



## **APPLIED DATA SCIENCE**

### **60% Individual Coursework**

**2023-24 Autumn**

**Student Name: AASUTOSH KUMAR VERMA**

**London Met ID: 22085760**

**College ID: NP01AI4S230020**

**Assignment Due Date: Monday, May 13, 2024**

**Assignment Submission Date: Sunday, May 12, 2024**

**Word Count: 1855**

*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

1. Data Understanding.....	3
2. Data Preparation.....	4
3. Data Analysis. ....	12
4. Data Exploration. ....	15
5. Conclusion.....	22
6. References.....	23

Figure 1 LOAD DATA.....	5
Figure 2 DROP COLUMN .....	6
Figure 3 NAN MISSING VALUE.....	7
Figure 4 NAN MISSING VALUE.....	7
Figure 5 NAN MISSING VALUE.....	8
Figure 6 CHECK DUPLICATE .....	8
Figure 7 CHECK DUPLICATE .....	9
Figure 8 UNIQUE VALUES.....	9
Figure 9 UNIQUE VALUES.....	10
Figure 10 UNIQUE VALUES.....	10
Figure 11 UNIQUE VALUES.....	11

Figure 12 COLUMN VALUE RENAME.....	12
Figure 13 BUSINESS ANALYTICS.....	13
Figure 14 BUSINESS ANALYTICS.....	13
Figure 15 CORRELATION .....	14
Figure 16 15 JOBS .....	16
Figure 17 GRAPH OF 15 JOBS .....	16
Figure 18 HIGHEST SALARY .....	17
Figure 19 HIGHEST SALARY GRAPH .....	17
Figure 20 SALARY EXPENDITURE LEVEL .....	18
Figure 21 GRAPH OF EXPENDITURE LEVEL.....	18
Figure 22 HISTOGRAM OF WORK YEAR .....	19
Figure 23 HISTOGRAM GRAPH OF WORK YEAR .....	20
Figure 24 BOX PLOT OF SALARY IN USD .....	20
Figure 25 BOX PLOT GRAPH OF SALARY IN USD .....	21
 Table 1 DATA TYPE .....	 4

## **ABSTRACT**

This document presents a comprehensive analysis of a given DataFrame, focusing on the extraction of meaningful insights and the application of analytical techniques. The analysis is question-driven, ensuring that the results are tailored to provide specific answers and solutions. The document covers various stages of the data analysis process, including data cleaning, exploration, visualization, and interpretation. It also discusses the application of statistical and machine learning methods to uncover patterns and trends in the data. The ultimate goal of this document is to enable readers to make data-driven decisions and strategies based on the findings

## **ACKNOWLEDGEMENTS**

I would like to express my deepest appreciation to my mentor, Dipshor Sir, who has the attitude and the substance of a genius. He continually and convincingly conveyed a spirit of adventure in regard to data analysis and an excitement in regard to teaching. Without his guidance and persistent help, this project would not have been possible. His guidance helped me in all the time of research and writing of this document. I could not have imagined having a better advisor and mentor for my project.

## INTRODUCTION

Welcome to our journey through a dataset that revolves around the fascinating world of data scientists. This dataset is like a treasure chest, filled with nuggets of information about their work years, experience levels, job titles, salaries, and more.

In the Data Understanding phase, we'll be like detectives, deciphering what each column in our dataset means. We'll also create a table that shows how complete our data is, by looking at the non-null count for each column.

Next, we'll roll up our sleeves for the Data Preparation phase. This is where we get our hands dirty with the data, loading it into a pandas DataFrame, trimming off unnecessary columns, dealing with those pesky missing values, checking for duplicate values, and giving the values in the experience level column a makeover for better understanding.

Then comes the Data Analysis phase, where we'll put on our statistician hats. We'll crunch numbers to calculate the sum, mean, standard deviation, skewness, and kurtosis of a chosen variable. We'll also look at how all the variables in our dataset interact with each other by calculating their correlations.

Finally, we'll dive deeper into the data in the Data Exploration phase. We'll find out which jobs are the top 15 in demand and illustrate this with a bar

graph. We'll also find out which job pays the highest salaries and represent this with another bar graph.

This document is like a map, guiding you through the dataset and the insights we can extract from it. It showcases the power of Python in handling and analyzing data, and the crucial role of data analysis in making informed decisions. To make the journey easier, we've included screenshots along the way.

## 1. Data Understanding.

### - Explanation what the dataset is about.

work\_year: This is the year the data was recorded, which is 2020 to 2023 for all entries.

experience\_level: This tells us the experience level of the employee. 'SE' stands for Senior, 'MI' for Mid-level, and 'EN' for Entry-level.

employment\_type: This indicates whether the employee is Full-Time (FT) or Contract (CT).

job\_title: This is the role or position of the employee in the company, such as Data Scientist, ML Engineer, etc.

salary: This is how much the employee earns in their local currency.

salary\_currency: This is the currency in which the employee is paid.

salary\_in\_usd: This is the employee's salary converted into US dollars for comparison.

employee\_residence: This is the country where the employee lives, represented by its two-letter country code.

remote\_ratio: This tells us what percentage of the time the employee works remotely. A value of 100 means they always work remotely, and 0 means they never do.

company\_location: This is the country where the company is located.



company\_size: This indicates the size of the company. 'S' stands for Small, 'M' for Medium, and 'L' for Large.

- **Making a table mentioning below information.**

Table 1 DATA TYPE

Column No.	Column Name	Non-Null Count
0	work_year	3755 non-null
1	experience_level	3755 non-null
2	employment_type	3755 non-null
3	job_title	3755 non-null
4	salary	3755 non-null
5	salary_currency	3755 non-null
6	salary_in_usd	3755 non-null
7	employee_residence	3755 non-null
8	remote_ratio	3755 non-null
9	company_location	3755 non-null
10	company_size	3755 non-null

## 2. Data Preparation.

- **Question & Answer along with screenshots.**
  - Write a python program to load data into pandas DataFrame

```
[198]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import os
from scipy.stats import skew, kurtosis

[200]: # Q1.write a python program to Load data into pandas DataFrame.
dataScience = pd.read_csv('dataScience.csv')
dataScience
```

	work_year	experience_level	employment_type	job_title	salary	salary_currency	salary_in_usd	employee_residence	remote_ratio	company_location	company_size
0	2023	SE	FT	Principal Data Scientist	80000	EUR	85847	ES	100	ES	L
1	2023	MI	CT	ML Engineer	30000	USD	30000	US	100	US	S
2	2023	MI	CT	ML Engineer	25500	USD	25500	US	100	US	S
3	2023	SE	FT	Data Scientist	175000	USD	175000	CA	100	CA	M
4	2023	SE	FT	Data Scientist	120000	USD	120000	CA	100	CA	M
...	...	...	...	...	...	...	...	...	...	...	...
3750	2020	SE	FT	Data Scientist	412000	USD	412000	US	100	US	M

Figure 1 LOAD DATA

➔ This above screenshots shows how to load csv data file in jupyter environment.

