'IL' 'GH' 'CO' 'SG' 'AU' 'SE' 'SI' 'MX' 'BR' 'PT' 'RU' 'TH' 'HR'
 'VN' 'EE' 'AM' 'BA' 'KE' 'GR' 'MK' 'LV' 'RO' 'PK' 'IT' 'MA' 'PL' 'AL'
 'AR' 'LT' 'AS' 'CR' 'IR' 'BS' 'HU' 'AT' 'SK' 'CZ' 'TR' 'PR' 'DK' 'BO'
 'PH' 'BE' 'ID' 'EG' 'AE' 'LU' 'MY' 'HN' 'JP' 'DZ' 'IQ' 'CN' 'NZ' 'CL'
 'MD' 'MT']
```

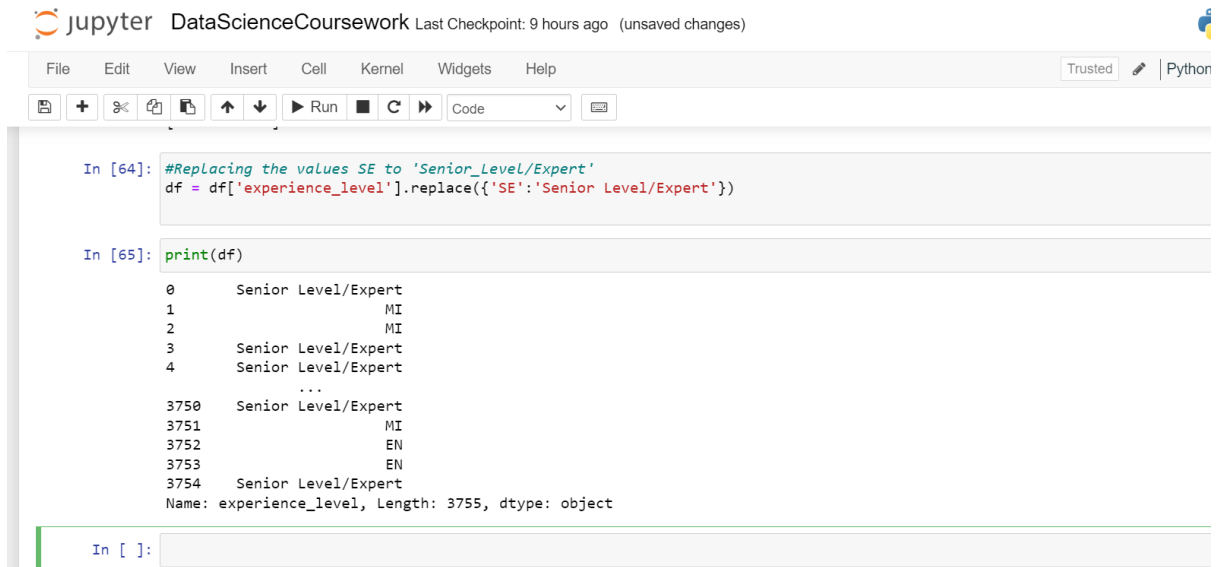
Figure 12: Unique Company Location

```
In [61]: #Unique values form column 'company_size'
print(df['company_size'].unique())

['L' 'S' 'M']
```

Figure 13: Unique Company Size

Renaming experience level column



Jupyter DataScienceCoursework Last Checkpoint: 9 hours ago (unsaved changes)

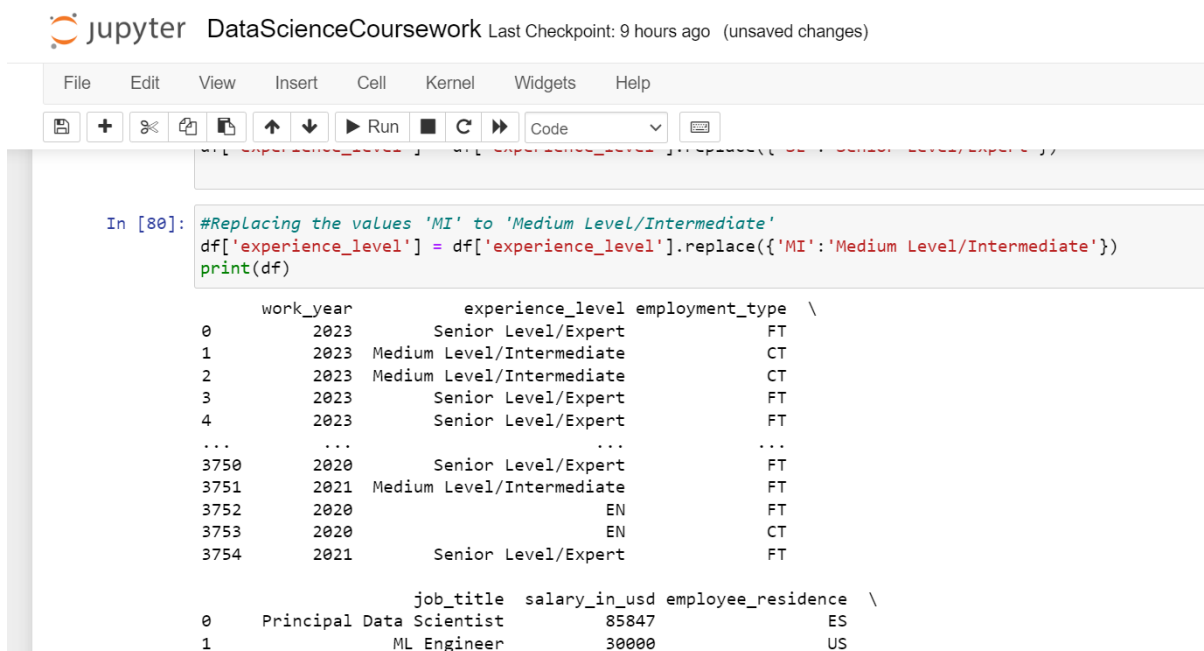
File Edit View Insert Cell Kernel Widgets Help Trusted Python

```
In [64]: #Replacing the values SE to 'Senior Level/Expert'
df = df['experience_level'].replace({'SE': 'Senior Level/Expert'})

In [65]: print(df)
0      Senior Level/Expert
1              MI
2              MI
3      Senior Level/Expert
4      Senior Level/Expert
...
