

In [167]:

```
1 #INTRODUCTION:
2 '''
3 In this project, we will delve into an intriguing dataset that encompasses
4 The dataset offers a multitude of information - from the company and job t
5 to the location, job type, experience level, and salary, among other eleme
6
7
8 This project forms a part of my data science portfolio, crafted with the i
9 My aim with this project is not just to unearth insightful trends about da
10 but also to demonstrate fundamental data science techniques - from preproc
11
12 '''
```

Out[167]: "\nIn this project, we will delve into an intriguing dataset that encompasses key aspects of data science job postings. \nThe dataset offers a multitude of information - from the company and job title, to the location, job type, experience level, and salary, among other elements.\n\nThe motivation behind this project is to unearth valuable insights into the current data science job market \nthat could potentially guide job seekers, recruiters, and companies alike.\nThese insights could reveal pivotal trends such as the most in-demand skills, \nsalary expectations across various geographies, and the types of companies hiring the most data scientists.\n\nTo accomplish this, we will initially carry out data preprocessing, \nwhich involves cleaning and preparing our data for analysis.\nFollowing this, we will conduct exploratory data analysis to scrutinize the underlying patterns and trends present in our data. \n\nStay tuned to uncover some intriguing insights from the ever-evolving world of data science occupations.\n\nThis project forms a part of my data science portfolio, crafted with the intention to showcase my data manipulation, analysis, and visualization skills. \nMy aim with this endeavor is not just to unearth insightful trends from the realm of data science jobs,\nbut also to exhibit a clear understanding and application of crucial data science techniques - from preprocessing raw data to extracting meaningful insights.\n\nAs we traverse through the project, you will witness the application of various data science methodologies, \nprogramming concepts, and libraries essential to the Python data science stack.\n\nLet's dive into the captivating world of data science jobs and uncover what lies beneath!\n\n"

In [168]:

```
1 import numpy as np
2 import matplotlib_inline
3 import matplotlib.pyplot as plt
4 import pandas as pd
5 import os
6 import math
7 import pathlib
8 from pathlib import Path
9 !pip install geonamescache
```

Requirement already satisfied: geonamescache in c:\users\matt\anaconda3\lib\site-packages (1.6.0)

In [169]:

```

1 #This code is used to easily read in the data path of the original .csv fi
2 data_path = Path(r"C:\Users\Matt\Desktop\Data_Science_Job_Project\data_sci
3 df = pd.read_csv(data_path, encoding='latin1')

```

In [170]:

```

1 #Displaying the raw data set
2 df

```

Out[170]:

	Company	Job Title	Location	Job Type	Experience level	Salary	Requirements
0	SGS	Clinical Data Analyst	Richardson, TX, United States	Full Time	Entry-level	48K+ *	Comput quality, Gen
1	Ocorian	AML/CFT & Data Analyst	Ebène, Mauritius	Full Time	Entry-level	48K+ *	Agile, Data management, Fi
2	Cricut	Machine Learning Engineer	South Jordan, UT, United States	Full Time	NaN	90K+ *	Agile, Architecture Sc
3	Bosch Group	Application Developer & Data Analyst	Nonantola, Italy	Full Time	Entry-level	48K+ *	Engineering, Industrial, Oracle, F
4	Publicis Groupe	Data Engineer Full time (Public Sector) USA	Arlington, VA, United States	Full Time	Mid-level	108K+	AWS, Science, Consi
...
3193	Western Digital	Data Scientist - New College Graduate	Biñan, Philippines	Full Time	Entry-level	39K+ *	APIs, Clustering, Computer Science
3194	Experian	Cloud Data Analyst	Heredia, Costa Rica	Full Time	Senior-level	92K+ *	AWS, Big Science, GCP
3195	Locus Robotics	Robotics Engineer, Sensors	Wilmington, MA, United States	Full Time	Senior-level	62K+ *	commerce, Engineering, Linux, Pyth
3196	ATB Financial	Data Scientist	Edmonton, Alberta, Canada	Full Time	Entry-level	39K+ *	Computer Science, Dat
3197	Shippeo	Senior Data Engineer	Paris, France	Full Time	Senior-level	115K+ *	Airflow, Architecture, BigQuery, (

3198 rows × 8 columns

In [171]:

```
1 #DATA CLEANING
2 #First, we will clean up the data set in various way for processing and an
```