- Write a python program to remove unnecessary columns i.e., **salary** and **salary currency**.

```

[206]: #Write python program to remove unnecessary columns i.e., salary and salary currency.
dataScience = dataScience.drop(columns=['salary', 'salary_currency'])

[208]: dataScience

[208]:
  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
...      ...                ...                ...      ...      ...                ...            ...                ...                ...
3750     2020                SE                FT    Data Scientist     412000                US            100                US                L
3751     2021                MI                FT  Principal Data Scientist     151000                US            100                US                L
3752     2020                EN                FT    Data Scientist     105000                US            100                US                S
3753     2020                EN                CT  Business Data Analyst     100000                US            100                US                L
3754     2021                SE                FT  Data Science Manager      94665                IN             50                IN                L

3755 rows x 9 columns

[210]: dataScience.info()

```

Figure 2 DROP COLUMN

➔ This above screenshots shows how to remove unwanted columns.

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

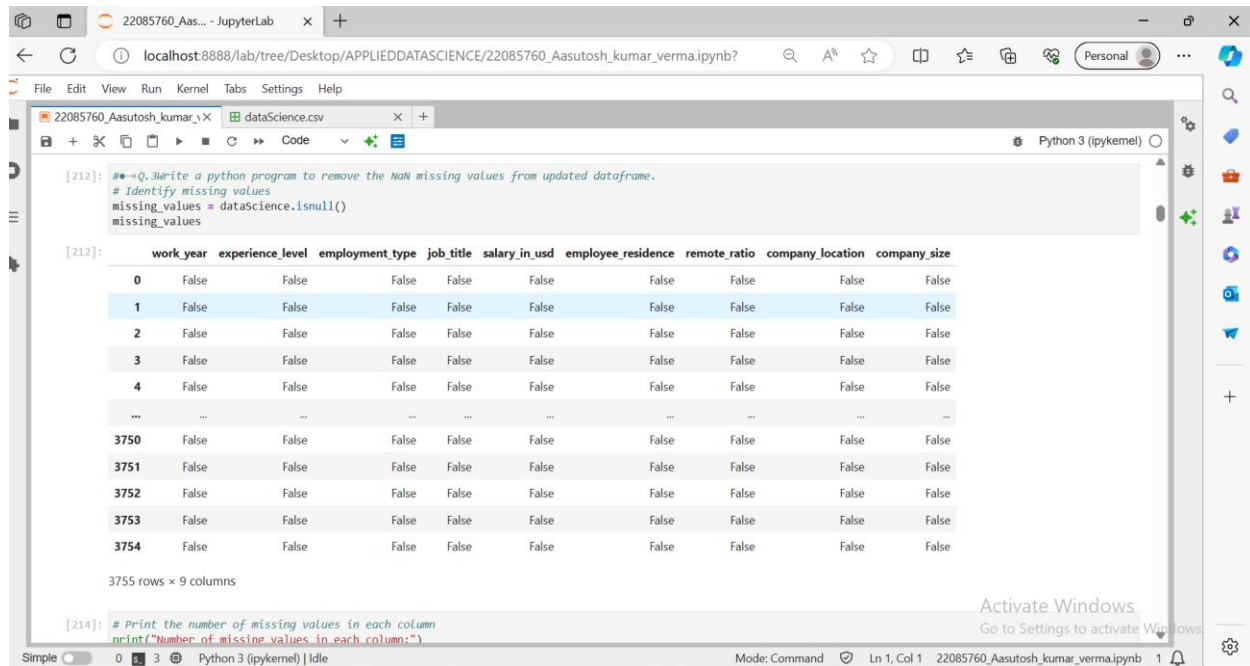


Figure 3 NAN MISSING VALUE

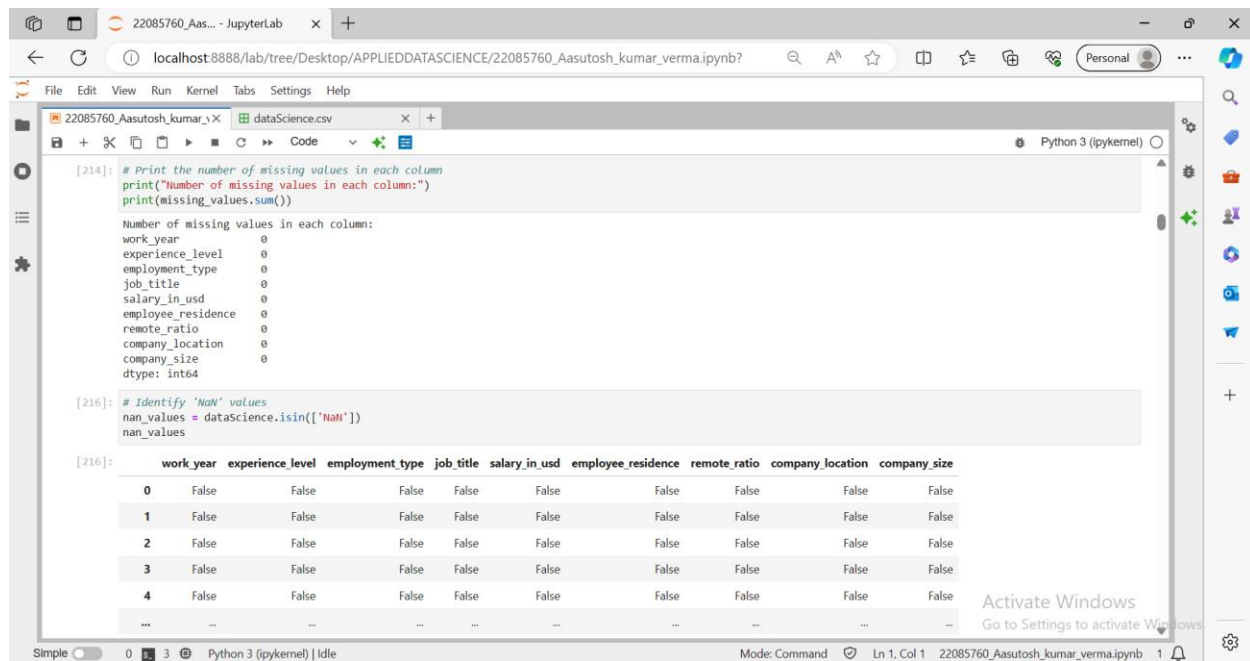


Figure 4 NAN MISSING VALUE

```

3755 rows x 9 columns

[218]: # Get rows with missing values
rows_with_nan = dataScience[missing_values.any(axis=1)]

# Print the rows with missing values
print(rows_with_nan)

Empty DataFrame
Columns: [work_year, experience_level, employment_type, job_title, salary_in_usd, employee_residence, remote_ratio, company_location, company_size]
Index: []

```

Figure 5 NAN MISSING VALUE

➔ This above screenshots shows how to check missing values with different ways as a part of data cleaning to get accurate value while analytical process.

- Write a python program to check duplicates value in the dataframe.

```

[220]: ##->Q.4write a python program to check duplicates value in the dataframe. ->
# Checking for duplicate rows
duplicate_rows = dataScience.duplicated()
duplicate_rows

[220]: 0      False
1      False
2      False
3      False
4      False
...
3750   False
3751   False
3752   False
3753   False
3754   False
Length: 3755, dtype: bool

[222]: # Print the number of duplicate rows
print("Number of duplicate rows:", duplicate_rows.sum())

Number of duplicate rows: 1171

[224]: # Get duplicate rows
duplicate_rows = dataScience[dataScience.duplicated()]

# Print the duplicate rows
print(duplicate_rows)

work_year  experience_level  employment_type  job_title \
115      2023              SE               FT  Data Scientist
