

```
# Importing Libraries
import ast
import pandas as pd
import seaborn as sns
from datasets import load_dataset
import matplotlib.pyplot as plt

# Loading Data
dataset = load_dataset('lukebarousse/data_jobs')
df = dataset['train'].to_pandas()

# Data Cleanup
df['job_posted_date'] = pd.to_datetime(df['job_posted_date'])
df['job_skills'] = df['job_skills'].apply(lambda x: ast.literal_eval(x) if pd.notna(x) else x)
```

/usr/local/lib/python3.12/dist-packages/huggingface_hub/utils/_auth.py:94: UserWarning:
The secret 'HF_TOKEN' does not exist in your Colab secrets.
To authenticate with the Hugging Face Hub, create a token in your settings tab (<https://huggingface.co/settings/tokens>), set it as secret in your Google Colab and
You will be able to reuse this secret in all of your notebooks.
Please note that authentication is recommended but still optional to access public models or datasets.
warnings.warn(
README.md: 3.25k/? [00:00<00:00, 237kB/s]

colab.research.google.com/drive/1t5p7y67hc1J7ZCyJdJ2mC9J5ULzpl_s-#scrollTo=LfqsnxDHsDkS

EDA intro Neha.ipynb

File Edit View Insert Runtime Tools Help

Commands + Code + Text Run all

You will be able to reuse this secret in all of your notebooks.
Please note that authentication is recommended but still optional to access public models or datasets.

warnings.warn(
README.md: 3.25k? [00:00<00:00, 237kB/s]
data_jobs.csv: 100% 231M/231M [00:02<00:00, 88.6MB/s]
Generating train split: 100% 785741/785741 [00:14<00:00, 101457.14 examples/s]

```
df_plot = df['job_title_short'].value_counts().to_frame()  
  
sns.set_theme(style='ticks')  
sns.barplot(data=df_plot, x='count', y='job_title_short', hue='count', palette='dark:b_r', legend=False)  
sns.despine()  
plt.title('Number of Jobs per Job Title')  
plt.xlabel('Number of Jobs')  
plt.ylabel('')  
plt.show()
```