In [172]:

```
1  #DATA CLEANING: 'Salary'
2  #We want to convert all salaries into a numerical value for comparison
3  #For simplicity, we will assume unless otherwise specified, that the data
4
5  #First, we remove all NaN salaries from the table:
6  df = df.dropna(subset=['Salary'])
7  df.sort_values("Salary", inplace = True)
8
9  #Printing the unique values to understand the data better
10 #print(df['Salary'].unique())
11 #Note that salaries are defined with + (ex: 50k+).
12 #For simplicity, we will round to the nearest thousand. So 50k -> 50000
13
14 #Note are GBP and EURO values. We will convert these to USD
15 #Note the exchange rates as of 7/3/23
16 #1 EUR = 1.09 USD
17 #1 GBP = 1.27 USD
18
19 #For processing, we set each value of the 'Salary' column to a string
20 df['Salary'] = df['Salary'].astype(str)
21
22 #We now process our data into a new column called 'Salary_USD'
23 #We will utilize regex expressions for processing
24 import re
25
26 #Defining function to process the 'Salary' column:
27 def parse_salary(salary):
28     salary = salary.strip()
29     salary_split = salary.split(' ')
30     #we will create a modifier value and multiply it to our result to get
31     modifier = 1
32     for item in salary_split:
33         if re.match(r'^\d', item):
34             numeric_value = re.findall(r'\d+', item)
35             if numeric_value:
36                 numeric_value = int(numeric_value[0])
37                 usd_salary = numeric_value*modifier*1000
38                 return usd_salary
39         if 'GBP' in item:
40             modifier = 1.27
41         if 'EUR' in item:
42             modifier = 1.09
43
44 df['Salary_USD'] = df['Salary'].apply(parse_salary)
```

```
C:\Users\Matt\AppData\Local\Temp\ipykernel_19976\3210654877.py:7: SettingWithCopyWarning:
```

A value is trying to be set on a copy of a slice from a DataFrame

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy (https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy)

```
df.sort_values("Salary", inplace = True)
```

```
C:\Users\Matt\AppData\Local\Temp\ipykernel_19976\3210654877.py:20: SettingWithCopyWarning:
```

A value is trying to be set on a copy of a slice from a DataFrame.

Try using `.loc[row_indexer,col_indexer] = value` instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy (https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy)

```
df['Salary'] = df['Salary'].astype(str)
```

```
C:\Users\Matt\AppData\Local\Temp\ipykernel_19976\3210654877.py:44: SettingWithCopyWarning:
```

A value is trying to be set on a copy of a slice from a DataFrame.

Try using `.loc[row_indexer,col_indexer] = value` instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy (https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy)

```
df['Salary_USD'] = df['Salary'].apply(parse_salary)
```

In [173]:

```
1 df
```

Out[173]:

	Company	Job Title	Location	Job Type	Experience level	Salary	Requirmen
2627	OPPO Research Center	Senior Research Scientist/Engineer -Video Unde...	Palo Alto, California, United States	Full Time	Mid-level	100K+	Caffe,Computer 5
113	Freeform	Software Engineer (Data Pipeline)	Los Angeles, CA	Full Time	Senior-level	100K+	Big Data Science
203	Vericast	Product Manager-Data Visualization & Analytics	Austin, TX, United States	Full Time	NaN	100K+	Business vis
478	Memora Health	Analyst, Reporting and Business Intelligence	Remote-U.S. Based	Full Time	Senior-level	100K+	Intelligence,Engineerin
2229	Zynga	Senior Gameplay AI Engineer - Star Wars	Canada	Full Time	Senior-level	100K+	Architecture,Engin
...	
2716	Monzo	Data Analyst, Borrowing	London	Full Time	Senior-level	GBP 55K	AWS,Banking,BigQuery
2492	Our Future Health	Data Manager - Airlock (18 month FTC)	London, England, United Kingdom	Full Time	Mid-level	GBP 55K+	Consul 27001,Pl
1384	iwoca	Senior Data Scientist - UK Ops	London, England, United Kingdom	Full Time	Mid-level	GBP 60K+	Bayesian,Eng Le
2536	Frazer-Nash Consultancy	Senior Data Architect	Burton upon Trent, England, United Kingdom	Full Time	Mid-level	GBP 65K+	Architecture,Data
1246	Audigent	Senior Data Engineer	New York / London / Remote	Full Time	Senior-level	GBP 80K+	Airflow,APIs,Architecture,A