123      2023              SE               FT  Analytics Engineer
152      2023              SE               FT  Data Engineer

```

Figure 6 CHECK DUPLICATE

22085760\_Aas... - JupyterLab

localhost:8888/lab/tree/Desktop/APPLIEDDATASCIENCE/22085760\_Aasutosh\_kumar\_verma.ipynb?

FileEditViewRunKernelTabbsSettingsHelp

22085760\_Aasutosh\_kumar\_v...dataScience.csv

work\_yearexperience\_levelemployment\_typejob\_title

1152023SEFTData Scientist

1232023SEFTAnalytics Engineer

1532023MIFTData Engineer

1542023MIFTData Engineer

1602023SEFTData Engineer

... ...

34392022MIFTData Scientist

34402022SEFTData Engineer

34412022SEFTData Engineer

35862021MIFTData Engineer

37092021MIFTData Scientist

salary\_in\_usdemployee\_residencerelease\_ratiocompany\_location

115150000US0US

123289800US0US

153100000US100US

15470000US100US

160115000US0US

... ...

343978000US100US

3440135000US100US

3441115000US100US

3586200000US100US

370990734DE50DE

company\_size

115M

123M

153M

154M

160M

Simple03Python 3 (ipykernel) | Idle

Mode: Command

Figure 7 CHECK DUPLICATE

➔ This above screenshots shows how to check duplicates value in a dataframe.

- Write a python program to see the unique values from all the columns in the dataframe.

```
[228]: unique_values = dataScience['experience_level'].unique()
       print(f"Unique values in 'experience_level': {unique_values}")

Unique values in 'experience_level': ['SE' 'MI' 'EN' 'EX']

[230]: unique_values1 = dataScience['job_title'].unique()
       print(f"Unique values in 'job_title': {unique_values1}")

Unique values in 'job_title': ['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'
'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']
```

Figure 8 UNIQUE VALUES

```

[246]: unique_values2 = dataScience['job_title'].unique()
       print(f"Unique values in 'job_title': {unique_values2}")

Unique values in 'job_title': ['Lead Data Scientist' 'Machine Learning Engineer' 'Data Scientist'
'Scientist' 'Data Analyst' 'Engineer' 'Analytics Engineer'
'Business Intelligence Developer' 'Data Engineer'
'Computer Vision Scientist' 'Data Architect' 'AI Engineer' 'Data Manager'
'Data Lead' 'Machine Learning Scientist' 'Big Data Engineer'
'Lead Data Analyst' 'Business Intelligence Data Engineer'
'Director of Data Science' 'Machine Learning Researcher' 'MLOps Engineer'
'Artificial Intelligence Scientist' 'Autonomous Vehicle Technician'
'Database Engineer' 'Financial Data Analyst'
'Artificial Intelligence Programmer' 'Deep Learning Scientist'
'Business Intelligence Analyst' 'Data Consultant'
'Business Intelligence Data Analyst' 'Data Science Lead'
'Big Data Architect' 'Computer Vision Engineer'
'Natural Language Processing Engineer' 'Machine Learning Developer'
'Computer Vision Researcher' 'Principal Machine Learning Engineer'
'Data Specialist' 'Tech Lead' 'Lead Data Architect'
'Machine Learning Lead' 'Lead Machine Learning Engineer' 'Data Developer'
'Lead Data Engineer']

