**LONDON METROPOLITAN UNIVERSITY**

**islington college**
(इस्लिङ्टन कलेज)

## CC5067NI

## 60% Individual Coursework

## 2023-24 Autumn

**Student Name: Amogh Man Bajracharya**

**London Met ID: 22067567**

**College ID: np0cp4a220064**

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

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

**Word Count: 1188**

Amogh Man Bajracharya

# Contents

## Table of Tables

## Table of Figures

## 1.  Data Understanding

This dataset is about information of different variables that could impact salaries such as experience level job title and many more. In this dataset we find out about employment type based on job title and their salary and details about the employee and company. We dive into data cleaning, data preparation, data analysis and data exploration to prepare and generate meaningful findings and draw the conclusions. The steps in this project we will follow will help systematically prepare the data to uncover the insights.

| S.no | Column Name | Description | Data Type |
|------|-------------|-------------|-----------|
| 1 | work_year | This column gives the work year of the employee | Int(64) |
| 2 | experience_level | This column give the experience level of the employee | object |
| 3 | employment_type | This column gives the employment type of the employee | object |
| 4 | job_title | This column state all the job titles of the employee | object |
| 5 | salary | This column states all the salary of given job title | Int(64) |
| 6 | salary_currency | This column describes which currency salary is received in | object |
| 7 | salary_in_usd | this column describes salary | Int(64) |

Amogh Man Bajracharya

| | | which is received in usd | |
|---|---|---|---|
| 8 | Employee_residence | This column describes the residence of the employee | object |
| 9 | Remote_ratio | This column describes whether the employee is working at site or remotely | Int(64) |
| 10 | Company_location | This column describes company's location which employee works in | object |
| 11 | Company_size | This column states company's size. | object |

*Table 1 Description of every column name of the dataset*

Amogh Man Bajracharya

## 2.  Data Preparation

Data preparation is the first step of cleaning and enriching raw data to help to make it ready for use in analytics and data science. Data preparation helps you to find prepare and use the prepared data faster. The idea behind data preparation is to change data into information which will be useful for data analysis. (Secoda, 2024)

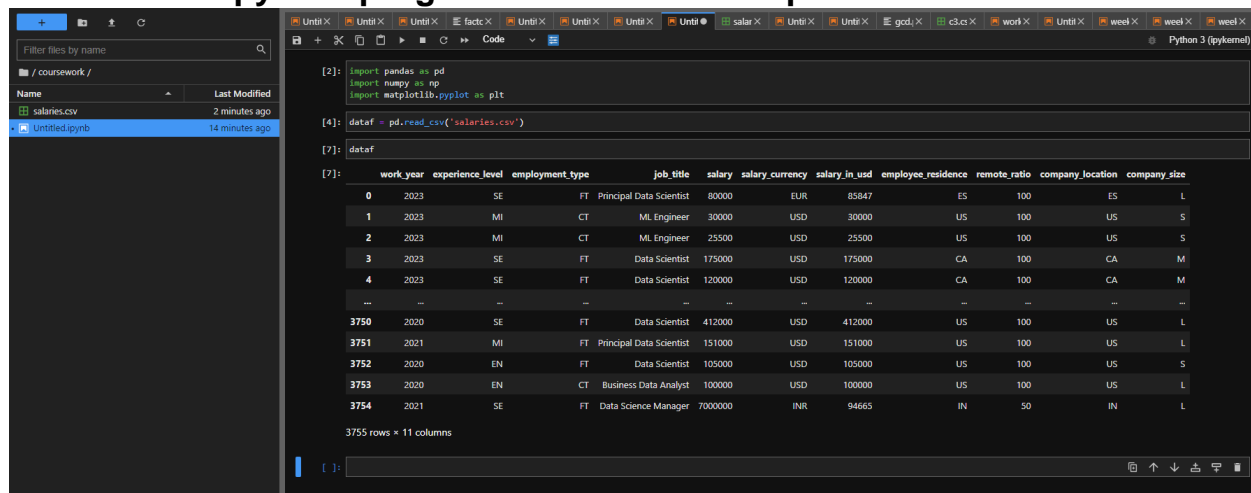### 2.1 Write a python program to load data into pandas DataFrame



*Figure 1 Loading data into pandas DataFrame*

This code imports pandas numpy and matplotlib dataframe and loads data into pandas dataframe making dataf as the name of dataframe.

