

```

1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 %matplotlib inline

```

```

1 df = pd.read_csv('/content/501+Case1+Dataset (4).csv')
2

```

```
1 df
```

| | Employee_Name | EmpID | Sex | GenderID | MaritalDesc | MarriedID | MaritalStatusID | DOB | State | Zip | ... | ManagerID | RecruitmentSource |
|-----|-------------------|-------|-----|----------|-------------|-----------|-----------------|------------|-------|------|-----|-----------|-------------------|
| 0 | Le, Binh | 10232 | F | 0 | Single | 0 | 0 | 06/14/87 | MA | 1886 | ... | 13.0 | Inde |
| 1 | Martin, Sandra | 10110 | F | 0 | Single | 0 | 0 | 11/07/1987 | MA | 2135 | ... | 10.0 | Google Sea |
| 2 | Myers, Michael | 10216 | M | 1 | Single | 0 | 0 | 04/18/80 | MA | 1550 | ... | 20.0 | Linke |
| 3 | Navathe, Kurt | 10079 | M | 1 | Single | 0 | 0 | 04/25/70 | MA | 2056 | ... | 13.0 | Inde |
| 4 | Sutwell, Barbara | 10209 | F | 0 | Single | 0 | 0 | 08/15/68 | MA | 2718 | ... | 16.0 | Inde |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 306 | Wilkes, Annie | 10204 | F | 0 | Divorced | 0 | 2 | 07/30/83 | MA | 1876 | ... | 19.0 | Google Sea |
| 307 | Demita, Carla | 10100 | F | 0 | Separated | 0 | 3 | 02/25/51 | MA | 2343 | ... | 18.0 | Google Sea |
| 308 | Lundy, Susan | 10096 | F | 0 | Widowed | 0 | 4 | 12/26/76 | MA | 2122 | ... | 22.0 | Linke |
| 309 | MacLennan, Samuel | 10191 | M | 1 | Widowed | 0 | 4 | 11/09/1972 | MA | 1938 | ... | 11.0 | Inde |
| 310 | Thibaud, Kenneth | 10268 | M | 1 | Widowed | 0 | 4 | 09/16/75 | MA | 2472 | ... | 39.0 | Ot |

311 rows × 33 columns

```

1 #1 Import the data & check the head, tail for it ?
2 df. head()

```

| | Employee_Name | EmpID | Sex | GenderID | MaritalDesc | MarriedID | MaritalStatusID | DOB | State | Zip | ... | ManagerID | RecruitmentSource |
|---|------------------|-------|-----|----------|-------------|-----------|-----------------|------------|-------|------|-----|-----------|-------------------|
| 0 | Le, Binh | 10232 | F | 0 | Single | 0 | 0 | 06/14/87 | MA | 1886 | ... | 13.0 | Inde |
| 1 | Martin, Sandra | 10110 | F | 0 | Single | 0 | 0 | 11/07/1987 | MA | 2135 | ... | 10.0 | Google Search |
| 2 | Myers, Michael | 10216 | M | 1 | Single | 0 | 0 | 04/18/80 | MA | 1550 | ... | 20.0 | LinkedIn |
| 3 | Navathe, Kurt | 10079 | M | 1 | Single | 0 | 0 | 04/25/70 | MA | 2056 | ... | 13.0 | Inde |
| 4 | Sutwell, Barbara | 10209 | F | 0 | Single | 0 | 0 | 08/15/68 | MA | 2718 | ... | 16.0 | Inde |

5 rows × 33 columns

```
1 df.tail
```

1 #2 Check the shape, size of the data ?