3009 rows × 9 columns



In [174]:

```

1  #DATA CLEANING: 'Location'
2
3
4  #I would like to investigate how the data is related to location
5  #To do this, I need to clean up the data set's 'Location' column.
6  #Ideally, I would like a single country of origin for each entry. For exam
7  #The new locational data will be stored in a seperate column called 'Count
8  # We will map any entry containing the word 'remote' to Remote
9
10 #We will start by removing all NAN values present in the "Location" column
11 df= df.dropna(axis = 'index', how = 'any', subset = ['Location'])
12
13 #We will then standardize all the locations to a format for processing:
14 #This is achieved by making all characters lowercase and removing any whit
15 df['Location'] = df['Location'].str.lower().str.strip()
16
17 #We now must go through the data set and standardize all the entries
18 #We will map each entry to a country name
19 #First, we will store each unique item in the 'Location' data set to an ar
20 locations = df['Location'].unique()
21
22 #to easily map our locations, we will utilize the GeonamesCache library
23 #this is a free repository of locational data that we can use to help clea
24
25 #we start by importing the library:
26 from geonamescache import GeonamesCache
27
28 #We then create a function to handle each location
29 gc = GeonamesCache()
30 def parse_location(location):
31
32     #If the location has remote, we set the location to "Remote"
33     if 'remote' in location:
34         return 'Remote'
35
36     #This section handles non-remote locations
37     else:
38         #splitting the location into an array of strings by using the comm
39         location_list = location.split(',')
40         location_list.sort()
41         for loc in location_list:
42             loc = loc.strip()
43             loc = loc.title()
44             #detects if country is in list -> sets value to country
45             if gc.get_countries_by_names().get(loc):
46                 return gc.get_countries_by_names().get(loc)['name']
47
48             #detects if US state is in list -> sets value to 'United State
49             if gc.get_us_states_by_names().get(loc):
50                 return 'United States'
51
52             #detects if US state is in list by state abbreviation -> sets
53             states = gc.get_us_states()
54             if loc.upper() in states:
55                 return 'United States'
56
57             #detects if city is in list-> sets value to country of origin

```



```

58         if gc.get_cities_by_name(loc):
59             city = gc.get_cities_by_name(loc)
60             unknown_key = list(city[0].keys())[0]
61             country_code = city[0][unknown_key]['countrycode'] # Access
62             c = gc.get_countries().get(country_code)
63             return c['name']
64
65 #Applying the function to our 'Location' column and printing the results
66 df['Country'] = df['Location'].apply(parse_location)
67
68 #We now have a unique location for each individual row:

```

In [175]:

```

1 #DATA ANALYSIS
2 #We now begin to analyze the data set in various different ways

