```
In [ ]: # Practicle 1 : Preprocessing of dataset
         # Preprocessing for data sciencentist salary dataset
         # 2021BIT023
In [1]:
         import pandas as pd
         import matplotlib.pyplot as plt
         import numpy as np
         df=pd.read_csv(ds_salaries1.csv)
In [2]:
                                                       Traceback (most recent call last)
         NameError
         Cell In[2], line 1
         ----> 1 df=pd.read_csv(ds_salaries1.csv)
                2 df
         NameError: name 'ds_salaries1' is not defined
         df=pd.read_csv("ds_salaries1.csv")
In [5]:
In [6]:
         df.head(5)
            work_year experience_level employment_type job_title
Out[6]:
                                                                   salary salary_currency salary_in_usd €
                                                        Principal
         0
                 2023
                                   SE
                                                    FT
                                                           Data
                                                                  0.00008
                                                                                    EUR
                                                                                              85847.0
                                                        Scientist
                                                             ML
         1
                                                                                              30000.0
                 2023
                                   MΙ
                                                                  30000.0
                                                                                    USD
                                                        Engineer
                                                             ML
         2
                 2023
                                   MΙ
                                                                  25500.0
                                                                                    USD
                                                                                              25500.0
                                                        Engineer
                                                           Data
         3
                 2023
                                   SE
                                                                 175000.0
                                                                                    USD
                                                                                             175000.0
                                                        Scientist
                                                           Data
         4
                 2023
                                   SE
                                                                 120000.0
                                                                                    USD
                                                                                             120000.0
                                                        Scientist
In [7]:
         df.tail(5)
```

| Out[7]: | V | work_year | experience_level | employment_ty | ре | job_title | e salary | salary_currency | salary_in_u | |
|---------|--|---|-------------------------------------|----------------|-----|------------------------------|-------------|-------------------|-------------|--|
| | 3750 | 2020 | SE | | FT | Data Scientis | | USD | 412000 | |
| | 3751 | 2021 | MI | | FT | Principa Data Scientis | a 151000.0 | USD | 151000 | |
| | 3752 | 2020 | EN | | FT | Data Scientis | | USD | 10500(| |
| | 3753 | 2020 | EN | | СТ | Busines Data Analys | a 100000.0 | USD | 100000 | |
| | 3754 | 2021 | SE | | FT | Data Science Manage | e 7000000.0 | INR | 94665 | |
| | | | | | | | | | • | |
| | df.shape # row, attributes | | | | | | | | | |
| | (3755, 11) | | | | | | | | | |
| | df1=df.copy() | | | | | | | | | |
|)]: | df1.he | ad(2) | | | | | | | | |
| 0]: | worl | k_year exp | erience_level er | nployment_type | joł | o_title | salary sala | ry_currency salar | y_in_usd er | |
| | 0 | 2023 | SE | FT | | ncipal Data 8 ientist | 0000.0 | EUR | 85847.0 | |
| | 1 | 2023 | MI | СТ | Eng | ML gineer ³ | 0000.0 | USD | 30000.0 | |
| | | | | | | | | | • | |
| 9]: | df1.is | null().sur | n() | | | | | | | |
| [9]: | employed job_ti salary salary salary employed remote | ence_level ment_type tle _currency _in_usd ee_resider _ratio y_locatior | 0 0 2 1 3 nce 0 1 | | | | | | | |
| 8]: | | series |)) # false for | r not null dat | a a | ani tru | e for NAN | value i.e for | null value | |

