

```
import pandas as pd
import matplotlib.pyplot as plt
import numpy as np
import seaborn as sns

df = pd.read_csv('data/sciencefile.csv')

In [2]:
df

Out[3]:
Unnamed: 0  work_year  experience_level  employment_type  job_title  salary_currency  salary_in_usd  employee_residence  remote_ratio  company_location  company_size
0          0          2020                MI              FT      Data Scientist      EUR      79833              DE              0              DE              L
1          1          2020                SE              FT      Machine Learning Scientist  USD      260000             JP              0              JP              S
2          2          2020                SE              FT      Big Data Engineer      GBP      109024             GB              50              GB              M
3          3          2020                MI              FT      Product Data Analyst      USD      20000              HN              0              HN              S
4          4          2020                SE              FT      Machine Learning Engineer      USD      150000             US              50              US              L
--          --          --                --              --              --              --              --              --              --              --
602         602         2022                SE              FT      Data Engineer      USD      154000             US              100              US              M
603         603         2022                SE              FT      Data Engineer      USD      126000             US              100              US              M
604         604         2022                SE              FT      Data Analyst      USD      129000             US              0              US              M
605         605         2022                SE              FT      Data Analyst      USD      150000             US              100              US              M
606         606         2022                MI              FT      AI Scientist      USD      200000             IN              100              US              L

607 rows x 12 columns

In [ ]:

In [4]: df.shape

Out[4]: (607, 12)

In [8]: df.drop(['Unnamed: 0'], 'salary', axis = 1, inplace = True)

In [9]: df

Out[9]:
work_year  experience_level  employment_type  job_title  salary_currency  salary_in_usd  employee_residence  remote_ratio  company_location  company_size
0          2020             MI              FT      Data Scientist      EUR      79833              DE              0              DE              L
1          2020             SE              FT      Machine Learning Scientist  USD      260000             JP              0              JP              S
2          2020             SE              FT      Big Data Engineer      GBP      109024             GB              50              GB              M
3          2020             MI              FT      Product Data Analyst      USD      20000              HN              0              HN              S
4          2020             SE              FT      Machine Learning Engineer      USD      150000             US              50              US              L
--          --          --                --              --              --              --              --              --              --
602         2022             SE              FT      Data Engineer      USD      154000             US              100              US              M
603         2022             SE              FT      Data Engineer      USD      126000             US              100              US              M
604         2022             SE              FT      Data Analyst      USD      129000             US              0              US              M
605         2022             SE              FT      Data Analyst      USD      150000             US              100              US              M
606         2022             MI              FT      AI Scientist      USD      200000             IN              100              US              L

607 rows x 10 columns

In [ ]:

In [15]: df1 = df.groupby('work_year')['salary_in_usd'].mean().round(2)
df1

Out[15]:
work_year
2020    95811.00
2021    99851.79
2022    124522.01
Name: salary_in_usd, dtype: float64

In [17]: df1.index

Out[17]: Index([2020, 2021, 2022], dtype='int64', name='work_year')

In [19]: df1.values

Out[19]: array([[2020., 95811.],
       [2021., 99851.],
       [2022., 124522.]])

In [19]: data = {
    'work_year': df1.index,
    'average_salary': df1.values
}
df1 = pd.DataFrame(data)
df1

Out[19]:
work_year  average_salary
0          2020    95811.00
1          2021    99851.79
2          2022    124522.01

In [21]: df1['average_salary'] = df1['average_salary']/1000 .round(2)

In [22]: df1

Out[22]:
work_year  average_salary
0          2020      95.81
1          2021      99.85
2          2022     124.52

In [34]: ax = df1.plot(kind = 'bar', x = 'work_year', y = 'average_salary', legend = 'True')
ax.bar.labels[ax.containers[0], labels = df1['average_salary'].map('{:2f}'.format)}
plt.subplot(2,1,1(top = 1,2))
ax.legend(['Average Salary'], loc = 'upper left')

plt.xlabel('Work Year')
plt.ylabel('Average Salary')
plt.title('Average salaries by years', color = 'red')
plt.show()

Average salaries by years

Average Salary
120
100
80
60
40
20
0
2020      2021      2022
Work Year

In [ ]:

In [38]: df2 = df.remote_ratio.value_counts()