```

Figure 9 UNIQUE VALUES

```

[248]: unique_values3 = dataScience['employee_residence'].unique()
       print(f"Unique values in 'employee_residence': {unique_values3}")

Unique values in 'employee_residence': ['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']

[250]: unique_values4 = dataScience['company_location'].unique()
       print(f"Unique values in 'company_location': {unique_values4}")

Unique values in 'company_location': ['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']

[252]: unique_values5 = dataScience['company_size'].unique()
       print(f"Unique values in 'company_size': {unique_values5}")

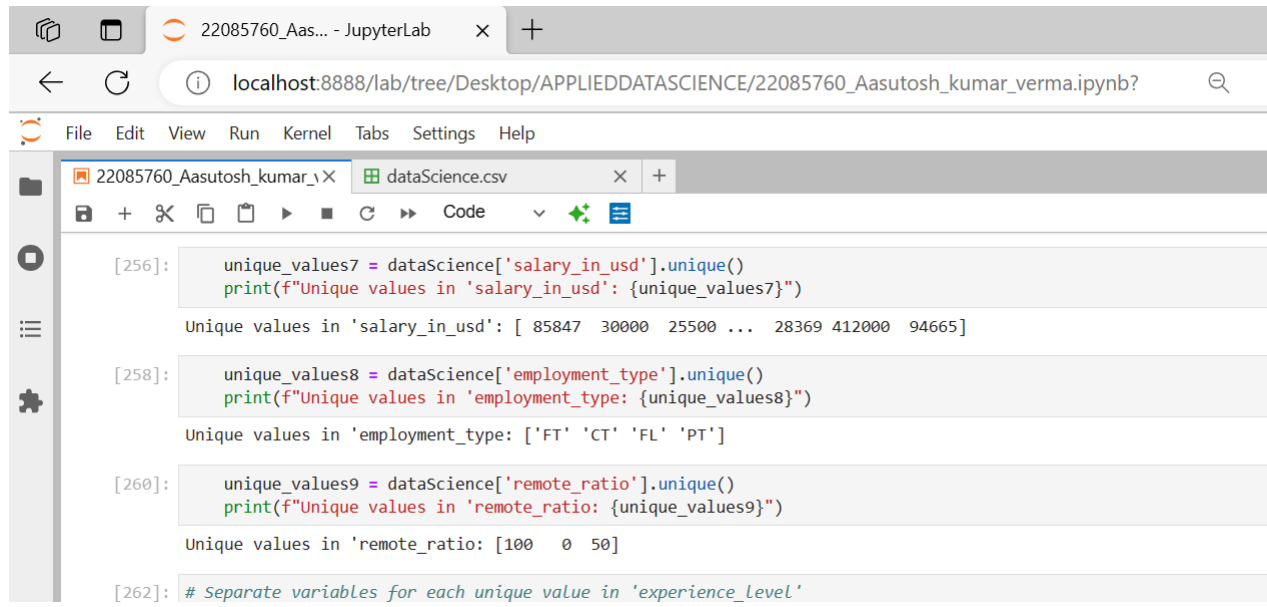
Unique values in 'company_size': ['L' 'S' 'M']

[254]: unique_values6 = dataScience['work_year'].unique()
       print(f"Unique values in 'work_year': {unique_values6}")

Unique values in 'work_year': [2023 2022 2020 2021]

```

Figure 10 UNIQUE VALUES



The screenshot shows a JupyterLab window with a browser tab at 'localhost:8888/lab/tree/Desktop/APPLIEDDATASCIENCE/22085760\_Aasutosh\_kumar\_verma.ipynb?'. The interface includes a menu bar (File, Edit, View, Run, Kernel, Tabs, Settings, Help) and a toolbar. The main area displays a notebook with the following code and output:

```
[256]: unique_values7 = dataScience['salary_in_usd'].unique()
       print(f"Unique values in 'salary_in_usd': {unique_values7}")
       Unique values in 'salary_in_usd': [ 85847  30000  25500 ... 28369 412000 94665]
```

```
[258]: unique_values8 = dataScience['employment_type'].unique()
       print(f"Unique values in 'employment_type': {unique_values8}")
       Unique values in 'employment_type': ['FT' 'CT' 'FL' 'PT']
```

```
[260]: unique_values9 = dataScience['remote_ratio'].unique()
       print(f"Unique values in 'remote_ratio': {unique_values9}")
       Unique values in 'remote_ratio': [100  0  50]
```

```
[262]: # Separate variables for each unique value in 'experience_level'
```

Figure 11 UNIQUE VALUES

➔ This above screenshots shows how to check unique values in a columns.

- Rename the experience level columns as below.

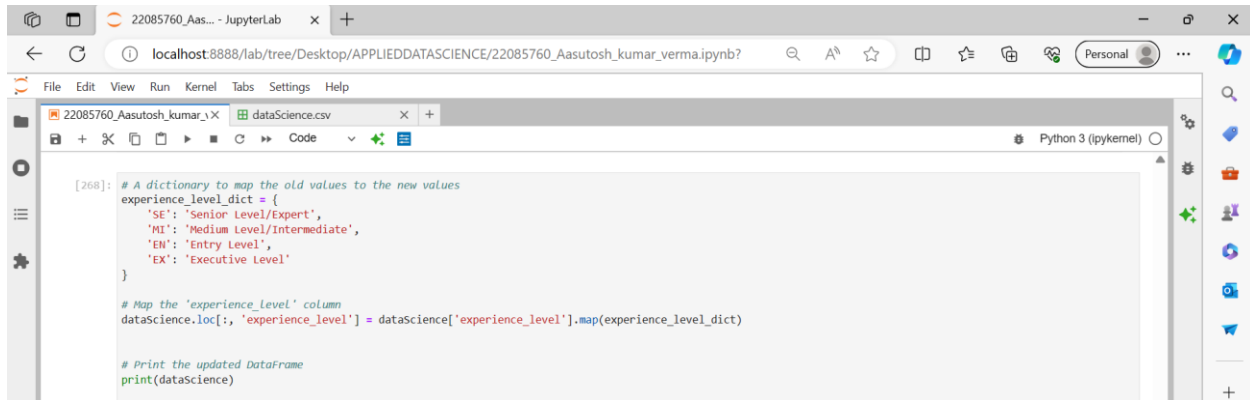
**SE – Senior Level/Expert**

**MI – Medium Level/Intermediate**

**EN – Entry Level**

**EX – Executive Level**





The screenshot shows a JupyterLab window with a browser tab titled '22085760\_Aas... - JupyterLab'. The address bar shows 'localhost:8888/lab/tree/Desktop/APPLIEDDATASCIENCE/22085760\_Aasutosh\_kumar\_verma.ipynb?'. The file explorer on the left shows a file named 'dataScience.csv'. The code editor displays the following Python code:

```
[268]: # A dictionary to map the old values to the new values
experience_level_dict = {
    'SE': 'Senior Level/Expert',
    'MI': 'Medium Level/Intermediate',
    'EN': 'Entry Level',
    'EX': 'Executive Level'
}

# Map the 'experience_level' column
dataScience.loc[:, 'experience_level'] = dataScience['experience_level'].map(experience_level_dict)

# Print the updated DataFrame
print(dataScience)
```

Figure 12 COLUMN VALUE RENAME

➔ This above screenshots shows how to rename column current value with new one.