```

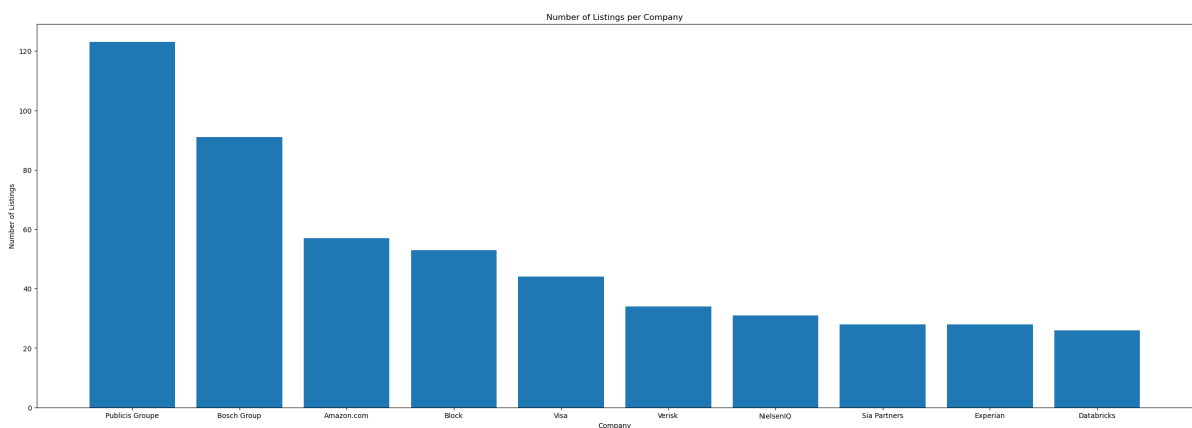
In [176]:

```

1 #creating a bar chart of companies vs Listing for top 10 most frequent com
2 company_counts_series = df['Company'].value_counts()[:10]
3 companies = company_counts_series.index
4 counts = company_counts_series.values
5
6 plt.figure(figsize=(30, 10))
7 plt.bar(companies, counts)
8 plt.xlabel('Company')
9 plt.ylabel('Number of Listings')
10 plt.title('Number of Listings per Company')

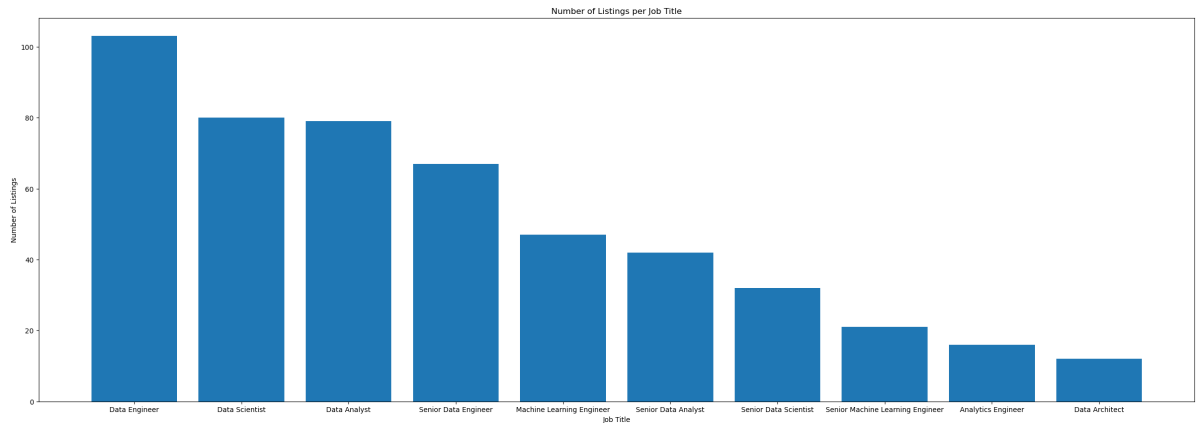
```

Out[176]: Text(0.5, 1.0, 'Number of Listings per Company')



```
In [177]: 1 #creating a bar chart of top 10 most frequent job titles
2 job_counts_series = df['Job Title'].value_counts()[:10]
3 jobs= job_counts_series.index
4 job_counts = job_counts_series.values
5
6 plt.figure(figsize=(30, 10))
7 plt.bar(jobs,job_counts)
8 plt.xlabel('Job Title')
9 plt.ylabel('Number of Listings')
10 plt.title('Number of Listings per Job Title')
```

Out[177]: Text(0.5, 1.0, 'Number of Listings per Job Title')



