

Project Name - Data Analyst Jobs (ML _ FA _ DA projects) (Part 1)

Project Type - Data Analysis

Industry - Unified Mentor

Contribution - Individual

Member Name - Hare Krishana Mishra

Task - 1

Project Summary -

Project Description:

This project analyzes over 2,000 job listings for Data Analyst roles to uncover industry trends, key skills, and salary patterns. Using data cleaning, exploratory data analysis (EDA), and machine learning, the project identifies factors influencing salary, such as company rating, size, sector, and required skills. A Random Forest Regressor model is trained to predict the average salary for given job attributes, helping job seekers and recruiters make data-driven career decisions.

Objective:

- To explore and visualize trends in Data Analyst job postings.
- To identify skills and company characteristics that significantly impact salaries.
- To develop a machine learning model that predicts average salaries based on job-related factors

Key Project Details:

Dataset: 2,253 job postings with details such as Salary Estimate, Company Rating, Location, Industry, and Job Description.

Data Cleaning: Removed duplicates, handled missing values, extracted numerical salary ranges, and standardized categorical values.

Feature Engineering:

Extracted technical skills (Python, Excel, SQL) from job descriptions.

Encoded categorical variables for machine learning compatibility.

Libraries & Tools:

Pandas, NumPy, Matplotlib, Seaborn, WordCloud, Scikit-learn (RandomForestRegressor, train_test_split, metrics), LabelEncoder

EDA Insights:

Top-paying sectors include Biotech & Pharmaceuticals, Real Estate, and Arts & Entertainment.

California locations dominate the highest average salaries.

Larger companies don't always pay more than smaller companies.

Model:

Random Forest Regressor trained to predict salaries.

Evaluation metrics: Mean Absolute Error (MAE) and R² Score.

Feature importance analysis to understand key drivers of salary.

Outcome:

A predictive model and visual insights that can guide job seekers, HR professionals, and analysts in understanding market trends.

Let's Begin:-

Data Collection

In []:

```
import pandas as pd
```

In []:

```
# Load the dataset
data = pd.read_csv("/content/DataAnalyst.csv")
```

In []:

```
# Inspect the dataset
data.head()
```

Out[]:

	Unnamed: 0	Job Title	Salary Estimate	Job Description	Rating	Company Name	Location	He
0	0	Data Analyst, Center on Immigration and Justic...	37K–66K (Glassdoor est.)	Are you eager to roll up your sleeves and harn...	3.2	Vera Institute of Justice\n3.2	New York, NY	Ne
1	1	Quality Data Analyst	37K–66K (Glassdoor est.)	Overview\n\nProvides analytical and technical ...	3.8	Visiting Nurse Service of New York\n3.8	New York, NY	Ne
2	2	Senior Data Analyst, Insights & Analytics Team...	37K–66K (Glassdoor est.)	We're looking for a Senior Data Analyst who ha...	3.4	Squarespace\n3.4	New York, NY	Ne

	Unnamed: 0	Job Title	Salary Estimate	Job Description	Rating	Company Name	Location	Headquarters
3	3	Data Analyst	37K–66K (Glassdoor est.)	Requisition NumberRR-0001939\nRemote:Yes\nWe c...	4.1	Celerity\n4.1	New York, NY	
4	4	Reporting Data Analyst	37K–66K (Glassdoor est.)	ABOUT FANDUEL GROUP\n\nFanDuel Group is a worl...	3.9	FanDuel\n3.9	New York, NY	Ne

In []:

```
print(data.info())
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 2253 entries, 0 to 2252
Data columns (total 16 columns):
#   Column                Non-Null Count  Dtype
---  -
0   Unnamed: 0            2253 non-null  int64
1   Job Title             2253 non-null  object
2   Salary Estimate       2253 non-null  object
3   Job Description       2253 non-null  object
4   Rating                2253 non-null  float64
5   Company Name          2252 non-null  object
6   Location              2253 non-null  object
7   Headquarters          2253 non-null  object
8   Size                  2253 non-null  object
9   Founded               2253 non-null  int64
10  Type of ownership     2253 non-null  object
11  Industry              2253 non-null  object
12  Sector               2253 non-null  object
13  Revenue               2253 non-null  object
14  Competitors          2253 non-null  object
15  Easy Apply            2253 non-null  object
dtypes: float64(1), int64(2), object(13)
memory usage: 281.8+ KB
None
```

Distribution of Company Ratings

In []:

```
import matplotlib.pyplot as plt
import seaborn as sns

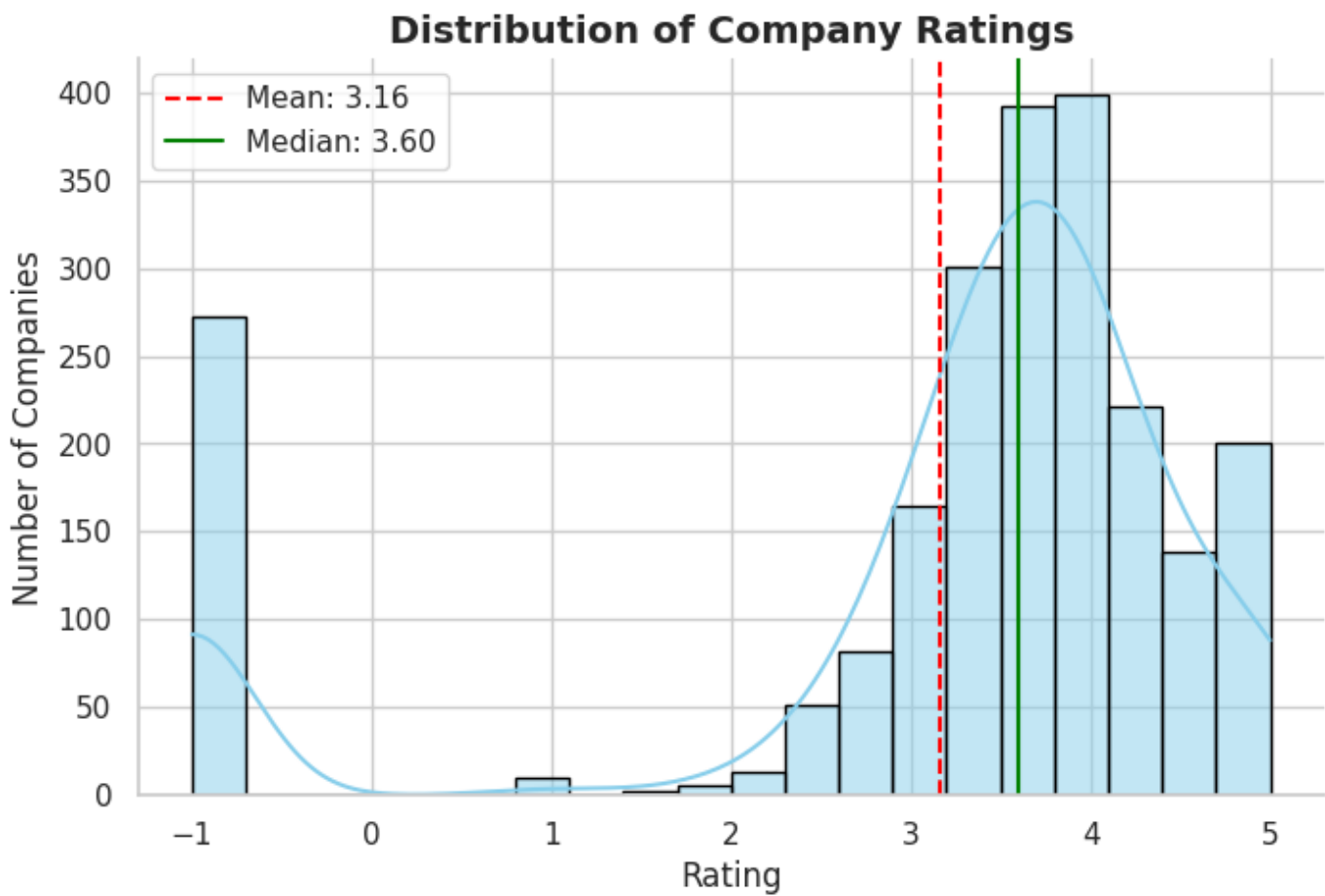
plt.figure(figsize=(8, 5))
sns.histplot(data=data, x='Rating', bins=20, kde=True, color='skyblue', edgecolor='black')

# Add mean and median lines
mean_rating = data['Rating'].mean()
median_rating = data['Rating'].median()
plt.axvline(mean_rating, color='red', linestyle='--', linewidth=1.5, label=f"Mean: {mean_rating}")
plt.axvline(median_rating, color='green', linestyle='-', linewidth=1.5, label=f"Median: {median_rating}")

# Styling
plt.title('Distribution of Company Ratings', fontsize=14, fontweight='bold')
plt.xlabel('Rating', fontsize=12)
```

```
plt.ylabel('Number of Companies', fontsize=12)
plt.legend()
sns.despine(top=True, right=True)

plt.show()
```



Exploratory Data Analysis (EDA)

```
In [ ]:
# Check for duplicates
print(f"Duplicate rows: {data.duplicated().sum()}")
```

Duplicate rows: 0

```
In [ ]:
# General statistics
data.describe(include='all')
```

Out[]:

	Unnamed: 0	Job Title	Salary Estimate	Job Description	Rating	Company Name	Location	Headquarters
count	2253.0000	2253	2253	2253	2253.000000	2252	2253	2253
unique	NaN	1272	90	2253	NaN	1513	253	483
top	NaN	Data Analyst	41K–78K (Glassdoor est.)	You.\n\nYou bring your body, mind, heart and s...	NaN	Staffigo Technical Services, LLC\n5.0	New York, NY	New York, NY

	Unnamed: 0	Job Title	Salary Estimate	Job Description	Rating	Company Name	Location	Headquarters
freq	NaN	405	57	1	NaN	58	310	206
mean	1126.0000	NaN	NaN	NaN	3.160630	NaN	NaN	NaN
std	650.5294	NaN	NaN	NaN	1.665228	NaN	NaN	NaN
min	0.0000	NaN	NaN	NaN	-1.000000	NaN	NaN	NaN
25%	563.0000	NaN	NaN	NaN	3.100000	NaN	NaN	NaN
50%	1126.0000	NaN	NaN	NaN	3.600000	NaN	NaN	NaN
75%	1689.0000	NaN	NaN	NaN	4.000000	NaN	NaN	NaN
max	2252.0000	NaN	NaN	NaN	5.000000	NaN	NaN	NaN

Visualization

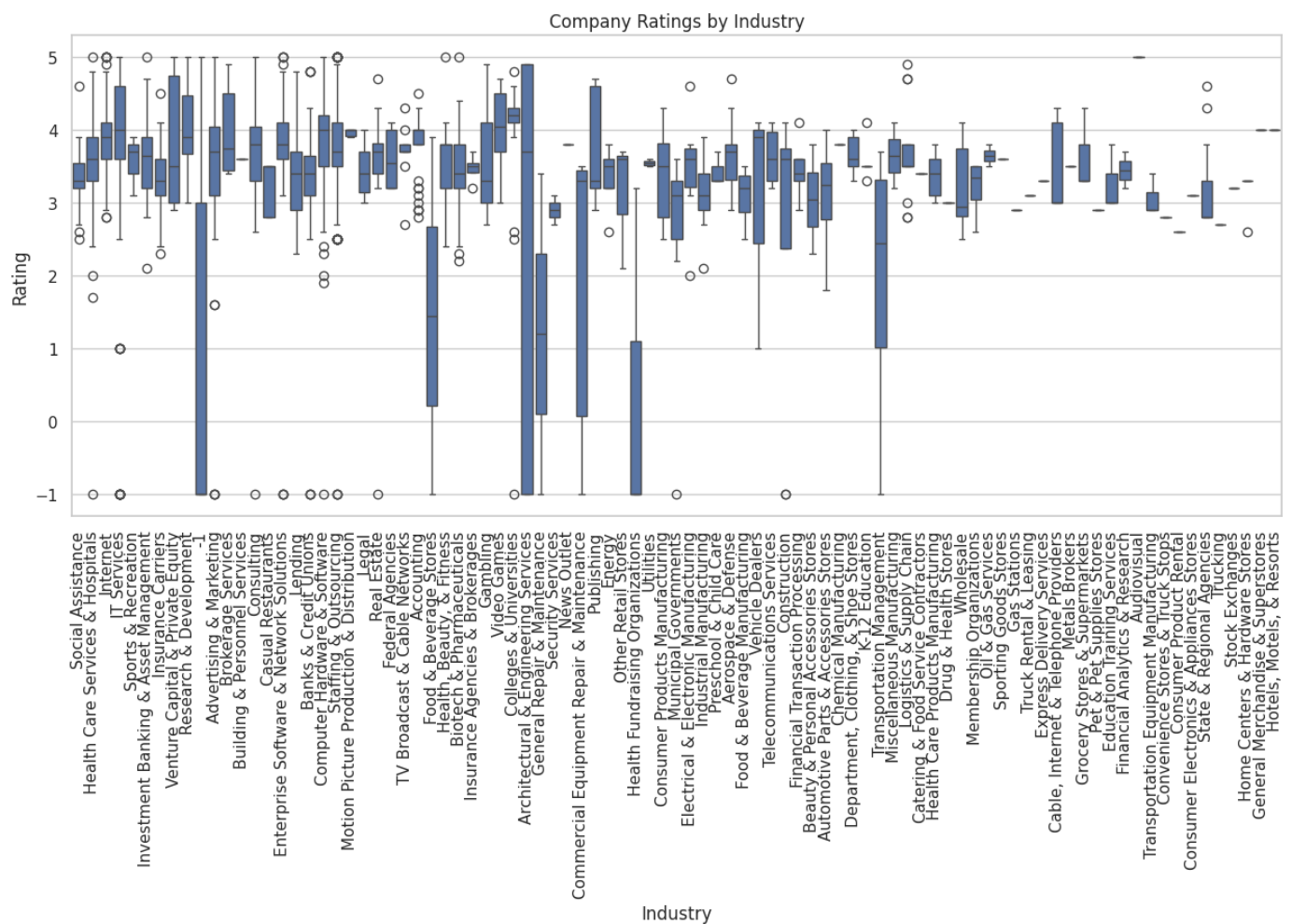
In []:

```
import matplotlib.pyplot as plt
import seaborn as sns
```