- Explain of output.

### 3. Data Analysis.

- Question note and explain of output information.
- Write a Python program to show summary statistics of sum, mean, standard deviation, skewness, and kurtosis of any chosen variable.

```

[270]: # Write a Python program to show summary statistics of sum, mean, standard deviation, skewness, and kurtosis of any chosen variable.

column = dataScience['salary_in_usd']

# Calculate summary statistics
sum_val = np.sum(column)
mean_val = np.mean(column)
std_dev = np.std(column)
skewness = skew(column)
kurtosis_val = kurtosis(column)

# Print the results
print(f"Sum: {sum_val}")

Sum: 344729580

[272]: print(f"Mean: {mean_val}")

Mean: 133409.28018575851

[274]: print(f"Standard Deviation: {std_dev}")

Standard Deviation: 67123.84519805372

[276]: print(f"Skewness: {skewness}")

Skewness: 0.6199567299104847

```

Figure 13 BUSINESS ANALYTICS

```

[278]: print(f"Kurtosis: {kurtosis_val}")

Kurtosis: 0.8230197549418379

```

Figure 14 BUSINESS ANALYTICS

### 3.1. Statistics of sum.

➔ 33729580 USD.

### 3.2. Statistics of Mean.

➔ 133409 USD. Its shows average salary in usd a employees recived.

### 3.3. Statistics of standard deviation.

➔ 67123.84519805372 USD. Its shows that the company provide salaries equally or not . so if it great means company need to ensure that the salaries is equally paid.

### 3.4. Statistics of skewness.

→ 0.6199567299104847. Its shows that there are number of employess that receive more salaries.

### 3.5. Statistics of Kurtosis.

→ 0.8230197549418379. Its shows that the most employees in the company have salaries that are close to the average(Mean) salary.

→ This above screenshots shows whether the analysis part or process is right or not by comparig the values received from analytical process using statistical calculations.

- Write a Python program to calculate and show correlation of all variables.

```
[284]: # only the numerical columns in the DataFrame
numerical_dataScience = dataScience[['work_year', 'salary_in_usd', 'remote_ratio']]

# Calculate the correlation matrix
correlation_matrix = numerical_dataScience.corr()

# Print the correlation matrix
print(correlation_matrix)
```

	work_year	salary_in_usd	remote_ratio
work_year	1.000000	0.236958	-0.219160
salary_in_usd	0.236958	1.000000	-0.084502
remote_ratio	-0.219160	-0.084502	1.000000

Figure 15 CORRELATION

- ➔ `work_year` and `work_year` (1.000000): This is the correlation of 'work\_year' with itself, which is always 1. It simply means each variable perfectly correlates with itself.
- ➔ `work_year` and `salary_in_usd` (0.236958): This shows a positive but weak relationship between the number of years worked ('work\_year') and the salary in USD. It suggests that as the number of years worked increases, the salary tends to increase slightly as well.
- ➔ `work_year` and `remote_ratio` (-0.219160): This shows a negative but weak relationship between the number of years worked ('work\_year') and the remote work ratio. It suggests that as the number of years worked increases, the remote work ratio tends to decrease slightly.
- ➔ `salary_in_usd` and `remote_ratio` (-0.084502): This shows a very weak negative relationship between salary and the remote work ratio. It suggests that there's almost no relationship between the salary and the remote work ratio.

#### 4. Data Exploration.

- **Question and explain of output information proper label of graphs.**
- Write a python program to find out top 15 jobs. Make a bar graph of sales as well.



Figure 16 15 JOBS

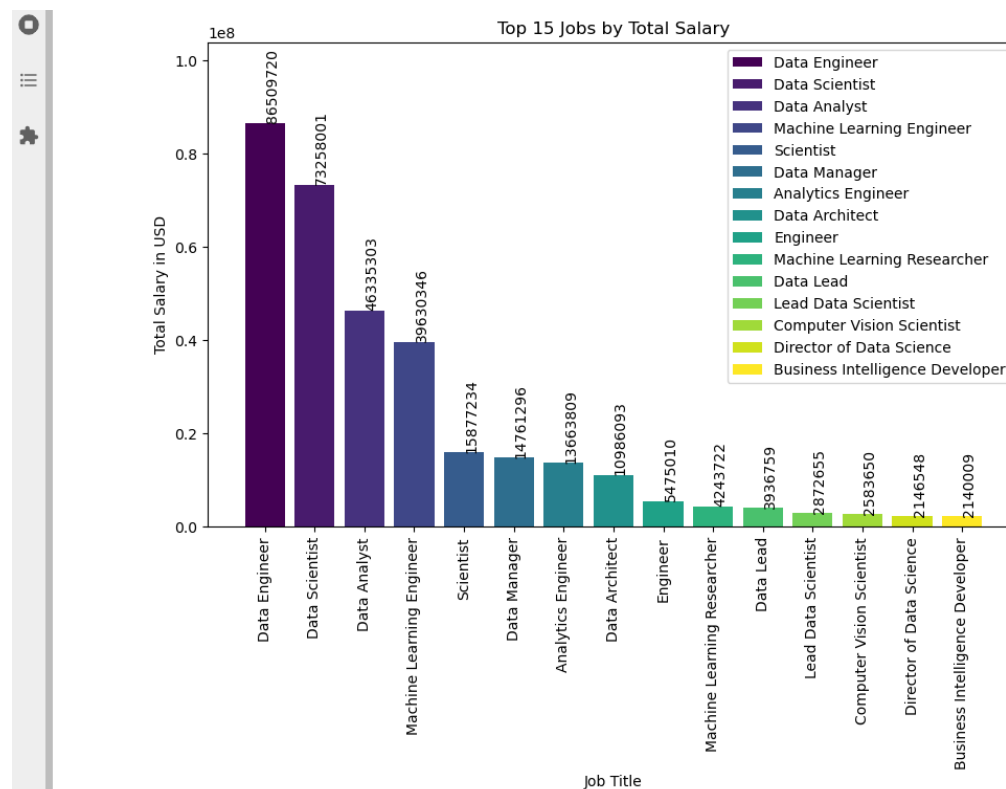


Figure 17 GRAPH OF 15 JOBS

➔ This screenshots shows that the data engineer holds higher demands out of 15 jobs.