Number of Jobs per Job Title

Data Analyst

Variables Terminal

Type here to search

SONY



EDA intro Neha.ipynb



File Edit View Insert Runtime Tools Help

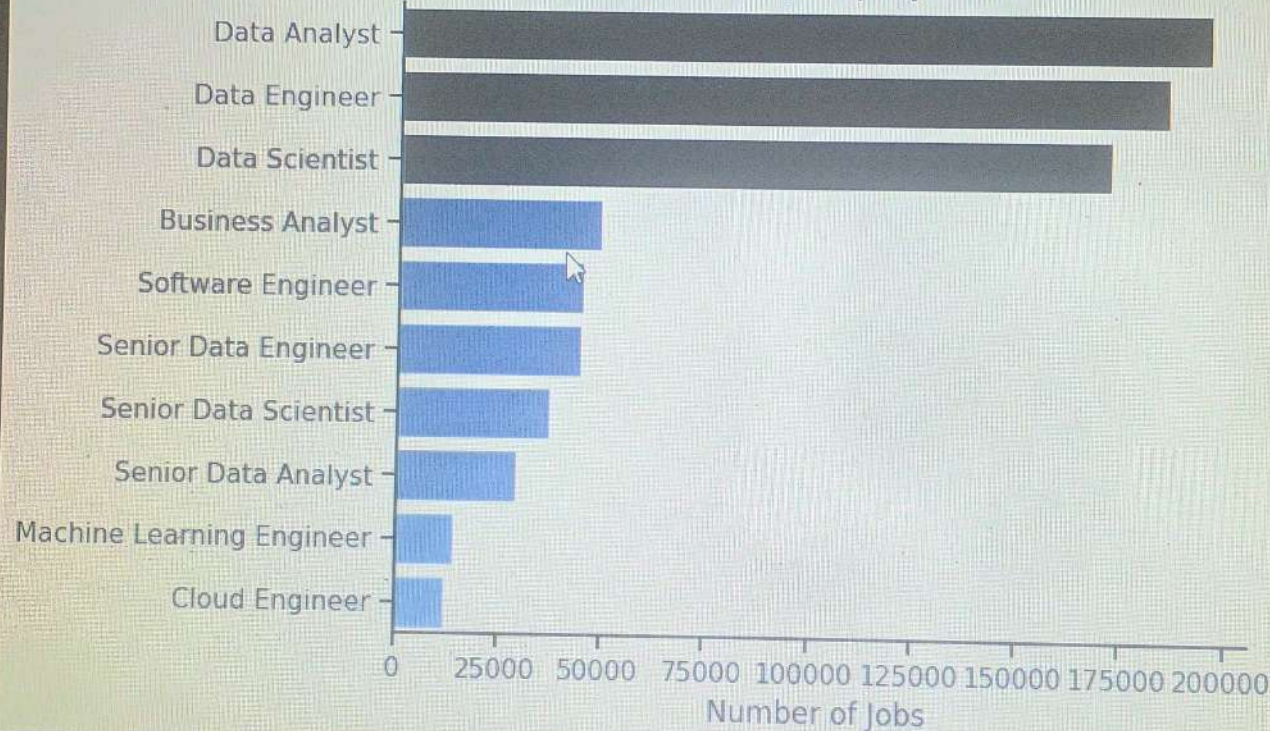
Commands

+ Code + Text

Run all



Number of Jobs per Job Title



Variables

Terminal



Type here to search



SONY

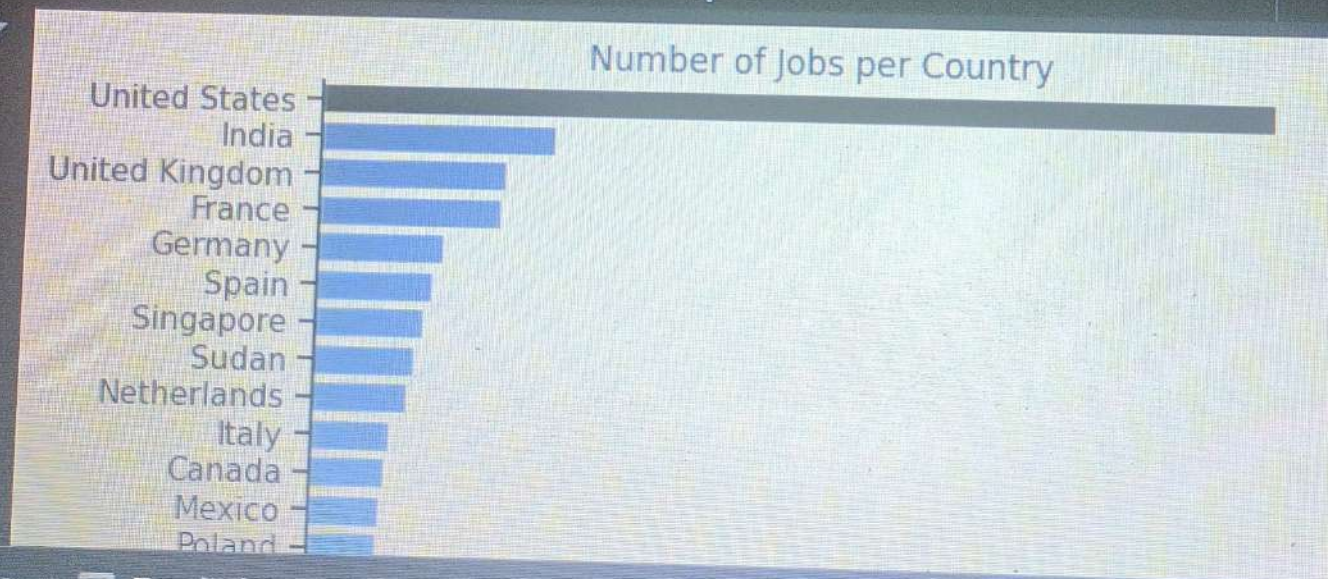
EDA intro Neha.ipynb

File Edit View Insert Runtime Tools Help

Commands + Code + Text ▶ Run all

```
df_plot = df['job_country'].value_counts().to_frame().head(20)
```

```
sns.set_theme(style='ticks')
sns.barplot(data=df_plot, x='count', y='job_country', hue='count', palette='dark:b_r', legend=False)
sns.despine()
plt.title('Number of Jobs per Country')
plt.xlabel('Number of Jobs')
plt.ylabel('')
plt.show()
```



Variables Terminal

Type here to search



SONY



EDA intro Neha.ipynb

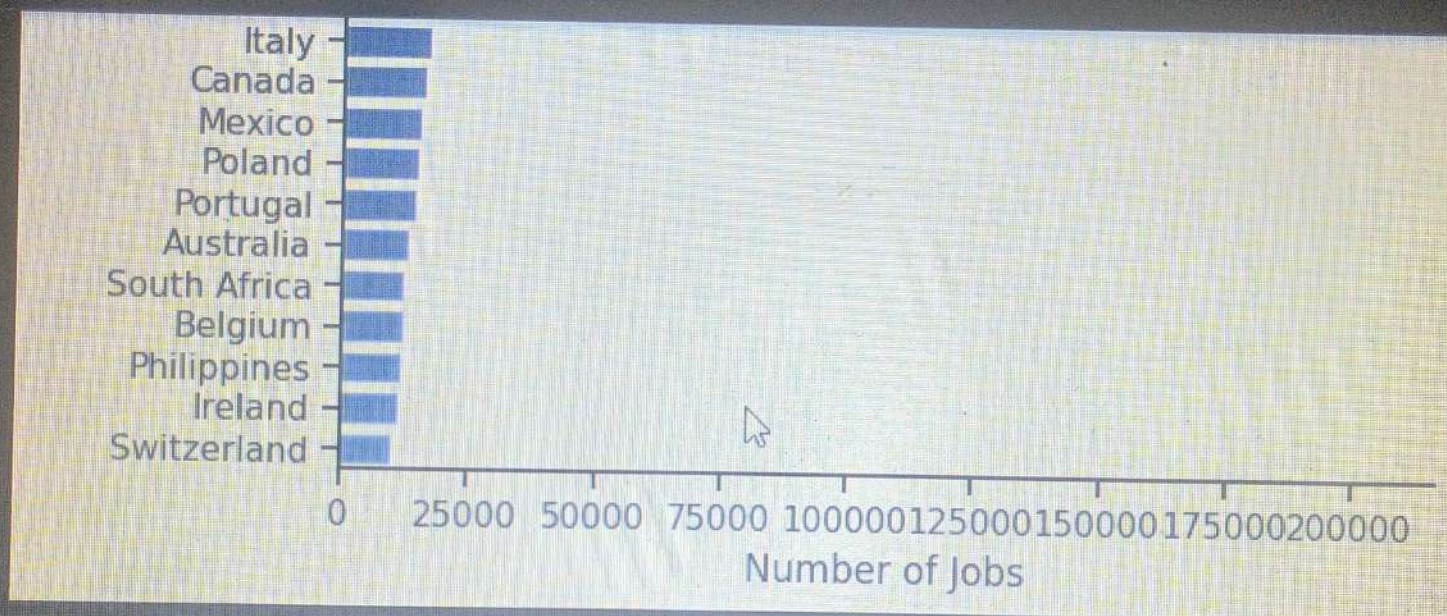


File Edit View Insert Runtime Tools Help

Commands

+ Code + Text

▶ Run all ▼



```
[4] df_plot = df['company_name'].value_counts().to_frame()[1:].head(20)

sns.set_theme(style='ticks')
sns.barplot(data=df_plot, x='count', y='company_name', hue='count', palette='dark:b_r', legend=False)
sns.despine()
plt.title('Number of Jobs per Company')
plt.xlabel('Number of Jobs')
plt.ylabel('')
plt.show()
```

