

A Complete End-to-End Analysis & Modeling Approach

Introduction:

Welcome to this comprehensive analysis of data science, AI, and ML job salaries in 2025! In this notebook, we'll dive deep into salary trends, build predictive models, and extract actionable insights for career planning in the tech industry.

Dataset overview: 133,349 entries with 11 columns spanning from 2020 to 2025, covering various job titles, experience levels, and geographical data.

Let's begin by importing the necessary libraries and setting up our environment.

```
In [7]: 1 # Import essential Libraries
2 import numpy as np
3 import pandas as pd
4 import matplotlib.pyplot as plt
5 import seaborn as sns
6 import plotly.express as px
7 import plotly.graph_objects as go
8 from plotly.subplots import make_subplots
9 import warnings
10
11 # ML Libraries
12 from sklearn.model_selection import train_test_split, cross_val_score, GridSearchCV
13 from sklearn.preprocessing import OneHotEncoder, StandardScaler, LabelEncoder
14 from sklearn.compose import ColumnTransformer
15 from sklearn.pipeline import Pipeline
16 from sklearn.ensemble import RandomForestRegressor, GradientBoostingRegressor
17 from sklearn.linear_model import LinearRegression, Ridge, Lasso
18 from sklearn.metrics import mean_squared_error, r2_score, mean_absolute_error
19 import xgboost as xgb
20 import lightgbm as lgb
21 import catboost as cb
22 from sklearn.svm import SVR
23
24 # For interactive visualizations
25 import plotly.io as pio
26 pio.templates.default = "plotly_white"
27
28 # Suppress warnings
29 warnings.filterwarnings('ignore')
30
31 # Set styling
32 sns.set(style="whitegrid")
33 plt.style.use('fivethirtyeight')
```

Data Loading and Initial Exploration

Let's load our dataset and take a first look:

```
In [11]: 1 pd.set_option("display.max_columns",None)
```

In [12]:

```
1 # Loading dataset
2 df=pd.read_csv("C:\\Users\\DILEEP V\\OneDrive\\Desktop\\Data_Science_Projects\\DataScience AI &
3 print(df)
```

	work_year	experience_level	employment_type	job_title	\
0	2025	SE	FT	Solutions Engineer	
1	2025	SE	FT	Solutions Engineer	
2	2025	MI	FT	Data Engineer	
3	2025	MI	FT	Data Engineer	
4	2025	EN	FT	Data Engineer	
...
133344	2020	SE	FT	Data Scientist	
133345	2021	MI	FT	Principal Data Scientist	
133346	2020	EN	FT	Data Scientist	
133347	2020	EN	CT	Business Data Analyst	
133348	2021	SE	FT	Data Scientist	
	salary	salary_currency	salary_in_usd	employee_residence	\
0	214000	USD	214000	US	
1	136000	USD	136000	US	
2	158800	USD	158800	AU	
3	139200	USD	139200	AU	
4	90000	USD	90000	US	
...
133344	412000	USD	412000	US	
133345	151000	USD	151000	US	
133346	105000	USD	105000	US	
133347	100000	USD	100000	US	
133348	7000000	INR	94665	IN	
	remote_ratio	company_location	company_size		
0	100	US	M		
1	100	US	M		
2	0	AU	M		
3	0	AU	M		
4	0	US	M		
...
133344	100	US	L		
133345	100	US	L		
133346	100	US	S		
133347	100	US	L		
133348	50	IN	L		

[133349 rows x 11 columns]

In [13]:

```

1 # Initial overview
2 print(f"Dataset Shape: {df.shape}")
3 print(f"\nMemory Usage: {df.memory_usage().sum() / 1024**2:.2f} MB")
4 print("\nFirst few rows:")
5 df.head()

```

Dataset Shape: (133349, 11)

Memory Usage: 11.19 MB

First few rows:

Out[13]:

	work_year	experience_level	employment_type	job_title	salary	salary_currency	salary_in_usd	employee_residence	re
0	2025	SE	FT	Solutions Engineer	214000	USD	214000		US
1	2025	SE	FT	Solutions Engineer	136000	USD	136000		US
2	2025	MI	FT	Data Engineer	158800	USD	158800		AU
3	2025	MI	FT	Data Engineer	139200	USD	139200		AU
4	2025	EN	FT	Data Engineer	90000	USD	90000		US

Dataset Information

In [14]:

```

1 # Display information about the dataset
2 df.info()

```

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 133349 entries, 0 to 133348
Data columns (total 11 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   work_year        133349 non-null  int64  
 1   experience_level 133349 non-null  object  
 2   employment_type   133349 non-null  object  
 3   job_title         133349 non-null  object  
 4   salary            133349 non-null  int64  
 5   salary_currency   133349 non-null  object  
 6   salary_in_usd    133349 non-null  int64  
 7   employee_residence 133349 non-null  object  
 8   remote_ratio      133349 non-null  int64  
 9   company_location  133349 non-null  object  
 10  company_size      133349 non-null  object  
dtypes: int64(4), object(7)
memory usage: 11.2+ MB

```

```
In [20]: 1 # Statistical summary of numerical columns
2 df.describe().T.style.background_gradient(cmap='Blues')
```

Out[20]:

	count	mean	std	min	25%	50%	75%
work_year	133349.000000	2024.358705	0.680788	2020.000000	2024.000000	2024.000000	2025.000000
salary	133349.000000	163283.322590	217386.021574	14000.000000	106020.000000	147000.000000	199000.000000
salary_in_usd	133349.000000	157617.272098	74288.363097	15000.000000	106000.000000	146206.000000	198000.000000
remote_ratio	133349.000000	20.905669	40.590044	0.000000	0.000000	0.000000	0.000000

Check for Missing Values