2 df.shape

```
(311, 33)
```

1 df .size

```
10263
```

```
Kenneth
```

1 df .info

```
<bound method DataFrame.info of
0      Le, Binh    10232  F      0      Single      0
1      Martin, Sandra 10110  F      0      Single      0
2      Myers, Michael 10216  M      1      Single      0
3      Navathe, Kurt 10079  M      1      Single      0
4      Sutwell, Barbara 10209  F      0      Single      0
..
306     Wilkes, Annie 10204  F      0      Divorced      0
307     Demita, Carla 10100  F      0      Separated      0
308     Lundy, Susan 10096  F      0      Widowed      0
309  MacLennan, Samuel 10191  M      1      Widowed      0
310     Thibaud, Kenneth 10268  M      1      Widowed      0

      MaritalStatusID      DOB State  Zip  ... ManagerID RecruitmentSource \
0      0      06/14/87  MA 1886  ...      13.0      Indeed
1      0      11/07/1987  MA 2135  ...      10.0      Google Search
2      0      04/18/80  MA 1550  ...      20.0      LinkedIn
3      0      04/25/70  MA 2056  ...      13.0      Indeed
4      0      08/15/68  MA 2718  ...      16.0      Indeed
..
306     2      07/30/83  MA 1876  ...      19.0      Google Search
307     3      02/25/51  MA 2343  ...      18.0      Google Search
308     4      12/26/76  MA 2122  ...      22.0      LinkedIn
309     4      11/09/1972  MA 1938  ...      11.0      Indeed
310     4      09/16/75  MA 2472  ...      39.0      Other

      LastPerformanceReview_Date PerformanceScore PerfScoreID EngagementSurvey \
0      01/08/2019      Fully Meets      3      4.10
1      1/14/2019      Fully Meets      3      4.50
2      1/22/2019      Fully Meets      3      4.10
3      2/25/2019      Fully Meets      3      5.00
4      1/31/2019      Fully Meets      3      3.40
..
306     02/06/2011      Fully Meets      3      3.60
307     05/06/2015      Fully Meets      3      4.62
308     06/10/2016      Fully Meets      3      4.65
309     04/01/2017      Fully Meets      3      3.08
310     7/14/2010      Fully Meets      3      4.10

      EmpSatisfaction SpecialProjectsCount DaysLateLast30 Absences
0      5      7      0      2
1      5      4      0      14
2      4      0      0      13
3      3      6      0      17
4      5      0      0      13
..
306     5      0      0      9
307     5      0      0      1
308     4      0      0      15
309     4      0      0      18
310     4      0      0      15
```

```
[311 rows x 33 columns]>
```

1 #3.How many columns have categorical features ?

2 len(df.select_dtypes(include="object").columns)

```
17
```

1 4 How many unique values are present in RaceDesc column ?

2 df['Racedesc'].unique()

```
File "<ipython-input-26-4a7ba73c9cdf>", line 1
  4 How many unique values are present in RaceDesc column?
    ^
```

```
SyntaxError: invalid syntax
```

```
1 #5 Check for mean, max .value, min .value, count, standard deviation of ManagerID
2 column ?
3 df['ManagerID'].describe()
```