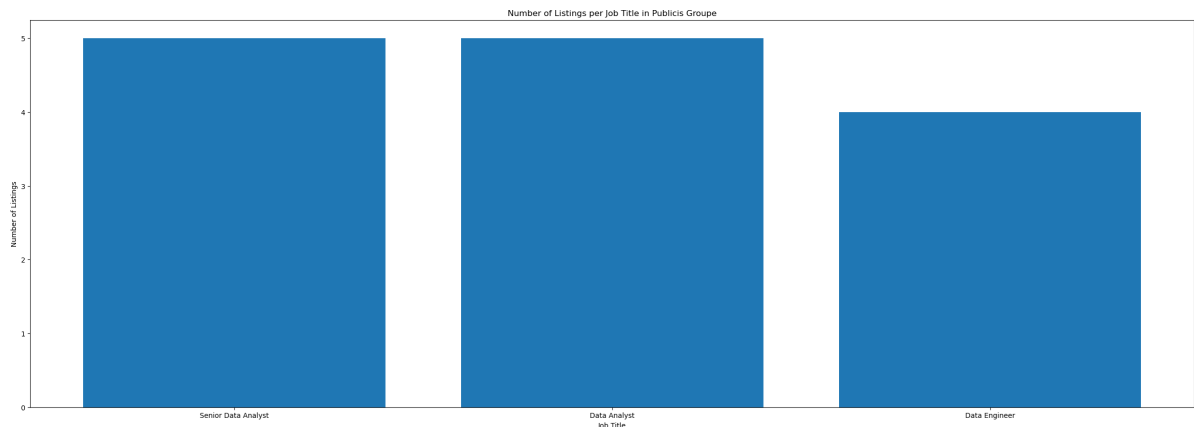
```
In [178]: 1 #We can see that 'Data Engineer', 'Data Scientist', and 'Data Analyst' pos
```

```

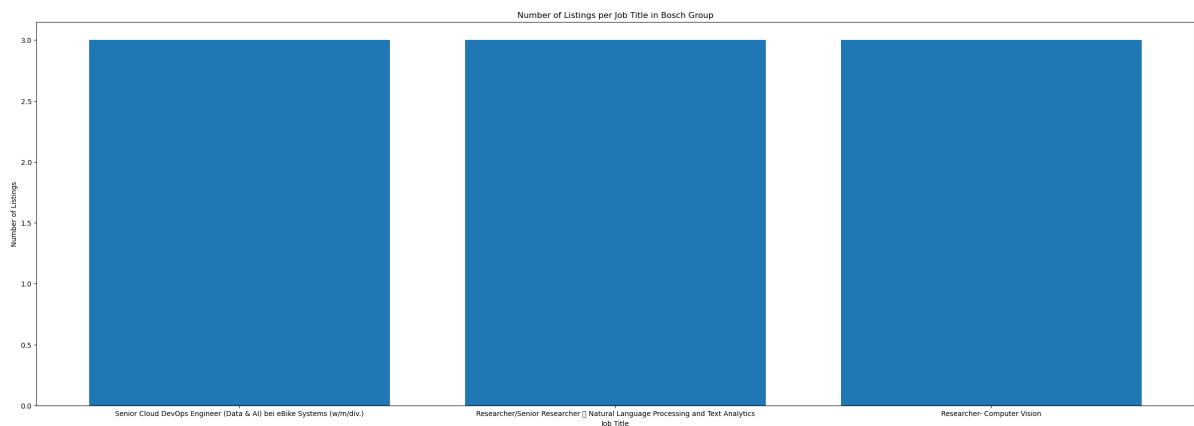
In [179]: 1 #we now want to categorize the most 3 frequent job titles for each of the
2
3 top_most_frequent_companies = companies[:3]
4 grouped_data = df.groupby(['Company'])
5 for x in top_most_frequent_companies:
6     company_grouped_data = grouped_data.get_group(x)['Job Title'].value_co
7     company_job = company_grouped_data.index
8     company_job_count = company_grouped_data.values
9     #plotting each bar graph
10    plt.figure(figsize=(30, 10))
11    plt.bar(company_job,company_job_count)
12    plt.xlabel('Job Title')
13    plt.ylabel('Number of Listings')
14    plt.title('Number of Listings per Job Title in ' + x)

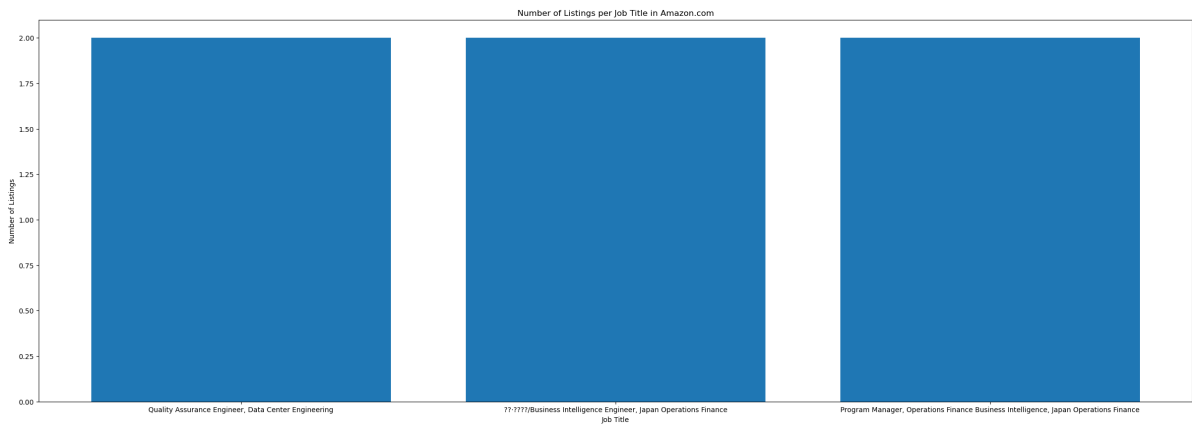
```