Ratings by Industry

In []:

```
plt.figure(figsize=(15, 6))
sns.boxplot(x='Industry', y='Rating', data=data)
plt.xticks(rotation=90)
plt.title("Company Ratings by Industry")
plt.show()
```

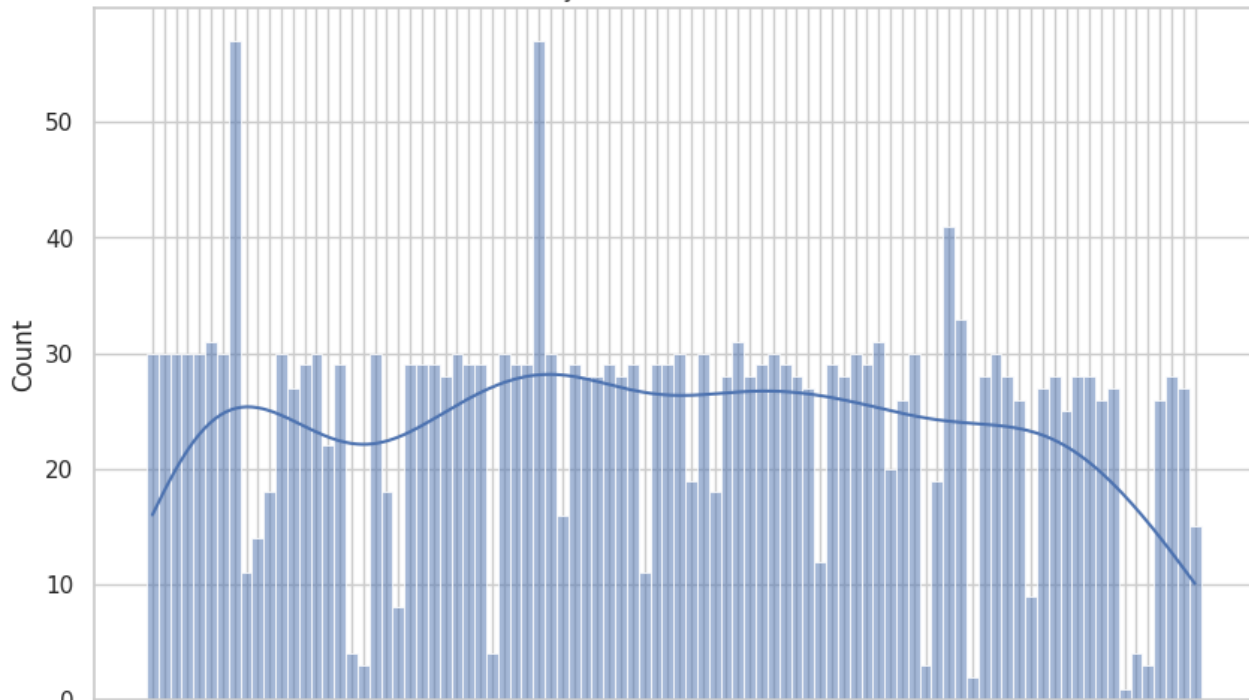


Salary Distribution

In []:

```
# Salary distribution
plt.figure(figsize=(10, 6))
sns.histplot(data['Salary Estimate'], kde=True, bins=20)
plt.title("Salary Estimate Distribution")
plt.xlabel("Salary")
plt.show()
```

Salary Estimate Distribution



```
34000 35000 36000 37000 38000 39000 40000 41000 42000 43000 44000 45000 46000 47000 48000 49000 50000 51000 52000 53000 54000 55000 56000 57000 58000 59000 60000 61000 62000 63000 64000 65000 66000 67000 68000 69000 70000 71000 72000 73000 74000 75000 76000 77000 78000 79000 80000 81000 82000 83000 84000 85000 86000 87000 88000 89000 90000 91000 92000 93000 94000 95000 96000 97000 98000 99000 100000 101000 102000 103000 104000 105000 106000 107000 108000 109000 110000 111000 112000 113000 114000 115000 116000 117000 118000 119000 120000 121000 122000 123000 124000 125000 126000 127000 128000 129000 130000 131000 132000 133000 134000 135000 136000 137000 138000 139000 140000
```

Salary

Data Cleaning

```
In [ ]:
```

```
# Check missing values
print(data.isnull().sum())
```

```
Unnamed: 0          0
Job Title          0
Salary Estimate    0
Job Description     0
Rating             0
Company Name       1
Location           0
Headquarters       0
Size               0
Founded            0
Type of ownership  0
Industry           0
Sector             0
Revenue            0
Competitors        0
Easy Apply         0
dtype: int64
```

```
In [ ]:
```

```
# Fill missing numerical values
data['Rating'].fillna(data['Rating'].median(), inplace=True)
```

```
/tmp/ipython-input-1749438820.py:2: FutureWarning:
```

A value is trying to be set on a copy of a DataFrame or Series through chained assignment using an inplace method.
The behavior will change in pandas 3.0. This inplace method will never work because the intermediate object on which we are setting values always behaves as a copy.

For example, when doing 'df[col].method(value, inplace=True)', try using 'df.method({col: value}, inplace=True)' or df[col] = df[col].method(value) instead, to perform the operation inplace on the original object.

```
In [ ]:
```

```
# Drop columns with > 30% missing data
threshold = len(data) * 0.3
data = data.dropna(thresh=threshold, axis=1)
```

```
In [ ]:
```

```
# Forward-fill categorical values
categorical_cols = ['Company Name', 'Industry', 'Sector', 'Type of ownership']
data[categorical_cols] = data[categorical_cols].fillna(method='ffill')
```

```
/tmp/ipython-input-3094392204.py:3: FutureWarning:
```

```
DataFrame.fillna with 'method' is deprecated and will raise in a future version. Use obj.ffill() or obj.bfill() instead.
```

Standardizing Data

```
In [ ]:
```

```
# Extract minimum salary
data['Min Salary'] = data['Salary Estimate'].str.extract(r'(\d+)').astype(float)
```

```
In [ ]:
```

```
# Extract maximum salary
data['Max Salary'] = data['Salary Estimate'].str.extract(r'-\s*(\d+)').astype(float)
```

```
In [ ]:
```

```
# Compute average salary
data['Avg Salary'] = (data['Min Salary'] + data['Max Salary']) / 2
```

```
In [ ]:
```

```
# Drop old salary column
data.drop('Salary Estimate', axis=1, inplace=True)
```

Feature Engineering

```
In [ ]:
```

```
# Extract keywords from Job Description
data['Python'] = data['Job Description'].str.contains('Python', case=False, na=False).astype(int)
data['Excel'] = data['Job Description'].str.contains('Excel', case=False, na=False).astype(int)
```

```
In [ ]:
```

```
# Create a tech skills score
data['Tech_Skills'] = data['Python'] + data['Excel']
```

Location Splits

```
In [ ]:
```



```
# Extract city and state from location
data['City'] = data['Location'].str.split(',', expand=True)[0]
data['State'] = data['Location'].str.split(',', expand=True)[1]
```

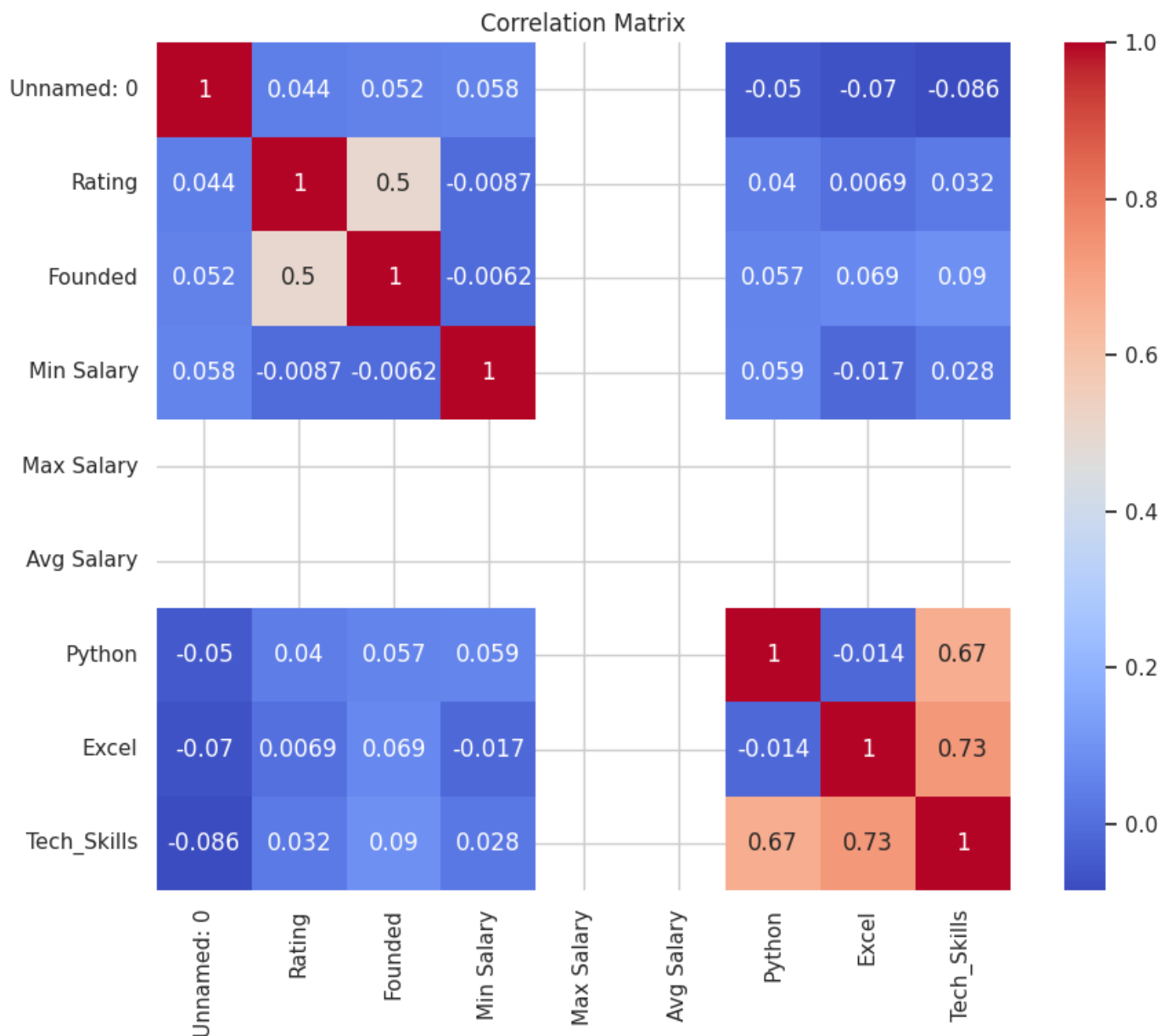
Statistics

Analyze relationships using correlation and significance tests.