3750   Senior Level/Expert
3751              MI
3752              EN
3753              EN
3754   Senior Level/Expert
Name: experience_level, Length: 3755, dtype: object
```

In []:

Figure 14: Replacing SE



Jupyter DataScienceCoursework Last Checkpoint: 9 hours ago (unsaved changes)

File Edit View Insert Cell Kernel Widgets Help

```
In [80]: #Replacing the values 'MI' to 'Medium Level/Intermediate'
df['experience_level'] = df['experience_level'].replace({'MI': 'Medium Level/Intermediate'})
print(df)
```

	work_year	experience_level	employment_type
0	2023	Senior Level/Expert	FT
1	2023	Medium Level/Intermediate	CT
2	2023	Medium Level/Intermediate	CT
3	2023	Senior Level/Expert	FT
4	2023	Senior Level/Expert	FT
...
3750	2020	Senior Level/Expert	FT
3751	2021	Medium Level/Intermediate	FT
3752	2020	EN	FT
3753	2020	EN	CT
3754	2021	Senior Level/Expert	FT

	job_title	salary_in_usd	employee_residence
0	Principal Data Scientist	85847	ES
1	ML Engineer	30000	US

Figure 15: Replacing MI

jupyter DataScienceCoursework Last Checkpoint: 9 hours ago (unsaved changes)

File Edit View Insert Cell Kernel Widgets Help

Run

```
In [81]: #Replacing the values of 'EN' to 'Entry Level'
df['experience_level'] = df['experience_level'].replace({'EN':'Entry Level'})
print(df['experience_level'])
```

```
0      Senior Level/Expert
1      Medium Level/Intermediate
2      Medium Level/Intermediate
3      Senior Level/Expert
4      Senior Level/Expert
...
3750     Senior Level/Expert
3751     Medium Level/Intermediate
3752              Entry Level
3753              Entry Level
3754     Senior Level/Expert
Name: experience_level, Length: 3755, dtype: object