```
In [21]: 1 # Check for missing values
2 missing_values = df.isnull().sum()
3 missing_percent = (missing_values / len(df)) * 100
4
5 missing_df = pd.DataFrame({
6     'Missing Values': missing_values,
7     'Percentage': missing_percent
8 })
9
10 print("Missing values check:")
11 missing_df[missing_df['Missing Values'] > 0]
```

Missing values check:

Out[21]:

Missing Values	Percentage
----------------	------------

Great! We have a clean dataset with no missing values. Let's proceed with our exploration

Exploratory Data Analysis

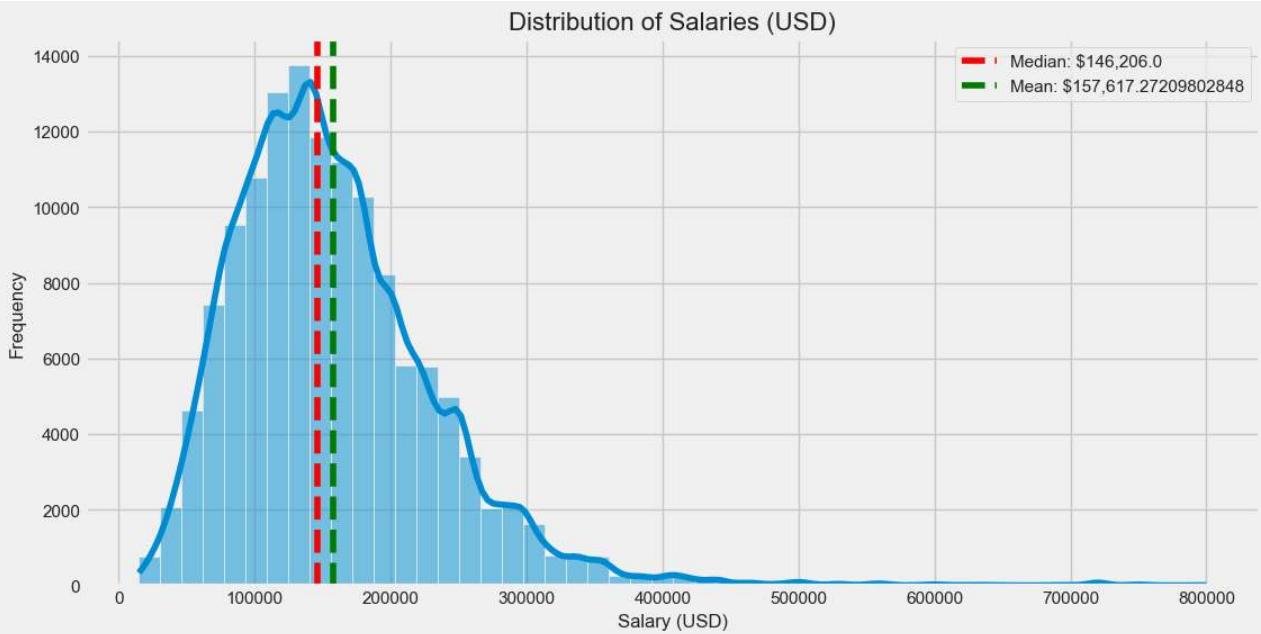
1. Salary Distribution

In [23]:

```

1 # Distribution of salaries in USD
2 plt.figure(figsize=(12, 6))
3 sns.histplot(df['salary_in_usd'], kde=True, bins=50)
4 plt.title('Distribution of Salaries (USD)', fontsize=16)
5 plt.xlabel('Salary (USD)', fontsize=12)
6 plt.ylabel('Frequency', fontsize=12)
7 plt.axvline(df['salary_in_usd'].median(), color='red', linestyle='--', label=f'Median: ${df["sa}
8 plt.axvline(df['salary_in_usd'].mean(), color='green', linestyle='--', label=f'Mean: ${df["sal}
9 plt.legend()
10 plt.show()
11
12
13 # Log-transformed salary distribution (to handle skewness)
14 plt.figure(figsize=(12, 6))
15 sns.histplot(np.log1p(df['salary_in_usd']), kde=True, bins=50)
16 plt.title('Log-Transformed Salary Distribution', fontsize=16)
17 plt.xlabel('Log(Salary+1)', fontsize=12)
18 plt.ylabel('Frequency', fontsize=12)
19 plt.show()

```

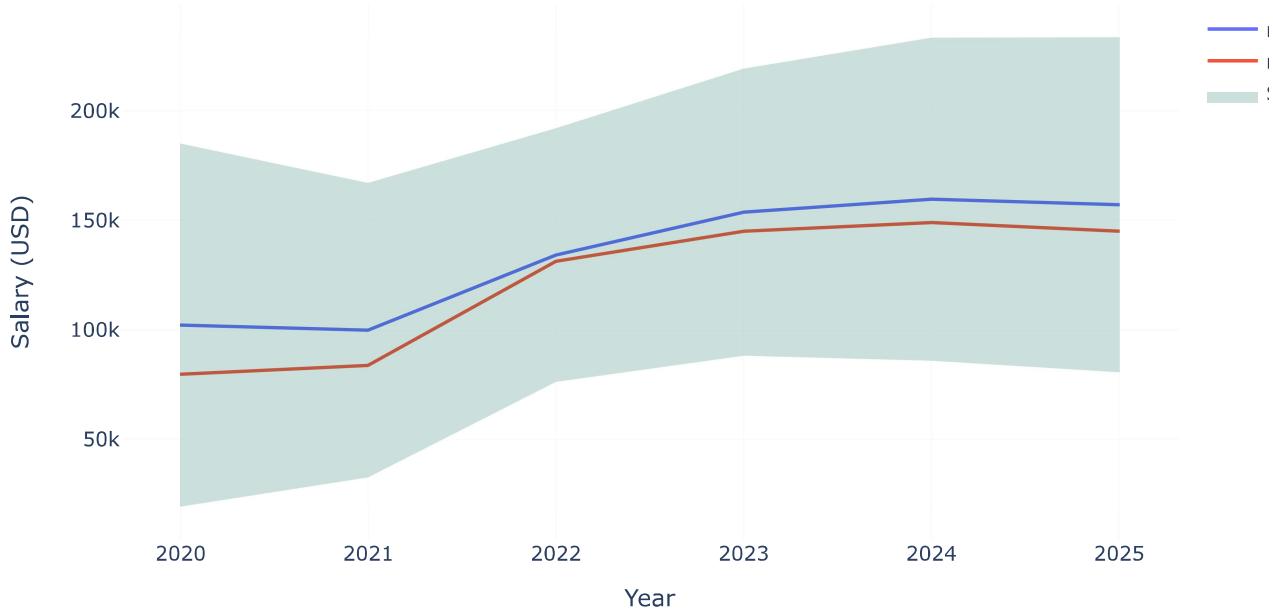


The salary distribution is right-skewed, with a significant number of high-paying outliers. This is typical for salary data. The log transformation gives us a more normal distribution, which will be useful for modeling.

2. Salary Trends Over Time