Amogh Man Bajracharya

## 2.2 Write a python program to remove unnecessary columns i.e., salary and salary currency.



*Figure 2 Removing Columns salary and salary currency*

This code removes unnecessary columns i.e. salary and currency which is repeated by salary in USD and company's location.

Amogh Man Bajracharya

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



*Figure 3 dropping any missing values*

This code of line removes the NaN missing values to remove it from the updated dataframe with no salary and currency column.

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



*Figure 4 checking duplicates in dataframe*

This code of line checks duplicates in dataframe



*Figure 5 droping duplicated values*

This code of line deletes duplicated values in the dataframe for better data consistency.

Amogh Man Bajracharya

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



*Figure 6 python program to see unique values of all the columns.*



*Figure 7 program to see unique values of all the columns.*

These codes of lines defines a function named unique_values_in_all_columns where all the unique_values are gathered from each columns and then later called to print each and every unique values with their respective columns.

Amogh Man Bajracharya

## 2.1 Rename the experience level columns as below.

```
[37]: value_rename_mapping = {
          'SE': 'Senior Level/Expert',
          'MI': 'Medium Level/Intermediate',
          'EN': 'Entry Level',
          'EX': 'Executive Level'
      }

[38]: dataf['experience_level'] = dataf['experience_level'].replace(value_rename_mapping)

[39]: dataf

[39]:
```

| | work_year | experience_level | employment_type | job_title | salary_in_usd | employee_residence | remote_ratio | company_location | company_size |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 2023 | Senior Level/Expert | FT | Principal Data Scientist | 85847 | ES | 100 | ES | L |
| 1 | 2023 | Medium Level/Intermediate | CT | ML Engineer | 30000 | US | 100 | US | S |
| 2 | 2023 | Medium Level/Intermediate | CT | ML Engineer | 25500 | US | 100 | US | S |
| 3 | 2023 | Senior Level/Expert | FT | Data Scientist | 175000 | CA | 100 | CA | M |
| 4 | 2023 | Senior Level/Expert | FT | Data Scientist | 120000 | CA | 100 | CA | M |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 3750 | 2020 | Senior Level/Expert | FT | Data Scientist | 412000 | US | 100 | US | L |
| 3751 | 2021 | Medium Level/Intermediate | FT | Principal Data Scientist | 151000 | US | 100 | US | L |
| 3752 | 2020 | Entry Level | FT | Data Scientist | 105000 | US | 100 | US | S |
| 3753 | 2020 | Entry Level | CT | Business Data Analyst | 100000 | US | 100 | US | L |
| 3754 | 2021 | Senior Level/Expert | FT | Data Science Manager | 94665 | IN | 50 | IN | L |

2584 rows × 9 columns

*Figure 8 renaming the experience level rows.*

This group of code renames experience level rows replacing and renaming every values in the column.

Amogh Man Bajracharya

# 3.  Data Analysis

Data analysis is a process of analysing, cleansing, manipulating, and modelling data in order to identify usable information and to draw conclusions which helps in decision-making. Data analysis is a process that uses a different approaches and methodologies to understand data from different sources in various formats, both structured and unstructured. (datacamp, 2023)

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



*Figure 9 showing summary statistics of sum, mean , sd, skewness and kurtosis*

These group of code defines the table's summary, statistics of sum, mean, Standard Deviation, skewness and kurtosis.

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



*Figure 10 Correlation of all the variables.*

These lines of code defines correlation of all the variables in the table.

Amogh Man Bajracharya

## 4. Data Cleaning

Data cleaning is the process of removing corrupted, inconsistent, or incomplete data in a dataset. It is a crucial step in the machine learning process to achieve and ensure that the data is accurate and consistent with minimum to none errors in the dataset. It is done to bring positive impact on performance of machine learning model. (Geek for Geeks, 2024)