In []:

```
# Keep only numeric columns
numeric_data = data.select_dtypes(include=['number'])

# Plot correlation matrix
plt.figure(figsize=(10, 8))
sns.heatmap(numeric_data.corr(), annot=True, cmap="coolwarm")
plt.title("Correlation Matrix")
plt.show()
```



Word Cloud: Most common words in Job Titles

```
from wordcloud import WordCloud
text = " ".join(data['Job Title'].astype(str))
wordcloud = WordCloud(width=800, height=400, background_color="white").generate(text)
plt.figure(figsize=(10,6))
plt.imshow(wordcloud, interpolation='bilinear')
plt.axis("off")
plt.title("Most Common Words in Job Titles")
plt.show()
```



```
plt.figure(figsize=(10,6))
sns.set_theme(style="whitegrid")

# Define custom colors for each ownership type
ownership_colors = {
    'Private': '#1f77b4',      # Blue
    'Public': '#ff7f0e',       # Orange
    'Nonprofit': '#2ca02c',    # Green
    'Subsidiary or Business Segment': '#d62728', # Red
    'Government': '#9467bd',   # Purple
    'Other': '#8c564b'         # Brown
}

# Map only colors that exist in the data
hue_order = [col for col in ownership_colors.keys() if col in data['Type of ownership']].
palette = {k: ownership_colors[k] for k in hue_order}

# KDE Plot
sns.kdeplot(
    data=data,
    x='Rating',
    hue='Type of ownership',
```

```

fill=True,
common_norm=False,
alpha=0.5,
palette=palette,
hue_order=hue_order
)

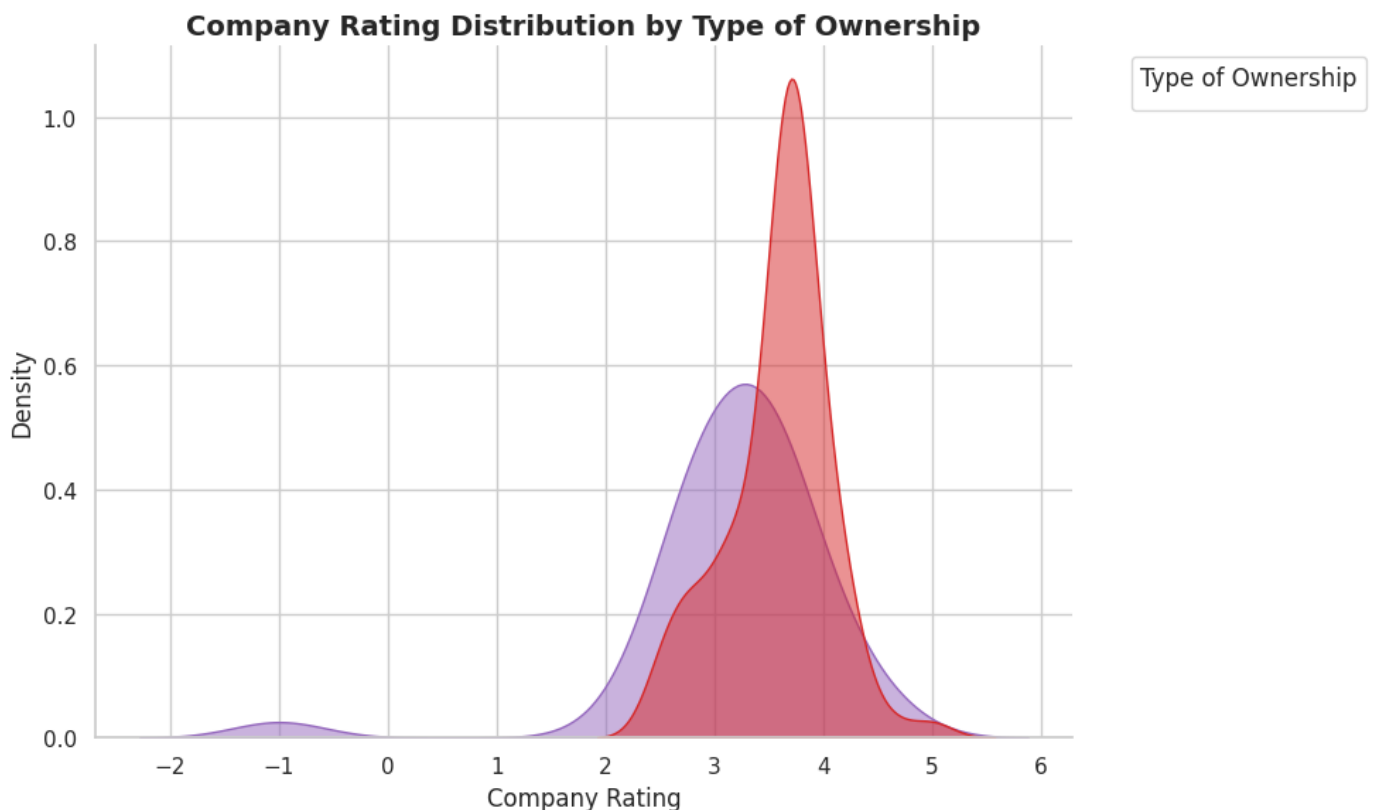
# Labels and title
plt.title("Company Rating Distribution by Type of Ownership", fontsize=14, fontweight='b')
plt.xlabel("Company Rating", fontsize=12)
plt.ylabel("Density", fontsize=12)
plt.legend(title='Type of Ownership', bbox_to_anchor=(1.05, 1), loc='upper left')

# Clean style
sns.despine()
plt.tight_layout()
plt.show()

```

/tmp/ipython-input-1431969422.py:34: UserWarning:

No artists with labels found to put in legend. Note that artists whose label start with an underscore are ignored when legend() is called with no argument.



Simple Bar Graph: Top 5 Job Titles

```

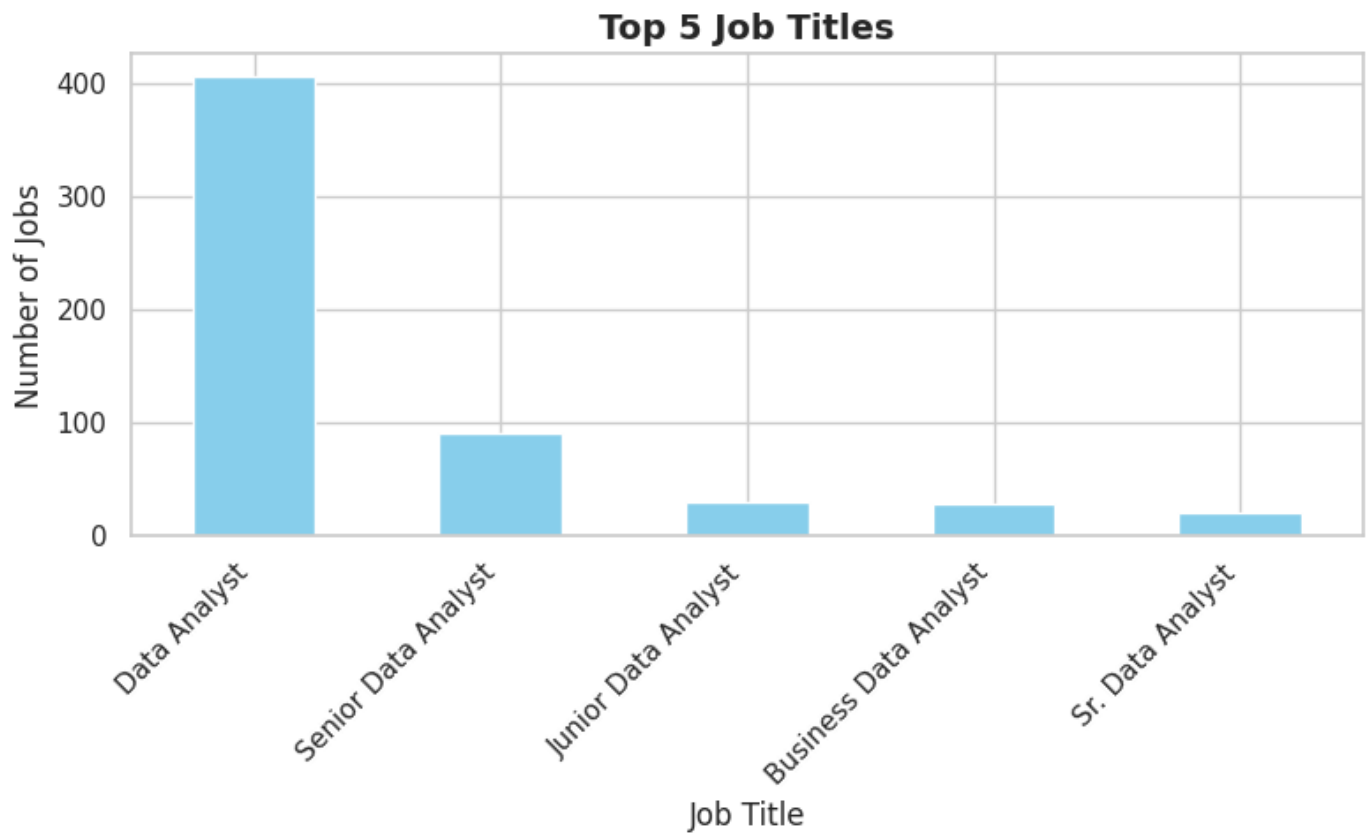
In [ ]:
top_jobs = data['Job Title'].value_counts().head(5)

plt.figure(figsize=(8,5))
top_jobs.plot(kind='bar', color='skyblue')

plt.title("Top 5 Job Titles", fontsize=14, fontweight='bold')
plt.xlabel("Job Title")
plt.ylabel("Number of Jobs")

```

```
plt.xticks(rotation=45, ha='right')
plt.tight_layout()
plt.show()
```



Model Development

Data Splitting

Split into features and target:

In []:

```
from sklearn.model_selection import train_test_split
# Define features and target
features = ['Rating', 'Tech_Skills', 'Size', 'Founded']
X = data[features]
y = data['Avg Salary']
# Train-test split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
```

Model Training

In []:

```
# =====
# Model Training for Data Analyst Jobs
# =====

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

```

from sklearn.model_selection import train_test_split
from sklearn.ensemble import RandomForestRegressor
from sklearn.metrics import mean_absolute_error, r2_score
from sklearn.preprocessing import LabelEncoder

# 1. Load Dataset
data = pd.read_csv("DataAnalyst.csv")

# 2. Basic Cleaning
data = data.drop_duplicates()
data = data.dropna(subset=["Salary Estimate", "Rating"]) # remove rows with no salary/r

# 3. Extract Min, Max, and Avg Salary
data['MinSalary'] = data['Salary Estimate'].str.extract(r'(\d+)').astype(float)
data['MaxSalary'] = data['Salary Estimate'].str.extract(r'(\d+)\D*$').astype(float)
data['AvgSalary'] = (data['MinSalary'] + data['MaxSalary']) / 2

# Drop old column
data.drop(['Salary Estimate', 'MinSalary', 'MaxSalary'], axis=1, inplace=True)

# 4. Encode categorical variables
categorical_cols = ['Company Name', 'Location', 'Size', 'Type of ownership', 'Industry',
for col in categorical_cols:
    if col in data.columns:
        le = LabelEncoder()
        data[col] = le.fit_transform(data[col].astype(str))

# 5. Extract skills from Job Description
data['Python'] = data['Job Description'].str.contains('Python', case=False, na=False).as
data['Excel'] = data['Job Description'].str.contains('Excel', case=False, na=False).asty
data['SQL'] = data['Job Description'].str.contains('SQL', case=False, na=False).astype(i
data['Tech_Skills'] = data['Python'] + data['Excel'] + data['SQL']

# 6. Feature selection
features = ['Rating', 'Company Name', 'Location', 'Size', 'Type of ownership', 'Industry
X = data[features]
y = data['AvgSalary']