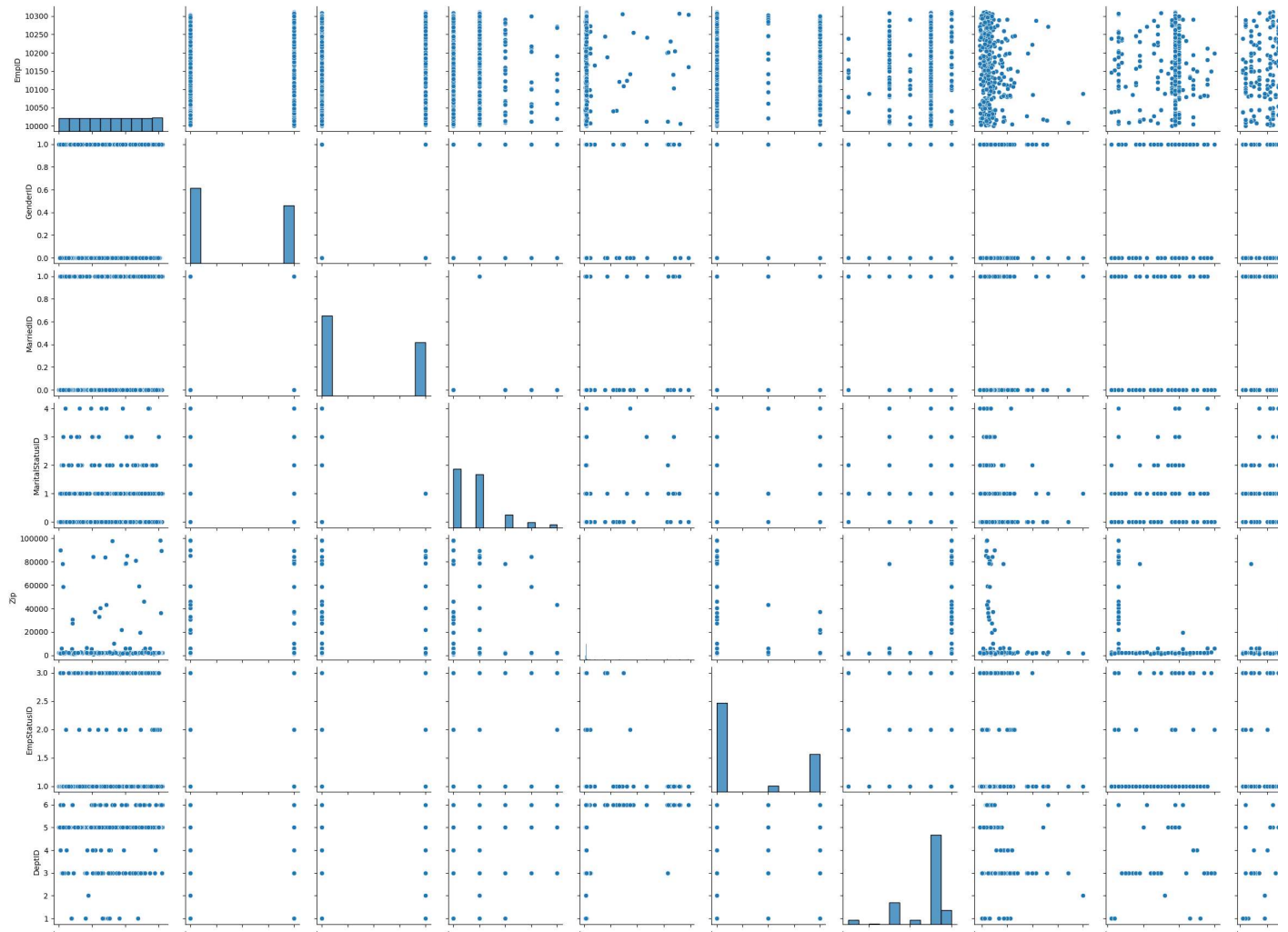
```
count    303.000000
mean      14.570957
std        8.078306
min         1.000000
25%       10.000000
50%       15.000000
75%       19.000000
max       39.000000
Name: ManagerID, dtype: float64
```

```
1 #6 Count the no of categorical , numerical columns ?
2 len(df.select_dtypes(exclude="object").columns)
```

```
16
```

```
1 # 7 Make a diversity report about the dataset?
2 import matplotlib.pyplot as plt
3 import seaborn as sns
4 plt.figure(figsize=(13,7))
5 sns.pairplot(data=df)
6
```

```
<seaborn.axisgrid.PairGrid at 0x7fe6339fc760>
<Figure size 1300x700 with 0 Axes>
```



1 #8 Which columns have correlation with each other, what will you interpret from it?

2 df.corr()

3

```
<ipython-input-5-2db1fe21e460>:2: FutureWarning: The default value of numeric_only in DataFrame.corr is deprecated. In a future version, df.corr()
```

```
1 #9 Add a new column which specifies the number of characters in Employee_Name
2 column?
3 df['employee_char'] = df['Employee_Name'].str.len()
4 df.head(2)
```

| | Employee_Name | EmpID | Sex | GenderID | MaritalDesc | MarriedID | MaritalStatusID | DOB | State | Zip | ... | RecruitmentSource | LastPerfc |
|---|----------------|-------|-----|----------|-------------|-----------|-----------------|------------|-------|------|-----|-------------------|-----------|
| 0 | Le, Binh | 10232 | F | 0 | Single | 0 | 0 | 06/14/87 | MA | 1886 | ... | Indeed | |
| 1 | Martin, Sandra | 10110 | F | 0 | Single | 0 | 0 | 11/07/1987 | MA | 2135 | ... | Google Search | |