*Figure 11 Before Data Cleaning*



*Figure 12 After Data Cleaning Machine Learning*



*Figure 13 Before Data Cleaning AI Programmer into AI Developer*



*Figure 14 After Data Cleaning AI Programmer into AI Developer*



*Figure 15 Before Data Cleaning Lead Data Scientist into Data Scientist Lead*



*Figure 16 After Data Cleaning Lead Data Scientist into Data Scientist Lead*

Amogh Man Bajracharya

## 5.  Data Exploration

Data exploration is one of the processes for machine learning which leads to reviewing of raw dataset that helps to figure out initial patterns for further analysis. As it is difficult to manage and review thousands of data elements to get proper analysis view of the dataset Data exploration helps to manage unstructured dataset and recognize patterns accordingly. (Qlik, 2024)

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



```
[93]: top_15_jobs = dataf['job_title'].value_counts().head(15)
      print(top_15_jobs)

job_title
Data Engineer                  598
Data Scientist                 538
Data Analyst                   396
Machine Learning Engineer      240
Analytics Engineer              91
Research Scientist              65
Data Architect                  64
Data Science Manager            52
Research Engineer               33
Applied Scientist               31
Machine Learning Scientist      26
Data Science Consultant         23
Data Manager                    23
Computer Vision Engineer        18
Data Analytics Manager          18
Name: count, dtype: int64
```

*Figure 17 python program to find out top 15 jobs.*

These lines of code print out top 15 jobs based on job title's frequency.



```
[108]: plt.figure(figsize=(12, 6))
       top_15_jobs.plot(kind='bar', color='red')
       plt.xlabel('Job Title')
       plt.ylabel('Frequency')
       plt.title('Top 15 Most Common Job Titles')
       plt.xticks(rotation=45, ha='right')
       plt.tight_layout()
       plt.show()
```

*Figure 18 Plotting bar graph of top 15 job titles.*

These lines of code plots top 15 most common job titles with proper x and y labelling and their title.

Amogh Man Bajracharya

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

```
[128]: top_10_highest_salary_jobs = dataf.groupby('job_title')['salary_in_usd'].mean().nlargest(10)
       print(top_10_highest_salary_jobs)

       job_title
       Data Science Tech Lead               375000.000000
       Cloud Data Architect                 250000.000000
       Data Lead                            212500.000000
       Data Analytics Lead                  211254.500000
       Principal Data Scientist             198171.125000
       Director of Data Science             195140.727273
       Principal Data Engineer              192500.000000
       Machine Learning Software Engineer   192420.000000
       Applied Scientist                    190342.580645
       Principal Machine Learning Engineer  190000.000000
       Name: salary_in_usd, dtype: float64
```

*Figure 19 python program to find out highest salaries based on job titles*

This line of code uses groupby() function to find out highest salaries of 10 job titles based on job titles.

```
[129]: plt.figure(figsize=(12, 6))
       top_10_highest_salary_jobs.plot(kind='bar', color='red')
       plt.xlabel('JOB')
       plt.ylabel('Salary')
       plt.title('Top 10 Jobs with the highest Salaries')
       plt.xticks(rotation=45, ha='right')
       plt.tight_layout()
       plt.show()
```

*Figure 20 bar graph to find out jobs with the highest salaries*

This bar graph plots  different top 10 job titles which has highest average salaries.

Amogh Man Bajracharya

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

```
[123]: experience_level_by_salary = dataf.groupby('experience_level')['salary_in_usd'].mean()
       print(experience_level_by_salary)

       experience_level
       Entry Level              72648.685185
       Executive Level         191078.208333
       Medium Level/Intermediate 101828.783133
       Senior Level/Expert     153897.435650
       Name: salary_in_usd, dtype: float64
```