C:\Users\Matt\anaconda3\lib\site-packages\IPython\core\events.py:89: UserWarning: Glyph 150 (\x96) missing from current font.
func(*args, **kwargs)



C:\Users\Matt\anaconda3\lib\site-packages\IPython\core\pylabtools.py:152: UserWarning: Glyph 150 (\x96) missing from current font.
fig.canvas.print_figure(bytes_io, **kw)





In [180]:

```
1 #We will now take a look at some descriptive statistics of the salaries in
2 df.describe()
```

Out[180]:

	Salary_USD
count	3009.000000
mean	90317.996012
std	42646.582032
min	30000.000000
25%	56000.000000
50%	77000.000000
75%	115000.000000
max	315000.000000

In [181]:

```
1 #From this, we can conclude that the average salary of our data set is aro
2 #Our minimum salary is around $30k USD and our maximum salary is $315K USD
```

```
In [182]: 1 #Let us explore the entries with the top 10 highest salaries in this data
          2 df.nlargest(10, 'Salary_USD')
```

Out[182]:

	Company	Job Title	Location	Job Type	Experience level	Salary	Requirement
868	Roblox	Senior/Principal ML Engineer, Content Understa...	san mateo, ca, united states	Full Time	Senior-level	315K+	Architecture,Comp
1422	Roblox	Senior Data Scientist - Creator Success	san mateo, ca, united states	Full Time	Senior-level	310K+	A/B testing,Airflow,
789	Roblox	Senior Data Scientist - Creator Success	san mateo, ca, united states	Full Time	Senior-level	310K+	A/B testing,Airflow,
		Senior Data	san mateo				

```
In [183]: 1 #We can see here that the company Roblox offers the highest salary compens
          2 #It is also interesting to note that these salaries are in the United Stat
```

```
In [184]: 1 #Lets explore the job distribution by country
          2 #First Lets see how many jobs there are per country:
          3 df['Country'].value_counts()
          4
          5
```

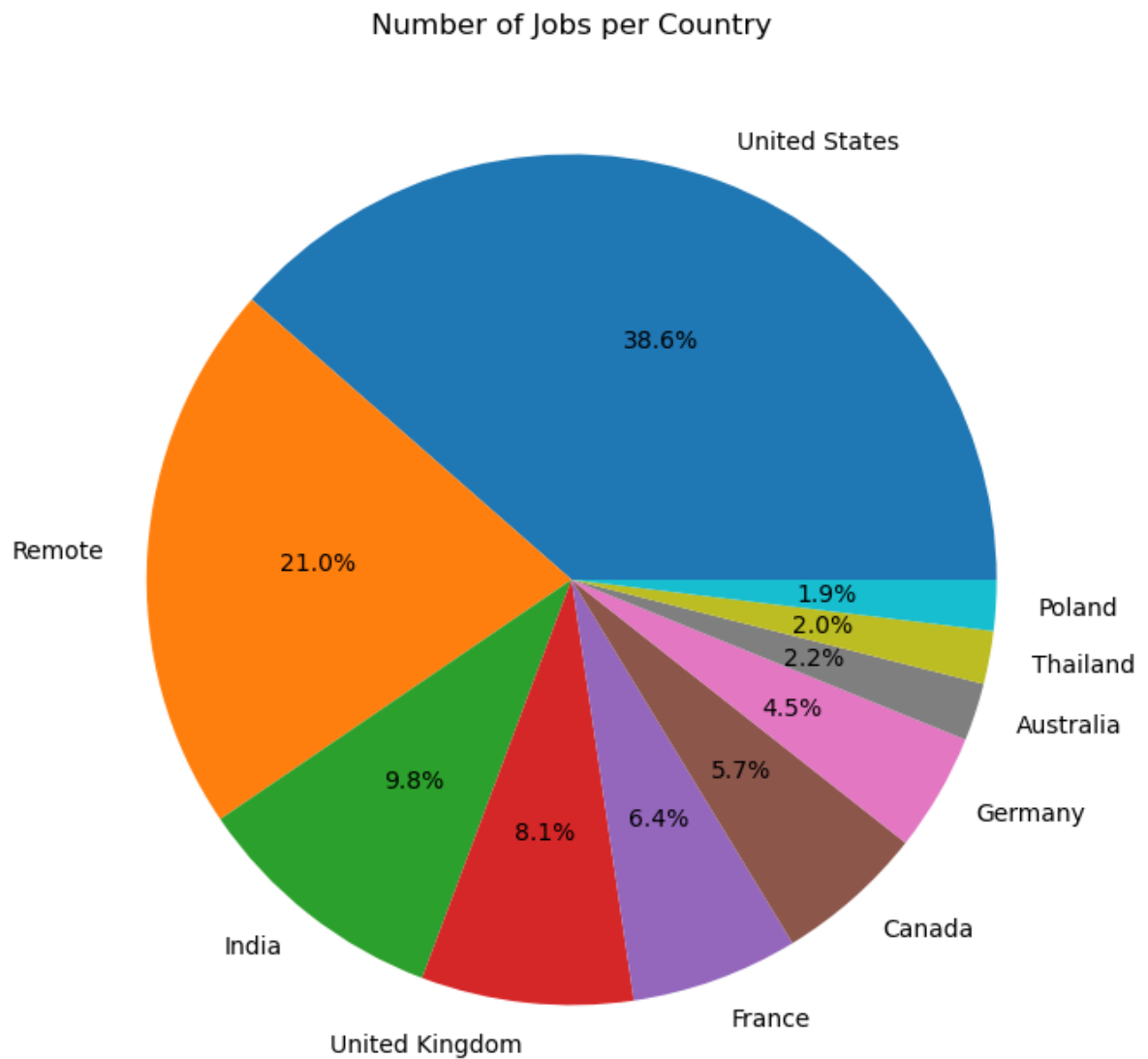
Out[184]:

United States	861
Remote	468
India	218
United Kingdom	180
France	142
...	
Chile	1
Niger	1
Jordan	1
Cyprus	1
Latvia	1

Name: Country, Length: 79, dtype: int64

```
In [185]: 1 #As we can see, a majority of the jobs lie in the United States
```

```
In [186]: 1 #Lets visualize this data as a pie chart
2 #For simplicity, we will limit our data to the top ten countries that hav
3
4 #Plotting the Pie Chart:
5 #extracts the number of jobs by country into a series
6 jobs_per_country = df['Country'].value_counts()[:10]
7 #extracts the countries
8 countries = jobs_per_country.index
9 #extracts the number of values
10 values = jobs_per_country.values
11
12 #plotting
13 plt.figure(figsize = (10,8))
14 plt.pie(values, labels = countries, autopct='%1.1f%%')
15 plt.title('Number of Jobs per Country')
16 plt.show()
```



In [187]: 1 *#We can see that the Remote work is the second highest source of Jobs for*

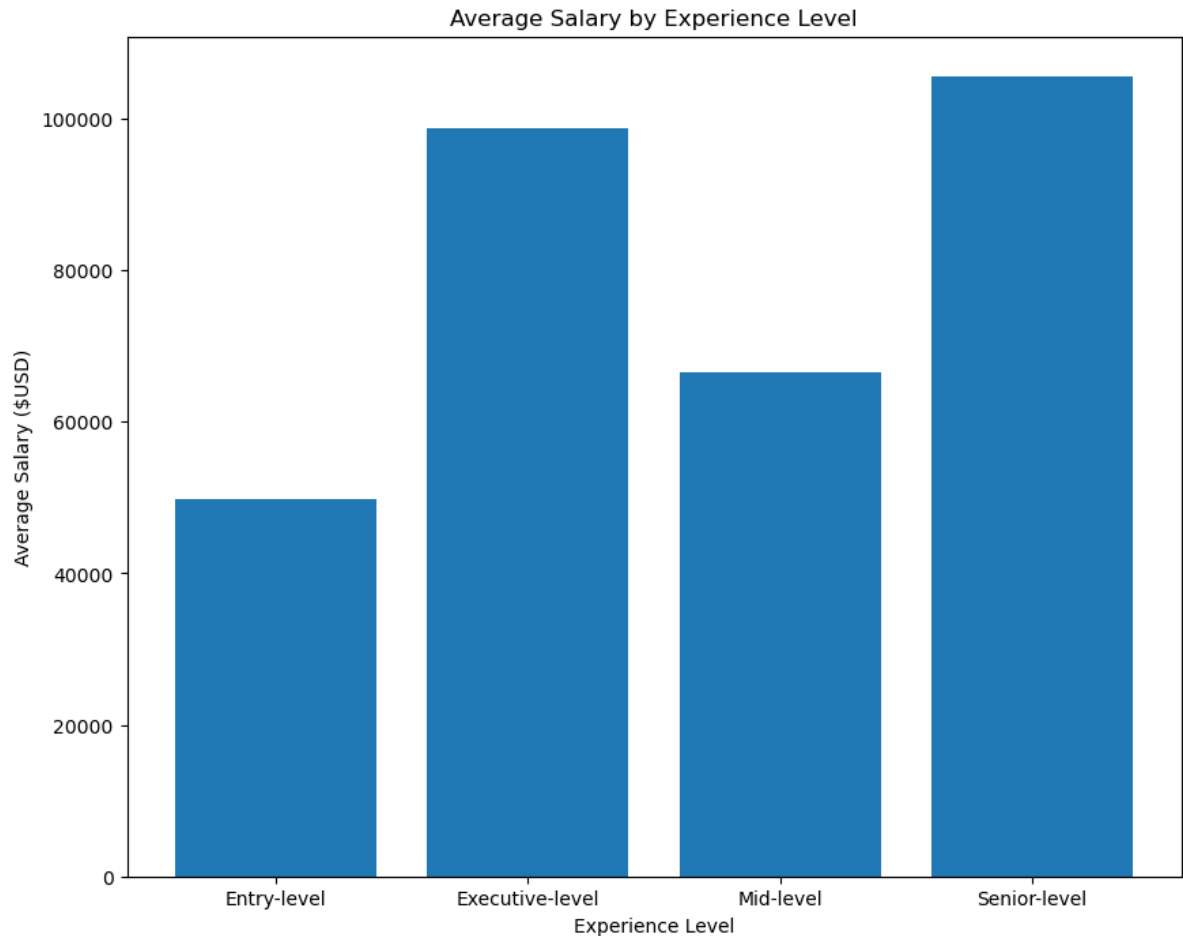
In [188]: 1 *#Now, Lets analyze salaries by experince level:*
 2 `print(df['Experience level'].value_counts())`
 3 `salary_group = df.groupby(['Experience level'])`
 4
 5
 6 *#Analyzing based on salary for each group:*
 7 `salary_group.describe()`

```
Senior-level      1825
Mid-level         447
Entry-level       383
Executive-level   122