2 rows × 34 columns

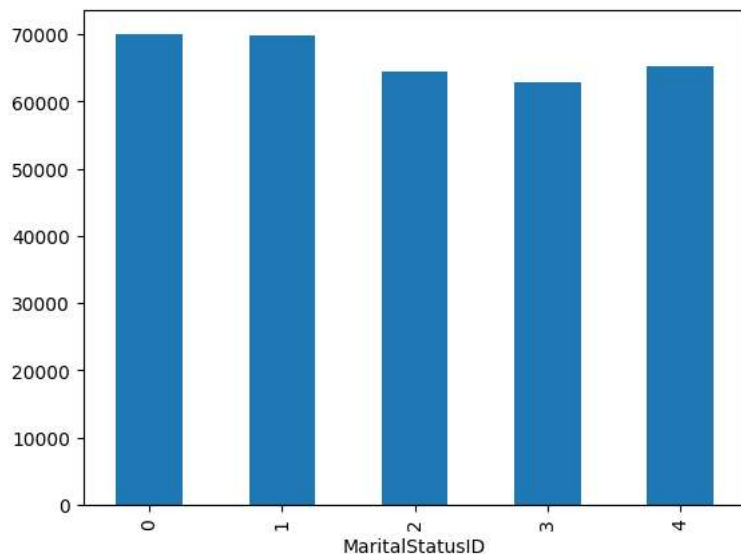
```
1 #10 Round up the values in EngagementSurvey column ?
2 df['EngagementSurvey'].round()
3
4
```

```
0      4.0
1      4.0
2      4.0
3      5.0
4      3.0
...
306    4.0
307    5.0
308    5.0
309    3.0
310    4.0
```

Name: EngagementSurvey, Length: 311, dtype: float64

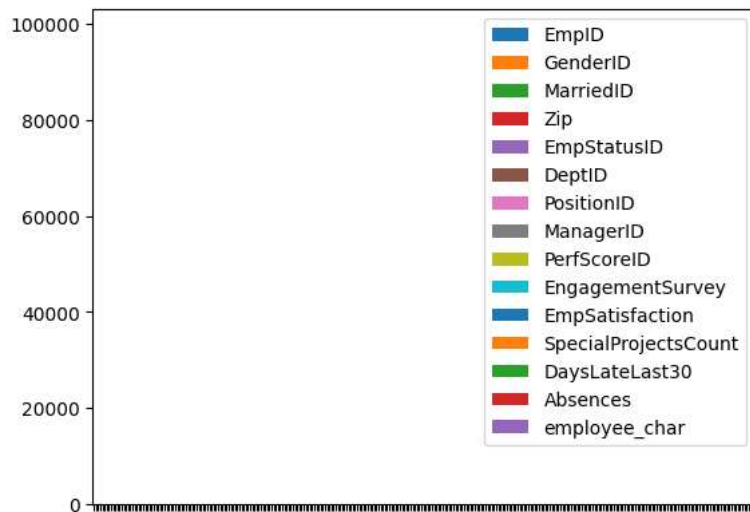
```
1 #11 a. Does marital status have any impact on salary?
2 df.groupby(['MaritalStatusID'])['Salary'].mean().plot(kind='bar')
```

<Axes: xlabel='MaritalStatusID'>



```
1 #11
2
3 df.groupby(['MaritalStatusID', 'Salary']).mean().plot(kind='bar')
```

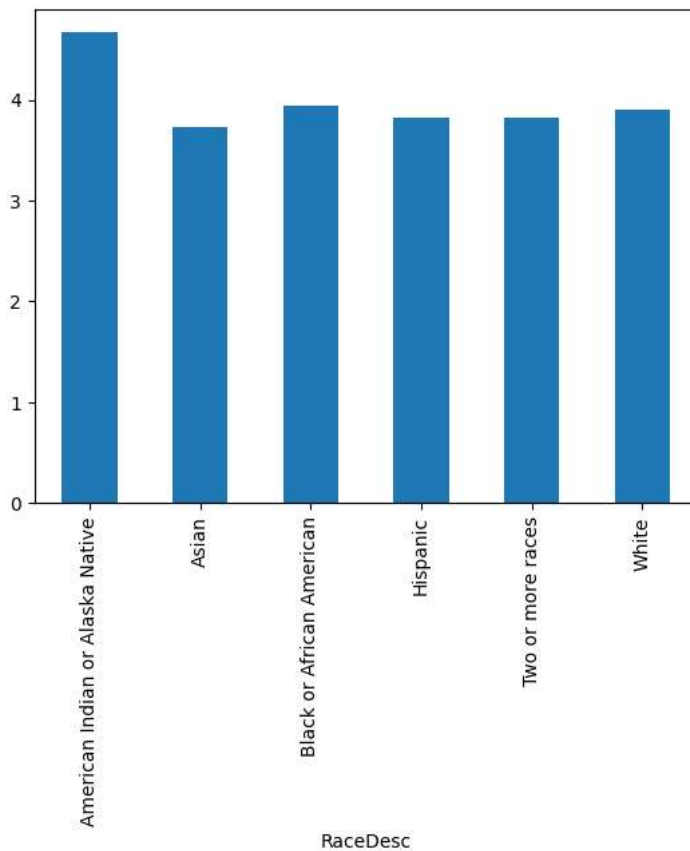
```
<ipython-input-48-d32890a7ede7>:3: FutureWarning: The default value of numeric_only in DataFrameGroupBy.mean is deprecated. In a future \
df.groupby(['MaritalStatusID','Salary']).mean().plot(kind='bar')
<Axes: xlabel='MaritalStatusID,Salary'>
```