- Which job has the highest salaries? Illustrate with bar graph.

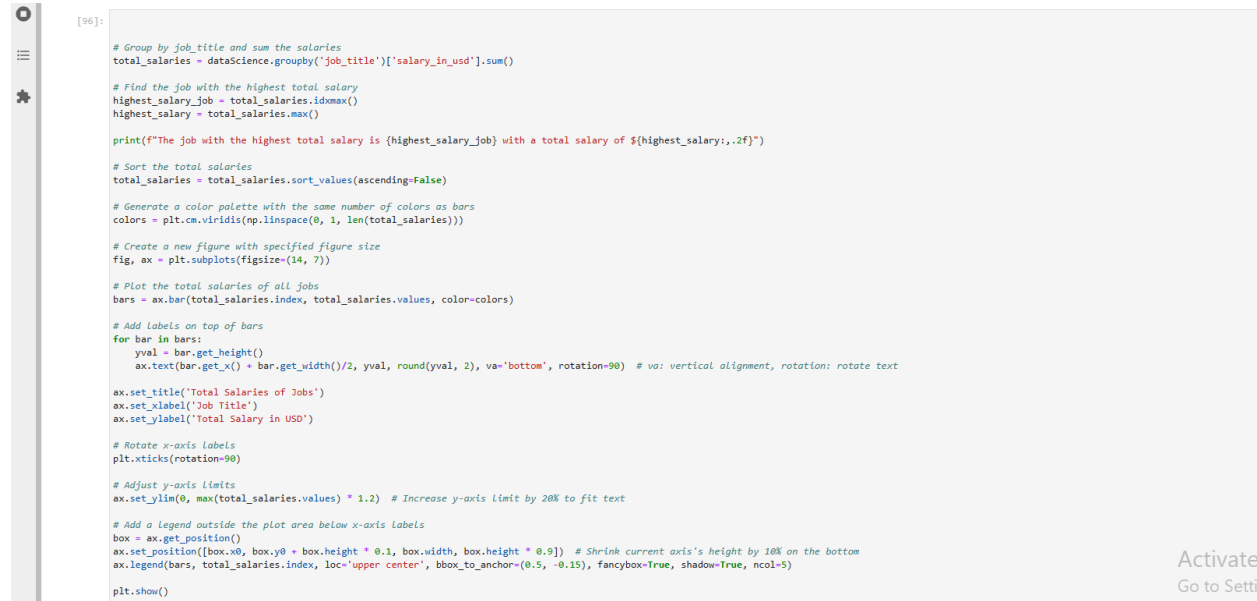


Figure 18 HIGHEST SALARY

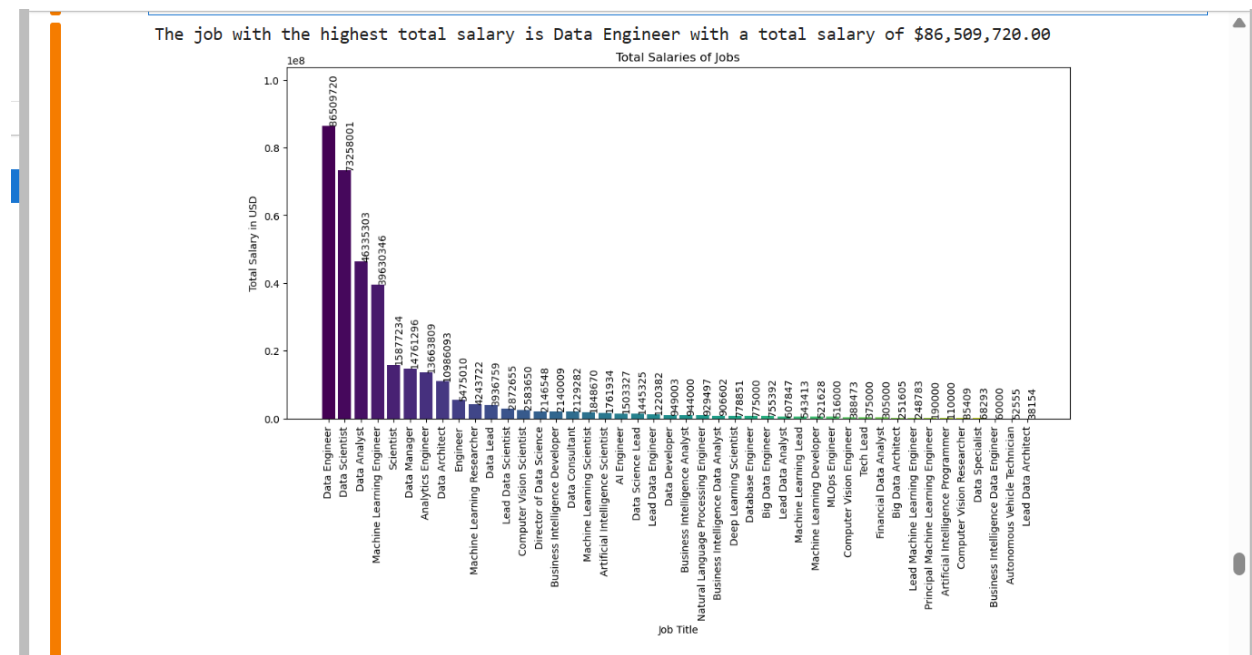


Figure 19 HIGHEST SALARY GRAPH



→ This screenshots shows that the data engineer holds the highest salary.

- Write a python program to find out salaries based on experience level. Illustrate it through bar graph.