# 7. Train/Test Split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42

# 8. Model Training
model = RandomForestRegressor(n_estimators=200, random_state=42)
model.fit(X_train, y_train)

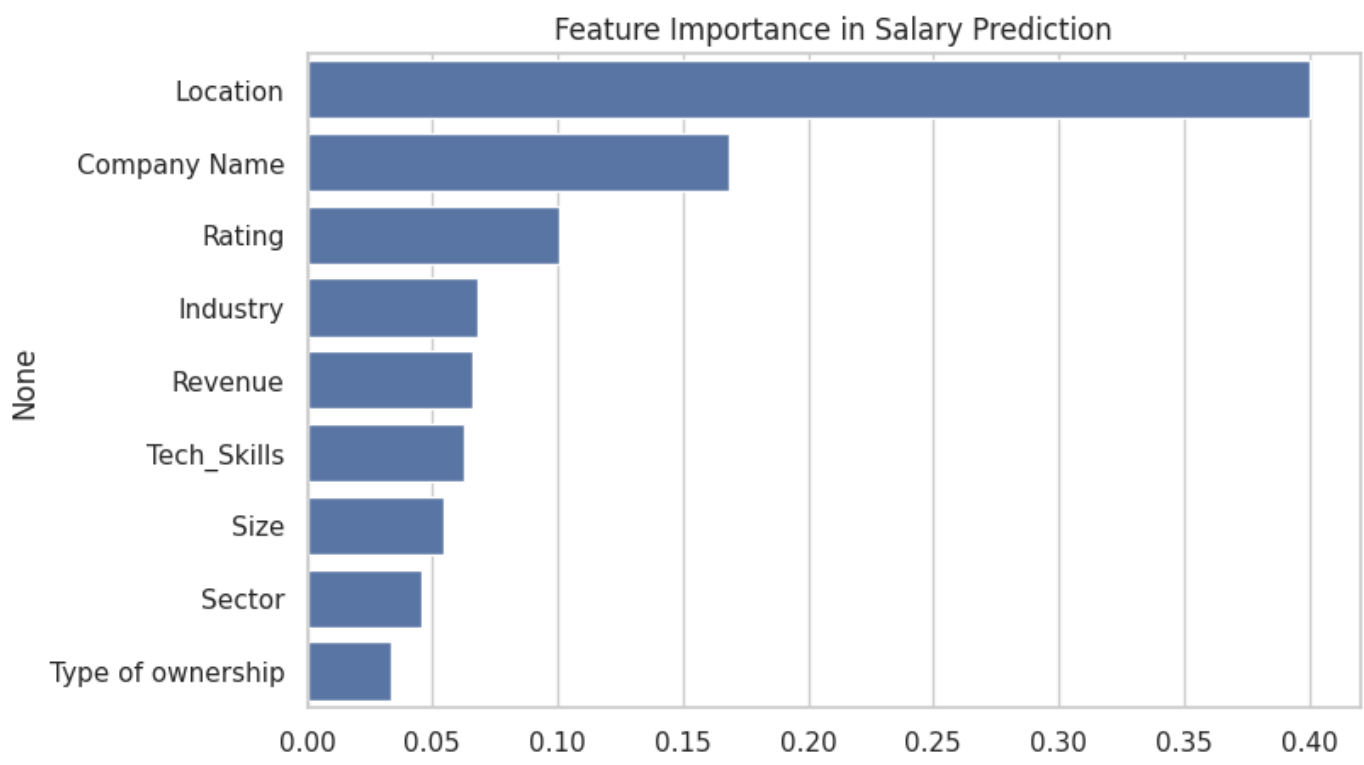
# 9. Predictions & Evaluation
y_pred = model.predict(X_test)
mae = mean_absolute_error(y_test, y_pred)
r2 = r2_score(y_test, y_pred)

print(f"MAE: {mae:.2f}")
print(f"R2 Score: {r2:.4f}")

# 10. Feature Importance
feat_importance = pd.Series(model.feature_importances_, index=features).sort_values(ascen
plt.figure(figsize=(8,5))
sns.barplot(x=feat_importance.values, y=feat_importance.index)
plt.title("Feature Importance in Salary Prediction")
plt.show()

```

MAE: 17.59
R² Score: 0.1343



Project Name - Data Analyst Jobs (ML _ FA _ DA projects) (Part 2)

Project Type - Data Analysis

Industry - Unified Mentor

Contribution - Individual

Member Name - Hare Krishana Mishra

Task - 2

Project Summary -

Project Description:

This project analyzes a dataset of over 2,000 job listings for Data Analyst positions collected from Glassdoor. It covers various attributes like salary estimates, location, company ratings, industry, job description, company size, and ownership type. The main goal is to uncover job market trends, evaluate salary ranges, and identify the factors that influence pay in the Data Analytics industry. It also includes predictive modeling to estimate salaries based on job and company features, enabling job seekers and recruiters to make informed decisions.

Objective:

- To analyze trends in Data Analyst job postings across different industries, sectors, and locations.
- To predict salary ranges based on job attributes like company rating, size, industry, and skills required.
- To provide actionable insights about company ratings, hiring patterns, and salary trends.
- To highlight top-paying sectors, industries, and locations for Data Analyst roles.

Key Project Details:

Dataset Source: Glassdoor job postings, >2,000 records, features like salary, location, company rating, job description, and ownership type.

Data Cleaning & Preprocessing: Removed duplicates, handled missing values, standardized column names, extracted salary ranges.

Exploratory Data Analysis (EDA):

Distribution of salaries, ratings, and company sizes.

Trends in job postings by industry, sector, and location.

Top industries and sectors hiring data analysts.

Libraries & Tools:

Pandas, NumPy, Matplotlib, Seaborn, Plotly (Express, Graph Objects, Subplots), Scikit-learn (RandomForestRegressor), Re, Warnings

Feature Engineering:

Extracted technical skills (Python, Excel) from job descriptions.

Created Tech_Skills score.

Split location into City and State.

Visualization Insights:

Top 10 job titles by count.

Average salary by job title, company size, and sector.

Salary trends by location.

Model Development:

Trained a Random Forest Regressor to predict average salary.

Features: Rating, Tech_Skills, Size, Founded.

Key Findings:

Highest salaries often in California-based locations.

Top-paying sectors: Biotech & Pharmaceuticals, Real Estate, Art, Entertainment & Recreation.

Private companies dominate hiring.

Deployment: Model can be deployed with Streamlit or Flask for interactive predictions.

Let's Begin:-

In []:

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import warnings
import re
import plotly.graph_objects as go
import plotly.express as px
from plotly.subplots import make_subplots
warnings.filterwarnings('ignore')
```

In []:

```
data_analyst_jobs = pd.read_csv('/content/DataAnalyst.csv')
```

In []:


```
data_analyst_jobs = data_analyst_jobs.drop('Unnamed: 0',axis=1)
data_analyst_jobs = data_analyst_jobs.drop('Founded', axis=1)
data_analyst_jobs = data_analyst_jobs.drop('Competitors',axis=1)
print(f'Number of rows:{data_analyst_jobs.shape[0]};Number ofcolumns:{data_analyst_jobs.
```

Number of rows:2253;Number ofcolumns:13; No of missingvalues:1

Dataset Overview

```
In [ ]:
data_analyst_jobs.head()
```

Out[]:

	Job Title	Salary Estimate	Job Description	Rating	Company Name	Location	Headquarters
0	Data Analyst, Center on Immigration and Justic...	37K–66K (Glassdoor est.)	Are you eager to roll up your sleeves and harn...	3.2	Vera Institute of Justice\n3.2	New York, NY	New York, NY
1	Quality Data Analyst	37K–66K (Glassdoor est.)	Overview\n\nProvides analytical and technical ...	3.8	Visiting Nurse Service of New York\n3.8	New York, NY	New York, NY
2	Senior Data Analyst, Insights & Analytics Team...	37K–66K (Glassdoor est.)	We're looking for a Senior Data Analyst who ha...	3.4	Squarespace\n3.4	New York, NY	New York, NY
3	Data Analyst	37K–66K (Glassdoor est.)	Requisition NumberRR-0001939\nRemote:Yes\nWe c...	4.1	Celerity\n4.1	New York, NY	McLean, VA
4	Reporting Data Analyst	37K–66K (Glassdoor est.)	ABOUT FANDUEL GROUP\n\nFanDuel Group is a worl...	3.9	FanDuel\n3.9	New York, NY	New York, NY

Quick view

```
In [ ]:
data_analyst_jobs.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 2253 entries, 0 to 2252
Data columns (total 13 columns):
#   Column                Non-Null Count  Dtype
---  -
0   Job Title              2253 non-null   object
1   Salary Estimate        2253 non-null   object
2   Job Description        2253 non-null   object
3   Rating                 2253 non-null   float64
4   Company Name           2252 non-null   object
5   Location               2253 non-null   object
6   Headquarters           2253 non-null   object
```

```
7   Size                2253 non-null    object
8   Type of ownership    2253 non-null    object
9   Industry             2253 non-null    object
10  Sector               2253 non-null    object
11  Revenue              2253 non-null    object
12  Easy Apply           2253 non-null    object
dtypes: float64(1), object(12)
memory usage: 228.9+ KB
```

Renaming Columns for Better Analysis

In []:

```
data_analyst_jobs.rename(columns={"Job Title": "job_title"},inplace=True)
```

In []:

```
data_analyst_jobs.rename(columns={"Salary Estimate":"salary_estimate"}, inplace=True)
data_analyst_jobs.rename(columns={"Job Description":"job_description"}, inplace=True)
data_analyst_jobs.rename(columns={"Company Name":"company_name"}, inplace=True)
data_analyst_jobs.rename(columns={"Location": "location"},inplace=True)
data_analyst_jobs.rename(columns={"Headquarters":"headquarters"}, inplace=True)
data_analyst_jobs.rename(columns={"Size": "size"},
inplace=True)
data_analyst_jobs.rename(columns={"Type of ownership":"type_of_ownership"}, inplace=True)
data_analyst_jobs.rename(columns={"Industry": "industry"},inplace=True)
data_analyst_jobs.rename(columns={"Sector": "sector"},inplace=True)
data_analyst_jobs.rename(columns={"Revenue": "revenue"},inplace=True)
data_analyst_jobs.rename(columns={"Easy Apply": "easy_apply"},inplace=True)
```

In []:

```
data_analyst_jobs.head()
```

Out[]:

	job_title	salary_estimate	job_description	Rating	company_name	location	headquarters
0	Data Analyst, Center on Immigration and Justic...	37K–66K (Glassdoor est.)	Are you eager to roll up your sleeves and harn...	3.2	Vera Institute of Justice\n3.2	New York, NY	New York,
1	Quality Data Analyst	37K–66K (Glassdoor est.)	Overview\n\nProvides analytical and technical ...	3.8	Visiting Nurse Service of New York\n3.8	New York, NY	New York,
2	Senior Data Analyst, Insights & Analytics Team...	37K–66K (Glassdoor est.)	We're looking for a Senior Data Analyst who ha...	3.4	Squarespace\n3.4	New York, NY	New York,
3	Data Analyst	37K–66K (Glassdoor est.)	Requisition NumberRR-0001939\n\nRemote:Yes\nWe c...	4.1	Celerity\n4.1	New York, NY	McLean,
4	Reporting Data Analyst	37K–66K (Glassdoor est.)	ABOUT FANDUEL GROUP\n\nFanDuel Group is a worl...	3.9	FanDuel\n3.9	New York, NY	New York,

Job Title

In []:

```
data_analyst_jobs['job_title'] =data_analyst_jobs['job_title'].replace(['Sr. Data Analys
```

In []:

```
data_analyst_jobs['job_title'] = data_analyst_jobs['job_title'].replace(['Data Analyst I
```

In []:

```
data_analyst_jobs['job_title'] =data_analyst_jobs['job_title'].replace(['Data Analyst II
```

In []:

```
# plot the most common types of jobs
to_plot = data_analyst_jobs.job_title.value_counts()[:5]
# ax = to_plot.plot(kind='bar',
color = sns.color_palette('Spectral')
to_plot
```

Out[]:

	count
job_title	
Data Analyst	405
Senior Data Analyst	120
Junior Data Analyst	58
Business Data Analyst	28
Data Quality Analyst	17

dtype: int64

Salary Estimate and Trends

In []:

```
## Changing Salary column to int for better calculation
data_analyst_jobs[['MinSalary', 'MaxSalary']] =data_analyst_jobs['salary_estimate'].str.
```

In []:

```
data_analyst_jobs['MinSalary'] =pd.to_numeric(data_analyst_jobs['MinSalary'])
data_analyst_jobs['MaxSalary'] =pd.to_numeric(data_analyst_jobs['MaxSalary'])
```

changing format to float

In []:

```
data_analyst_jobs['MinSalary'] =data_analyst_jobs['MinSalary'].astype(float)
data_analyst_jobs['MaxSalary'] =data_analyst_jobs['MaxSalary'].astype(float)
data_analyst_jobs['average_salary'] =(data_analyst_jobs['MaxSalary'] +
data_analyst_jobs['MinSalary']) / 2
#drop salary estimate(unuseful column)
data_analyst_jobs.drop(['salary_estimate', 'MinSalary','MaxSalary'], axis=1, inplace=True
```