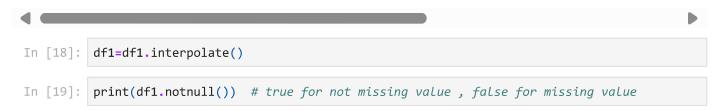
```
experience level
                                                employment type
                                                                  job title
                work year
                                                                              salary \
          0
                     False
                                        False
                                                           False
                                                                       False
                                                                               False
          1
                     False
                                        False
                                                           False
                                                                       False
                                                                               False
          2
                     False
                                        False
                                                           False
                                                                       False
                                                                               False
          3
                     False
                                        False
                                                           False
                                                                       False
                                                                               False
          4
                     False
                                        False
                                                           False
                                                                       False
                                                                               False
                       . . .
                                           . . .
                                                             . . .
                                                                         . . .
                                                                                  . . .
          3750
                     False
                                        False
                                                           False
                                                                       False
                                                                               False
                     False
                                        False
                                                                       False
                                                                               False
          3751
                                                           False
          3752
                     False
                                        False
                                                           False
                                                                       False
                                                                               False
          3753
                     False
                                        False
                                                           False
                                                                       False
                                                                               False
          3754
                     False
                                        False
                                                           False
                                                                       False
                                                                               False
                 salary_currency
                                   salary_in_usd
                                                   employee_residence remote_ratio
          0
                           False
                                            False
                                                                 False
                                                                                 False
                                                                                False
          1
                           False
                                            False
                                                                 False
          2
                           False
                                            False
                                                                 False
                                                                                 False
          3
                           False
                                            False
                                                                 False
                                                                                 False
                           False
                                                                                 False
          4
                                            False
                                                                 False
                              . . .
                                              . . .
                                                                    . . .
                                                                                   . . .
          . . .
          3750
                           False
                                            False
                                                                 False
                                                                                 False
          3751
                           False
                                            False
                                                                 False
                                                                                False
                                                                                 False
          3752
                           False
                                            False
                                                                 False
          3753
                           False
                                            False
                                                                 False
                                                                                 False
          3754
                                                                 False
                                                                                 False
                           False
                                            False
                company_location
                                    company_size
          0
                            False
                                            False
          1
                            False
                                            False
          2
                            False
                                            False
          3
                            False
                                            False
          4
                            False
                                            False
                               . . .
                                              . . .
          3750
                            False
                                            False
          3751
                            False
                                            False
          3752
                            False
                                            False
          3753
                             False
                                            False
          3754
                             False
                                            False
          [3751 rows x 11 columns]
          df1.dropna(inplace=True)
In [12]:
          df1.isnull().sum()
In [13]:
                                  0
          work year
Out[13]:
          experience level
                                  0
                                  0
          employment_type
                                  0
          job_title
                                  0
          salary
          salary_currency
                                  0
          salary_in_usd
                                  0
          employee_residence
                                  0
                                  0
          remote_ratio
                                  0
          company_location
                                  0
          company size
          dtype: int64
          # in salary_currency we replice salary from INR to USD
In [15]:
          df repl=df1.replace({'INR':'USD'})
```

df_repl.head(3)

```
Out[15]:
             work_year experience_level employment_type job_title
                                                                   salary salary_currency salary_in_usd er
                                                         Principal
          0
                  2023
                                    SE
                                                                  80000.0
                                                                                    EUR
                                                                                              85847.0
                                                     FT
                                                            Data
                                                          Scientist
                                                              ML
          1
                                                                  30000.0
                                                                                    USD
                                                                                              30000.0
                  2023
                                    MΙ
                                                         Engineer
                                                              ML
                                                                  25500.0
          2
                  2023
                                    MI
                                                                                    USD
                                                                                              25500.0
                                                         Engineer
          df intrplt=df1.interpolate()
In [16]:
          df intrplt.head(2)
             work_year experience_level employment_type job_title
                                                                   salary salary_currency salary_in_usd er
Out[16]:
                                                         Principal
          0
                  2023
                                    SE
                                                     FT
                                                            Data
                                                                  80000.0
                                                                                    EUR
                                                                                              85847.0
                                                          Scientist
                                                              ML
          1
                  2023
                                    MI
                                                                  30000.0
                                                                                    USD
                                                                                              30000.0
                                                         Engineer
          df1=df1.infer object()
In [17]:
          AttributeError
                                                        Traceback (most recent call last)
          ~\AppData\Local\Temp\ipykernel_6456\2477449844.py in ?()
          ---> 1 df1=df1.infer object()
          ~\anaconda3\Lib\site-packages\pandas\core\generic.py in ?(self, name)
             5985
                                and name not in self._accessors
                                and self._info_axis._can_hold_identifiers_and_holds_name(name)
             5986
             5987
                           ):
             5988
                                return self[name]
                            return object.__getattribute__(self, name)
          -> 5989
          AttributeError: 'DataFrame' object has no attribute 'infer_object'
          df2=df1.infer_objects()
In [19]:
          df2
In [20]:
```

| Out[20]: | | work_year | experience_level | employment_type | job_title | salary | salary_currency | salary_in_u |
|----------|------|-----------|------------------|-----------------|--------------------------------|-----------|-----------------|-------------------|
| | 0 | 2023 | SE | FT | Principal Data Scientist | 80000.0 | EUR | 8584 ⁷ |
| | 1 | 2023 | MI | СТ | ML Engineer | 30000.0 | USD | 30000 |
| | 2 | 2023 | МІ | СТ | ML Engineer | 25500.0 | USD | 2550(|
| | 3 | 2023 | SE | FT | Data Scientist | 175000.0 | USD | 17500(|
| | 4 | 2023 | SE | FT | Data Scientist | 120000.0 | USD | 12000(|
| | ••• | | | | ••• | ••• | | |
| | 3750 | 2020 | SE | FT | Data Scientist | 412000.0 | USD | 412000 |
| | 3751 | 2021 | MI | FT | Principal Data Scientist | 151000.0 | USD | 151000 |
| | 3752 | 2020 | EN | FT | Data Scientist | 105000.0 | USD | 10500(|
| | 3753 | 2020 | EN | СТ | Business Data Analyst | 100000.0 | USD | 100000 |
| | 3754 | 2021 | SE | FT | Data Science Manager | 7000000.0 | INR | 9466! |