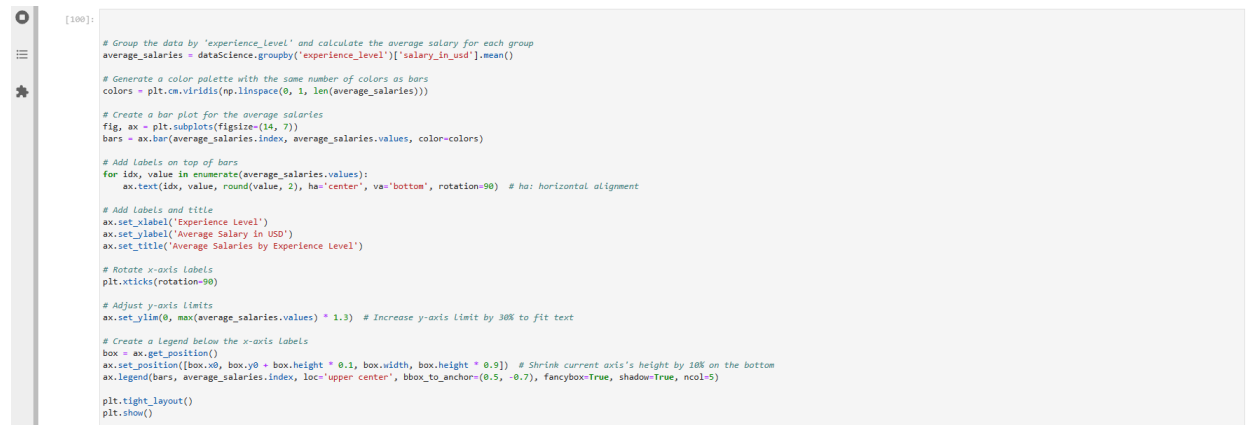


Figure 20 SALARY EXPENDITURE LEVEL

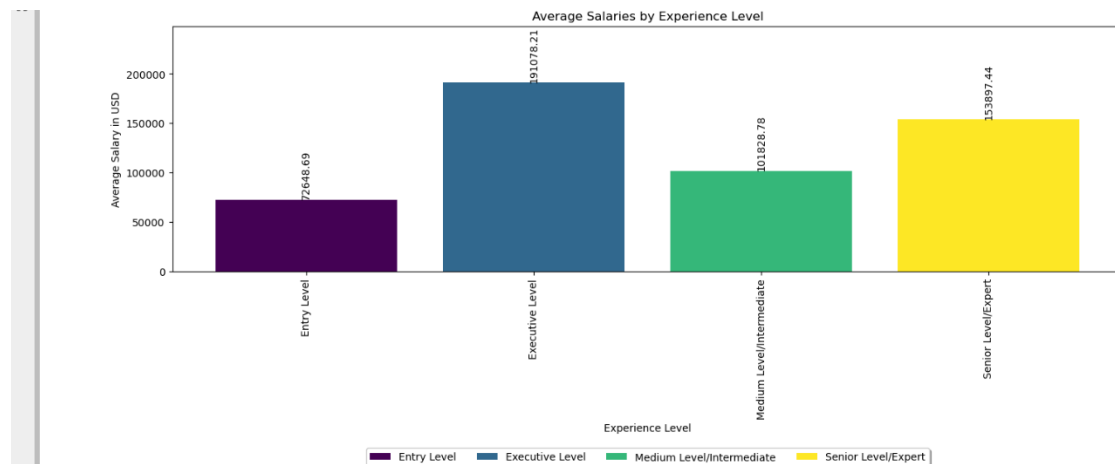


Figure 21 GRAPH OF EXPENDITURE LEVEL

Activate Wind  
Go to Settings to ac

➔ The screenshots shows that the employees who earns more salaries belongs to executive level.

- Write a Python program to show histogram and box plot of any chosen different variables. Use proper labels in the graph.

```
[103]: # Create a histogram for 'work_year'
counts, bins, patches = plt.hist(dataScience['work_year'], bins = range(2020, 2024), edgecolor='black')

# Add frequency annotations on top of each bar
for count, bin, patch in zip(counts, bins, patches):
    plt.annotate(f'{int(count)}', (bin, count), xytext=(0, 5), textcoords='offset points', ha='center')

plt.title('Work Year Histogram')
plt.xlabel('Work Year')
plt.ylabel('Frequency')
plt.show()
```