{ } Variables

Terminal



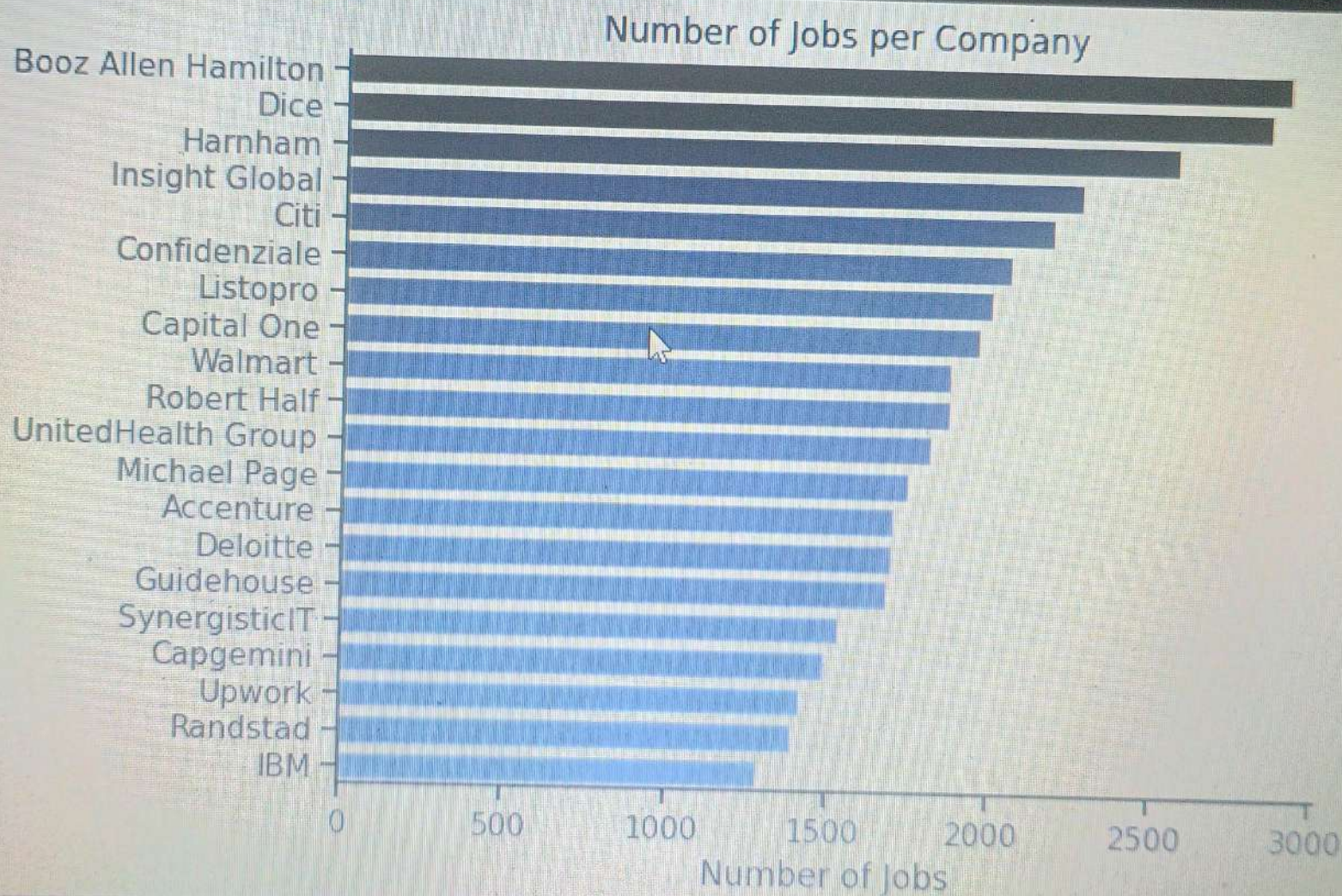
Type here to search



CO EDA intro Neha.ipynb ☆ ☁

File Edit View Insert Runtime Tools Help

Commands + Code + Text ▶ Run all ▼



{ } Variables Terminal

Type here to search





EDA intro Neha.ipynb



File Edit View Insert Runtime Tools Help

Commands

+ Code

+ Text

▶ Run all ▼



0s



```
dict_column = {  
    'job_work_from_home': 'Work from Home Offered',  
    'job_no_degree_mention': 'Degree Requirement',  
    'job_health_insurance': 'Health Insurance Offered'  
}
```

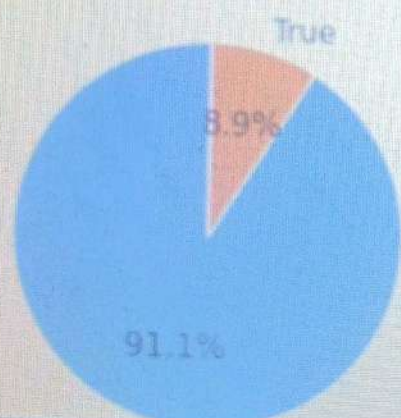
```
fig, ax = plt.subplots(1, 3, figsize=(11, 3.5))
```

```
for i, (column, title) in enumerate(dict_column.items()):  
    ax[i].pie(df[column].value_counts(), labels=['False', 'True'], autopct='%1.1f%%', startangle=90)  
    ax[i].set_title(title)
```

```
plt.show()
```