df2

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

In [40]: values = df2.to_list()

In [46]: values

Out[46]: [381, 127, 99]

In [48]: labels = ['Fully Remote', 'No Remote', 'Partially Remote']

In [48]: s = plt.bar(labels, values, width = 0.5, color = 'red')
plt.bar.labels[labels = values]
plt.title('Employees in remote positions')
plt.ylabel('Counts')
plt.show()

Employees in remote positions

Counts
400
350
300
250
200
150
100
50
0
Fully Remote  No Remote  Partially Remote

In [ ]:

In [52]: df3 = df.company_size.value_counts()

In [53]: df3

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

In [53]: df3.index.to_list()

Out[53]: ['M', 'L', 'S']

In [56]: df3.to_list()

Out[56]: [326, 198, 83]

In [58]: values2 = df3.values

In [59]: values2

Out[59]: array([326, 198, 83], dtype=int64)

In [ ]:

In [63]: labels_for_company = ['Medium', 'Large', 'Small']

In [67]: plt.figure(figsize = (13,6))
plt.pie(x = values2, labels = None, autopct = '%1.1f%', shadow = True)
plt.axis('equal')
plt.legend(labels = labels_for_company, loc = 'upper right')
plt.title('Company Size')
plt.show()

Company Size

53.7%
32.6%
13.7%
Medium
Large
Small

In [70]: df4 = df.job_title.value_counts()[1:5]

In [71]: df4

Out[71]:
job_title
Data Scientist      143
Data Engineer       132
Data Analyst         97
Machine Learning Engineer   41
Research Scientist    16
Name: count, dtype: int64

In [ ]:

In [76]: plt.figure(figsize = (12,8))
#Plot bar(x = df1.index, height = df1.values)
sns.barplot(x = df1.index, y = df1.values, palette = 'viridis')
label_bar = Top 5 Job Titles
plt.ylabel('Count')
plt.xticks(rotation = -20)
plt.show()

C:\Users\Boschke\envs\anaconda\lib\site-packages\seaborn\_oldcore.py:1498: FutureWarning: is_categorical_dtype is deprecated and will be removed in a future version. Use isinstance(dtype, CategoricalDtype) instead
  if pd.api.types.is_categorical_dtype(vector):
C:\Users\Boschke\envs\anaconda\lib\site-packages\seaborn\_oldcore.py:1498: FutureWarning: use_inf_as_na option is deprecated and will be removed in a future version. Use isinstance(dtype, CategoricalDtype) instead
  with pd.option_context(mode.use_inf_as_na, True):
C:\Users\Boschke\envs\anaconda\lib\site-packages\seaborn\_oldcore.py:1498: FutureWarning: is_categorical_dtype is deprecated and will be removed in a future version. Use isinstance(dtype, CategoricalDtype) instead
  if pd.api.types.is_categorical_dtype(vector):
C:\Users\Boschke\envs\anaconda\lib\site-packages\seaborn\_oldcore.py:1498: FutureWarning: use_inf_as_na option is deprecated and will be removed in a future version. Use isinstance(dtype, CategoricalDtype) instead
  with pd.option_context(mode.use_inf_as_na, True):
C:\Users\Boschke\envs\anaconda\lib\site-packages\seaborn\_oldcore.py:1498: FutureWarning: is_categorical_dtype is deprecated and will be removed in a future version. Use isinstance(dtype, CategoricalDtype) instead
  if pd.api.types.is_categorical_dtype(vector):
C:\Users\Boschke\envs\anaconda\lib\site-packages\seaborn\_oldcore.py:1498: FutureWarning: use_inf_as_na option is deprecated and will be removed in a future version. Use isinstance(dtype, CategoricalDtype) instead
  with pd.option_context(mode.use_inf_as_na, True):

Top 5 Job Titles

Count
140
120
100
80
60
40
20
0
Data Scientist  Data Engineer  Data Analyst  Machine Learning Engineer  Research Scientist
job_title

In [ ]:

In [79]: df

Out[79]:
work_year  experience_level  employment_type  job_title  salary_currency  salary_in_usd  employee_residence  remote_ratio  company_location  company_size
0          2020             MI              FT      Data Scientist      EUR      79833              DE              0              DE              L
1          2020             SE              FT      Machine Learning Scientist  USD      260000             JP              0              JP              S
2          2020             SE              FT      Big Data Engineer      GBP      109024             GB              50              GB              M
3          2020             MI              FT      Product Data Analyst      USD      20000              HN              0              HN              S
4          2020             SE              FT      Machine Learning Engineer      USD      150000             US              50              US              L
--          --          --                --              --              --              --              --              --              --
602         2022             SE              FT      Data Engineer      USD      154000             US              100              US              M
603         2022             SE              FT      Data Engineer      USD      126000             US              100              US              M
604         2022             SE              FT      Data Analyst      USD      129000             US              0              US              M
605         2022             SE              FT      Data Analyst      USD      150000             US              100              US              M
606         2022             MI              FT      AI Scientist      USD      200000             IN              100              US              L

607 rows x 10 columns

In [ ]:

In [83]: s = df3[df3['company_size'] == 'S']
M = df3[df3['company_size'] == 'M']
L = df3[df3['company_size'] == 'L']
labels = ['Medium', 'Large', 'Small']
sal_mean = [s['salary_in_usd'].mean(), M['salary_in_usd'].mean(), L['salary_in_usd'].mean()]