```

In []:

Figure 16: Replacing EN

jupyter DataScienceCoursework Last Checkpoint: 9 hours ago (autosaved)

File Edit View Insert Cell Kernel Widgets Help

Run

```
In [82]: #Replacing the values of 'EX' to 'Executive Level'
df['experience_level'] = df['experience_level'].replace({'EX':'Executive Level'})
print(df['experience_level'])
```

```
0      Senior Level/Expert
1      Medium Level/Intermediate
2      Medium Level/Intermediate
3      Senior Level/Expert
4      Senior Level/Expert
...
3750     Senior Level/Expert
3751     Medium Level/Intermediate
3752              Entry Level
3753              Entry Level
3754     Senior Level/Expert
Name: experience_level, Length: 3755, dtype: object
```

Figure 17: Replacing EX

Data Analysis

Python Program to show summary statistics of sum, mean, standard deviation, skewness, and kurtosis of any chosen variable

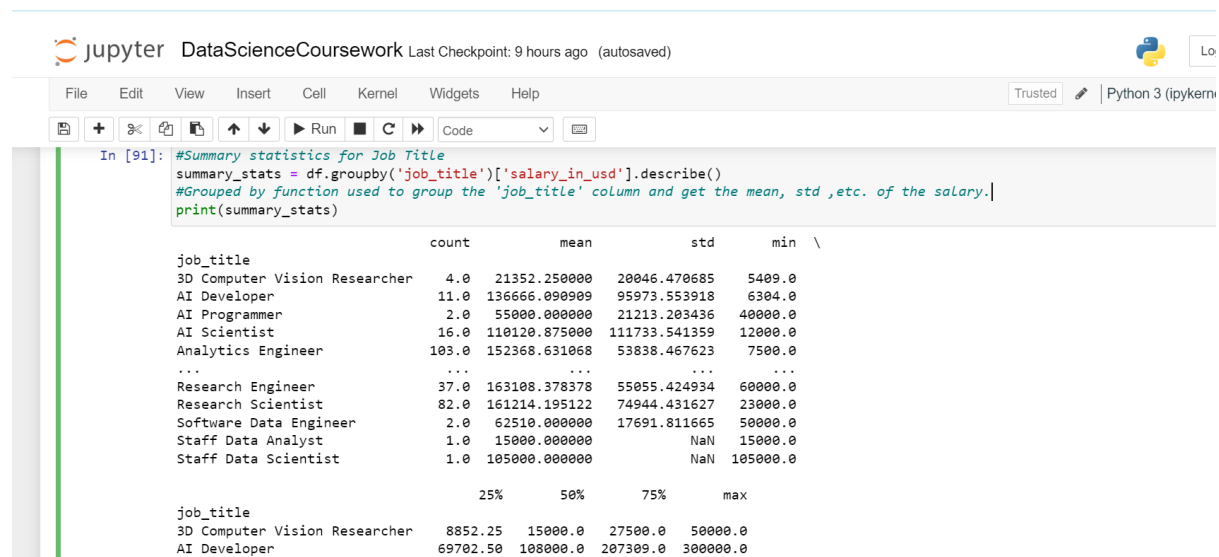


Figure 18: Summary Statistics

Summary Stats are being shown here using the `.describe()` method of the dataframe object. The dataframe is first grouped the on 'job title' and then we calculated the summary statistics of the 'salary' based on 'job title'

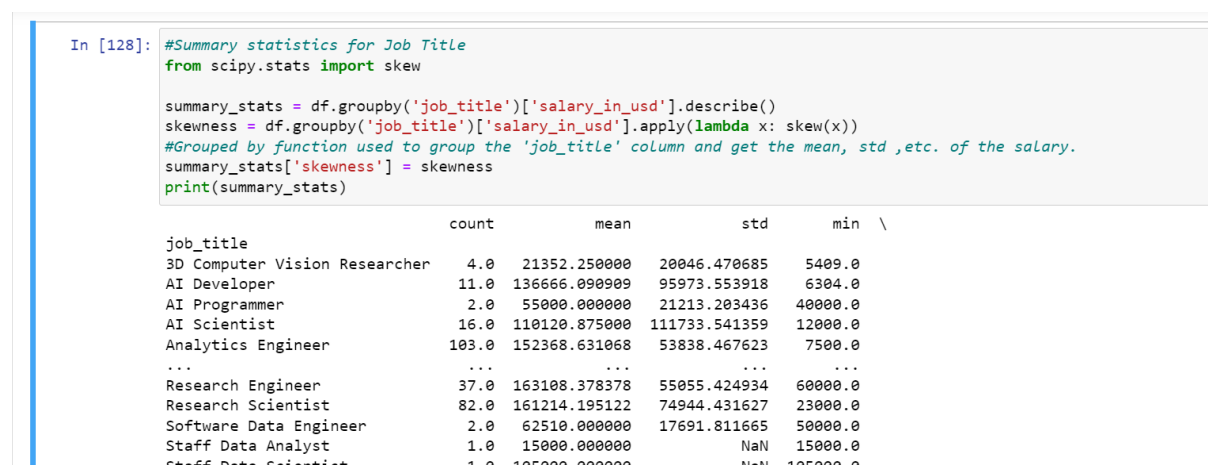


Figure 19: Skewness

Here I have used the scipy libraries' skew() method to get the skew of the data.

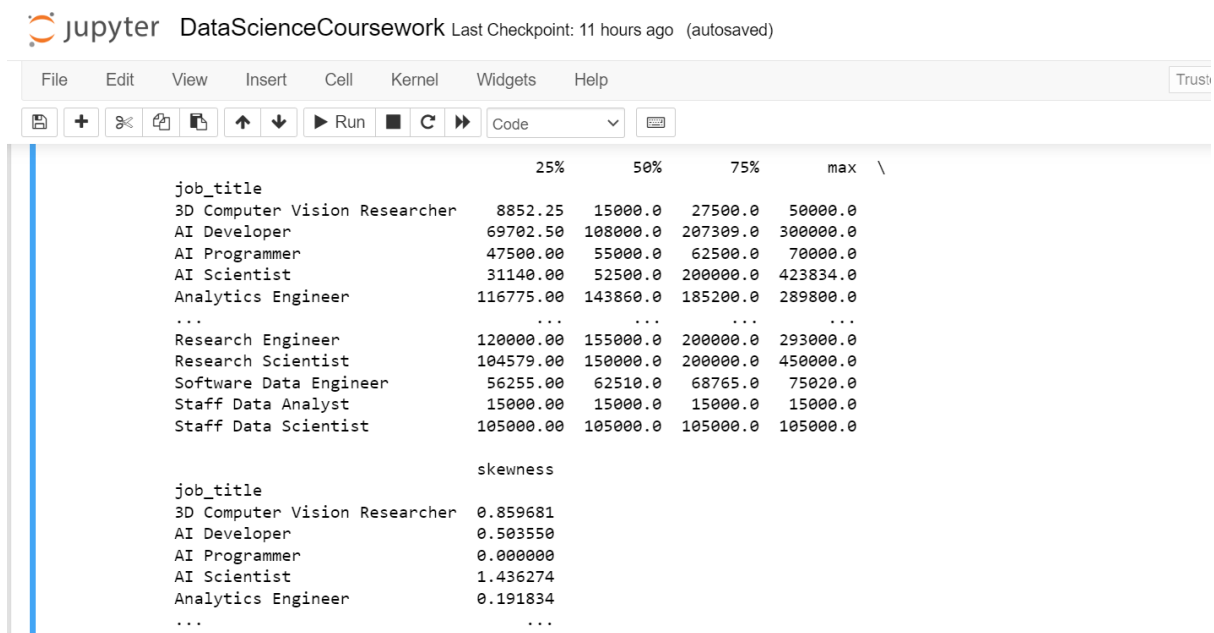


Figure 19.1: Skewness

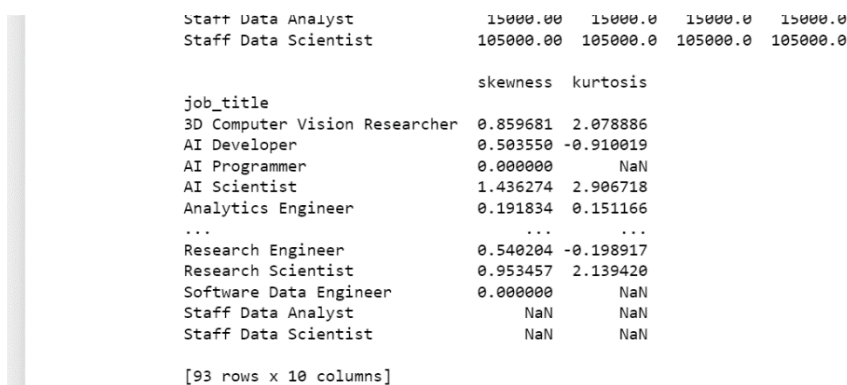


Figure 20: Kurtosis

Skewness values for the job title 'Staff Data Analyst' and 'Staff Data Scientist' is NaN as there is only 1 person for each of the job title. Standard Deviation is NaN as well.