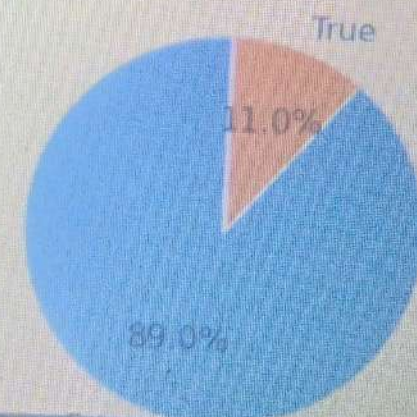
Work from Home Offered



Degree Requirement



Health Insurance Offered



Variables

Terminal



Type here to search





EDA intro Neha.ipynb



File Edit View Insert Runtime Tools Help

Commands

+ Code + Text

▶ Run all ▼



22s



Importing Libraries

```
import ast
import pandas as pd
import seaborn as sns
from datasets import load_dataset
import matplotlib.pyplot as plt
```



Loading Data

```
dataset = load_dataset('lukebarousse/data_jobs')
df = dataset['train'].to_pandas()
```

Data Cleanup

```
df['job_posted_date'] = pd.to_datetime(df['job_posted_date'])
df['job_skills'] = df['job_skills'].apply(lambda x: ast.literal_eval(x) if pd.notna(x) else x)
```

0s



```
df_DA_US = df[(df['job_country'] == 'United States') & (df['job_title_short'] == 'Data Analyst')]
```

0s



```
df_plot = df_DA_US['job_location'].value_counts().head(10).to_frame()
```

```
sns.set_theme(style='ticks')
sns.barplot(data=df_plot, x='count', y='job_location', hue='count', palette='dark:b_r', legend=False)
sns.despine()
plt.title('Counts of Job Locations for Data Analyst in the US')
plt.xlabel('Number of Jobs')
plt.ylabel('')
plt.show()
```



Variables



Terminal



Type here to search



SONY



[9] # rewrite the above with a for loop



EDA intro Neha.ipynb

File Edit View Insert Runtime Tools Help

Commands + Code + Text ▶ Run all

```
# rewrite the above with a for loop
dict_column = {
    'job_work_from_home': 'Work from Home Offered',
    'job_no_degree_mention': 'Degree Requirement',
    'job_health_insurance': 'Health Insurance Offered'
}

fig, ax = plt.subplots(1, 3)
fig.set_size_inches((12, 5))

for i, (column, title) in enumerate(dict_column.items()):
    ax[i].pie(df_DA_US[column].value_counts(), labels=['False', 'True'], autopct='%1.1f%%', startangle=90)
    ax[i].set_title(title)

# plt.suptitle('Benefit Analysis of Data Jobs', fontsize=16)
plt.show()
```



