

Importing Needed Libraries:

In [1]: 1 !pip install country_converter

```
Collecting country_converter
  Downloading country_converter-1.0.0-py3-none-any.whl (44 kB)
Requirement already satisfied: pandas>=1.0 in c:\users\titan rafi\anaconda3_python\lib\site-packages (from country_converter) (1.4.2)
Requirement already satisfied: numpy>=1.18.5 in c:\users\titan rafi\anaconda3_python\lib\site-packages (from pandas>=1.0->country_converter) (1.21.5)
Requirement already satisfied: pytz>=2020.1 in c:\users\titan rafi\anaconda3_python\lib\site-packages (from pandas>=1.0->country_converter) (2021.3)
Requirement already satisfied: python-dateutil>=2.8.1 in c:\users\titan rafi\anaconda3_python\lib\site-packages (from pandas>=1.0->country_converter) (2.8.2)
Requirement already satisfied: six>=1.5 in c:\users\titan rafi\anaconda3_python\lib\site-packages (from python-dateutil>=2.8.1->pandas>=1.0->country_converter) (1.16.0)
Installing collected packages: country-converter
Successfully installed country-converter-1.0.0
```

In [2]: 1 import pandas as pd
2 import numpy as np
3 import country_converter as coco
4 import matplotlib.pyplot as plt
5 import seaborn as sns
6 import plotly.express as px
7 import plotly.figure_factory as ff
8 import plotly.graph_objects as go
9 import warnings
10 warnings.filterwarnings('ignore')
11 pd.options.display.max_columns = None

In [3]: 1 df = pd.read_csv("C:/Users/Titan Rafi/Dropbox/PC/Desktop/ML Projects - Self/ds_salaries.csv")
2 df.head()

Out[3]:

	work_year	experience_level	employment_type	job_title	salary	salary_currency	salary_in_usd	employee_residence
0	2023	SE	FT	Principal Data Scientist	80000	EUR	85847	ES
1	2023	MI	CT	ML Engineer	30000	USD	30000	US
2	2023	MI	CT	ML Engineer	25500	USD	25500	US
3	2023	SE	FT	Data Scientist	175000	USD	175000	CA
4	2023	SE	FT	Data Scientist	120000	USD	120000	CA

In [5]: 1 df.shape

Out[5]: (3755, 11)

In [7]: 1 df.isnull().sum()

Out[7]: work_year 0
experience_level 0
employment_type 0
job_title 0
salary 0
salary_currency 0
salary_in_usd 0
employee_residence 0
remote_ratio 0
company_location 0
company_size 0
dtype: int64

In [8]: 1 df.duplicated().sum()

Out[8]: 1171

In [9]: 1 # drop the duplicated rows
2 df.drop_duplicates(inplace=True)

In [10]: 1 df.shape

Out[10]: (2584, 11)

In [11]: 1 df.info()

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 2584 entries, 0 to 3754
Data columns (total 11 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   work_year        2584 non-null   int64  
 1   experience_level 2584 non-null   object  
 2   employment_type  2584 non-null   object  
 3   job_title         2584 non-null   object  
 4   salary            2584 non-null   int64  
 5   salary_currency  2584 non-null   object  
 6   salary_in_usd    2584 non-null   int64  
 7   employee_residence 2584 non-null   object  
 8   remote_ratio      2584 non-null   int64  
 9   company_location  2584 non-null   object  
 10  company_size      2584 non-null   object  
dtypes: int64(4), object(7)
memory usage: 242.2+ KB
```

In [12]: 1 df.columns

Out[12]: Index(['work_year', 'experience_level', 'employment_type', 'job_title',
'salary', 'salary_currency', 'salary_in_usd', 'employee_residence',
'remote_ratio', 'company_location', 'company_size'],
dtype='object')

In [13]: 1 # drop 2 columns: 'salary' and 'salary_currency'
2 df.drop(['salary'], axis=1, inplace=True)
3
4 # rename salary_in_usd to salary
5 df.rename(columns={'salary_in_usd': 'salary'}, inplace=True)
6 # Inflation rates

```
In [14]: 1 # Inflation rates
2 us_inflation_rates = {2019: 0.0181, 2020: 0.0123, 2021: 0.0470, 2022: 0.065}
3 global_inflation_rates = {2019: 0.0219, 2020: 0.0192, 2021: 0.0350, 2022: 0.088}
4
5 # Function to adjust salary
6 def adjust_salary(row):
7     year = row['work_year']
8     original_salary = row['salary']
9     currency = row['salary_currency']
10
11     if year == 2023:
12         return original_salary
13
14     adjusted_salary = original_salary
15     for y in range(year, 2023):
16         if currency == 'USD':
17             inflation_rate = us_inflation_rates[y]
18         else:
19             inflation_rate = global_inflation_rates[y]
20
21     adjusted_salary *= (1 + inflation_rate)
22
23     return adjusted_salary
24
25 # Apply the function to the dataset
26 df['adjusted_salary'] = df.apply(adjust_salary, axis=1)
```