```
In [24]: 1 # Yearly salary trends
2 yearly_stats = df.groupby('work_year')['salary_in_usd'].agg(['mean', 'median', 'std']).reset_index()
3
4 fig = px.line(yearly_stats, x='work_year', y=['mean', 'median'],
5                 title='Data Science Salary Trends (2020-2025)',
6                 labels={'value': 'Salary (USD)', 'work_year': 'Year', 'variable': 'Metric'},
7                 template='plotly_white')
8
9 fig.update_layout(legend_title_text='', hovermode='x unified',
10                     width=900, height=500)
11
12 # Add range for standard deviation
13 fig.add_trace(go.Scatter(
14     x=np.concatenate([yearly_stats['work_year'], yearly_stats['work_year'][::-1]]),
15     y=np.concatenate([yearly_stats['mean'] + yearly_stats['std'],
16                      (yearly_stats['mean'] - yearly_stats['std'])[::-1]]),
17     fill='toself',
18     fillcolor='rgba(0,100,80,0.2)',
19     line=dict(color='rgba(255,255,255,0)'),
20     name='Standard Deviation'
21 ))
22
23 fig.show()
```

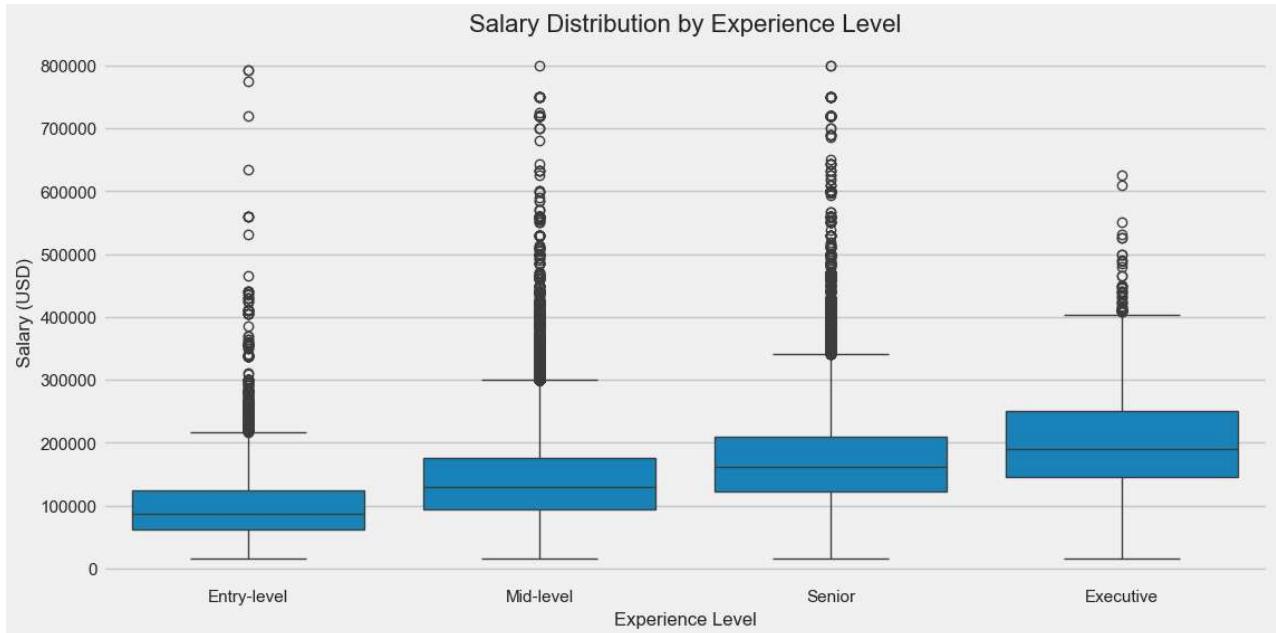
Data Science Salary Trends (2020-2025)

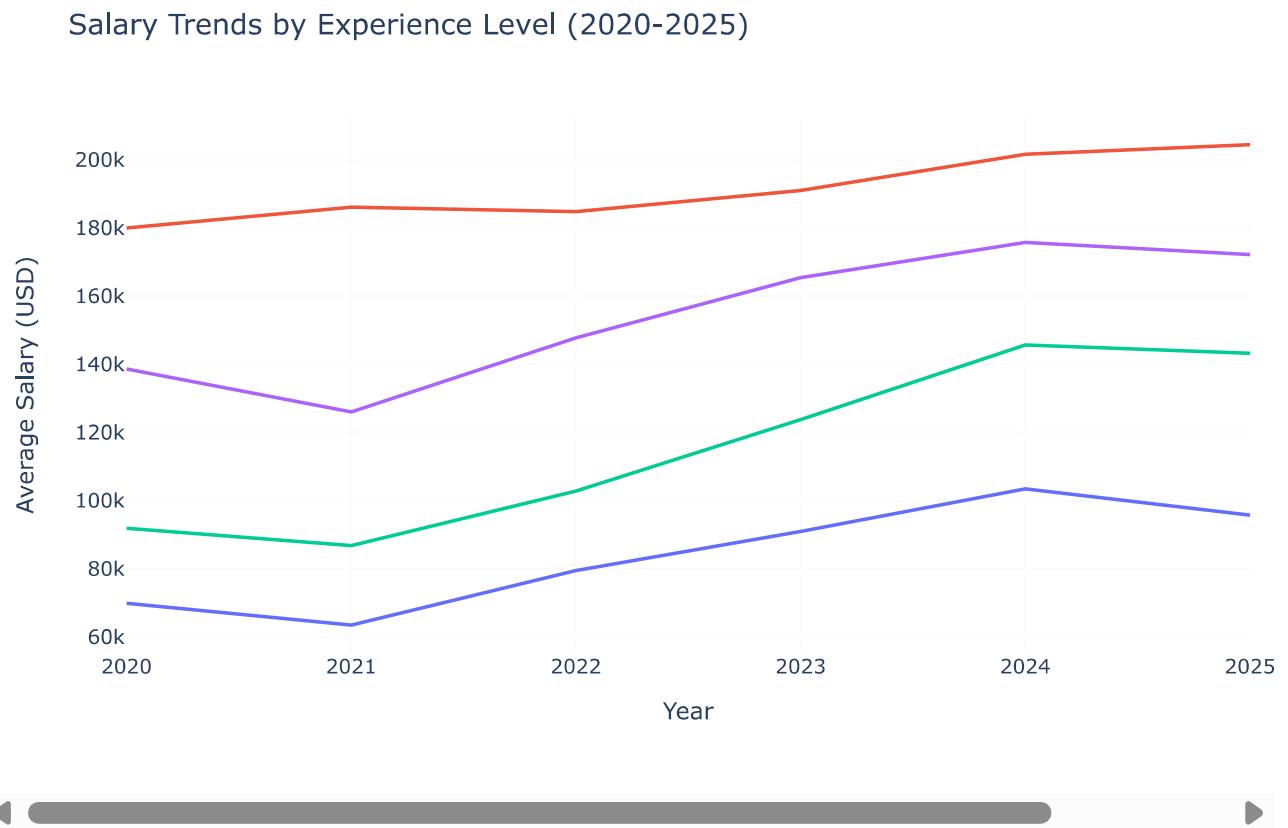


We can observe a steady increase in both mean and median salaries from 2020 to 2025, with the gap between them widening slightly, indicating increasing inequality in the field.

3. Experience Level Analysis