1 ##11 (b) Does RaceDesc have any impact on EmpSatisfaction?

```
2 df.groupby(['RaceDesc'])['EmpSatisfaction'].mean().plot(kind='bar')
```

<Axes: xlabel='RaceDesc'>



1 ##11 (3)

1 ## CASE 2

```
1 db = pd.read_csv('//content/502+Case2+Dataset.csv')
```

```
1 db.columns
```

```
Index(['Unnamed: 0', 'work_year', 'experience_level', 'employment_type',
       'job_title', 'salary', 'salary_currency', 'salary_in_usd',
```

```
'employee_residence', 'remote_ratio', 'company_location',
'company_size'],
dtype='object')
```

```
1 db.head(5)
```

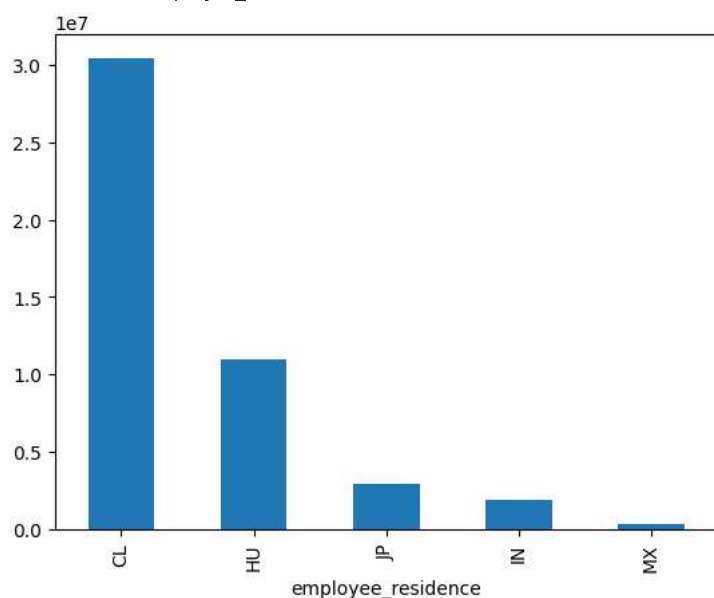
| | Unnamed: 0 | work_year | experience_level | employment_type | job_title | salary | salary_currency | salary_in_usd | employee_res |
|---|------------|-----------|------------------|-----------------|----------------------------|--------|-----------------|---------------|--------------|
| 0 | 0 | 2020 | MI | FT | Data Scientist | 70000 | EUR | 79833 | |
| 1 | 1 | 2020 | SE | FT | Machine Learning Scientist | 260000 | USD | 260000 | |
| 2 | 2 | 2020 | SE | FT | Big Data Engineer | 85000 | GBP | 109024 | |
| 3 | 3 | 2020 | MI | FT | Product Data Analyst | 20000 | USD | 20000 | |
| 4 | 4 | 2020 | SE | FT | Machine Learning Engineer | 150000 | USD | 150000 | |

```
1 ##1 How does where you live affect your salary?
```

```
2
```

```
3 db.groupby(['employee_residence'])['salary'].mean().sort_values(ascending=False).head().plot(kind='bar')
```

<Axes: xlabel='employee_residence'>



```
1 #2 How has the demand of the jobs been throughout the years?
```

```
2 db['work_year'].value_counts()
```

```
3
```

```
2022    318
2021    217
2020     72
Name: work_year, dtype: int64
```

```
1 ##3 How common is to work remote?
```

```
2 db['remote_ratio'].value_counts()
```

```
3
```

```
100    381
0       127
50      99
Name: remote_ratio, dtype: int64
```

```
1 ##4 What the highest paying jobs with entry level as well as for senior level experienced?
```

```
2
```

```
3 db[(db['experience_level'] == "EN") | (db['experience_level'] == 'SE')]
```

```
4
```

```
5 db.groupby(['experience_level'])['salary'].max()
```

```
6
```

```

experience_level
EN      4450000
EX      6000000
MI      3040000
SE      7000000
Name: salary, dtype: int64