```
In [15]: 1 df.experience_level.value_counts()
```

```
Out[15]: SE    1554
          MI    664
          EN    270
          EX     96
Name: experience_level, dtype: int64
```

```
In [16]: 1 df['experience_level'] = df['experience_level'].replace({
2     'SE': 'Senior',
3     'EN': 'Junior',
4     'EX': 'Executive',
5     'MI': 'Mid-level',
6 })
```

```
In [17]: 1 converted_country = coco.convert(names=df['employee_residence'], to="name")
2 df['employee_residence'] = converted_country
```

```
In [18]: 1 converted_location = coco.convert(names=df['company_location'], to="name")
2 df['company_location'] = converted_location
```

```
In [19]: 1 df['remote_ratio'].value_counts()
```

```
Out[19]: 100    1211
          0      1186
          50     187
Name: remote_ratio, dtype: int64
```

```
In [21]: 1 df['job_type'] = df['remote_ratio'].replace({0: 'onsite', 50: 'hybrid', 100:'remote'})
```

In [22]:

```

1 #df.job_title.unique()
2 def assign_broader_category(job_title):
3     data_scientist = ['Principal Data Scientist', 'Data Scientist', 'Applied Scientist', 'Applied Data Scientist', 'Data Science Manager', 'Director of Data Science', 'Lead Data Scientist', 'Data Science Lead', 'Data Science Consultant', 'Head of Data Science', 'Data Science Engineer', 'Data Science Tech Lead', 'Data Scientist Lead', 'Product Data Scientist', 'Staff Data Scientist']
4 ]
5     data_engineer = [ 'Data Modeler', 'Data Strategist', 'Data Engineer', 'Data Architect', 'Head of Data', 'Data Manager', 'Data Operations Engineer', 'Azure Data Engineer', 'Big Data Engineer', 'Cloud Database Engineer', 'Marketing Data Engineer', 'Data Lead', 'Data Infrastructure Engineer', 'Software Data Engineer', 'Data Specialist', 'BI Data Engineer', 'BI Developer', 'Big Data Architect', 'Cloud Data Engineer', 'Data Operations Analyst', 'Power BI Developer', 'Principal Data Architect', 'Cloud Data Architect', 'Lead Data Engineer', 'Principal Data Engineer']
6     machine_learning= [ 'ML Engineer', 'Machine Learning Engineer', 'Applied Machine Learning Engineer', 'Machine Learning Researcher', 'Machine Learning Scientist', 'MLOps Engineer', 'AI Scientist', 'AI Developer', 'Applied Machine Learning Scientist', 'AI Programmer', 'Deep Learning Researcher', 'Machine Learning Infrastructure Engineer', 'Deep Learning Engineer', 'Machine Learning Software Engineer', 'Machine Learning Research Engineer', 'NLP Engineer', 'Machine Learning Developer', 'Principal Machine Learning Engineer', 'Machine Learning Manager', 'Lead Machine Learning Engineer', 'Head of Machine Learning', 'Machine Learning Analyst' ]
7 ]
8     data_analyst = [ 'Data Analyst', 'Analytics Engineer', 'Data Quality Analyst', 'Compliance Data Analyst', 'Data Analytics Manager', 'Business Data Analyst', 'Staff Data Analyst', 'Lead Data Analyst', 'Financial Data Analyst', 'BI Analyst', 'Product Data Analyst', 'Data Analytics Lead', 'Data Analytics Specialist', 'BI Data Analyst', 'Insight Analyst', 'Data Analytics Engineer', 'Data Analytics Consultant', 'Marketing Data Analyst', 'Principal Data Analyst', 'Finance Data Analyst' ]
9     if job_title in data_engineer:
10         return "Data Engineer"
11     elif job_title in data_scientist:
12         return "Data Scientist"
13     elif job_title in machine_learning:
14         return "ML Engineer"
15     elif job_title in data_analyst:
16         return "Data Analyst"
17     else:
18         return "Other"
19
20
21
22 df['job_category'] = df['job_title'].apply(assign_broader_category)
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42

```

In [23]:

```

1 df['employment_type'] = df['employment_type'].replace({
2     'FL': 'Freelancer',
3     'CT': 'Contractor',
4     'FT' : 'Full-time',
5     'PT' : 'Part-time'
6 })
7 df['company_size'] = df['company_size'].replace({
8     'S': 'Small',
9     'M': 'Medium',
10    'L' : 'Large',
11 })

```

In [24]: 1 df.head()

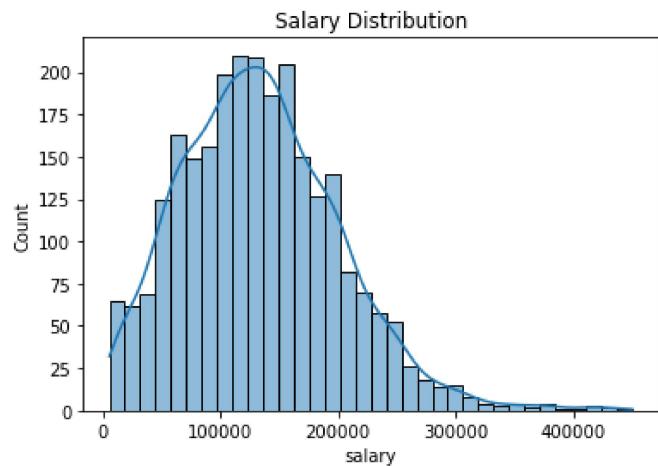
Out[24]:

	work_year	experience_level	employment_type	job_title	salary_currency	salary	employee_residence	remote_ratio
0	2023	Senior	Full-time	Principal Data Scientist	EUR	85847	Spain	100
1	2023	Mid-level	Contractor	ML Engineer	USD	30000	United States	100
2	2023	Mid-level	Contractor	ML Engineer	USD	25500	United States	100
3	2023	Senior	Full-time	Data Scientist	USD	175000	Canada	100
4	2023	Senior	Full-time	Data Scientist	USD	120000	Canada	100

Salary Distribution:

In [25]: 1 ax = sns.histplot(df['salary'], kde = True)
2 ax.set_title('Salary Distribution')

Out[25]: Text(0.5, 1.0, 'Salary Distribution')



Experience Level:

In [26]:

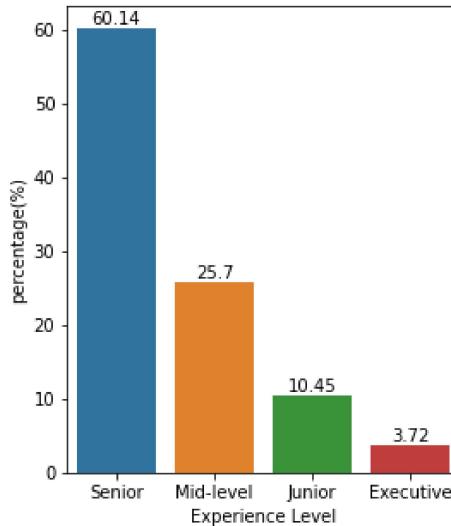
```

1 ex_level = df['experience_level'].value_counts().to_frame()
2 ex_level.columns = ['count']
3 ex_level['percentage(%)'] = round(100 * df['experience_level'].value_counts(normalize= True),2)
4
5 fig = plt.figure(figsize=(9,5))
6 ax1 = plt.subplot(121)
7 ax1.axis('off')
8 bbox =[0,0,0.8,0.8]
9 ax1.table(cellText = ex_level.values, rowLabels=ex_level.index, colLabels=ex_level.columns, borderpad=0, cellLoc='center', loc='center',bbox=bbox)
10 ax2 = fig.add_subplot(122)
11 ax2 = sns.barplot(data = ex_level, x=ex_level.index,y=ex_level['percentage(%)'])
12 ax2.bar_label(ax2.containers[0])
13 # plt.xticks(rotation= 30)
14 plt.xlabel('Experience Level')
15 plt.suptitle('Experience Level count')
16 plt.show()

```

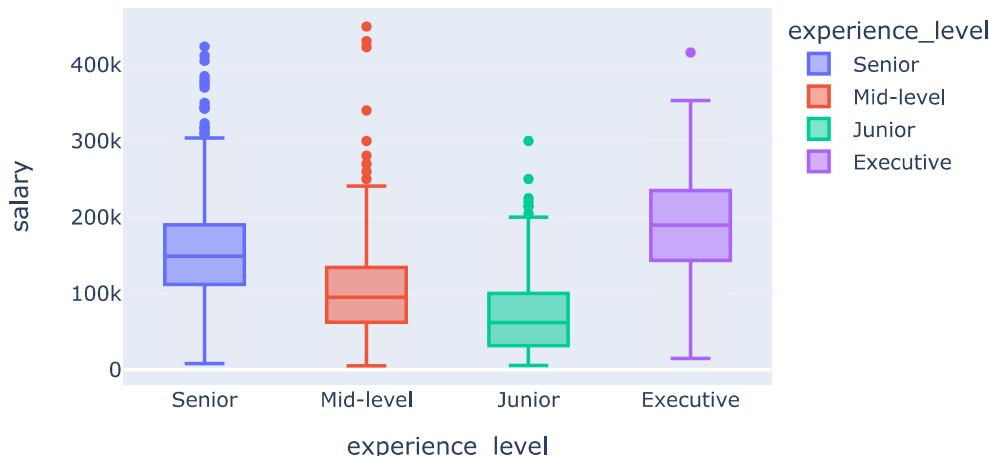
Experience Level count

	count	percentage(%)
Senior	1554.0	60.14
Mid-level	664.0	25.7
Junior	270.0	10.45
Executive	96.0	3.72