```
In [25]: 1 # Experience Level salary comparison
2 plt.figure(figsize=(12, 6))
3 sns.boxplot(x='experience_level', y='salary_in_usd', data=df, order=['EN', 'MI', 'SE', 'EX'])
4 plt.title('Salary Distribution by Experience Level', fontsize=16)
5 plt.xlabel('Experience Level', fontsize=12)
6 plt.ylabel('Salary (USD)', fontsize=12)
7 plt.xticks(ticks=[0, 1, 2, 3], labels=['Entry-level', 'Mid-level', 'Senior', 'Executive'])
8 plt.show()
9
10 # Mean salary by experience level over time
11 exp_time = df.groupby(['work_year', 'experience_level'])['salary_in_usd'].mean().reset_index()
12 exp_time['experience_level'] = exp_time['experience_level'].replace({
13     'EN': 'Entry-level', 'MI': 'Mid-level', 'SE': 'Senior', 'EX': 'Executive'
14 })
15
16 fig = px.line(exp_time, x='work_year', y='salary_in_usd', color='experience_level',
17                 title='Salary Trends by Experience Level (2020-2025)',
18                 labels={'salary_in_usd': 'Average Salary (USD)', 'work_year': 'Year'},
19                 template='plotly_white')
20
21 fig.update_layout(width=900, height=500, hovermode='x unified')
22 fig.show()
```





As expected, experience level has a significant impact on salary, with executives earning substantially more than other levels. The growth rate for executive salaries appears to be steeper, indicating increasing premiums for leadership roles.

4. Job Title Analysis

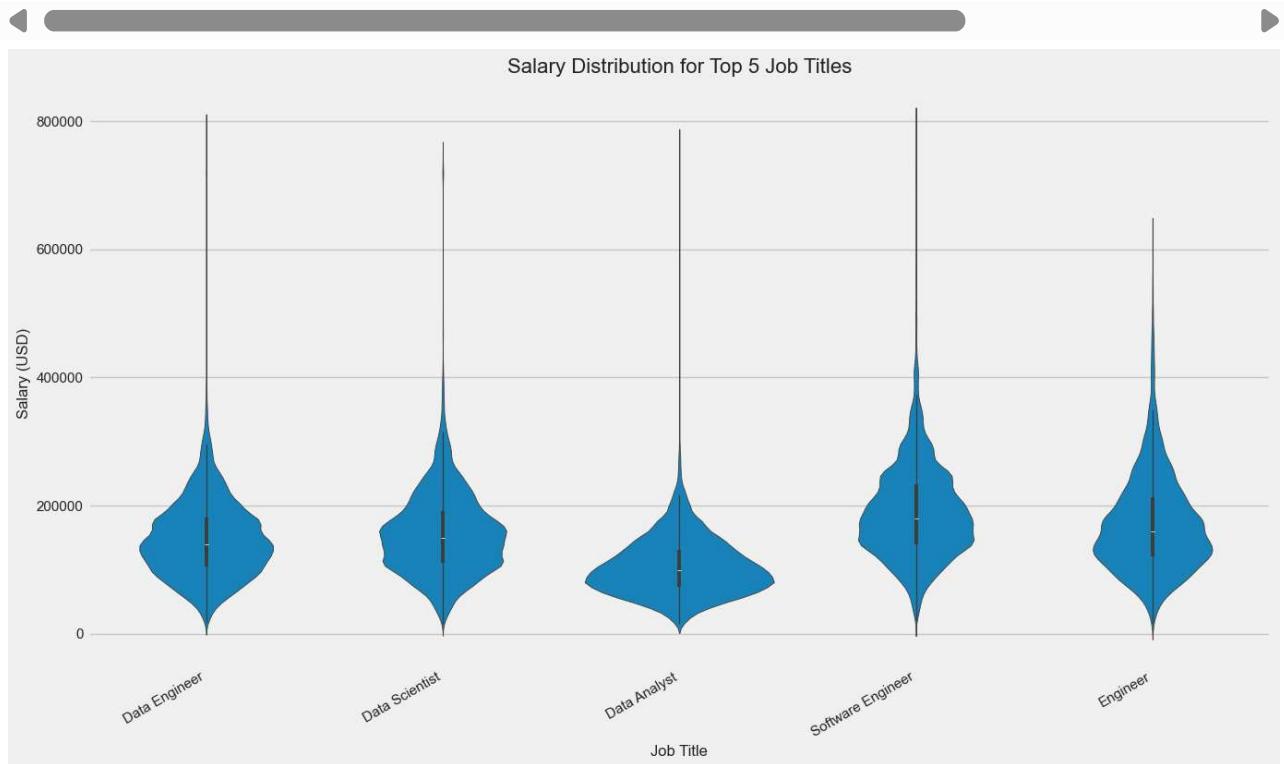
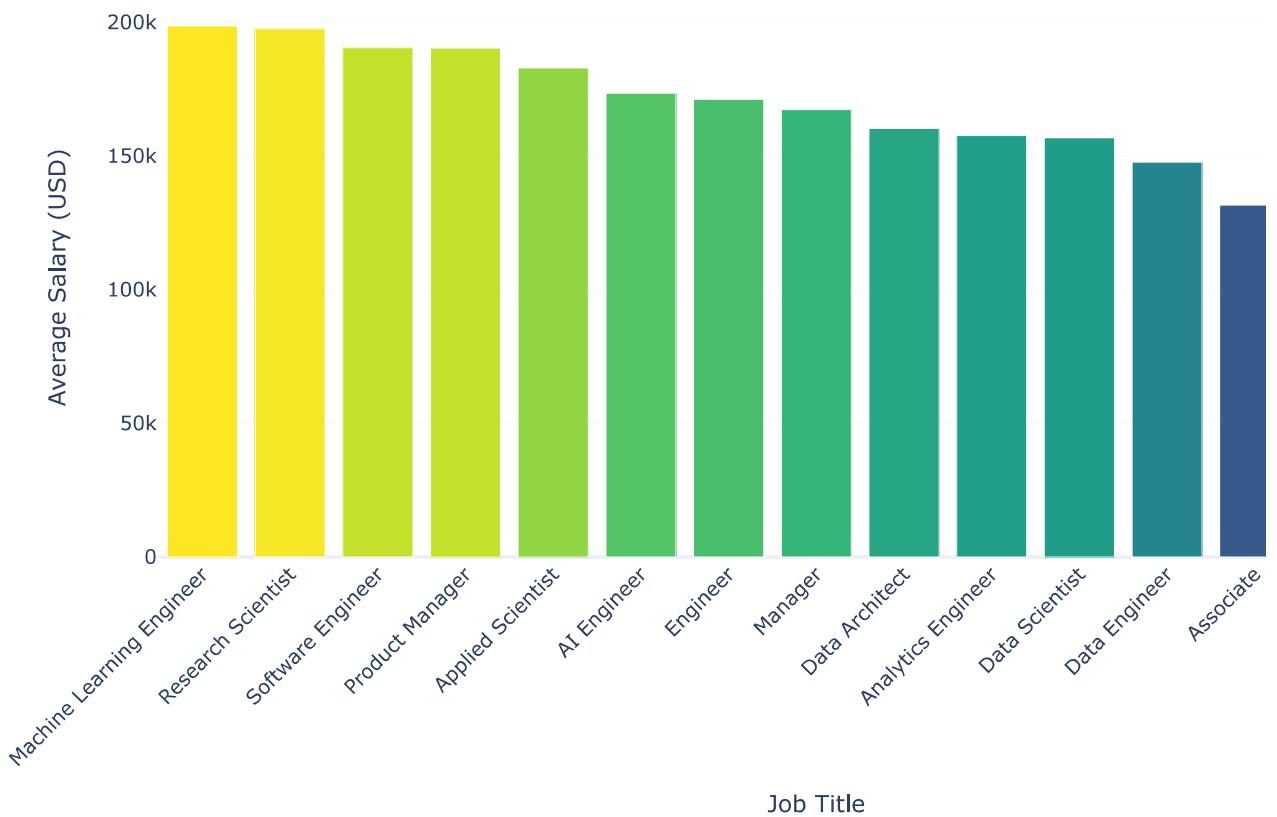
```
In [27]: 1 # Top 15 job titles by count
2 top_jobs = df['job_title'].value_counts().head(15)
3 print("Distribution of top job titles:")
4 print(top_jobs)
5
6 # Mean salary for top job titles