```

```

1 ##5 What company size hire the most?
2 db['company_size'].value_counts().head()
3

```

```

M      326
L      198
S       83
Name: company_size, dtype: int64

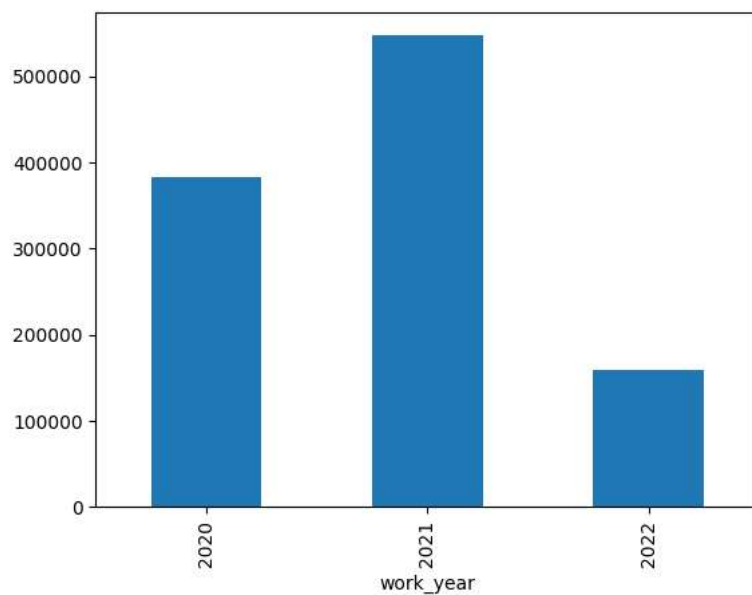
```

```

1 ##6 How has average salary changed throughout the years?
2 db.groupby(['work_year'])['salary'].mean().plot(kind='bar')
3

```

<Axes: xlabel='work_year'>



```

1 ##7 What are most popular roles in Data Science ?
2
3 db['job_title'].value_counts()

```

```

Data Scientist      143
Data Engineer       132
Data Analyst        97
Machine Learning Engineer  41
Research Scientist  16
Data Science Manager  12
Data Architect      11
Big Data Engineer   8
Machine Learning Scientist  8
Principal Data Scientist  7
AI Scientist        7
Data Science Consultant  7
Director of Data Science  7
Data Analytics Manager  7
ML Engineer         6
Computer Vision Engineer  6
BI Data Analyst     6
Lead Data Engineer  6
Data Engineering Manager  5
Business Data Analyst  5
Head of Data        5
Applied Data Scientist  5
Applied Machine Learning Scientist  4
Head of Data Science  4
Analytics Engineer  4
Data Analytics Engineer  4