```
In [27]: 1 px.box(df, x= 'experience_level', y = 'salary',
2           color= 'experience_level',title= 'Salary Distribution by Experience Level',
3           width=600, height=400)
```

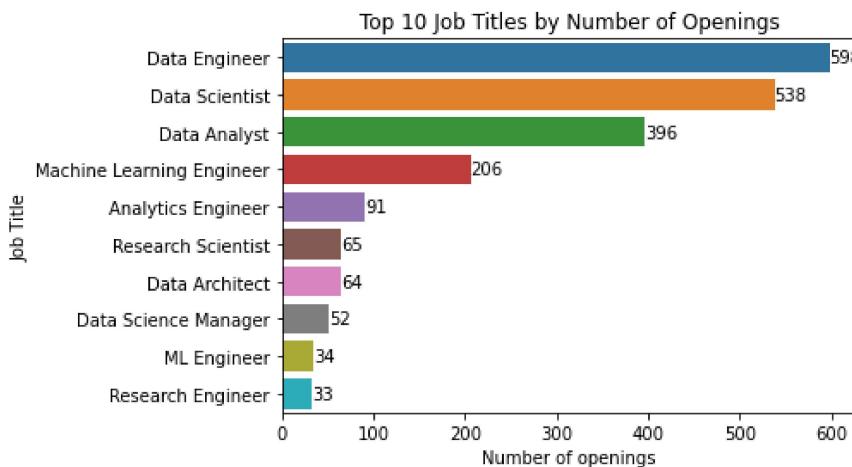
Salary Distribution by Experience Level



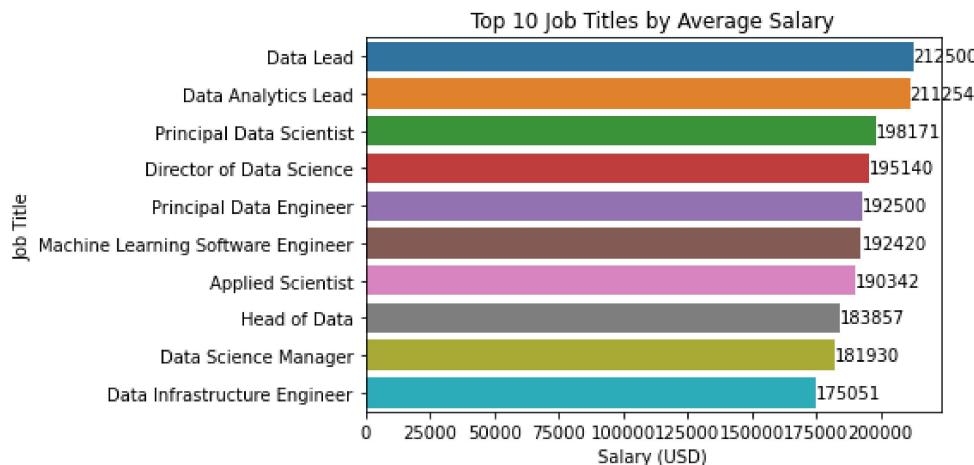
Which Job Roles are Most Popular and Most paid in Data Science:

```
In [28]: 1 print('There are {} job titles in this dataset.'.format(len(df['job_title'].value_counts())))
There are 93 job titles in this dataset.
```

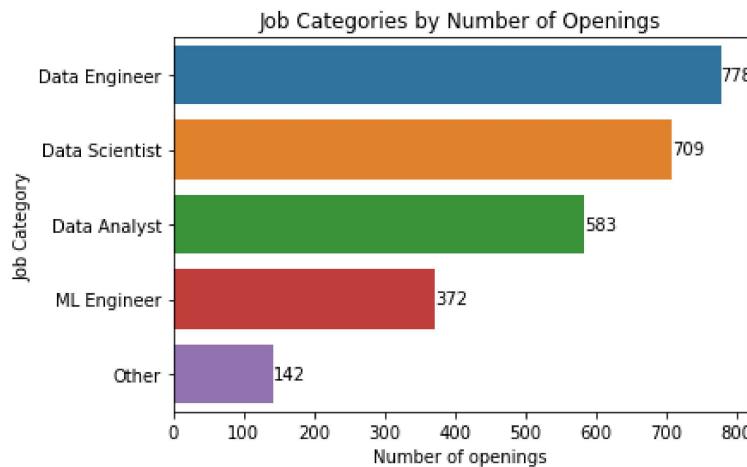
```
In [29]: 1 top10_count = df['job_title'].value_counts()[:10].to_frame()
2 top10_count.columns = ['count']
3
4 plt.figure(figsize=(6,4))
5
6 ax= sns.barplot(data = top10_count, y = top10_count.index, x= top10_count['count'])
7 plt.ylabel('Job Title')
8 plt.xlabel('Number of openings')
9 plt.title('Top 10 Job Titles by Number of Openings')
10 ax.bar_label(ax.containers[0], fmt = '%d')
11 plt.show()
```