Write a python program to calculate and show the correlation of all variables

```
In [130]: correlation = df.corr()
print(correlation)
```

	work_year	salary_in_usd	remote_ratio
work_year	1.00000	0.228290	-0.236430
salary_in_usd	0.22829	1.000000	-0.064171
remote_ratio	-0.23643	-0.064171	1.000000

C:\Users\ashish\AppData\Local\Temp\ipykernel_32764\1693706327.py:1: FutureWarning: The def

Figure 21: Correlation

From this output, we can see that work_year and salary_in_usd has a positive correlation of 0.22, which means that an increase in the year would likely mean an increase in salary.

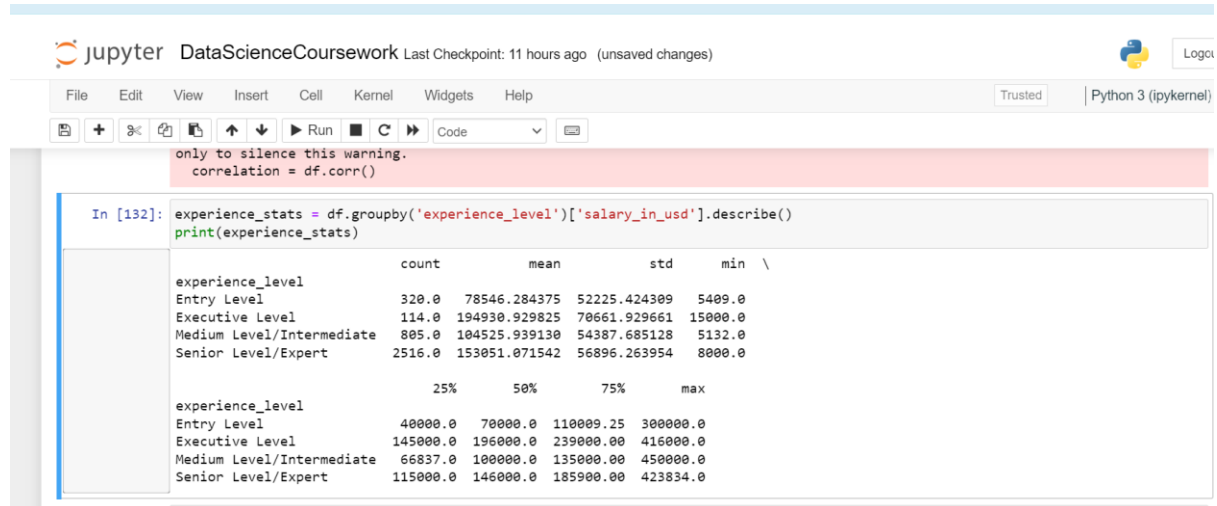


Figure 22: Correlation Non-Numeric

We can observe a correlation between the non-numeric data (experience_level) and the salary variable. Experience level 'Executive' earns the highest, followed by 'Senior Level/Expert', followed by 'Medium Level/Intermediate', followed by 'Entry Level'.

Data Exploration

Write a python program to find out the top 15 jobs. Make a bar graph of sales as well.



Figure 23: Top 15 Jobs

Here we have grouped the data by 'job title' and then selected the 'salary_in_usd' column which we then calculated the mean for.

The `plt.figure(figsize = (10,6))` creates a new figure with a size of 10 (width) by 6 (height).

We then used the `top15 Series` value to plot the bar. With its x-label being defined using the method `.xlabel()` and y-label being defined using the `.ylabel()`

The `.xticks()` method rotates the label for better readability.

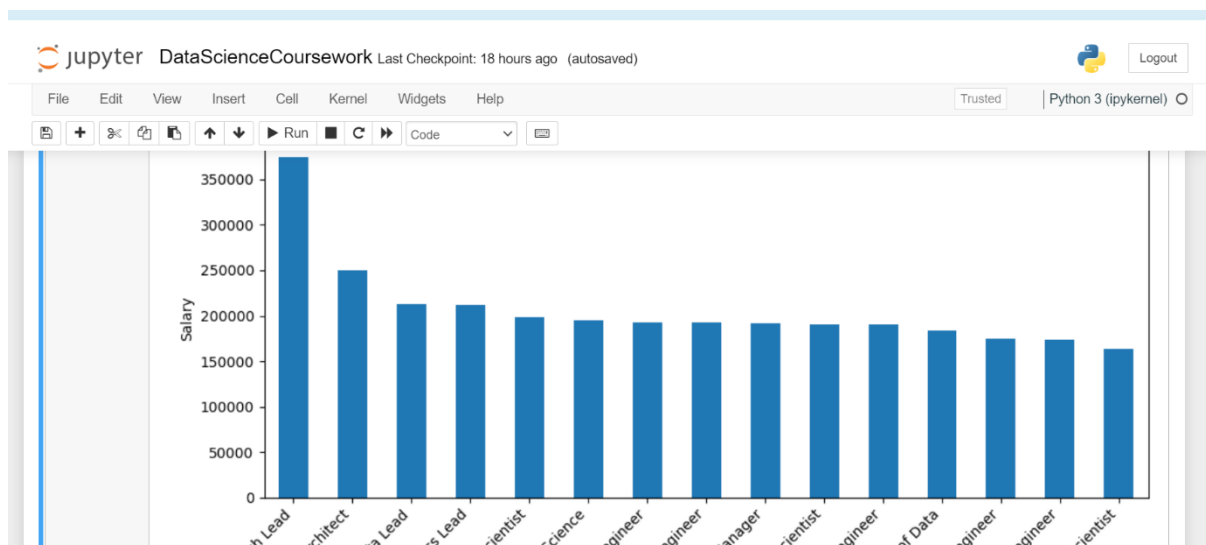


Figure 23.1: Top 15 Jobs

The figure that was shown when executing the code, here we can see top 15 jobs on the basis of the jobs salary mean value.