Name: Experience level, dtype: int64
```

Out[188]:

								Salary_USD
	count	mean	std	min	25%	50%	75%	max
Experience level								
Entry-level	383.0	49795.665796	21816.940402	30000.0	39000.0	44000.0	52000.0	167000.0
Executive-level	122.0	98803.278689	39876.594627	57000.0	73000.0	81000.0	102000.0	230000.0
Mid-level	447.0	66465.033557	22365.372417	30000.0	54500.0	65000.0	73000.0	207000.0
Senior-level	1825.0	105580.953425	42877.321482	30000.0	69000.0	110000.0	129000.0	315000.0

```
In [189]: 1 #lets plot the average salaries for each group in a bar chart:
2 average_salary = salary_group['Salary_USD'].mean()
3 plt.figure(figsize=(10, 8))
4 plt.bar(average_salary.index, average_salary.values)
5 plt.xlabel('Experience Level')
6 plt.ylabel('Average Salary ($USD)')
7 plt.title('Average Salary by Experience Level')
8 plt.xticks(rotation=0)
9 plt.show()
```




```
In [190]: 1 #Lets visualize the groupings for the entry_level positions using a
2
3 entry_level_filt = df['Experience level'] == 'Entry-level'
4 entry_level = df[entry_level_filt]
5
6 plt.figure(figsize=(8, 6))
7 plt.boxplot(entry_level['Salary_USD'])
8 plt.title('Box and Whisker Plot - Entry-level Salaries')
9 plt.xlabel('Positions')
10 plt.ylabel('Salary_USD')
11 plt.xticks([1], ['Entry'])
12 plt.show()
13
```