```
In [30]: 1 top_avg_salary = df.groupby('job_title')['salary'].mean().round(2).sort_values(ascending=False)
2 role_count = df['job_title'].value_counts()
3 top10_avg_salary = top_avg_salary[role_count[role_count > 1].index].sort_values(ascending = False)
4
5 plt.figure(figsize=(6,4))
6
7 ax = sns.barplot(data = top10_avg_salary, y=top10_avg_salary.index,x=top10_avg_salary['salary'])
8 plt.ylabel('Job Title')
9 plt.xlabel('Salary (USD)')
10 plt.title('Top 10 Job Titles by Average Salary')
11 ax.bar_label(ax.containers[0], fmt = '%d')
12 plt.show()
```



```
In [31]: 1 temp = df['job_category'].value_counts().to_frame()
2 temp.columns = ['count']
3
4 plt.figure(figsize=(6,4))
5
6 ax= sns.barplot(data = temp, y = temp.index, x= temp['count'])
7 plt.ylabel('Job Category')
8 plt.xlabel('Number of openings')
9 plt.title('Job Categories by Number of Openings')
10 ax.bar_label(ax.containers[0], fmt = '%d')
11 plt.show()
```

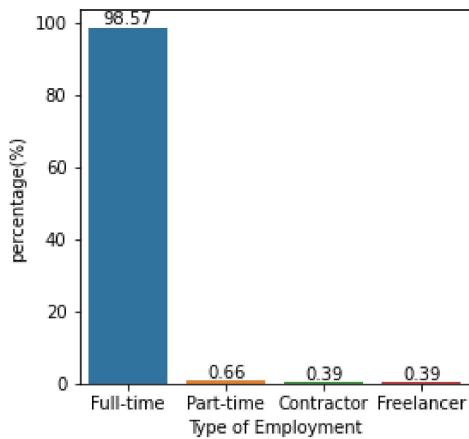


Which type of Employment is Lucrative?

```
In [32]: 1 employment = df['employment_type'].value_counts().to_frame()
2 employment.columns = ['count']
3 employment['percentage(%)'] = round(100 * df['employment_type'].value_counts(normalize= True))
4
5 fig = plt.figure(figsize=(9,4))
6 ax1 = plt.subplot(121)
7 ax1.axis('off')
8 bbox =[0,0,0.8,0.8]
9 ax1.table(cellText = employment.values, rowLabels=employment.index, colLabels=employment.col
10 ax2 = fig.add_subplot(122)
11 ax2 = sns.barplot(data = employment, x=employment.index,y=employment['percentage(%)'])
12 ax2.bar_label(ax2.containers[0])
13 # plt.xticks(rotation= 30)
14 plt.xlabel('Type of Employment')
15 plt.suptitle('Type of Employment count')
16 plt.show()
```

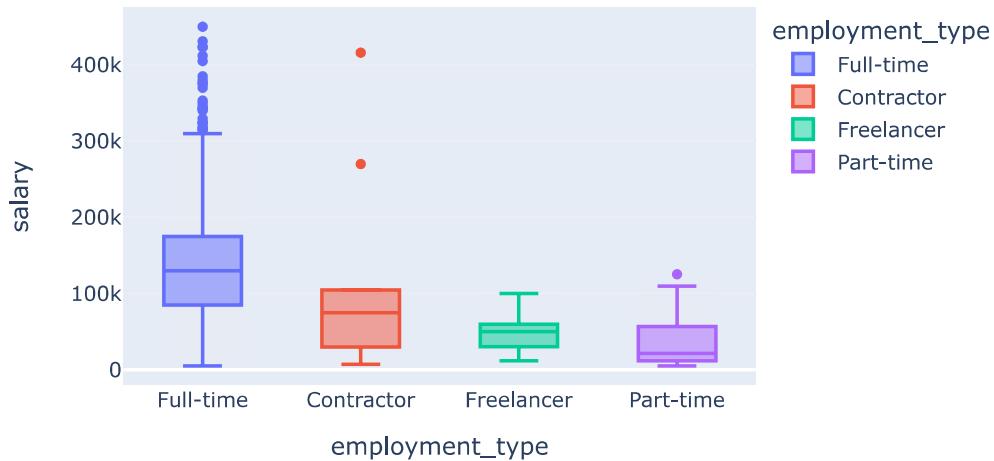
Type of Employment count

	count	percentage(%)
Full-time	2547.0	98.57
Part-time	17.0	0.66
Contractor	10.0	0.39
Freelancer	10.0	0.39



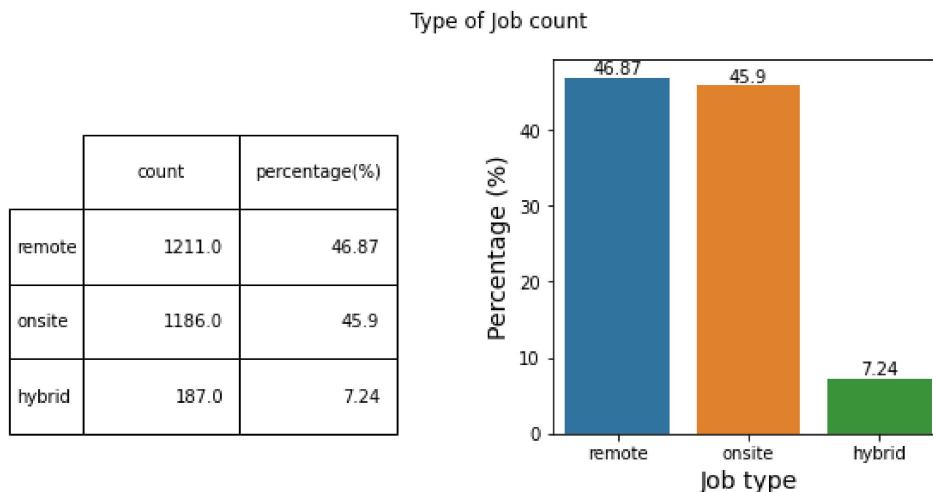
```
In [33]: 1 px.box(df,x='employment_type',y='salary',color='employment_type',
2           title= 'Data Sceince Salaries by the type of employee',
3           width=600, height=400)
```

Data Sceince Salaries by the type of employee



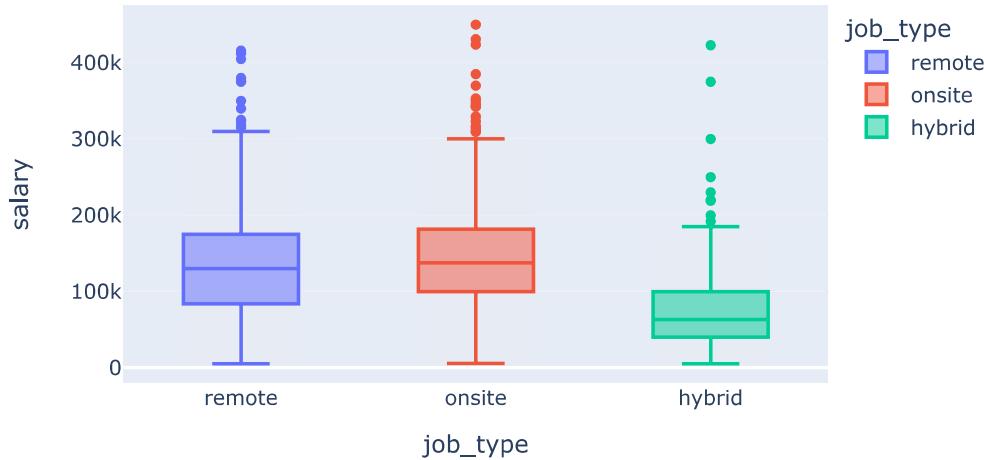
Which type of job has more Average Salary?