Which job has the highest salary? Illustrate with bar graph

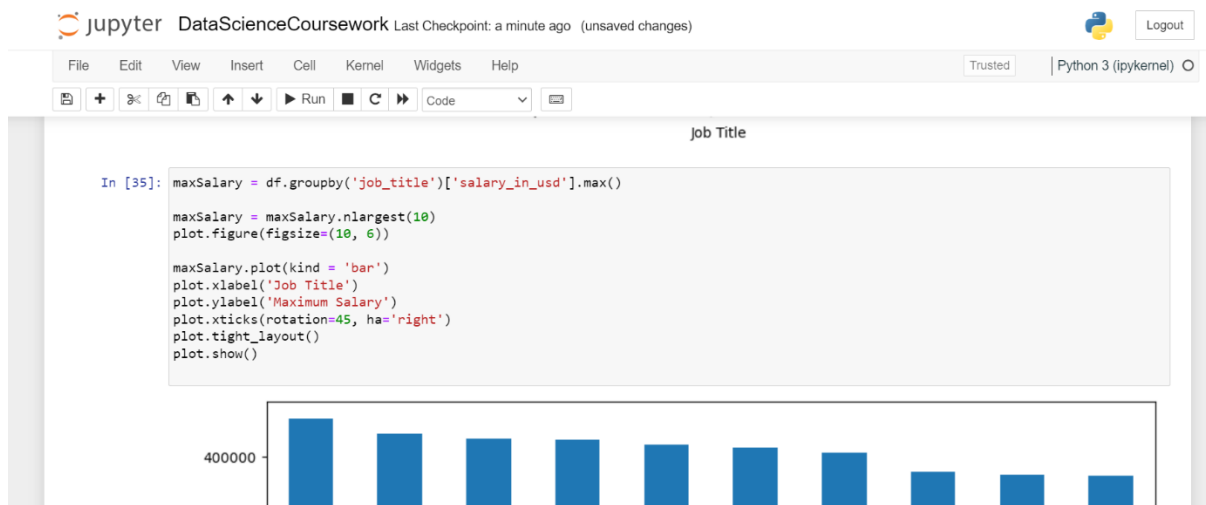


Figure 24: Highest Salary

We grouped the data by 'job title' and then selected the 'salary_in_usd' column which we then identified the maximum value.

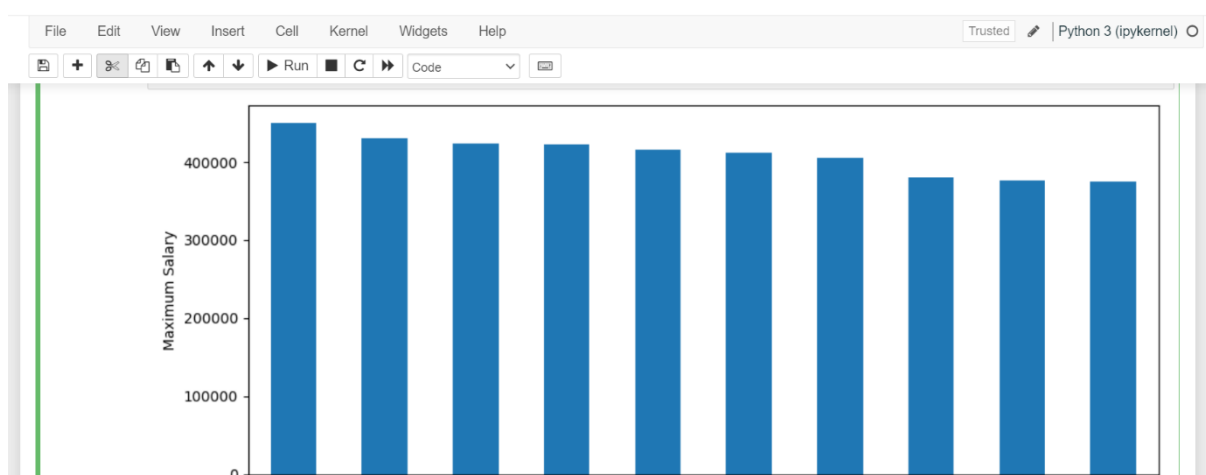


Figure 24.1: Highest Salary

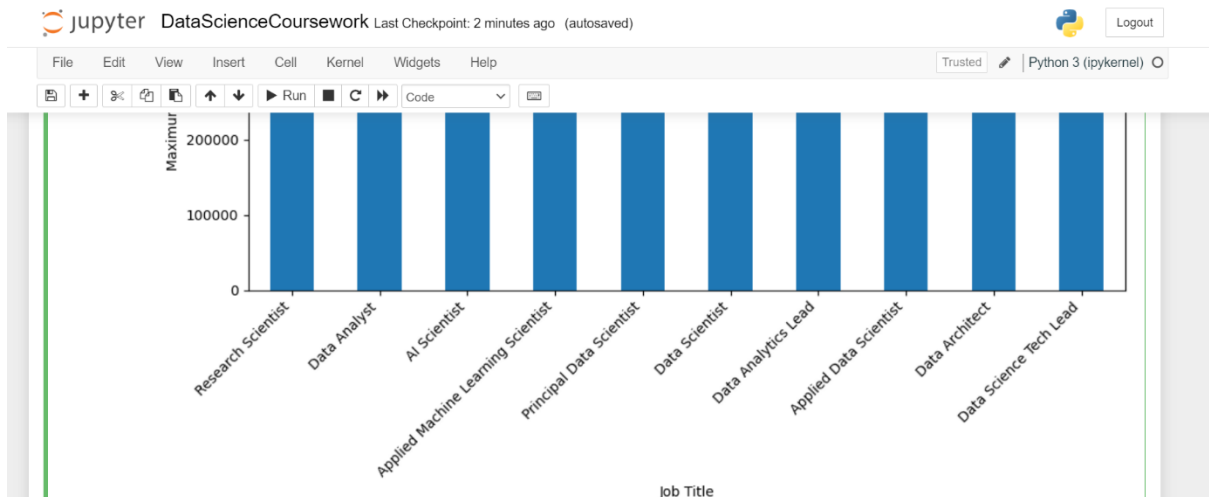


Figure 24.2: Highest Salary

From the bar graph, it shows that the job 'Research Scientist' has the highest salaries

Python program to find out salaries based on experience level.