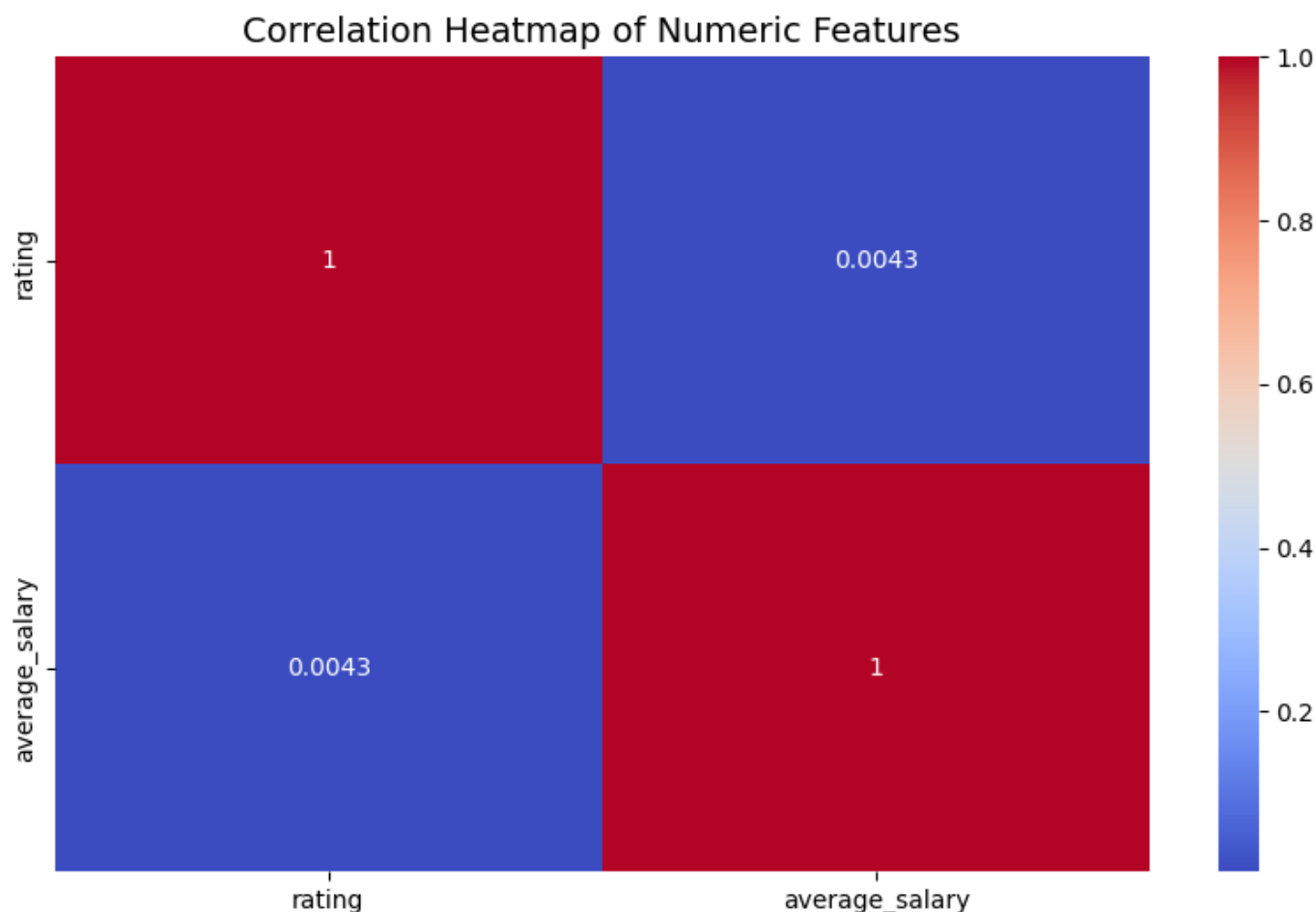
Graphs and Chart Visualizations

Correlation Heatmap of Numerical Features in Data Analyst Jobs Dataset

In []:

```
# Clean up column names (similar to your PDF renaming step)
data_analyst_jobs.columns = data_analyst_jobs.columns.str.strip().str.lower().str.replace(' ', '_')

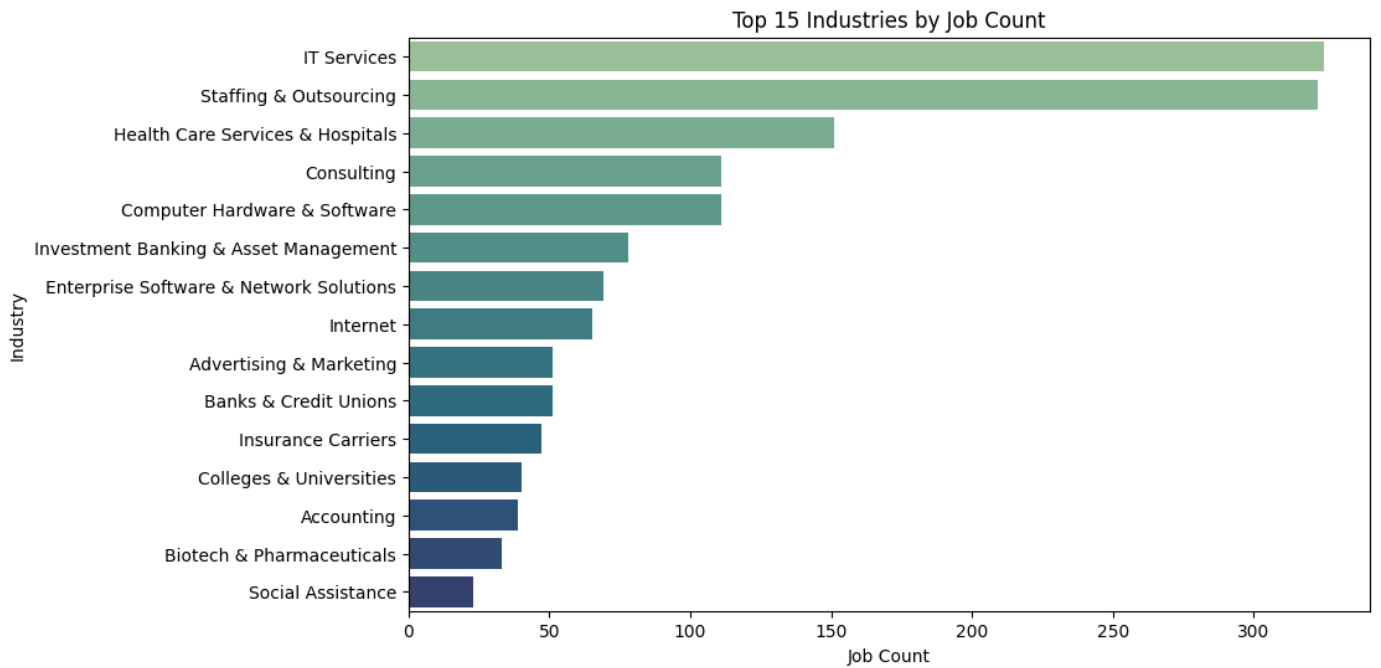
# --- 1. Correlation heatmap for numeric columns only ---
plt.figure(figsize=(10, 6))
sns.heatmap(data_analyst_jobs.select_dtypes(include='number').corr(), annot=True, cmap='magma')
plt.title("Correlation Heatmap of Numeric Features", fontsize=14)
plt.show()
```



Top 15 Industries Hiring Data Analysts

In []:

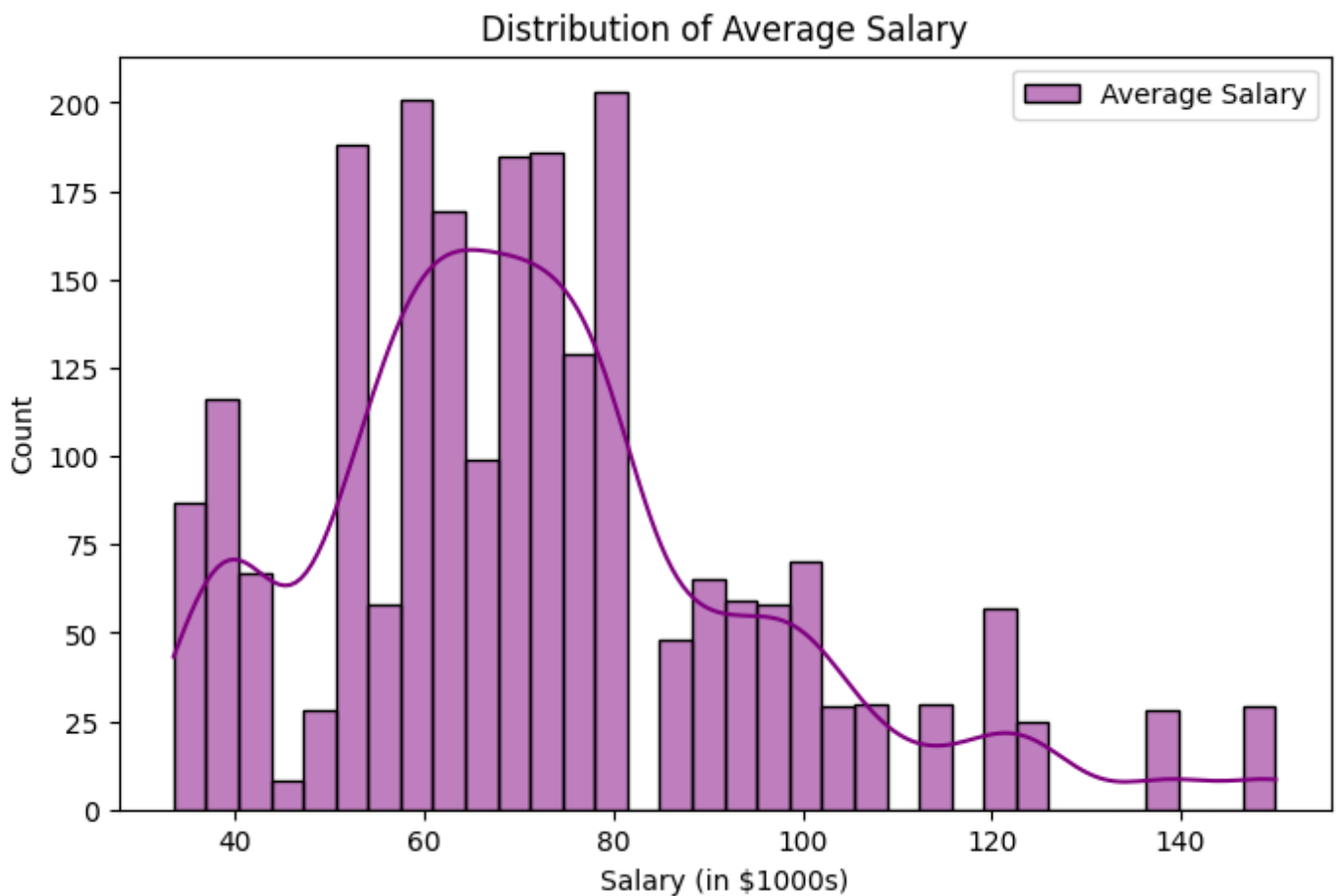
```
# --- 3. Top industries by job count ---
top_industries = data_analyst_jobs[data_analyst_jobs['industry'] != '-1']['industry'].value_counts().head(15)
plt.figure(figsize=(10, 6))
sns.barplot(y=top_industries.index, x=top_industries.values, palette='crest')
plt.title("Top 15 Industries by Job Count")
plt.xlabel("Job Count")
plt.ylabel("Industry")
plt.show()
```



Salary Distribution for Data Analyst Roles

In []:

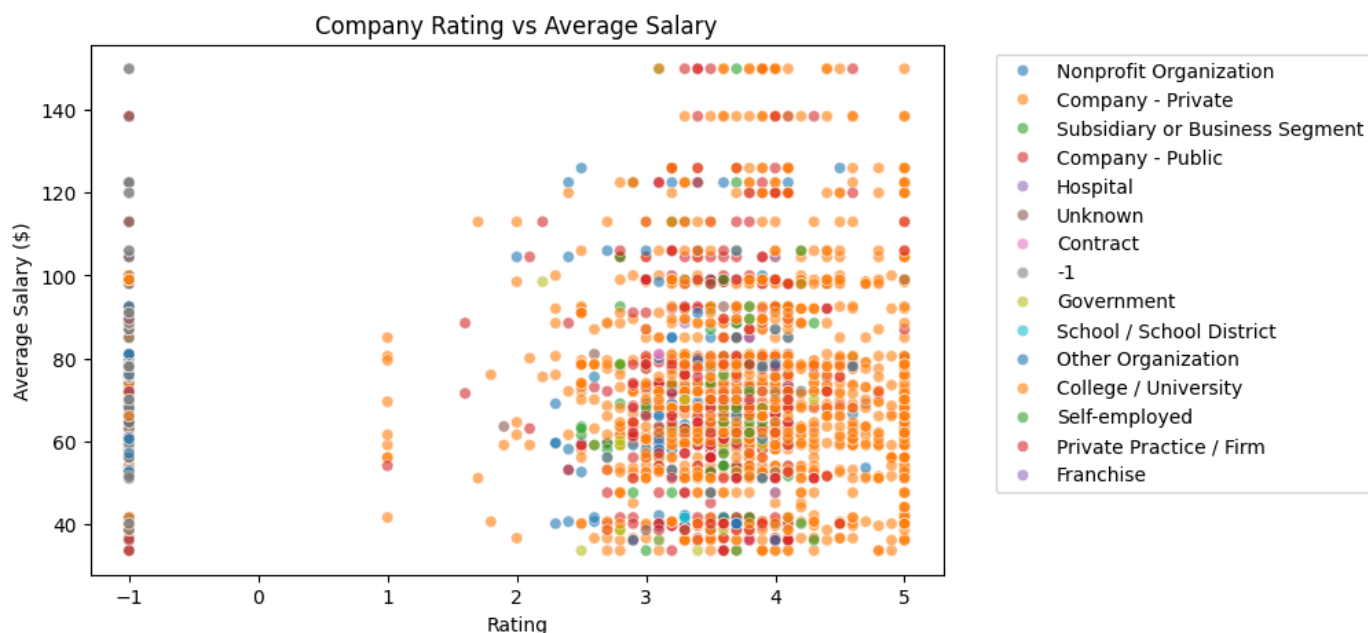
```
plt.figure(figsize=(8, 5))
sns.histplot(data_analyst_jobs['average_salary'], color='purple', label='Average Salary')
plt.legend()
plt.title('Distribution of Average Salary')
plt.xlabel('Salary (in $1000s)')
plt.show()
```