```

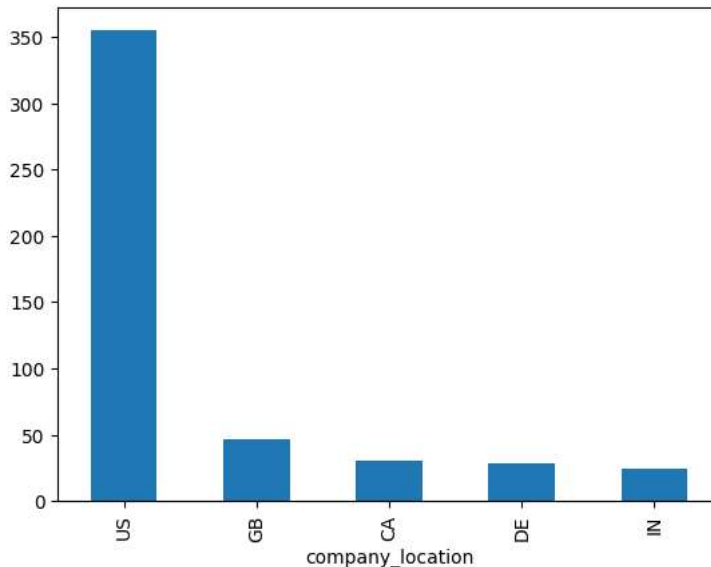

| | |
|--|---|
| Machine Learning Developer | 3 |
| Machine Learning Infrastructure Engineer | 3 |
| Lead Data Scientist | 3 |
| Computer Vision Software Engineer | 3 |
| Lead Data Analyst | 3 |
| Data Science Engineer | 3 |
| Principal Data Engineer | 3 |
| Principal Data Analyst | 2 |
| ETL Developer | 2 |
| Product Data Analyst | 2 |
| Director of Data Engineering | 2 |
| Financial Data Analyst | 2 |
| Cloud Data Engineer | 2 |
| Lead Machine Learning Engineer | 1 |
| NLP Engineer | 1 |
| Head of Machine Learning | 1 |
| 3D Computer Vision Researcher | 1 |
| Data Specialist | 1 |
| Staff Data Scientist | 1 |
| Big Data Architect | 1 |
| Finance Data Analyst | 1 |
| Marketing Data Analyst | 1 |
| Machine Learning Manager | 1 |
| Data Analytics Lead | 1 |

Name: job_title, dtype: int64

1 ##8 Which country hire the most people in Data Science?

2 db.groupby(['company_location'])['job_title'].count().sort_values(ascending=False).head().plot(kind='bar')

<Axes: xlabel='company_location'>



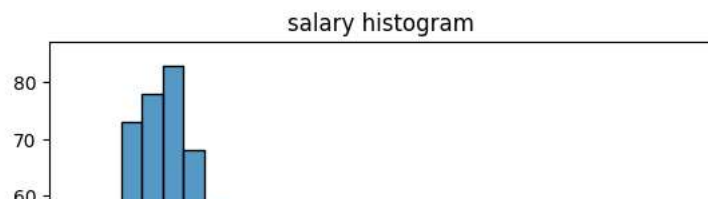
1 ##9 What is the distribution of Salaries?

2

3 sns.histplot(db['salary_in_usd'])

4 plt.title('salary histogram')

5 plt.show()



1 ##10. How much can you expect depending on your years of experience?

2 db.groupby(['experience_level'])['salary'].mean()

3

```
experience_level
EN    264622.454545
EX    427072.115385
MI    480617.690141
SE    213949.353571
Name: salary, dtype: float64
```



1 ##11. Which year do people prefer to stay at home the most?

2 db.groupby(['remote_ratio'])['work_year'].value_counts()

```
remote_ratio  work_year
0             2022      78
              2021      34
              2020      15
50            2021      66
              2020      21
              2022      12
100           2022     228
              2021     117
              2020      36
Name: work_year, dtype: int64
```

1 ##12. Which country has the highest pay?

2 db.groupby(['company_location'])['salary'].max().head(5)

3

4

```
company_location
AE    120000
AS    1335000
AT     80000
AU    150000
BE     75000
Name: salary, dtype: int64
```

1 ##13. Which job title has the highest pay?

2 db.groupby(['job_title'])['salary'].max().head().sort_values(ascending=False)

```
job_title
AI Scientist          1335000
Applied Machine Learning Scientist  423000
3D Computer Vision Researcher      400000
Applied Data Scientist      380000
Analytics Engineer      205300
Name: salary, dtype: int64
```

1 ##14 . Is freelancing is worth or not?

2 db.groupby(['employment_type'])['salary'].agg(['mean', 'max', 'min', 'sum'])



| | mean | max | min | sum |
|-----------------|---------------|----------|-------|-----------|
| employment_type | | | | |
| CT | 184000.000000 | 416000 | 29000 | 920000 |
| FL | 48000.000000 | 100000 | 12000 | 192000 |
| FT | 331124.622449 | 30400000 | 4000 | 194701278 |
| PT | 85476.000000 | 400000 | 8760 | 854760 |

```

1 ### CASE - 1 PART A
2
3 import pandas as pd
4 import numpy as np
5 import matplotlib.pyplot as plt
6 import seaborn as sns
7 %matplotlib inline
8
9
10

```

```

1 ##1 1. Write a program to display "Welcome " if a number entered by user is a multiple of five otherwise pri
2
3 number = int(input("Enter a number: "))
4
5 if number % 5 == 0:
6     print("Welcome")
7 else:
8     print("Bye")

```

```

Enter a number: 4
Bye

```

```

1 ##2 2. Write a program to find the largest number out of three numbers excepted from user ?
2 num1 = int(input("enter the first number:"))
3 num2 = int(input("enter the second number:"))
4 num3 = int(input("enter the third number:"))
5 largest = max(num1, num2, num3)
6 print("the largest number is :",largest)

```

```

enter the first number:45
enter the second number:98
enter the third number:98
the largest number is : 98