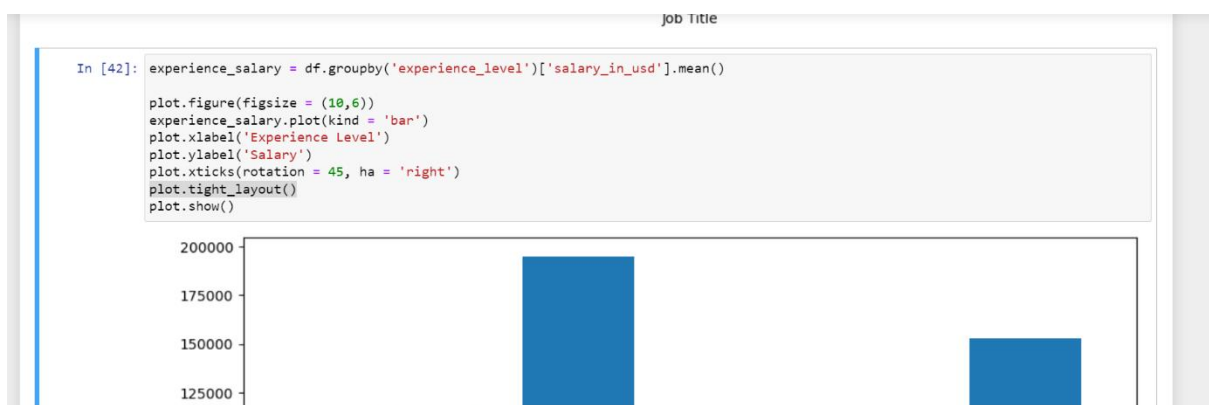


Figure 25: Experience Level Salary

Here the dataframe is grouped by 'experience_level' and then it's 'salary_in_usd' column selected, then the mean of the selected column is calculated for each group.

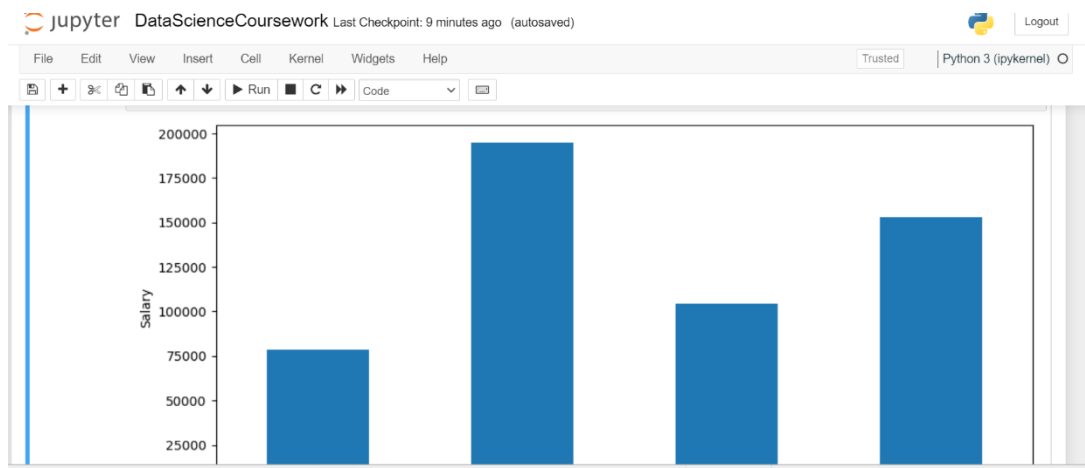


Figure 25.1: Experience Level Salary

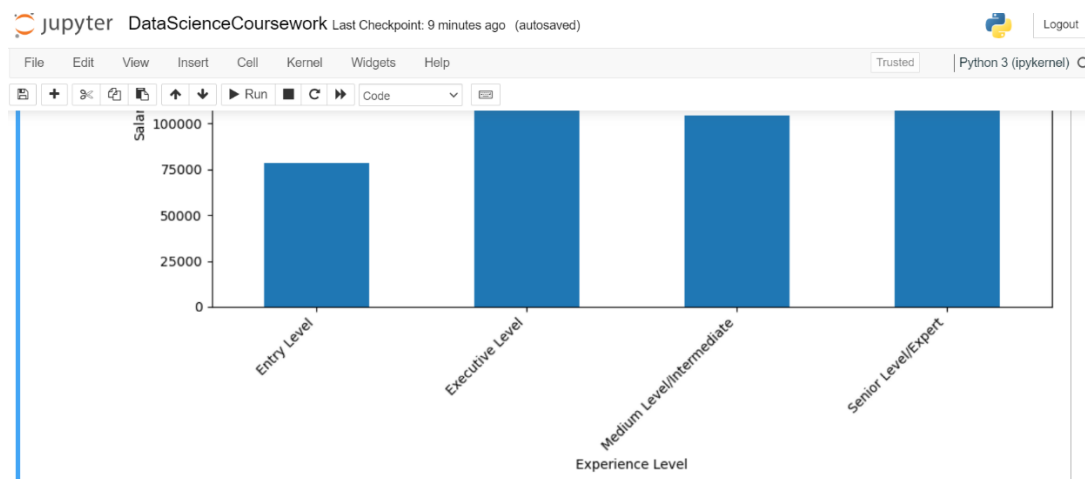


Figure 25.2: Experience Level Salary

This bar chart shows that people with experience 'Executive Level' is paid the highest, followed