Relationship Between Company Rating and Average Salary by Ownership Type

In []:

```
# --- 5. Company rating vs Average Salary scatter plot ---
plt.figure(figsize=(8, 5))
sns.scatterplot(data=data_analyst_jobs, x='rating', y='average_salary', alpha=0.6, hue='')
plt.title('Company Rating vs Average Salary')
plt.xlabel('Rating')
plt.ylabel('Average Salary ($)')
plt.legend(bbox_to_anchor=(1.05, 1), loc='upper left')
plt.show()
```



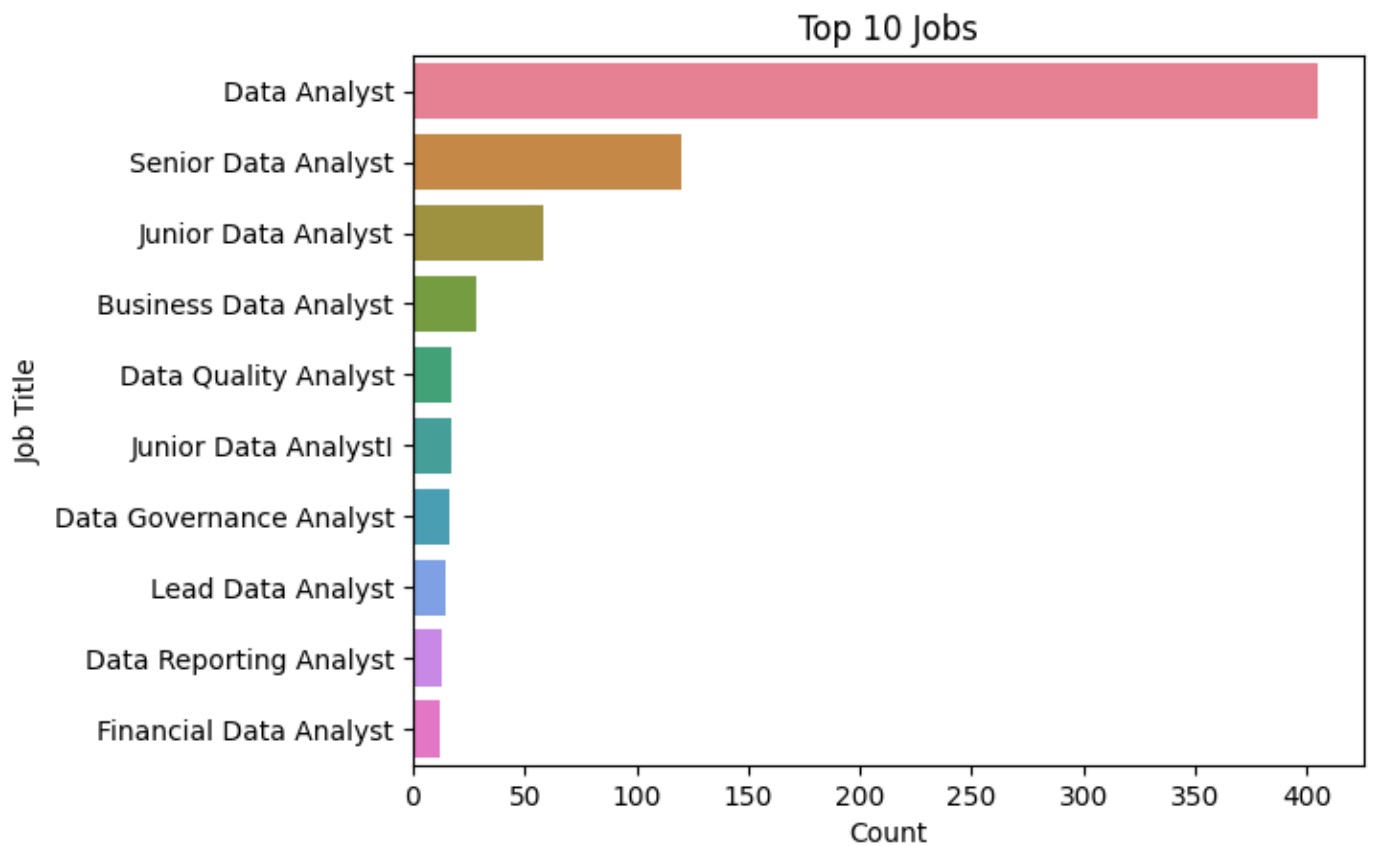
Top 10 Most Common Data Analyst Job Titles

In []:

```
top_jobs = data_analyst_jobs['job_title'].value_counts().head(10)

# Create the bar plot with custom colors
sns.barplot(
    x=top_jobs.values,
    y=top_jobs.index,
    palette=sns.color_palette("husl", len(top_jobs)) # "husl" gives different bright co
)

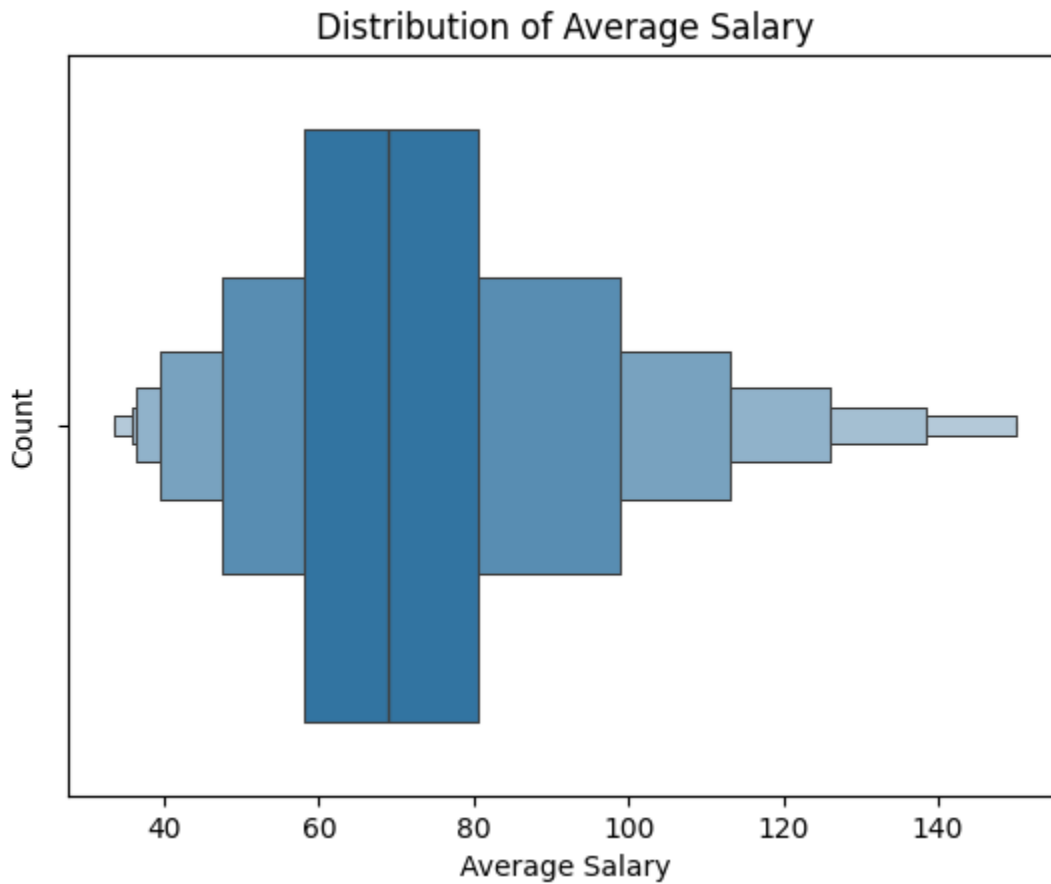
plt.xlabel('Count')
plt.ylabel('Job Title')
plt.title('Top 10 Jobs')
plt.show()
```



Average Salary

In []:

```
# Average Salary
sns.boxenplot(data=data_analyst_jobs, x='average_salary')
plt.xlabel('Average Salary')
plt.ylabel('Count')
plt.title('Distribution of Average Salary')
plt.show()
```



Top 10 Data Analyst Job Titles by Average Salary

In []:

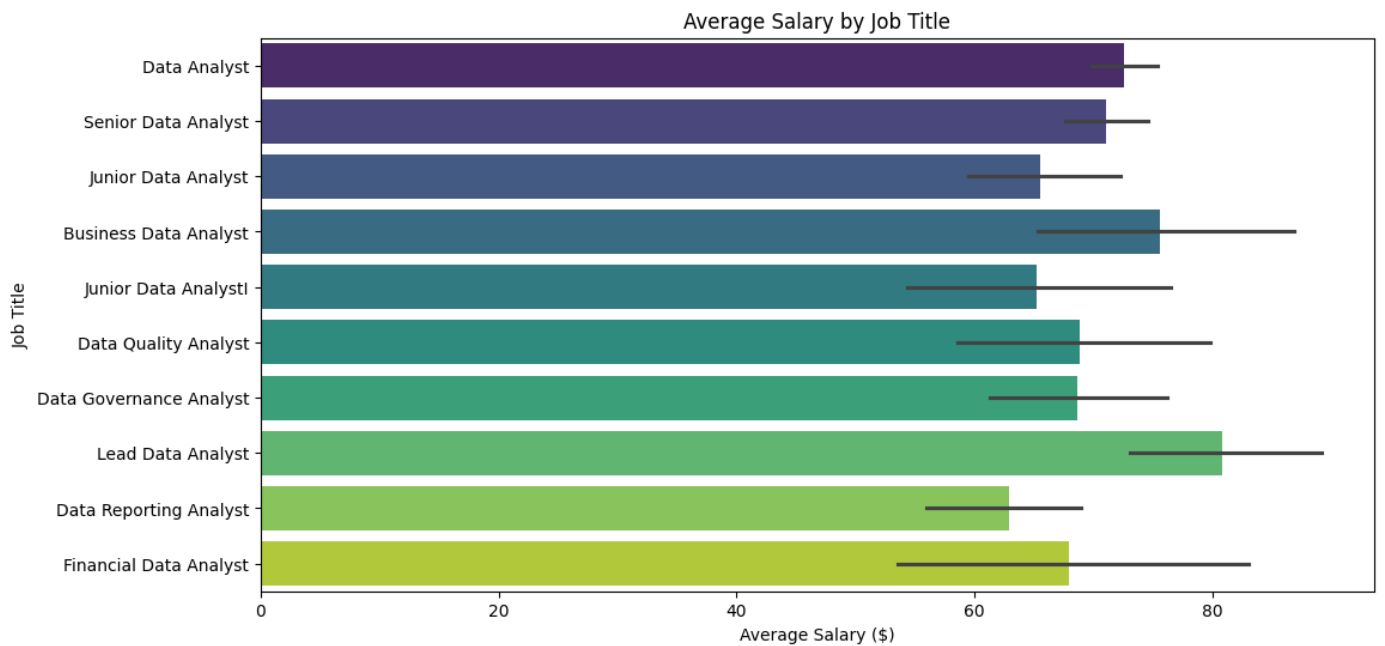
```
import seaborn as sns
import matplotlib.pyplot as plt

# Sort data by salary
data_analyst_jobs_sorted = data_analyst_jobs.sort_values(by='average_salary', ascending=

# Top 10 job titles by frequency
top_10_jobs = data_analyst_jobs_sorted['job_title'].value_counts().head(10).index

plt.figure(figsize=(12, 6))
sns.barplot(
    x='average_salary',
    y='job_title',
    data=data_analyst_jobs_sorted,
    orient='h',
    order=top_10_jobs,
    palette=sns.color_palette("viridis", len(top_10_jobs)) # nice gradient palette
)

plt.xlabel('Average Salary ($)')
plt.ylabel('Job Title')
plt.title('Average Salary by Job Title')
plt.show()
```