*Figure 21 python program to find out salaries based on experience levels.*

```
[127]: plt.figure(figsize=(12, 6))
       experience_level_by_salary.plot(kind='bar', color='red')
       plt.xlabel('Experience Level')
       plt.ylabel('Salary')
       plt.title('Average Salaries based on experience levels')
       plt.xticks(rotation=45, ha='right')
       plt.tight_layout()
       plt.show()
```
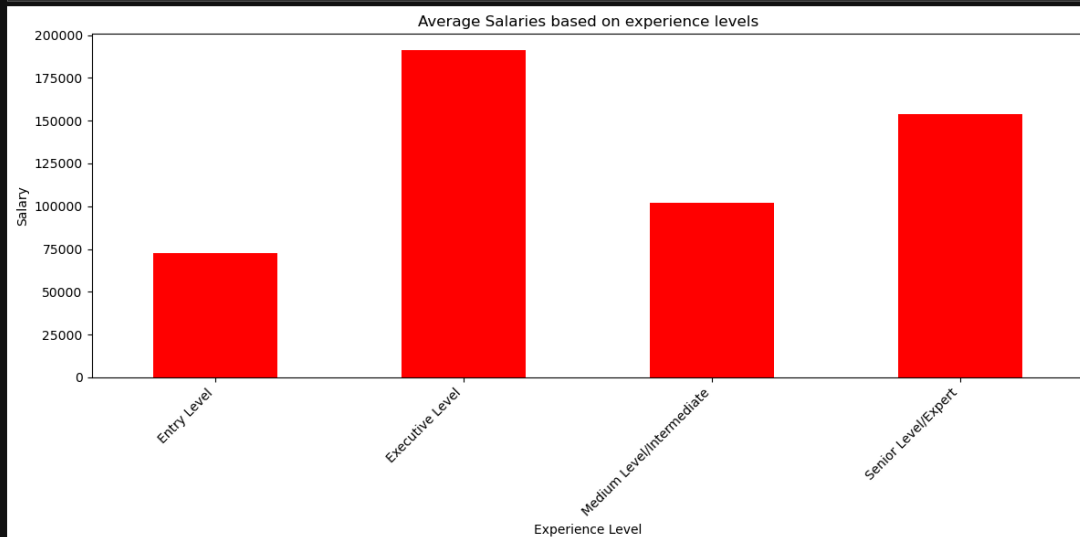


*Figure 22 bar graph to plot average salaries based on experience levels*

This code of line plots average salaries based on experience levels of the employees in the company.

Amogh Man Bajracharya

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

```
[173]: salaries_given_based_year = dataf.groupby('salary_in_usd')['work_year'].mean()
       plt.figure(figsize=(12, 6))
       experience_level_by_salary.plot(kind='hist')
       plt.xlabel('Experience Level')
       plt.ylabel('Salary')
       plt.title('number of salaries given based on year ')
       plt.xticks(rotation=45, ha='right')
       plt.tight_layout()
       plt.hist(experience_level_by_salary,edgecolor='black', linewidth=1.5)
       plt.show()
```
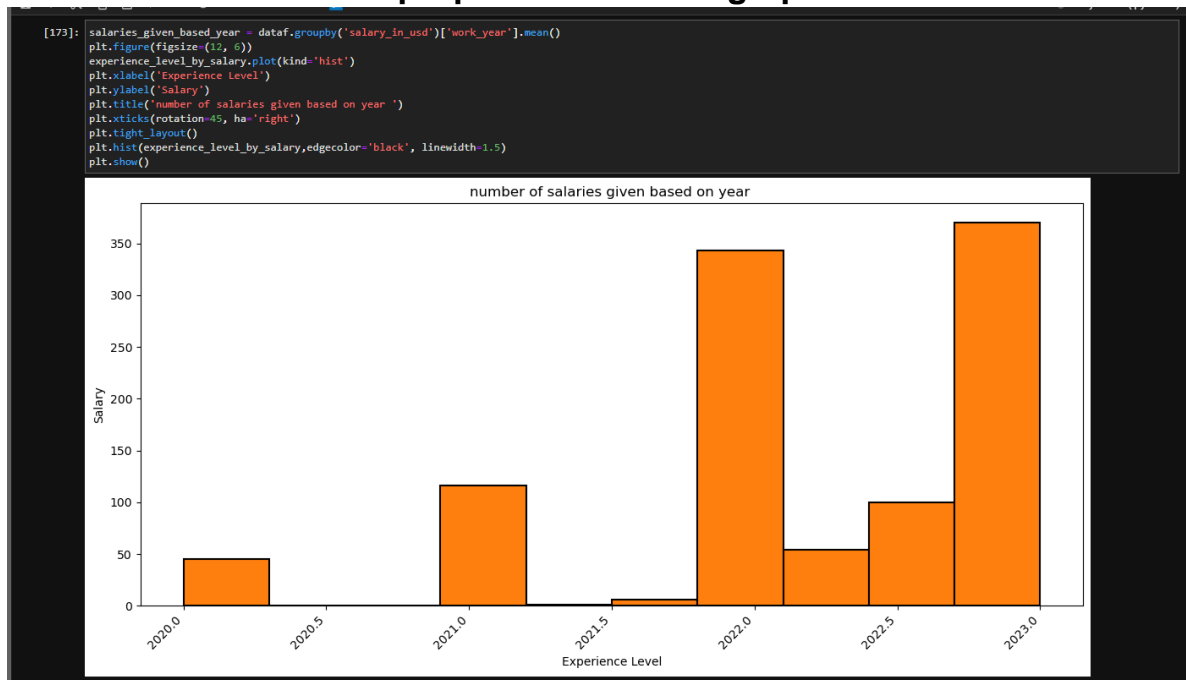