```
In [34]: 1 jobType = df['job_type'].value_counts().to_frame()
2 jobType.columns = ['count']
3 jobType['percentage(%)'] = round(100 * df['job_type'].value_counts(normalize= True),2)
4
5 fig = plt.figure(figsize=(9,4))
6 ax1 = plt.subplot(121)
7 ax1.axis('off')
8 bbox =[0,0,0.8,0.8]
9 ax1.table(cellText = jobType.values, rowLabels=jobType.index, colLabels=jobType.columns,bbox=bbox)
10 ax2 = fig.add_subplot(122)
11 ax2 = sns.barplot(data = jobType, x=jobType.index,y=jobType['percentage(%)'])
12 ax2.bar_label(ax2.containers[0])
13 # plt.xticks(rotation= 30)
14 plt.xlabel('Job type', fontdict={'fontsize':14})
15 plt.ylabel('Percentage (%)', fontdict={'fontsize':14})
16
17 plt.suptitle('Type of Job count')
18 plt.show()
```



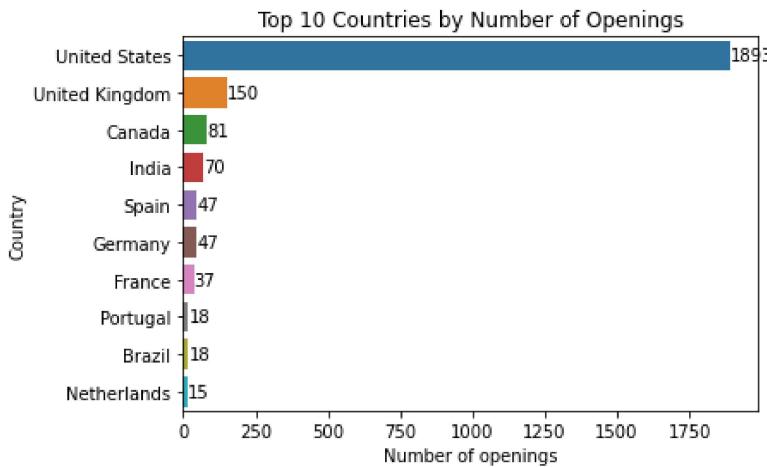
```
In [35]: 1 px.box(df,x='job_type',y='salary',color='job_type',
2           title='Data Science Salaries by the Type of Job',width=600, height=400)
```

Data Science Salaries by the Type of Job

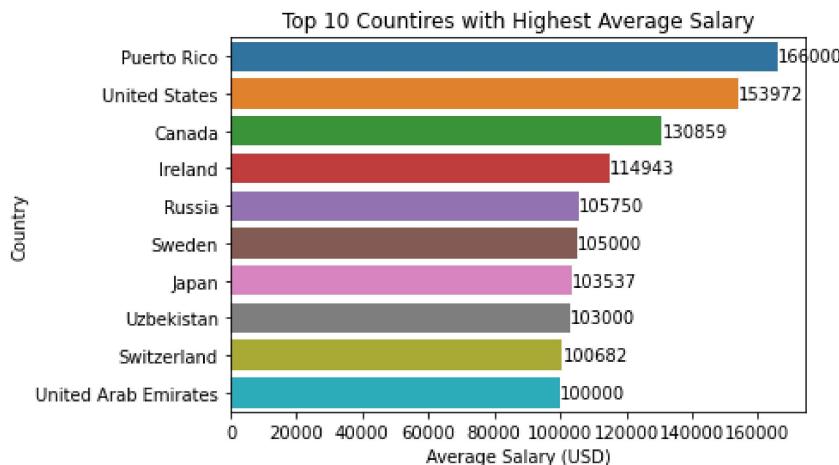


Company Location:

```
In [36]: 1 top10_count = df['company_location'].value_counts()[:10].to_frame()
2 top10_count.columns = ['count']
3
4 plt.figure(figsize=(6,4))
5
6 ax= sns.barplot(data = top10_count, y = top10_count.index, x= top10_count['count'])
7 plt.ylabel('Country')
8 plt.xlabel('Number of openings')
9 plt.title('Top 10 Countries by Number of Openings')
10 ax.bar_label(ax.containers[0], fmt = '%d')
11 plt.show()
```

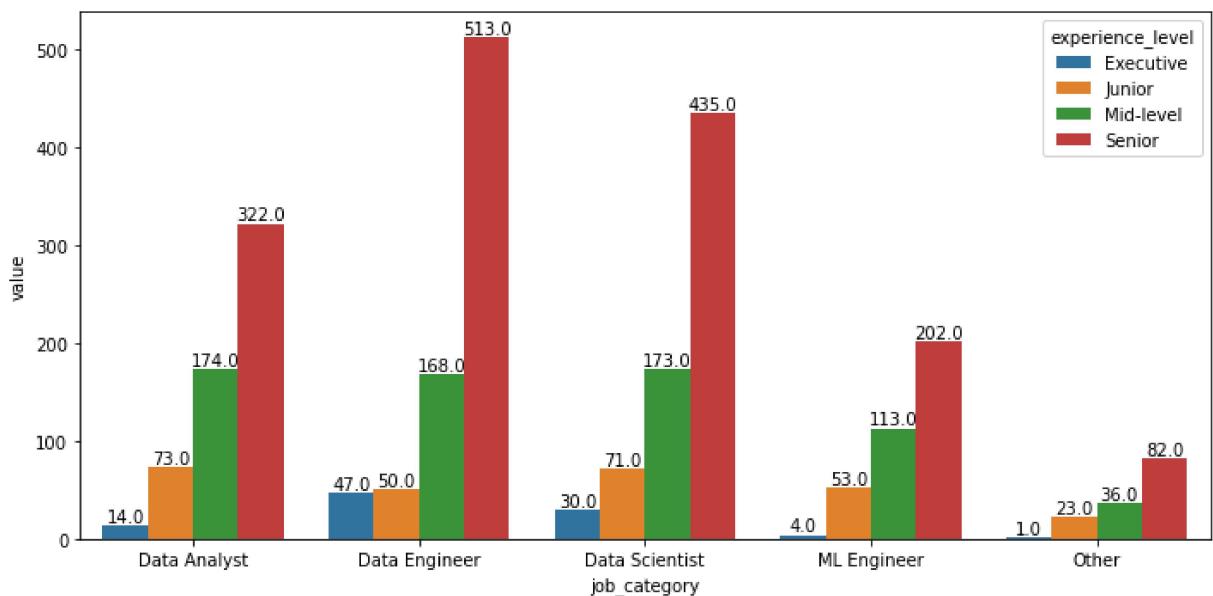