3751 rows × 11 columns



```
experience level
                                                 employment type
                                                                    job title
                                                                                salary \
                 work year
          0
                      True
                                          True
                                                             True
                                                                         True
                                                                                  True
          1
                      True
                                          True
                                                             True
                                                                         True
                                                                                  True
          2
                      True
                                          True
                                                             True
                                                                         True
                                                                                  True
          3
                      True
                                          True
                                                             True
                                                                         True
                                                                                  True
          4
                      True
                                          True
                                                             True
                                                                         True
                                                                                  True
                       . . .
                                           . . .
                                                              . . .
                                                                           . . .
                                                                                   . . .
          3750
                      True
                                          True
                                                             True
                                                                         True
                                                                                  True
                                                                         True
                                                                                  True
          3751
                      True
                                          True
                                                             True
          3752
                      True
                                          True
                                                             True
                                                                         True
                                                                                  True
          3753
                      True
                                          True
                                                             True
                                                                         True
                                                                                  True
          3754
                      True
                                          True
                                                             True
                                                                         True
                                                                                  True
                                                    employee_residence
                 salary_currency
                                    salary_in_usd
                                                                          remote_ratio
          0
                                                                    True
                             True
                                              True
                                                                                   True
          1
                             True
                                              True
                                                                    True
                                                                                   True
          2
                             True
                                              True
                                                                    True
                                                                                   True
          3
                             True
                                              True
                                                                    True
                                                                                   True
                             True
          4
                                              True
                                                                    True
                                                                                   True
                                                                     . . .
          . . .
                              . . .
                                               . . .
                                                                                    . . .
          3750
                             True
                                              True
                                                                    True
                                                                                   True
          3751
                             True
                                              True
                                                                    True
                                                                                   True
                                              True
                                                                    True
                                                                                   True
          3752
                             True
          3753
                             True
                                              True
                                                                    True
                                                                                   True
                                                                                   True
          3754
                             True
                                              True
                                                                    True
                 company_location
                                     company_size
          0
                              True
                                              True
          1
                              True
                                              True
          2
                              True
                                              True
          3
                              True
                                              True
          4
                              True
                                              True
                               . . .
                                               . . .
          3750
                              True
                                              True
          3751
                              True
                                              True
          3752
                              True
                                              True
          3753
                              True
                                              True
          3754
                              True
                                              True
          [3751 rows x 11 columns]
          df1.notnull().sum()
In [20]:
          work_year
                                   3751
Out[20]:
          experience level
                                   3751
                                   3751
          employment_type
          job_title
                                   3751
          salary
                                   3751
          salary_currency
                                   3751
          salary in usd
                                   3751
          employee residence
                                   3751
          remote_ratio
                                   3751
          company_location
                                   3751
                                   3751
          company size
          dtype: int64
          df2=df.copy()
In [21]:
          df2_filna=df2.fillna(0)
In [23]:
```

```
In [24]: #filter data
filter_data=df2[df['salary']>=81000]
filter_data
```

| Out[24]: | | work_year | experience_level | employment_type | job_title | salary | salary_currency | salary_in_u |
|----------|------|-----------|------------------|-----------------|--------------------------------|-----------|-----------------|-------------|
| | 3 | 2023 | SE | FT | Data Scientist | 175000.0 | USD | 175000 |
| | 4 | 2023 | SE | FT | Data Scientist | 120000.0 | USD | 120000 |
| | 5 | 2023 | SE | FT | Applied Scientist | 222200.0 | USD | 22220(|
| | 6 | 2023 | SE | FT | Applied Scientist | 136000.0 | USD | 136000 |
| | 7 | 2023 | SE | FT | Data Scientist | 219000.0 | USD | 21900(|
| | ••• | | ••• | | | ••• | ••• | |
| 3 | 3750 | 2020 | SE | FT | Data Scientist | 412000.0 | USD | 41200(|
| 8 | 3751 | 2021 | MI | FT | Principal Data Scientist | 151000.0 | USD | 151000 |
| 3 | 3752 | 2020 | EN | FT | Data Scientist | 105000.0 | USD | 105000 |
| 3 | 3753 | 2020 | EN | СТ | Business Data Analyst | 100000.0 | USD | 100000 |
| 3 | 3754 | 2021 | SE | FT | Data Science Manager | 7000000.0 | INR | 94665 |

3082 rows × 11 columns

In [37]: # dict2lst=df2.to_dict(orient='list')
print(dict2lst)