Variables Terminal

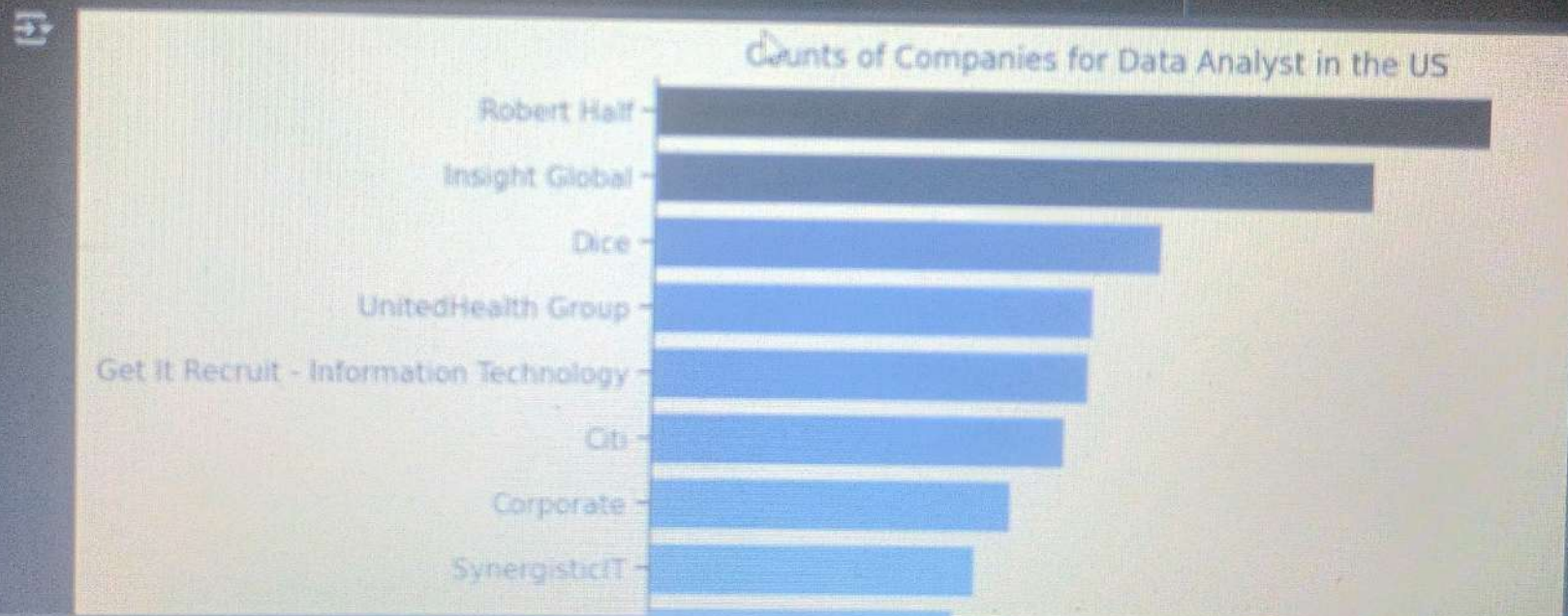
EDA intro Neha.ipynb

File Edit View Insert Runtime Tools Help

Commands + Code + Text Run all

```
df_plot = df_DA_US['company_name'].value_counts().head(10).to_frame()

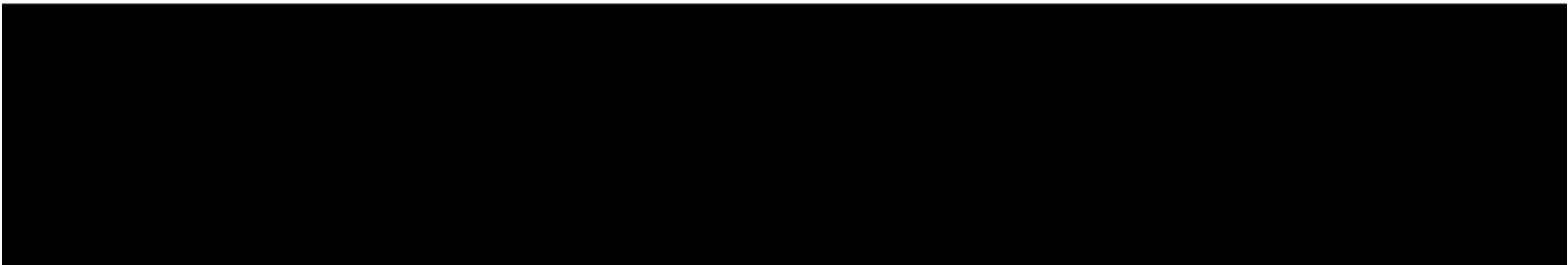
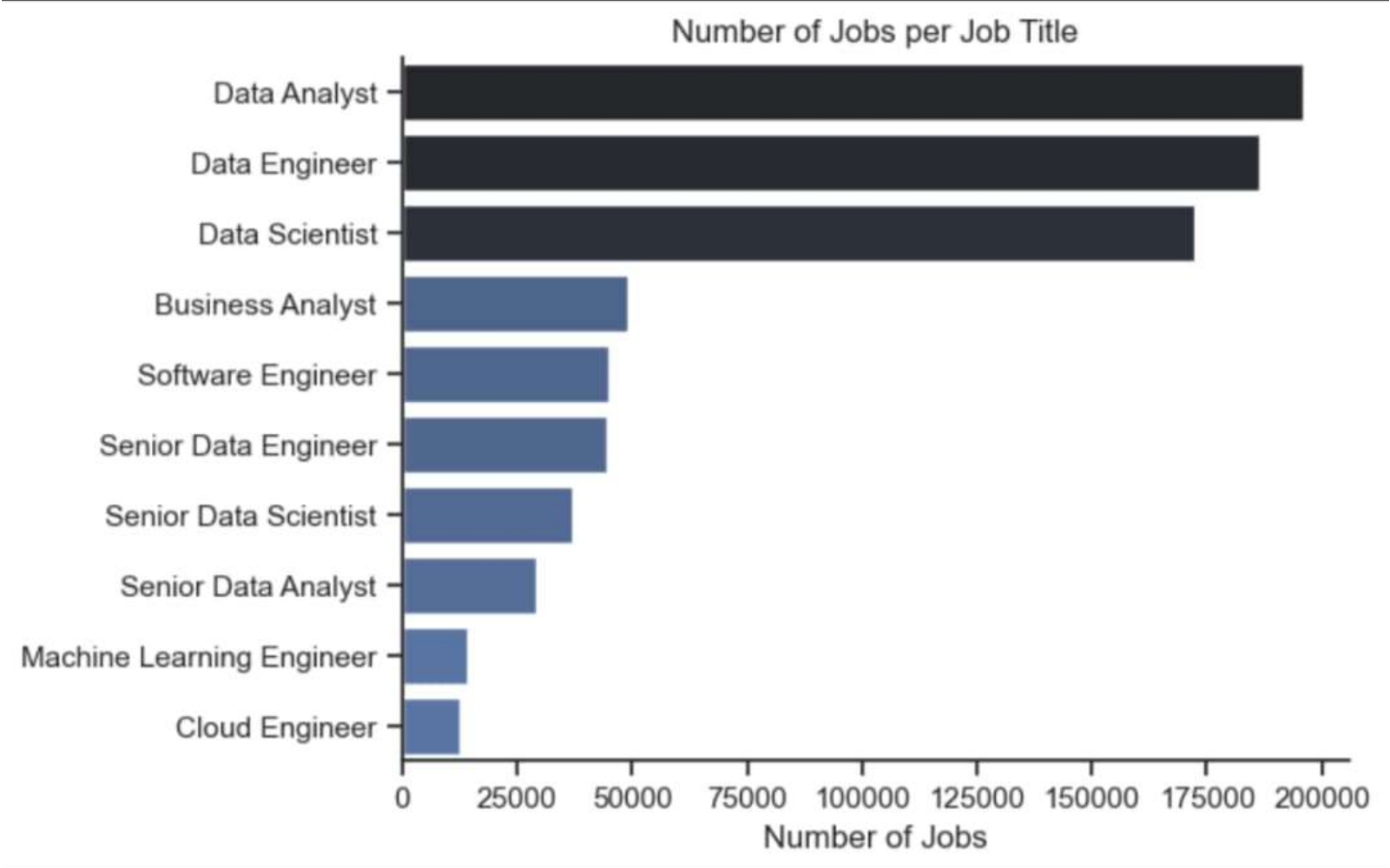
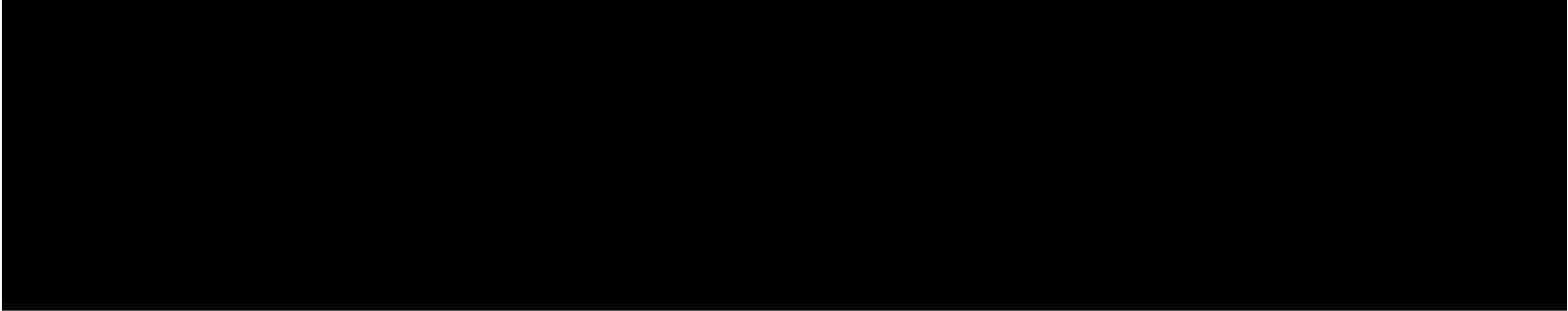
sns.set_theme(style='ticks')
sns.barplot(data=df_plot, x='count', y='company_name', hue='count', palette='dark:b_r', legend=False)
sns.despine()
plt.title('Counts of Companies for Data Analyst in the US')
plt.xlabel('Number of Jobs')
plt.ylabel('')
plt.show()
```