In [84]: sal_mean

Out[84]: [79632.67468979519, 114905.46625746871, 119242.99494949495]

In [86]: label_change = np.round(1/100 * (s.sal_mean, 2))
label_change = list(map(str, label_change))
label_change = [x + '%' for x in label_change]
label_change

Out[86]: ['77.63K', '114.91K', '119.24K']

In [89]: s = plt.bar(labels, sal_mean)
plt.bar.labels[labels, alpha = 0.8]
plt.title('Distribution of Salary by Company Size')
plt.xlabel('Size of company')
plt.ylabel('Salary in usd')
plt.show()

Distribution of Salary by Company Size

Salary in usd
120000
100000
80000
60000
40000
20000
0
Medium  Large  Small
Size of company

In [ ]:

In [93]: sns.histplot(s['salary_in_usd'], label = 'Small', kde = True)
sns.histplot(M['salary_in_usd'], label = 'Medium', kde = True)
sns.histplot(L['salary_in_usd'], label = 'Large', kde = True)
plt.xlabel('Distribution of Salary by Company Size')
plt.show()

C:\Users\Boschke\envs\anaconda\lib\site-packages\seaborn\_oldcore.py:1498: FutureWarning: is_categorical_dtype is deprecated and will be removed in a future version. Use isinstance(dtype, CategoricalDtype) instead
  if pd.api.types.is_categorical_dtype(vector):
C:\Users\Boschke\envs\anaconda\lib\site-packages\seaborn\_oldcore.py:1498: FutureWarning: use_inf_as_na option is deprecated and will be removed in a future version. Use isinstance(dtype, CategoricalDtype) instead
  with pd.option_context(mode.use_inf_as_na, True):
C:\Users\Boschke\envs\anaconda\lib\site-packages\seaborn\_oldcore.py:1498: FutureWarning: is_categorical_dtype is deprecated and will be removed in a future version. Use isinstance(dtype, CategoricalDtype) instead
  if pd.api.types.is_categorical_dtype(vector):
C:\Users\Boschke\envs\anaconda\lib\site-packages\seaborn\_oldcore.py:1498: FutureWarning: use_inf_as_na option is deprecated and will be removed in a future version. Use isinstance(dtype, CategoricalDtype) instead
  with pd.option_context(mode.use_inf_as_na, True):
C:\Users\Boschke\envs\anaconda\lib\site-packages\seaborn\_oldcore.py:1498: FutureWarning: is_categorical_dtype is deprecated and will be removed in a future version. Use isinstance(dtype, CategoricalDtype) instead
  if pd.api.types.is_categorical_dtype(vector):
C:\Users\Boschke\envs\anaconda\lib\site-packages\seaborn\_oldcore.py:1498: FutureWarning: use_inf_as_na option is deprecated and will be removed in a future version. Use isinstance(dtype, CategoricalDtype) instead
  with pd.option_context(mode.use_inf_as_na, True):

Distribution of Salary by Company Size

Count
50
40
30
20
10
0
0 100000 200000 300000 400000 500000 600000
Salary in usd

In [ ]:

In [96]: df

Out[96]:
work_year  experience_level  employment_type  job_title  salary_currency  salary_in_usd  employee_residence  remote_ratio  company_location  company_size
0          2020             MI              FT      Data Scientist      EUR      79833              DE              0              DE              L
1          2020             SE              FT      Machine Learning Scientist  USD      260000             JP              0              JP              S
2          2020             SE              FT      Big Data Engineer      GBP      109024             GB              50              GB              M
3          2020             MI              FT      Product Data Analyst      USD      20000              HN              0              HN              S
4          2020             SE              FT      Machine Learning Engineer      USD      150000             US              50              US              L
--          --          --                --              --              --              --              --              --              --
602         2022             SE              FT      Data Engineer      USD      154000             US              100              US              M
603         2022             SE              FT      Data Engineer      USD      126000             US              100              US              M
604         2022             SE              FT      Data Analyst      USD      129000             US              0              US              M
605         2022             SE              FT      Data Analyst      USD      150000             US              100              US              M
606         2022             MI              FT      AI Scientist      USD      200000             IN              100              US              L

607 rows x 10 columns

In [ ]:

In [97]: df5 = df.experience_level.value_counts()

In [100]: df5

Out[100]:
experience_level
SE      80
MI      213
EN      68
EX      26
Name: count, dtype: int64

In [103]: df5.index

Out[103]: Index(['SE', 'MI', 'EN', 'EX'], dtype='object', name='experience_level')

In [106]: df5.index.to_list()

Out[106]: ['SE', 'MI', 'EN', 'EX']

In [ ]:

In [102]: df['experience_level'].replace('SE', 'Senior', inplace = True)

In [107]: exp_map = {
    'SE': 'Senior',
    'MI': 'Midlevel',
    'EN': 'Entry level',
    'EX': 'Executive'
}

In [108]: df['experience_level'].replace(exp_map, inplace = True)

In [ ]:

In [109]: labels = df5.index.to_list()

Out[109]: ['SE', 'MI', 'EN', 'EX']

In [115]: values = df5.values
values

Out[115]: array([80, 213, 68, 26], dtype=int64)

In [121]: plt.figure(figsize = (6,12))
plt.pie(x = values, labels = labels, autopct='%1.2f%', shadow = 2)
plt.title('Experience Level')
plt.show()

Experience Level

46.13%
4.28%
14.50%
35.09%
SE
EX
EN
MI
```