*Figure 23 histogram to plot number of salaries which was given based on the year.*

These code of lines plots the number of salaries which was given every year from 2020 to 2023 in the histogram.

```
[174]: plt.boxplot(salaries_given_based_year)
       plt.title('Box Plot of number of salaries given in each year')
       plt.ylabel('Year')
       plt.show()
```
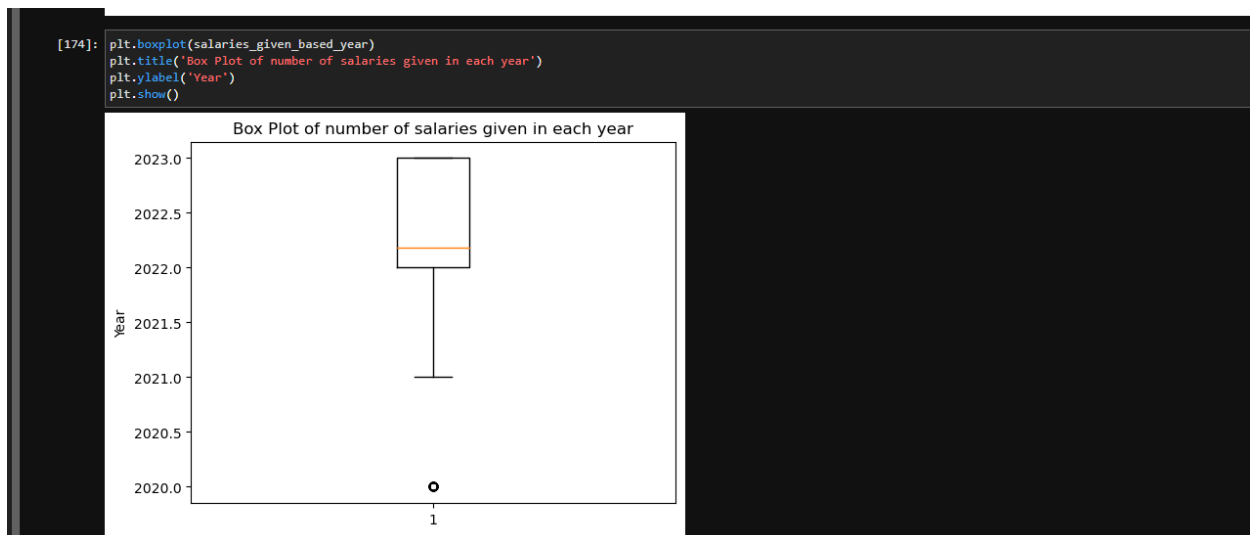


*Figure 24 box plot to plot the number of salaries given in each year.*

These codes of lines plots in box plot to show the number of salaries given to the employees each year from 2020 to 2023.

13

Amogh Man Bajracharya

## 6. References

datacamp. (2023, july). *What is Data Analysis? An Expert Guide With Examples*. Retrieved from datacamp: https://www.datacamp.com/blog/what-is-data-analysis-expert-guide

Geek for Geeks. (2024, January 25). *ML | Overview of Data Cleaning*. Retrieved from Geek for Geeks: https://www.geeksforgeeks.org/data-cleansing-introduction/

Qlik. (2024). *Data Exploration*. Retrieved from Qilk: https://www.qlik.com/us/data-analytics/data-exploration

Secoda. (2024). *What is Data Preperation?* Retrieved from Secoda: https://www.secoda.co/glossary/data-preperation-definition

Amogh Man Bajracharya

Amogh Man Bajracharya