7 top_jobs_salary = df[df['job_title'].isin(top_jobs.index)].groupby('job_title')['salary_in_usd']
8
9 fig = px.bar(x=top_jobs_salary.index, y=top_jobs_salary.values,
10               labels={'x': 'Job Title', 'y': 'Average Salary (USD)'},
11               title='Average Salary by Top Job Titles',
12               color=top_jobs_salary.values, color_continuous_scale='Viridis')
13
14 fig.update_layout(xaxis_tickangle=-45, width=1000, height=600)
15 fig.show()
16
17 # Salary distribution for top 5 job titles
18 plt.figure(figsize=(14, 8))
19 top5_jobs = top_jobs.index[:5]
20 sns.violinplot(x='job_title', y='salary_in_usd', data=df[df['job_title'].isin(top5_jobs)])
21 plt.title('Salary Distribution for Top 5 Job Titles', fontsize=16)
22 plt.xlabel('Job Title', fontsize=12)
23 plt.ylabel('Salary (USD)', fontsize=12)
24 plt.xticks(rotation=30, ha='right')
25 plt.tight_layout()
26 plt.show()
```

Distribution of top job titles:

job_title	count
Data Scientist	17314
Software Engineer	15007
Data Engineer	14868
Data Analyst	12381
Engineer	9456
Machine Learning Engineer	8205
Manager	6679
Analyst	4364
Research Scientist	3202
Product Manager	2230
Applied Scientist	2205
Associate	2073
Data Architect	2024
Analytics Engineer	1895
AI Engineer	1620

Name: count, dtype: int64

Average Salary by Top Job Titles

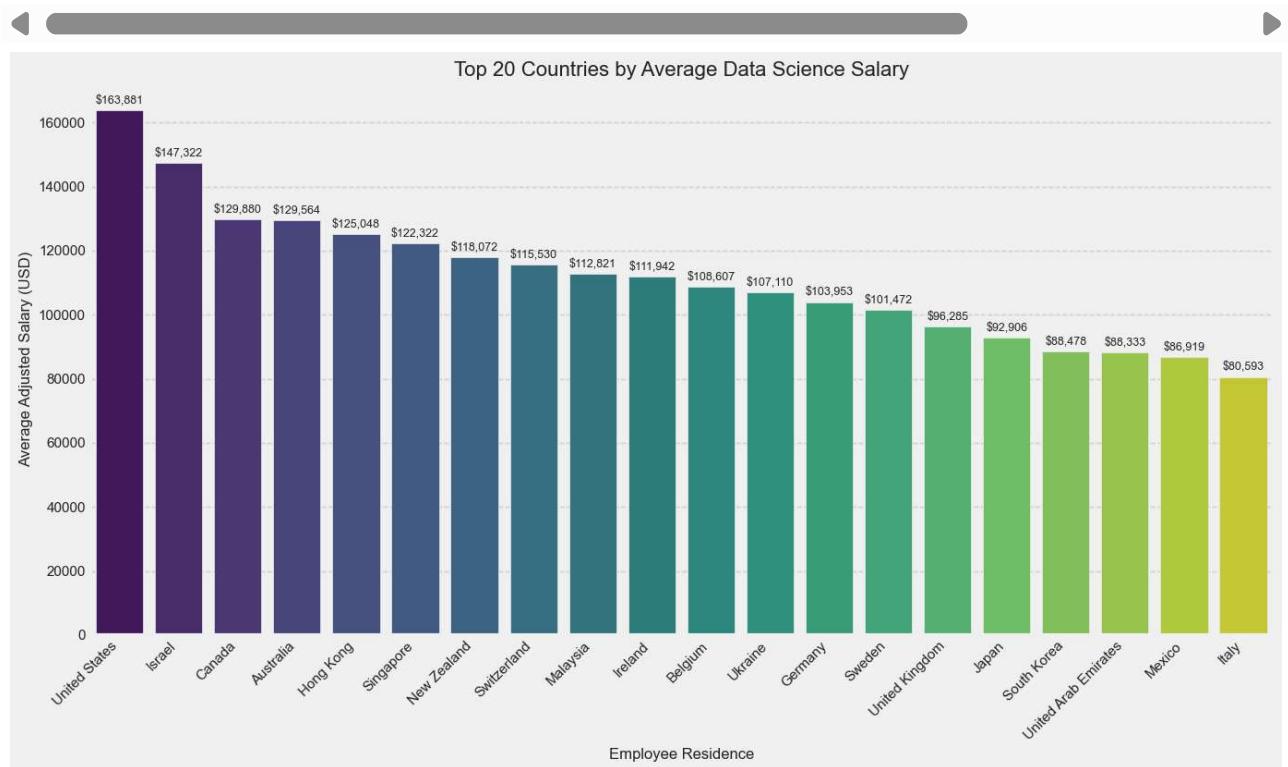
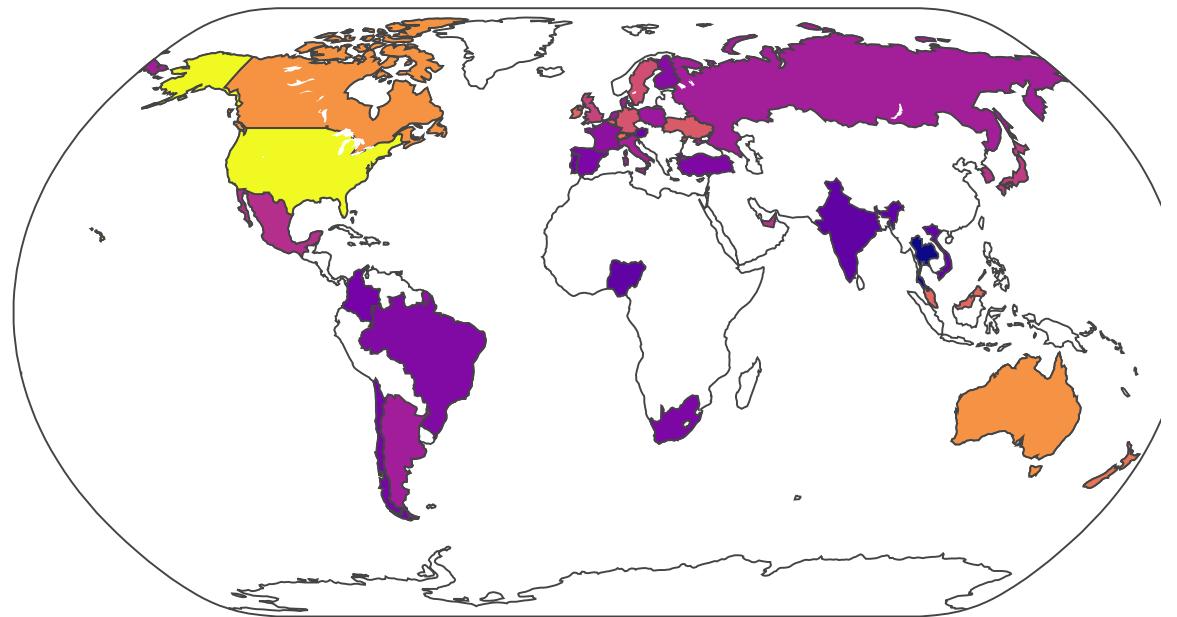


Among the most common job titles, Machine Learning Engineers and Research Scientists tend to command the highest salaries, while Data Analysts typically earn less. Software Engineers show the widest salary distribution, reflecting the diverse range of roles under this title..

5. Geographical Analysis