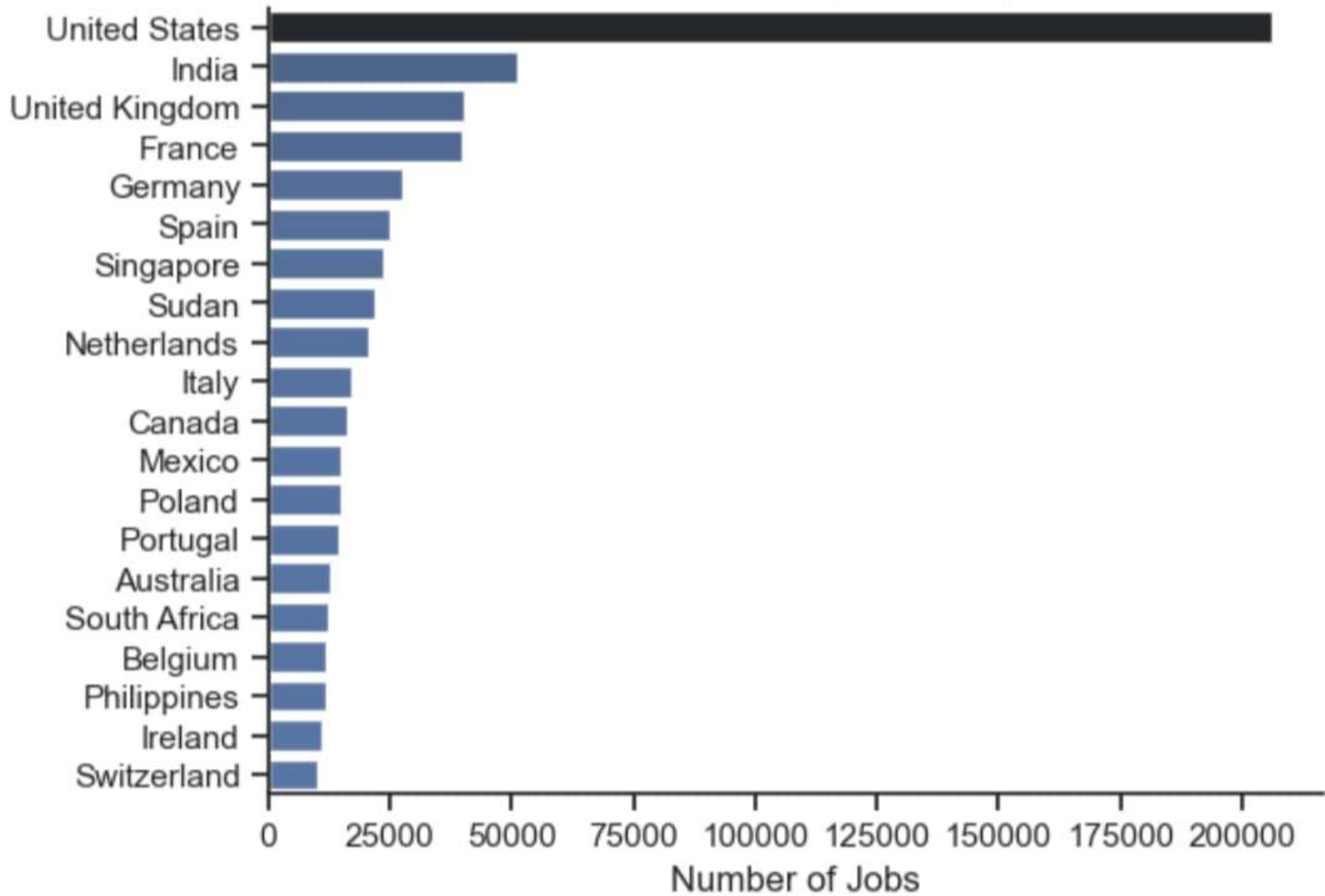
Variables Terminal

Type here to search

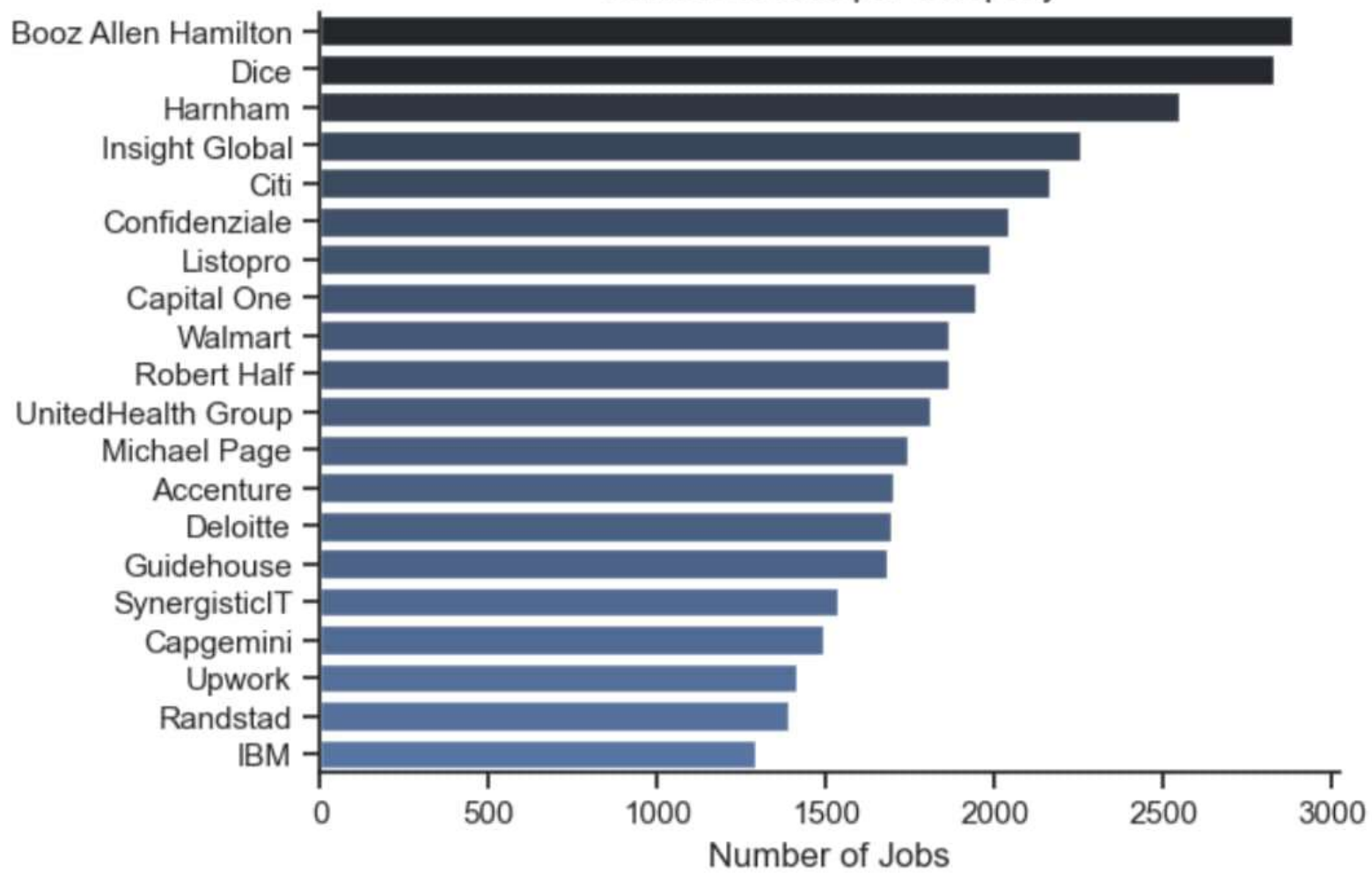




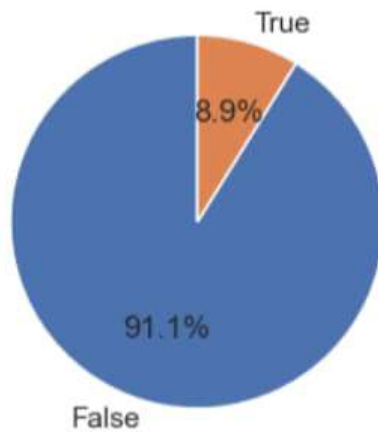
Number of Jobs per Country



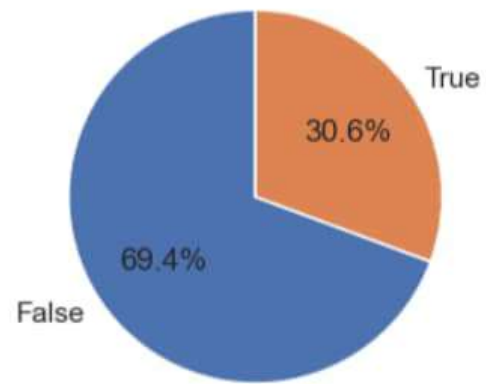