by 'Senior Level/Expert', followed by 'Medium Level / Intermediate', followed by 'Entry Level'

Python program to show histogram and boxplot of any chosen different variables

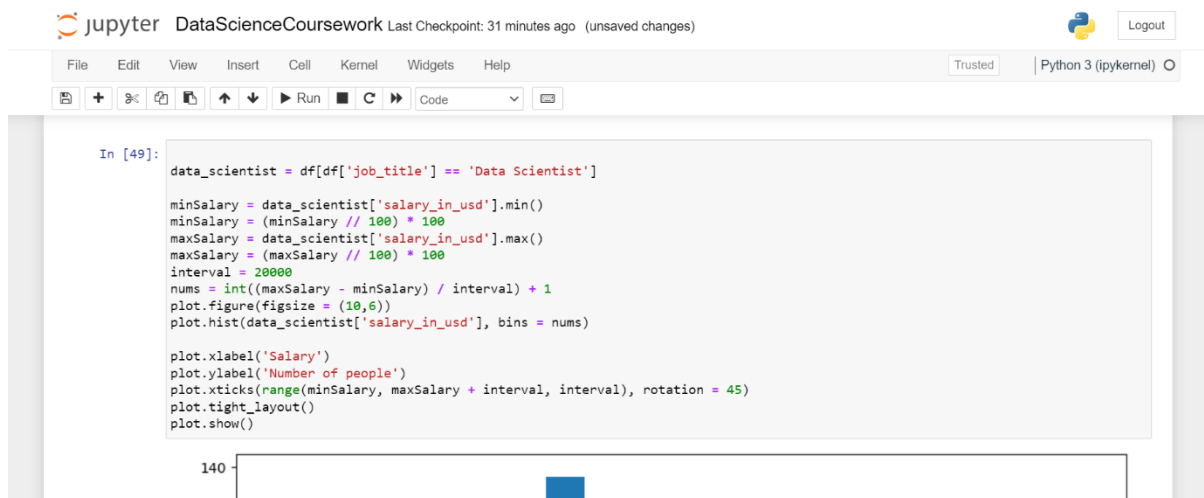


Figure 26: Histogram

Here the statement creates a histogram based on salary as the interval and the height of the histogram representing the number of people in the job 'Data Scientist' having an earning within that specific interval.