Top Locations Based on Average Salary

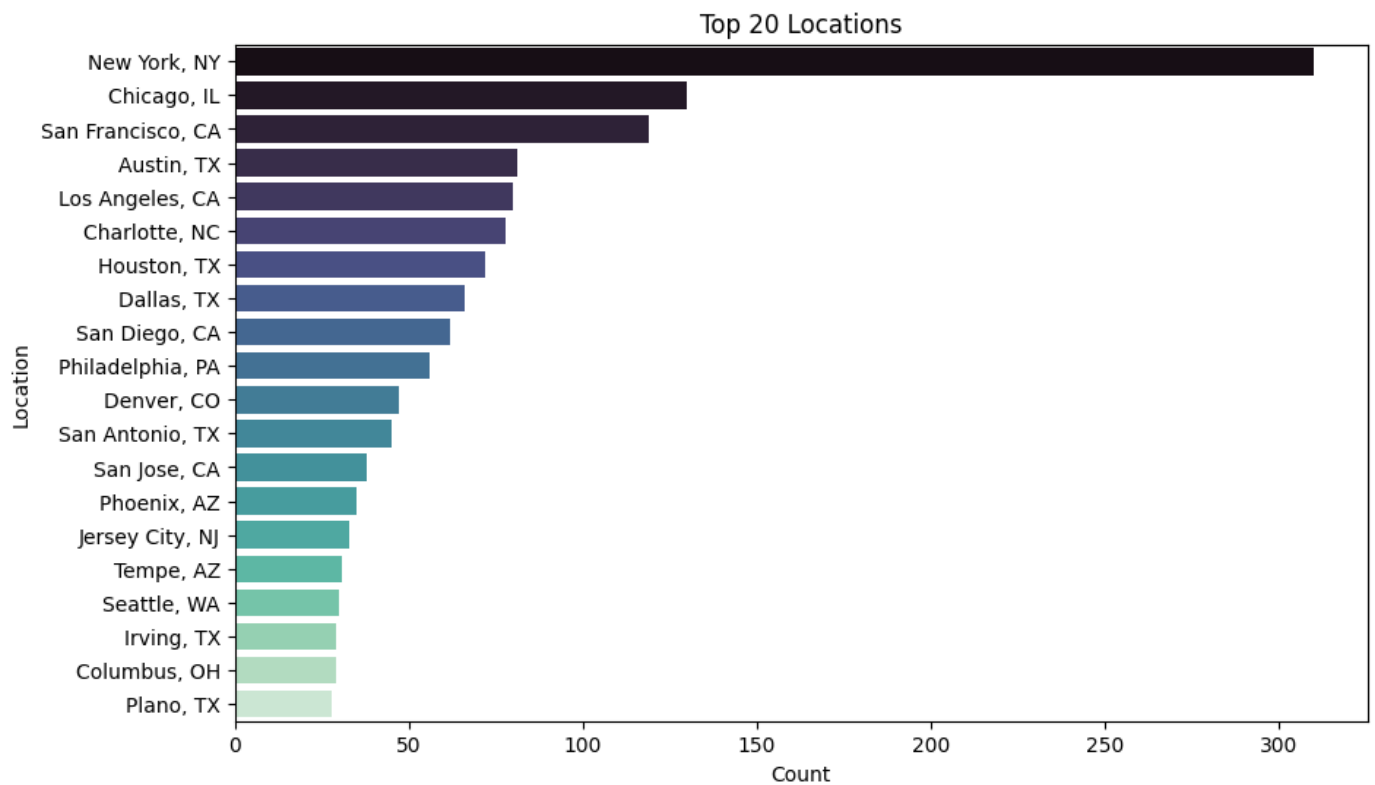
In []:

```
import seaborn as sns
import matplotlib.pyplot as plt

# Top work locations
top_locations = data_analyst_jobs['location'].value_counts().head(20)

plt.figure(figsize=(10, 6))
sns.barplot(
    x=top_locations.values,
    y=top_locations.index,
    palette=sns.color_palette("mako", len(top_locations)) # gradient colors
)

plt.xlabel('Count')
plt.ylabel('Location')
plt.title('Top 20 Locations')
plt.show()
```



Salary Trends by Location

In []:

```
job_location = data_analyst_jobs.groupby('location')['average_salary'].mean().reset_index
top_10 = job_location.sort_values(by = "average_salary", ascending=False).head(10)
```

In []:

```
fig = px.bar(top_10, x='average_salary', y='location', orientation='h', title='Salary Trends by Location')
fig.update_layout(xaxis_title='AVG Salary (USD)', yaxis_title='Location', showlegend = False)
fig.show()
```

Distribution of Company Sizes for Data Analyst Roles

In []:

```
import seaborn as sns
import matplotlib.pyplot as plt

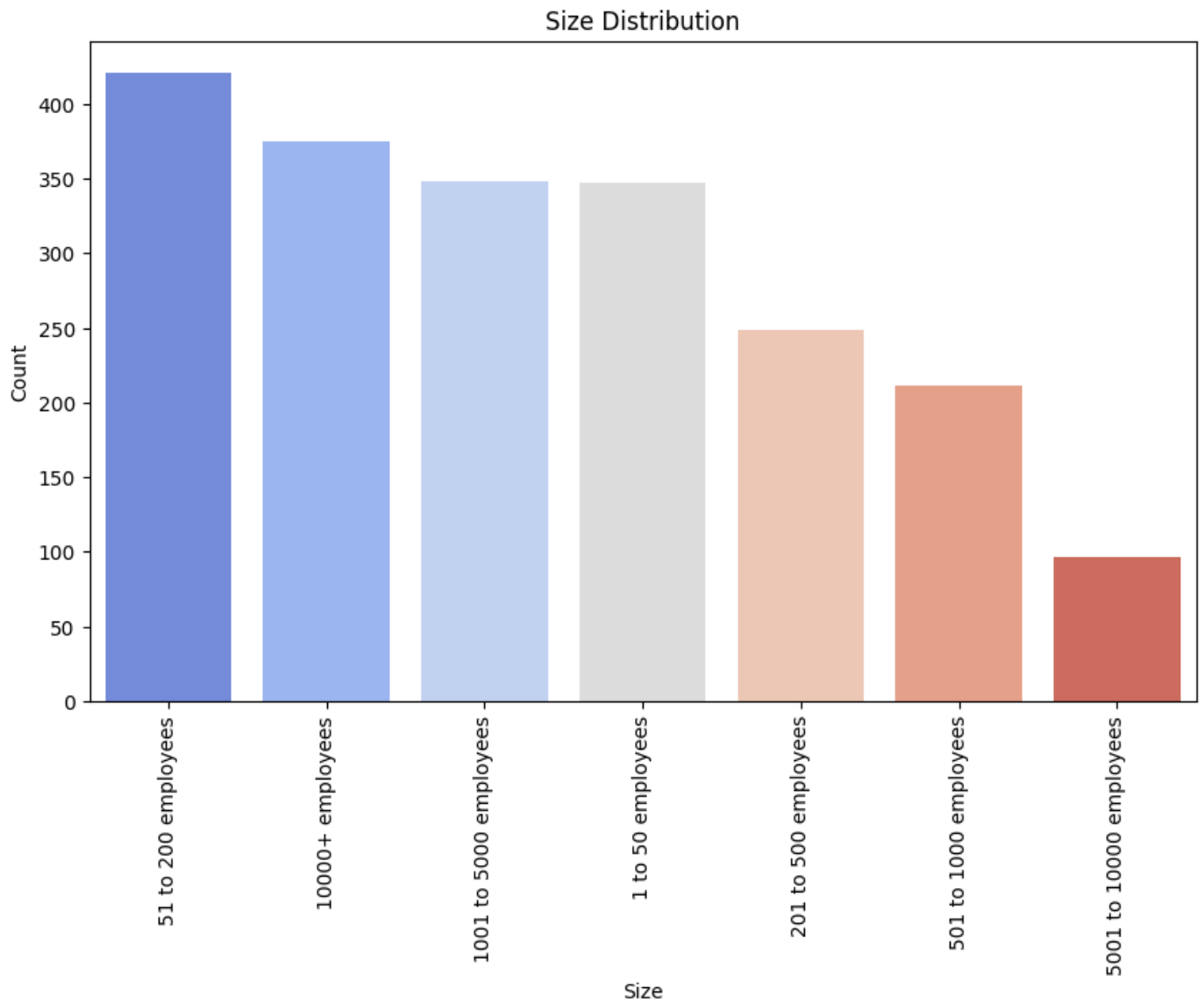
filtered_size = data_analyst_jobs[
    (data_analyst_jobs['size'] != '-1') &
    (data_analyst_jobs['size'] != 'Unknown')
]

data_analyst_jobs_size = filtered_size['size'].value_counts().head(20)

plt.figure(figsize=(10, 6))
sns.barplot(
    x=data_analyst_jobs_size.index,
    y=data_analyst_jobs_size.values,
    palette=sns.color_palette("coolwarm", len(data_analyst_jobs_size)) # gradient
)

plt.xlabel('Size')
plt.ylabel('Count')
```

```
plt.title('Size Distribution')
plt.xticks(rotation=90)
plt.show()
```



Average Salary by Company Size

```
In [ ]:
data_analyst_jobs_filtered = data_analyst_jobs[(data_analyst_jobs['size'] != '-1') & (data_analyst_jobs_sizeXsalary = data_analyst_jobs_filtered.groupby('size')['average_salary
```

```
In [ ]:
# Sort the DataFrame by 'AverageSalary' in descending order
data_analyst_jobs_sizeXsalary = data_analyst_jobs_sizeXsalary.sort_values(by='average_sal
```

```
In [ ]:
import seaborn as sns
import matplotlib.pyplot as plt
import numpy as np

# Normalize salary values to a 0-1 range for color mapping
norm = plt.Normalize(
    data_analyst_jobs_sizeXsalary['average_salary'].min(),
    data_analyst_jobs_sizeXsalary['average_salary'].max()
```

```

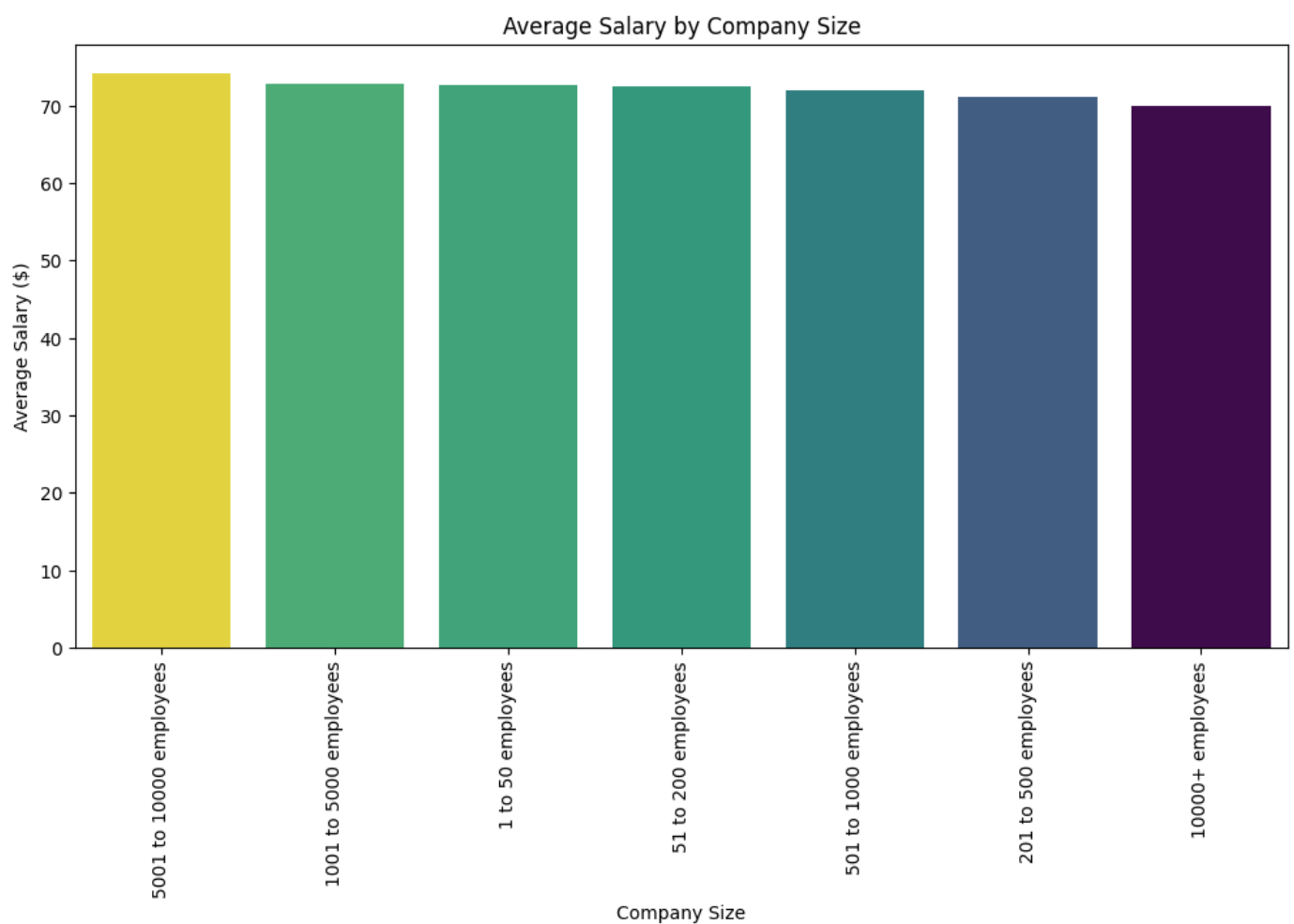
)

# Choose a colormap (e.g., viridis, plasma, coolwarm)
colors = plt.cm.viridis(norm(data_analyst_jobs_sizeXsalary['average_salary']))

plt.figure(figsize=(12, 6))
sns.barplot(
    x='size',
    y='average_salary',
    data=data_analyst_jobs_sizeXsalary,
    palette=colors
)

plt.xlabel('Company Size')
plt.ylabel('Average Salary ($)')
plt.title('Average Salary by Company Size')
plt.xticks(rotation=90)
plt.show()