```
In [37]: 1 top_avg_salary = df.groupby('company_location')['salary'].mean().round(2).sort_values(ascending=True)
2 role_count = df['company_location'].value_counts()
3 top10_avg_salary = top_avg_salary[role_count[role_count > 1].index].sort_values(ascending = True)
4
5 plt.figure(figsize=(6,4))
6
7 ax = sns.barplot(data = top10_avg_salary, y=top10_avg_salary.index,x=top10_avg_salary['salary'])
8 plt.ylabel('Country')
9 plt.xlabel('Average Salary (USD)')
10 plt.title('Top 10 Countries with Highest Average Salary')
11 ax.bar_label(ax.containers[0], fmt = '%d')
12 plt.show()
```



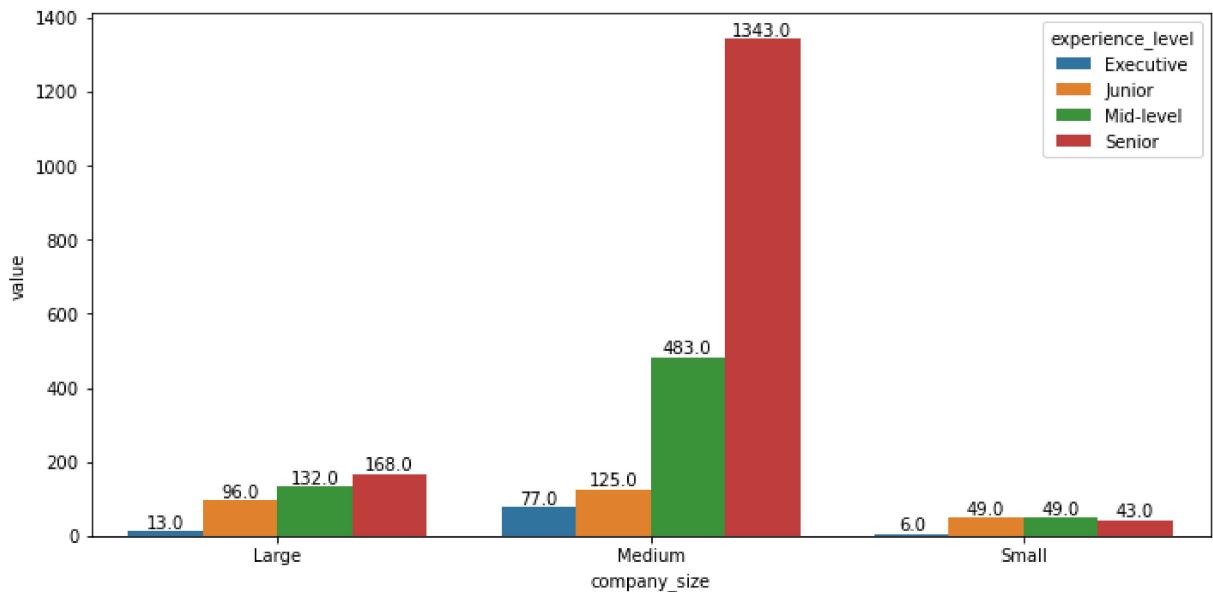
Experience level and Job Category:

```
In [38]: 1 exlevel_type = df.groupby(['experience_level','job_category']).size().reset_index()
2 exlevel_type.columns = ['experience_level','job_category', 'value']
3 plt.figure(figsize=(10,5))
4 ax = sns.barplot(data = exlevel_type, x=exlevel_type['job_category'], y=exlevel_type['value'],
5                   hue=exlevel_type['experience_level'])
6 for bars in ax.containers:
7     ax.bar_label(bars, fmt='%.1f')
8 plt.tight_layout()
9 plt.show()
```



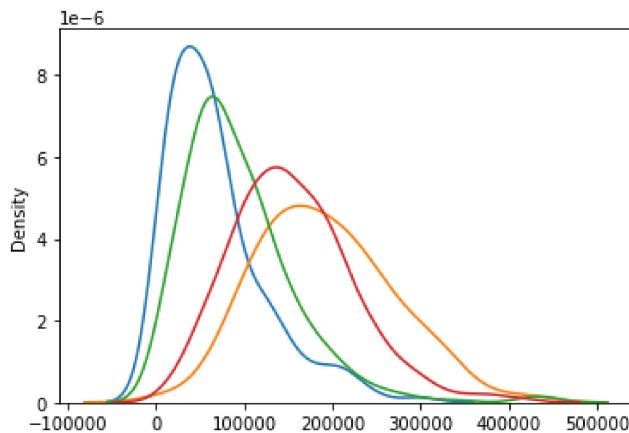
Experience level and Company Size:

```
In [39]: 1 exlevel_type = df.groupby(['experience_level','company_size']).size().reset_index()
2 exlevel_type.columns = ['experience_level','company_size', 'value']
3 plt.figure(figsize=(10,5))
4 ax = sns.barplot(data = exlevel_type, x=exlevel_type['company_size'], y=exlevel_type['value']
5                   hue=exlevel_type['experience_level'])
6 for bars in ax.containers:
7     ax.bar_label(bars, fmt='%.1f')
8 plt.tight_layout()
9 plt.show()
```



```
In [40]: 1 exlevel_salary = df.groupby(['experience_level','salary']).size()
2 entry_salary = exlevel_salary['Junior'].reset_index()
3 executive_salary = exlevel_salary['Executive'].reset_index()
4 mid_salary = exlevel_salary['Mid-level'].reset_index()
5 senior_salary = exlevel_salary['Senior'].reset_index()
6 sns.distplot(entry_salary[['salary']], hist=False)
7 sns.distplot(executive_salary[['salary']], hist=False)
8 sns.distplot(mid_salary[['salary']], hist=False)
9 sns.distplot(senior_salary[['salary']], hist=False)
```

Out[40]: <AxesSubplot:ylabel='Density'>



In []:

1