```
In [29]: 1 # First, let's create adjusted_salary column (assuming it's based on salary_in_usd)
2 df['adjusted_salary'] = df['salary_in_usd'] # You can modify this if there's a specific adjustment
3
4 # Create mapping from ISO-2 to country names for visualization
5 iso2_to_name = {
6     'US': 'United States', 'GB': 'United Kingdom', 'DE': 'Germany', 'FR': 'France',
7     'CA': 'Canada', 'IN': 'India', 'AU': 'Australia', 'ES': 'Spain', 'BR': 'Brazil',
8     'NL': 'Netherlands', 'JP': 'Japan', 'CH': 'Switzerland', 'IT': 'Italy',
9     'SG': 'Singapore', 'SE': 'Sweden', 'MX': 'Mexico', 'FI': 'Finland', 'DK': 'Denmark',
10    'PL': 'Poland', 'PT': 'Portugal', 'NZ': 'New Zealand', 'IE': 'Ireland',
11    'HK': 'Hong Kong', 'RU': 'Russia', 'BE': 'Belgium', 'IL': 'Israel',
12    'UA': 'Ukraine', 'TR': 'Turkey', 'AE': 'United Arab Emirates', 'ZA': 'South Africa',
13    'CO': 'Colombia', 'AR': 'Argentina', 'CL': 'Chile', 'AT': 'Austria', 'MY': 'Malaysia',
14    'NG': 'Nigeria', 'VN': 'Vietnam', 'KR': 'South Korea', 'TH': 'Thailand'
15    # Add more mappings as needed
16 }
17
18 # Calculate average salary by employee residence
19 avg_salary_by_residence = df.groupby('employee_residence')['adjusted_salary'].mean().reset_index()
20
21 # Add country names for mapping
22 avg_salary_by_residence['country_name'] = avg_salary_by_residence['employee_residence'].map(iso2_to_name)
23 avg_salary_by_residence = avg_salary_by_residence.dropna(subset=['country_name']) # Drop unmapped entries
24
25 # Create the choropleth map
26 fig2 = px.choropleth(avg_salary_by_residence,
27                       locations='country_name',
28                       locationmode='country names',
29                       color='adjusted_salary',
30                       hover_name='country_name',
31                       hover_data={'employee_residence': True, 'adjusted_salary': ':,.0f'},
32                       color_continuous_scale=px.colors.sequential.Plasma,
33                       title='Average Salary by Employee Residence',
34                       labels={'adjusted_salary': 'Average Adjusted Salary'},
35                       projection='natural earth')
36
37 fig2.update_layout(width=1000, height=600)
38 fig2.show()
39
40 # Create a sorted bar chart for top 20 countries
41 top_countries = avg_salary_by_residence.sort_values('adjusted_salary', ascending=False).head(20)
42
43 plt.figure(figsize=(14, 8))
44 chart = sns.barplot(x='country_name', y='adjusted_salary', data=top_countries,
45                      palette='viridis', order=top_countries['country_name'])
46 plt.title('Top 20 Countries by Average Data Science Salary', fontsize=16)
47 plt.xlabel('Employee Residence', fontsize=12)
48 plt.ylabel('Average Adjusted Salary (USD)', fontsize=12)
49 plt.xticks(rotation=45, ha='right')
50 plt.grid(axis='y', linestyle='--', alpha=0.7)
51
52 # Add salary values on top of bars
53 for i, bar in enumerate(chart.patches):
54     chart.text(bar.get_x() + bar.get_width()/2.,
55                bar.get_height() + 2000,
56                f'${top_countries["adjusted_salary"].iloc[i]:,.0f}',
57                ha='center', fontsize=9)
58
59 plt.tight_layout()
60 plt.show()
```

Average Salary by Employee Residence

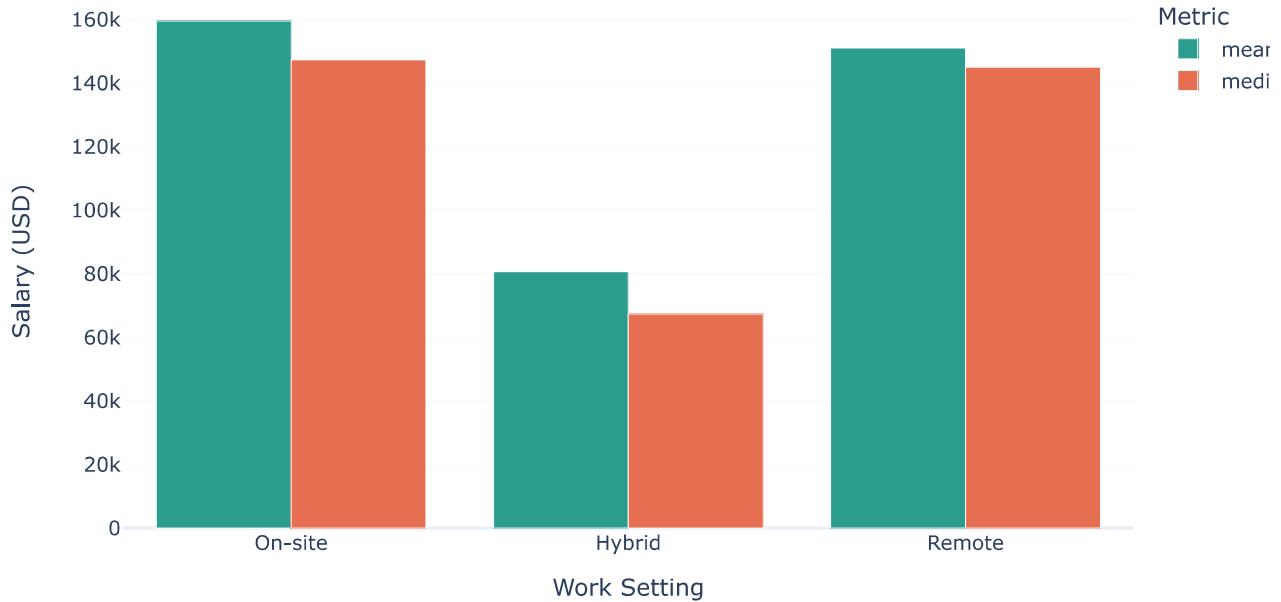


The US, Switzerland, and Israel lead in data science salaries. Interestingly, professionals working for companies outside their country of residence tend to earn more on average, highlighting the benefits of remote work arrangements with companies in high-paying markets.

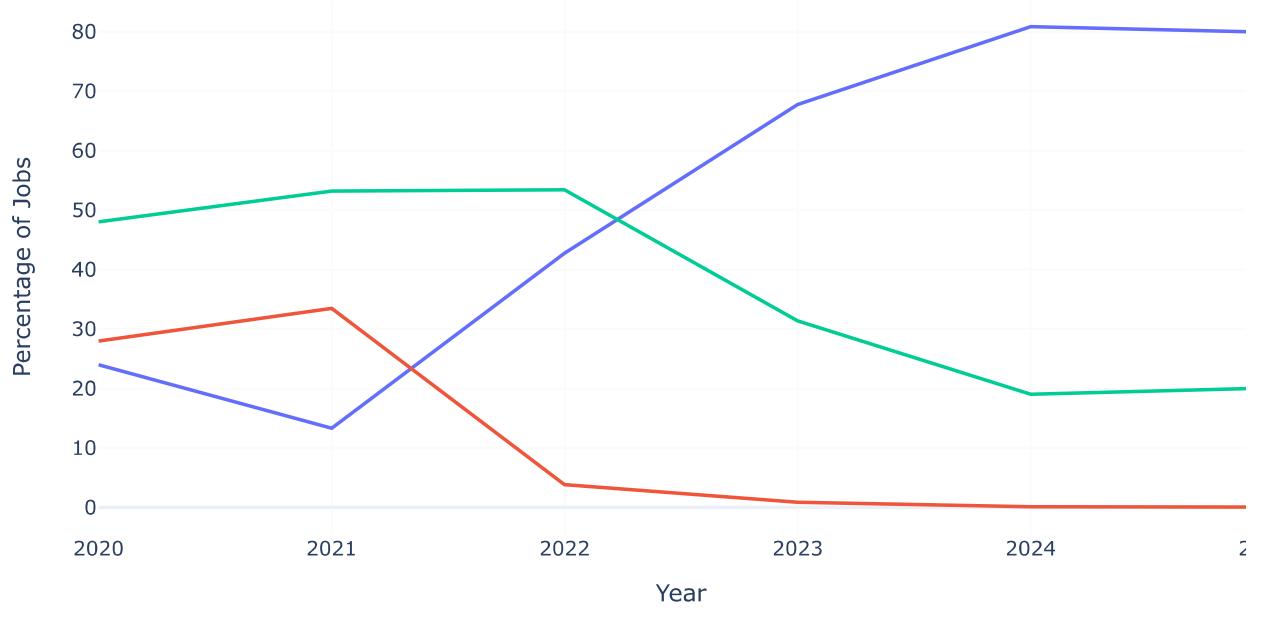
6. Remote Work Analysis