```

```

1 ##3. Write a program to accept a number from 1 to 7 and display the name of the day like 1 for sunday , 2
2 day = int(input("Enter a number from 1 to 7: "))
3
4 if day == 1:
5     print("Sunday")
6 elif day == 2:
7     print("Monday")
8 elif day == 3:
9     print("Tuesday")
10 elif day == 4:
11     print("Wednesday")
12 elif day == 5:
13     print("Thursday")
14 elif day == 6:
15     print("Friday")
16 elif day == 7:
17     print("Saturday")
18 else:
19     print("Invalid input. Please enter a number from 1 to 7.")
20
21

```

```

Enter a number from 1 to 7: 4
Wednesday

```

```

1 #4. 4. Accept any city from the user and display monument of that city ?
2 city = input("Enter a city: ")
3

```

```

4 if city == "Delhi":
5     print("Monument: Red Fort")
6 elif city == "Agra":
7     print("Monument: Taj Mahal")
8 elif city == "Jaipur":
9     print("Monument: Jal Mahal")
10 else:
11     print("Monument not found for the given city.")

```

```

Enter a city: jaipur
Monument not found for the given city.

```

```

1 # 5. Write a program to accept two numbers and mathematical operators and perform
2 #operation accordingly
3 #Like:
4 #Enter First Number: 7
5 #Enter Second Number : 9
6 #Enter operator : + (you can use different operators as well)
7 #Your Answer is : 16
8 First_number = int(input('First_number: '))
9 Second_number = int(input('Second_number: '))
10 operator = input('Enter operator (+,-,*,/,%)')
11 if operator == '+':
12     print(First_number+Second_number)
13 elif operator == '-':
14     print(First_number-Second_number)
15 elif operator == '*':
16     print(First_number*Second_number)
17 elif operator == '/':
18     print(First_number%Second_number)
19 else:
20     print('Invalid')
21
22
23

```

```

First_number: 20
Second_number: 10
Enter operator (+,-,*,/,%)+
Invalid

```

```

1 #6 Check if the input is Leap Year , write a function
2 #We add a Leap Day on February 29, almost every four years.
3 #The leap day is an extra, or intercalary day and we add it to the shortest month of the year,
4 #February.
5 #In the Gregorian calendar three criteria must be taken into account to identify leap years:
6 #1. The year can be evenly divided by 4, is a leap year, unless:
7 #2. The year can be evenly divided by 100, it is NOT a leap year, unless:
8 #3. The year is also evenly divisible by 400. Then it is a leap year.
9 #Examples : This means that in the Gregorian calendar, the years 2000 and 2400 are leap
10 #years, while 1800, 1900, 2100, 2200, 2300 and 2500 are NOT leap years.
11 #What you have to do?
12 #You are given the year, and you have to write a function to check if the year is leap or not.
13 #Note that you have to complete the function and remaining code is given as template.
14 #You can use a variable as input with a fix value of user input
15 Year = int(input('Enter your year here: '))
16 if Year%400==0:
17     print('Leap Year')
18 elif Year%100==0:
19     print('No Leap Year')
20 elif Year%4==0:
21     print('Leap Year')
22 else:
23     print('No Leap Year')

```

Enter your year here: 2220
Leap Year

```
1 #7 # Rock Paper Scissors , if else or while ?
2 #Make a two-player Rock-Paper-Scissors game. (Hint: Ask for player plays (using input),
3 #compare them, print out a message of congratulations to the winner, and ask if the players
4 #want to start a new game)
5 #Remember the rules:
6 #1. Rock beats scissors
7 #2. Scissors beats paper
8 #3. Paper beats rock
9 #If you don't know what is this game , go watch it :)
10
11 Player_1 = str(input())
12 Player_2 = str(input())
13 if Player_1==Player_2:
14     print('Tie')
15 elif Player_1 == "Rock" and Player_2 == "Scissors":
16     print('Player_1 wins')
17 elif Player_1 == "Rock" and Player_2 == "Paper":
18     print('Player_2 wins')
19 elif Player_1 == 'Paper' and Player_2 == 'Rock':
20     print('Player_2 wins')
21 elif Player_1 == "paper" and Player_2 == "Scissors":
22     print('Player_2 wins')
23 elif Player_1 == "Scissors" and Player_2 == "Rock":
24     print('Player_1 wins')
25 else:
26     print('Player_1 wins')

Rock
paper
Player_1 wins
```

1