Figure 22 HISTOGRAM OF WORK YEAR



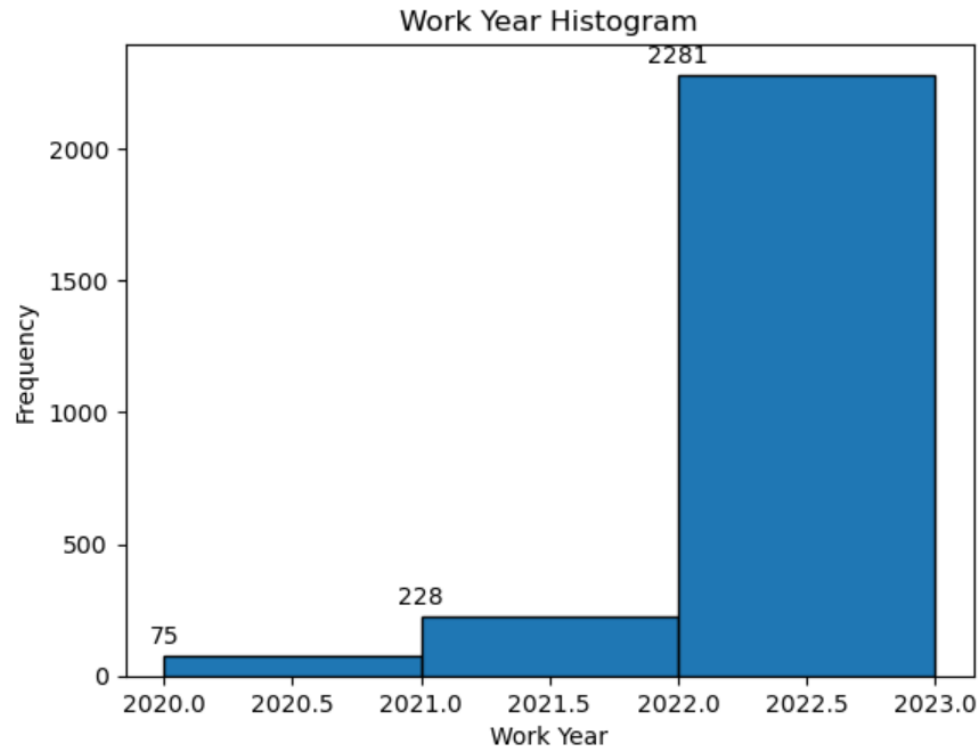


Figure 23 HISTOGRAM GRAPH OF WORK YEAR

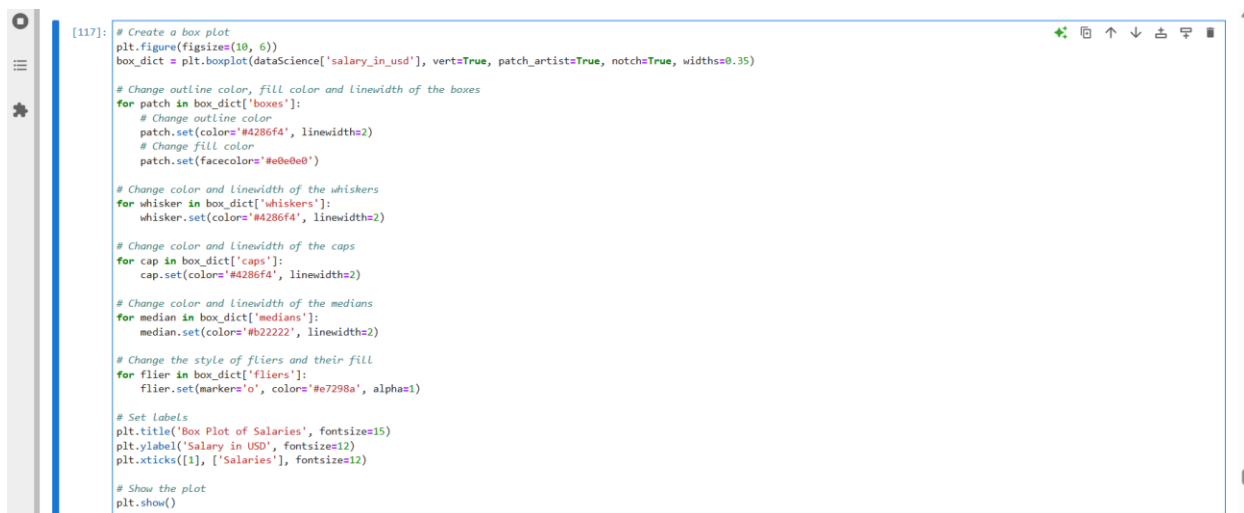


Figure 24 BOX PLOT OF SALARY IN USD

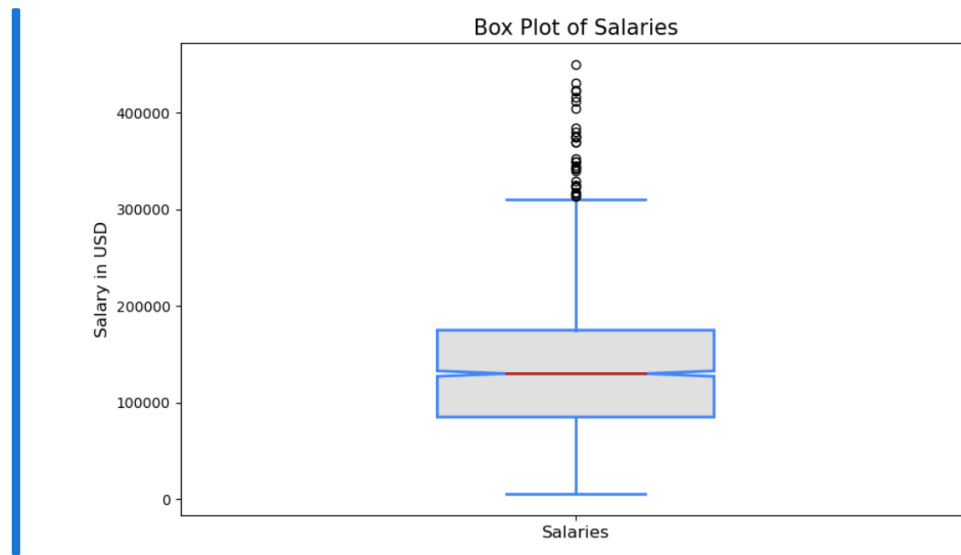


Figure 25 BOX PLOT GRAPH OF SALARY IN USD

- ➔ The screenshots of histogram shows that the increase of employees who started work between 2022 to 2023 year.
- ➔ FOR BOX PLOT:-
- ➔ Median (Middle Line): This is the middle salary. Half of the people earn less than this and half earn more.
- ➔ Box (Q1 to Q3): This is where the middle 50% of salaries fall. If the box is skewed, it means more people earn either towards the lower end (box skewed towards top) or higher end (box skewed towards bottom).
- ➔ Whiskers (Lines extending from the box): These represent the overall range of salaries, excluding extreme outliers. If the whiskers are long, it means there's a large spread in what people earn.
- ➔ Outliers (Dots): These are salaries that are unusually low or high compared to the rest. If there are many dots, it means there are many people with unusually high or low salaries.
- ➔ Notches: If the notches of two boxes do not overlap, this suggests a statistically significant difference between the medians.

## 5. Conclusion.

This project has highlighted the significance of the Python language in data analysis. While some believe Python may not be suitable for large-scale data analysis due to limitations in processing multiple lines of code simultaneously, this is a misconception. Python is indeed capable of handling extensive datasets with the appropriate use of libraries.

It is acknowledged that C, often referred to as the 'mother of all languages,' can be utilized alongside Python to optimize performance, particularly in global interpretation scenarios. Python was designed to simplify coding, and while it's essential to critically evaluate information from scientific sources, news articles, and lectures, dismissing Python's capabilities based on unverified claims would be misguided. Recent rumors of Google disbanding its Python team are unfounded and should not influence our perception of the language's reliability.

The assertion that Python's security can be compromised by C is an oversimplification. Security vulnerabilities are a concern for any language and are more dependent on the implementation rather than the language itself. It's important to recognize that while C developers possess the knowledge to create sophisticated AI and robotics, this expertise is not exclusive to them, nor does it diminish Python's contributions to these fields.

While Python does incorporate features from other languages, this interoperability is one of its strengths, allowing for a more comprehensive approach to problem-solving. If one's goal is to excel in data analysis or machine learning, Python remains a strong contender due to its extensive libraries and supportive community.

In terms of data processing, this project has implemented data cleaning techniques such as removing duplicates, renaming values, and segregating unique values into separate variables for efficient access. These steps are crucial for achieving accurate results and are detailed in the code section of the report. parts as well (Mckinney, 2017).

## 6. References

Mckinney, W., 2017. *Python for Data Analysis*. 2nd ed. Mumbai: SHROFF.

Islington Library, Data science and analytical Books .