Number of Jobs per Company



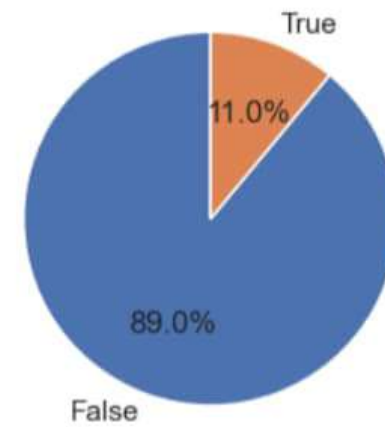
Work from Home Offered



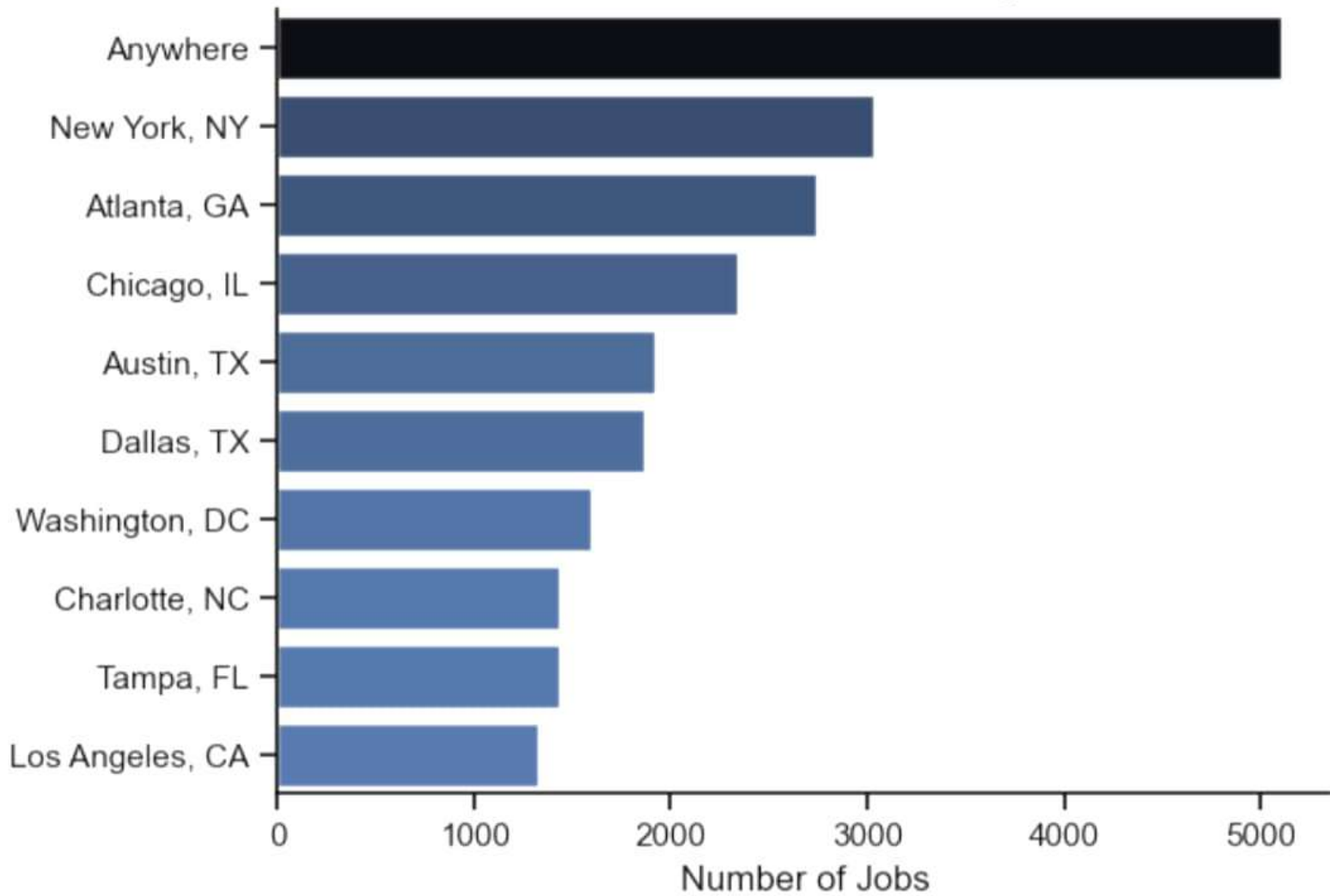
Degree Requirement



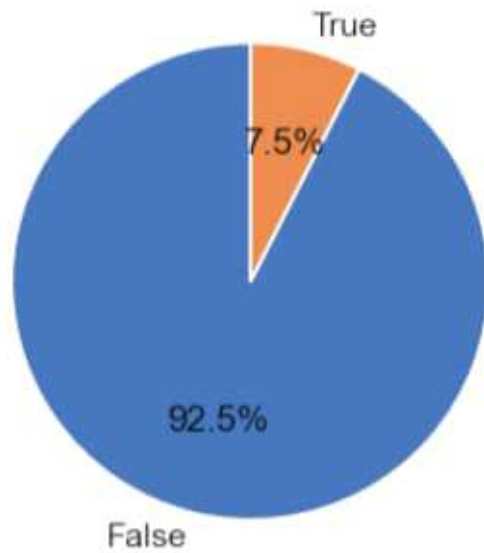
Health Insurance Offered