The $\text{minSalary} = (\text{minSalary} // 100) * 100$ is used to round down the value to the nearest hundred.

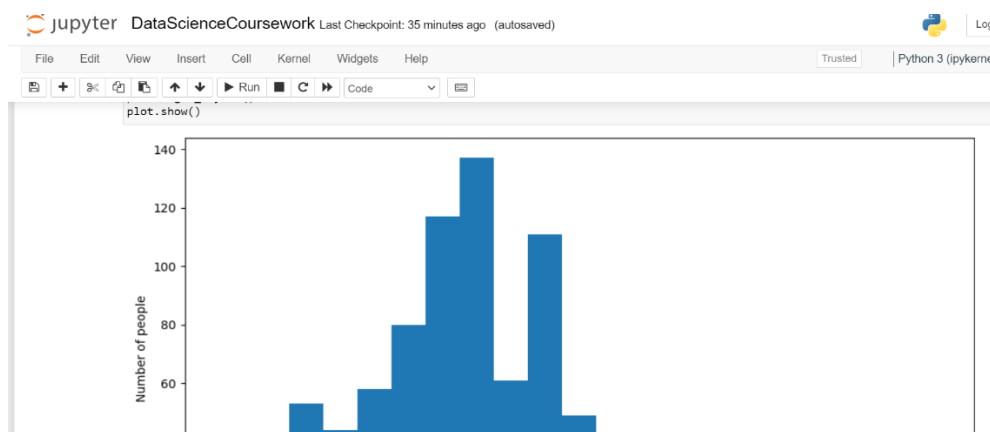


Figure 26.1: Histogram

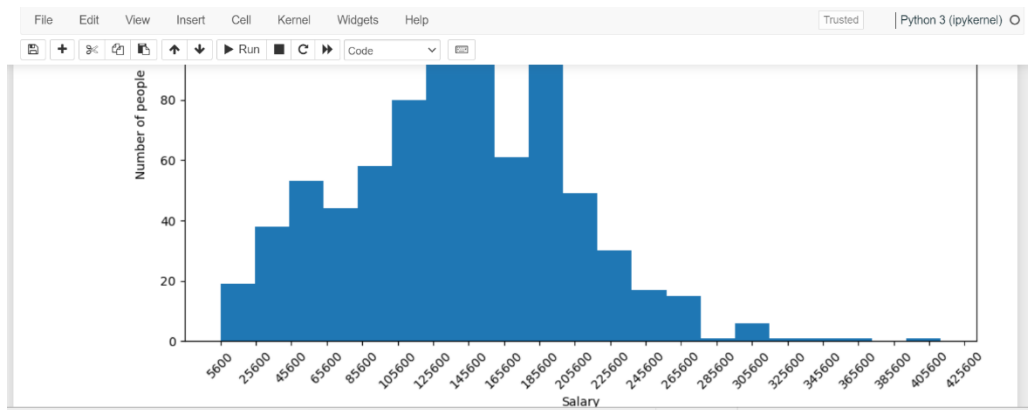


Figure 26.2: Histogram

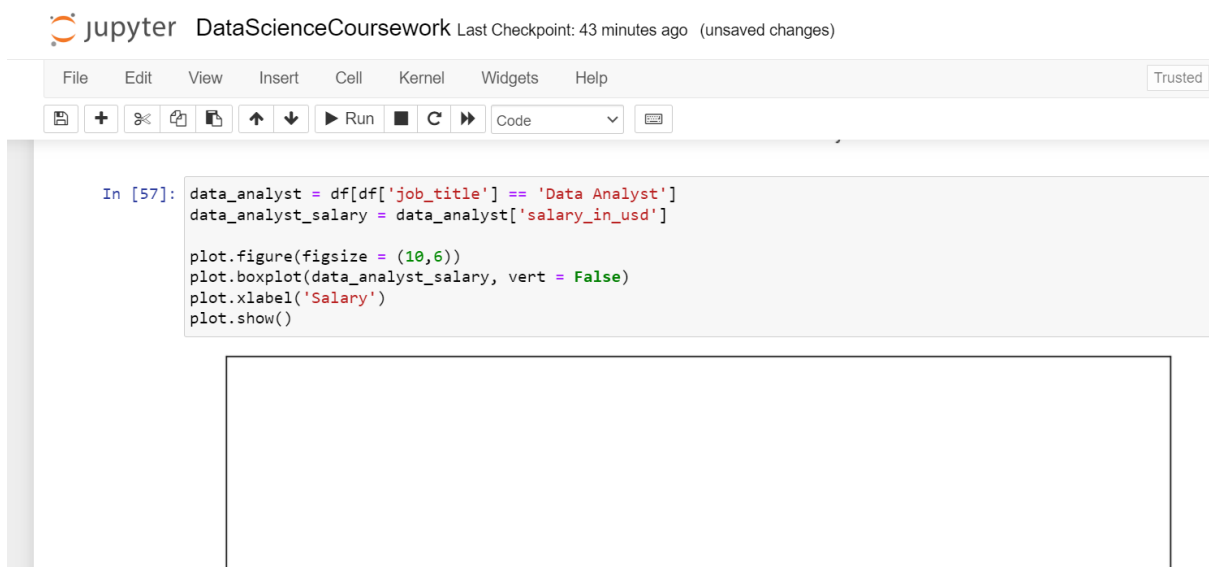


Figure 27: Box Plot

This code creates a box plot of salary of people with the job title 'Data Analyst'

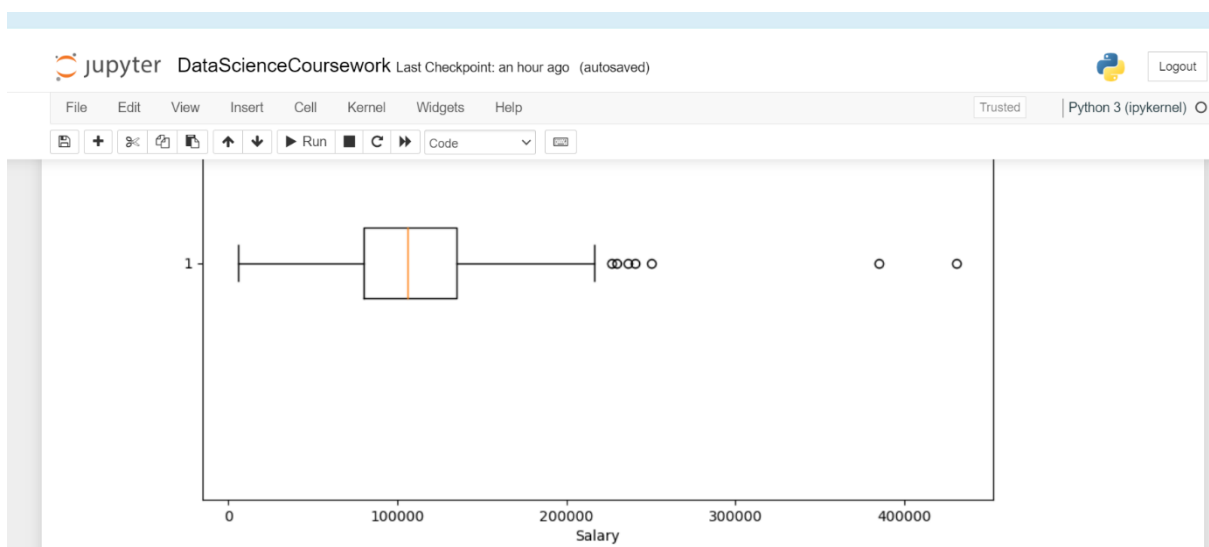


Figure 27.1: Box Plot

We can see the Q1, Q3, median, min, max, and outliers from this boxplot.