```



Distribution of Company Rating

In []:

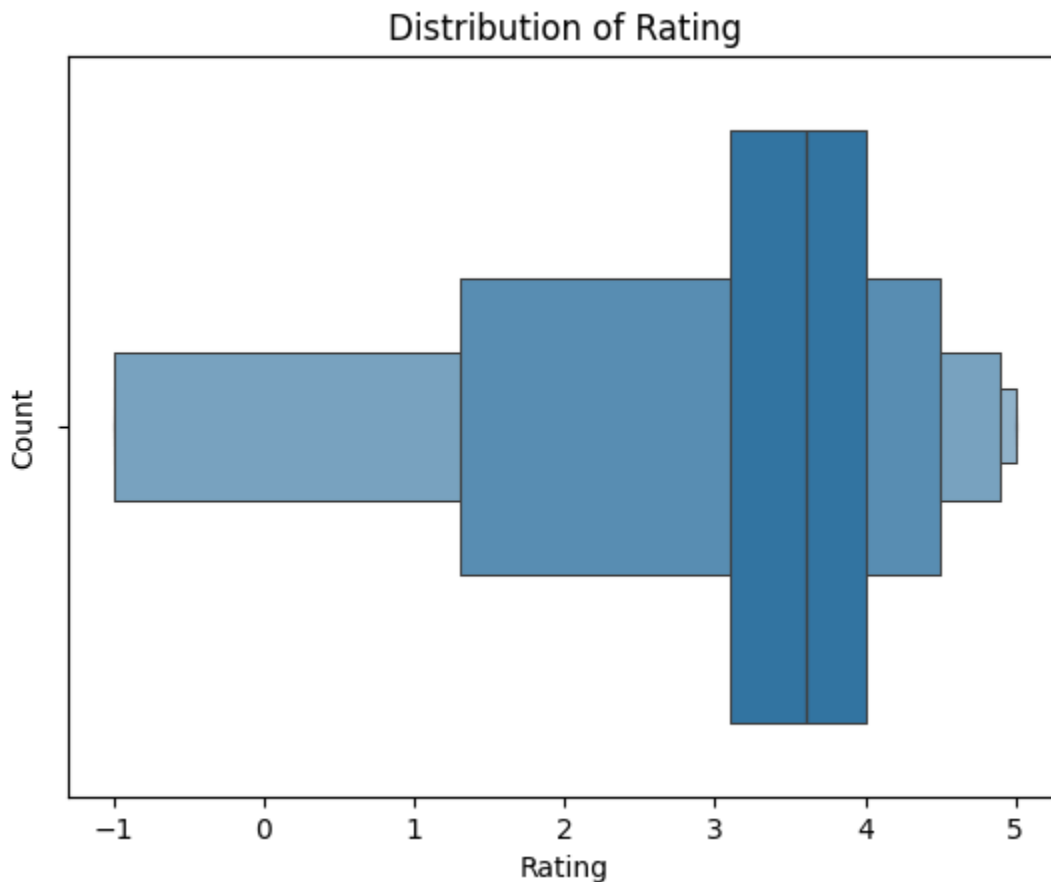
```

import matplotlib.pyplot as plt
import seaborn as sns

sns.boxenplot(data=data_analyst_jobs, x='rating')
plt.xlabel('Rating')
plt.ylabel('Count')

```

```
plt.title('Distribution of Rating')
plt.show()
```



Type of Ownership

In []:

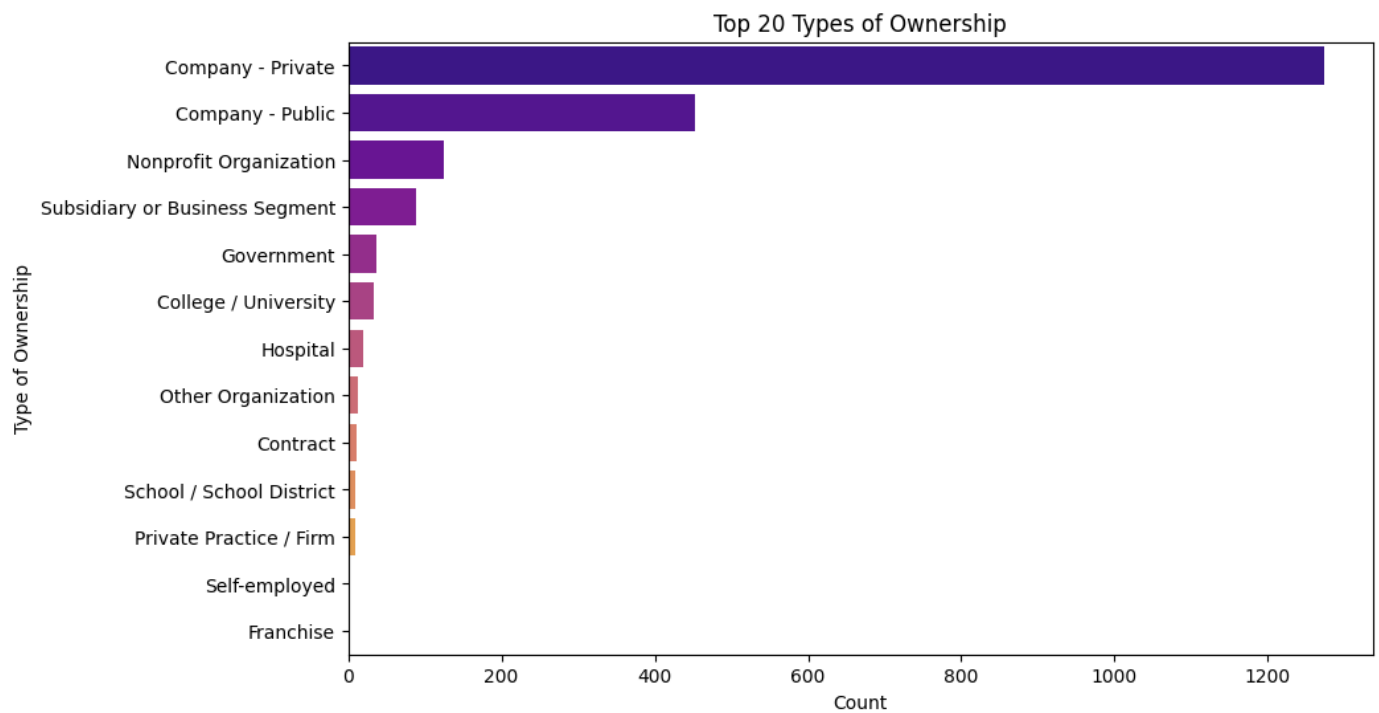
```
import seaborn as sns
import matplotlib.pyplot as plt

TOP = data_analyst_jobs[
    (data_analyst_jobs['type_of_ownership'] != '-1') &
    (data_analyst_jobs['type_of_ownership'] != 'Unknown')
]

TOP = TOP['type_of_ownership'].value_counts().head(20)

plt.figure(figsize=(10, 6))
sns.barplot(
    x=TOP.values,
    y=TOP.index,
    palette=sns.color_palette("plasma", len(TOP)) # colorful gradient
)

plt.xlabel('Count')
plt.ylabel('Type of Ownership')
plt.title('Top 20 Types of Ownership')
plt.show()
```



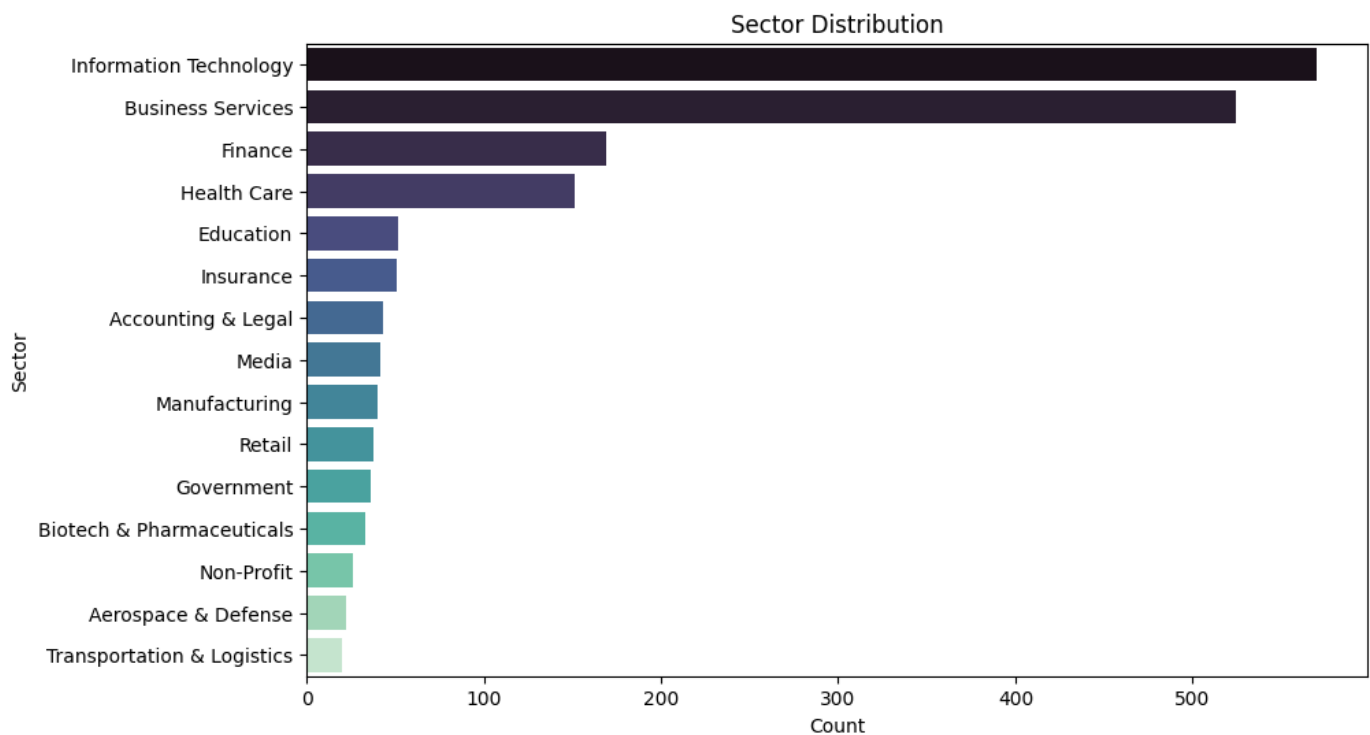
Top 15 Sectors Employing Data Analysts

In []:

```
data_analyst_jobs_sector = data_analyst_jobs[
    data_analyst_jobs['sector'] != '-1'
][['sector']].value_counts().head(15)

plt.figure(figsize=(10, 6))
sns.barplot(
    x=data_analyst_jobs_sector.values,
    y=data_analyst_jobs_sector.index,
    palette=sns.color_palette("mako", len(data_analyst_jobs_sector)) # gradient
)

plt.xlabel('Count')
plt.ylabel('Sector')
plt.title('Sector Distribution')
plt.show()
```



Average Salary by Sector

In []:

```
average_salary_by_sector = data_analyst_jobs[data_analyst_jobs['sector'] != '-1'].groupby(
average_salary_by_sector = average_salary_by_sector.sort_values(by='average_salary', ascen
```

In []:

```
# Normalize salary values for color mapping
norm = plt.Normalize(
    average_salary_by_sector['average_salary'].min(),
    average_salary_by_sector['average_salary'].max()
)
colors = plt.cm.viridis(norm(average_salary_by_sector['average_salary']))

plt.figure(figsize=(12, 6))
sns.barplot(
    x='sector',
    y='average_salary',
    data=average_salary_by_sector,
    palette=colors
)

plt.xticks(rotation=90)
plt.xlabel('Sector')
plt.ylabel('Average Salary (Thousands Dollars)')
plt.title('Average Salary by Sector')
plt.show()
```

Average Salary by Sector