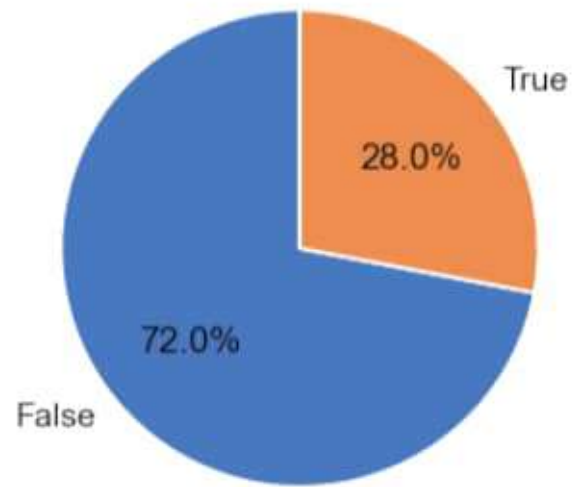
Counts of Job Locations for Data Analyst in the US



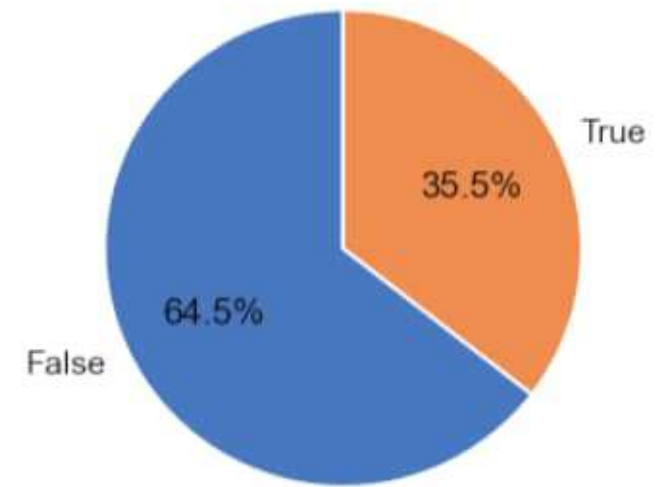
Work from Home Offered



Degree Requirement



Health Insurance Offered



Counts of Companies for Data Analyst in the US