```
In [30]: 1 # Salary by remote ratio
2 remote_salary = df.groupby('remote_ratio')['salary_in_usd'].agg(['mean', 'median', 'count']).reset_index()
3 remote_salary['remote_ratio'] = remote_salary['remote_ratio'].map({0: 'On-site', 50: 'Hybrid', 100: 'Remote'})
4
5 fig = px.bar(remote_salary, x='remote_ratio', y=['mean', 'median'],
6               barmode='group', title='Salary by Remote Work Ratio',
7               labels={'value': 'Salary (USD)', 'remote_ratio': 'Work Setting', 'variable': 'Metric'},
8               color_discrete_sequence=['#2a9d8f', '#e76f51'])
9
10 fig.update_layout(width=800, height=500)
11 fig.show()
12
13 # Remote ratio trends over time
14 remote_time = df.groupby(['work_year', 'remote_ratio']).size().reset_index(name='count')
15 total_per_year = remote_time.groupby('work_year')['count'].sum().reset_index()
16 remote_time = remote_time.merge(total_per_year, on='work_year', suffixes=('', '_total'))
17 remote_time['percentage'] = (remote_time['count'] / remote_time['count_total']) * 100
18 remote_time['remote_ratio'] = remote_time['remote_ratio'].map({0: 'On-site', 50: 'Hybrid', 100: 'Remote'})
19
20 fig = px.line(remote_time, x='work_year', y='percentage', color='remote_ratio',
21                title='Remote Work Trends (2020-2025)',
22                labels={'percentage': 'Percentage of Jobs', 'work_year': 'Year'},
23                template='plotly_white')
24
25 fig.update_layout(width=900, height=500, hovermode='x unified')
26 fig.show()
```

Salary by Remote Work Ratio



Remote Work Trends (2020-2025)



Fully remote positions tend to offer higher salaries on average. We can also observe a significant shift toward remote work after 2020 (likely due to the COVID-19 pandemic), with the trend stabilizing around 2023-2024 but still showing a higher remote work percentage compared to pre-2020 levels.

7. Company Size Impact

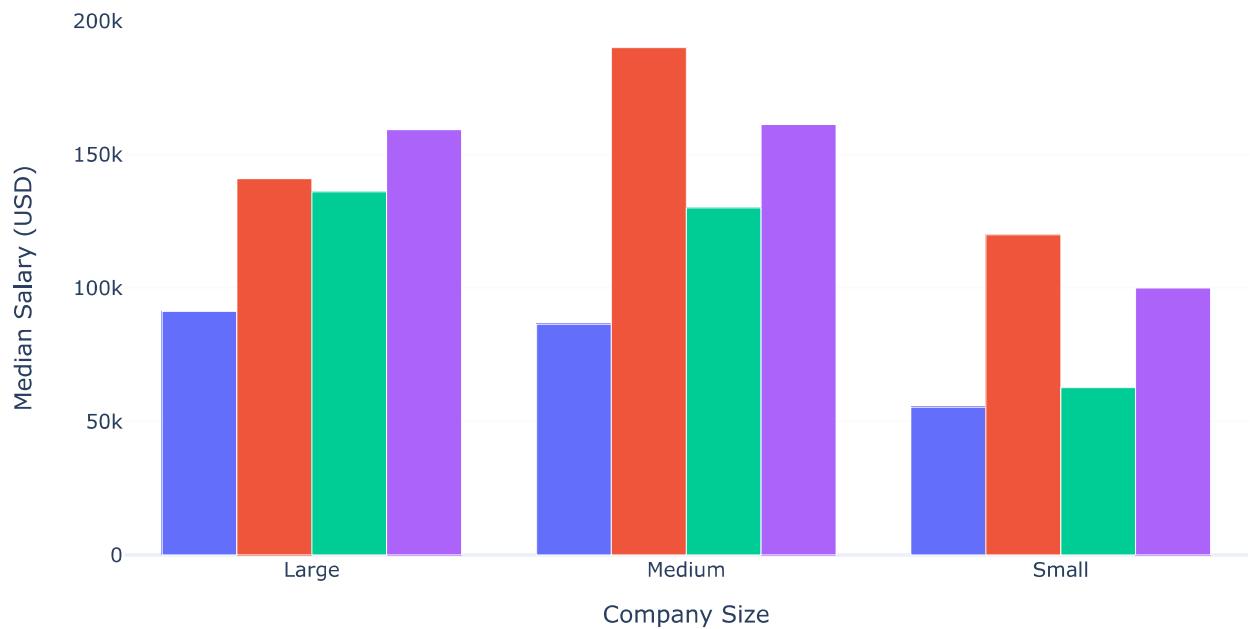
In [31]:

```

1 # Salary by company size
2 company_salary = df.groupby(['company_size', 'experience_level'])['salary_in_usd'].median().reset_index()
3 company_salary['company_size'] = company_salary['company_size'].map({'S': 'Small', 'M': 'Medium',
4 company_salary['experience_level'] = company_salary['experience_level'].map({
5     'EN': 'Entry-level', 'MI': 'Mid-level', 'SE': 'Senior', 'EX': 'Executive'
6 })
7
8 fig = px.bar(company_salary, x='company_size', y='salary_in_usd', color='experience_level',
9             barmode='group', title='Median Salary by Company Size and Experience Level',
10            labels={'salary_in_usd': 'Median Salary (USD)', 'company_size': 'Company Size'},
11            template='plotly_white')
12
13 fig.update_layout(width=900, height=500)
14 fig.show()

```

Median Salary by Company Size and Experience Level

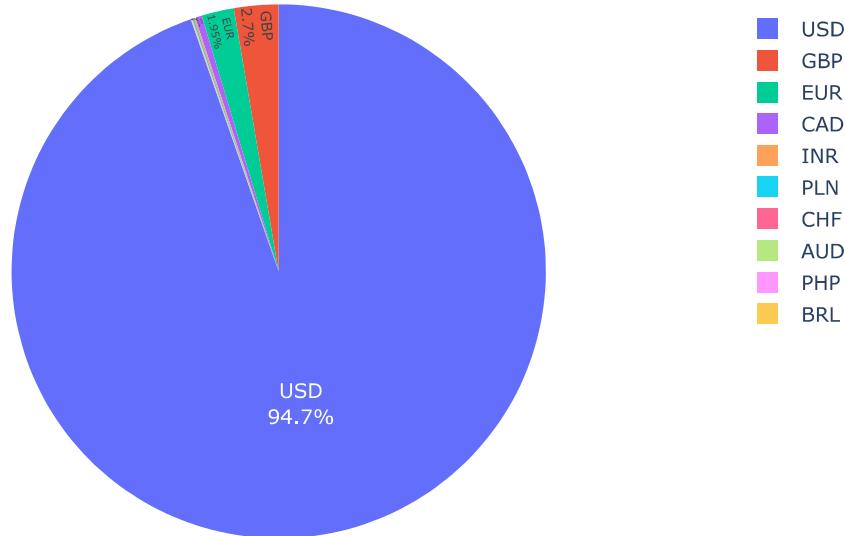


Larger companies generally offer higher salaries across all experience levels, with the gap most pronounced at the executive level. This reflects the greater resources and revenue of larger organizations

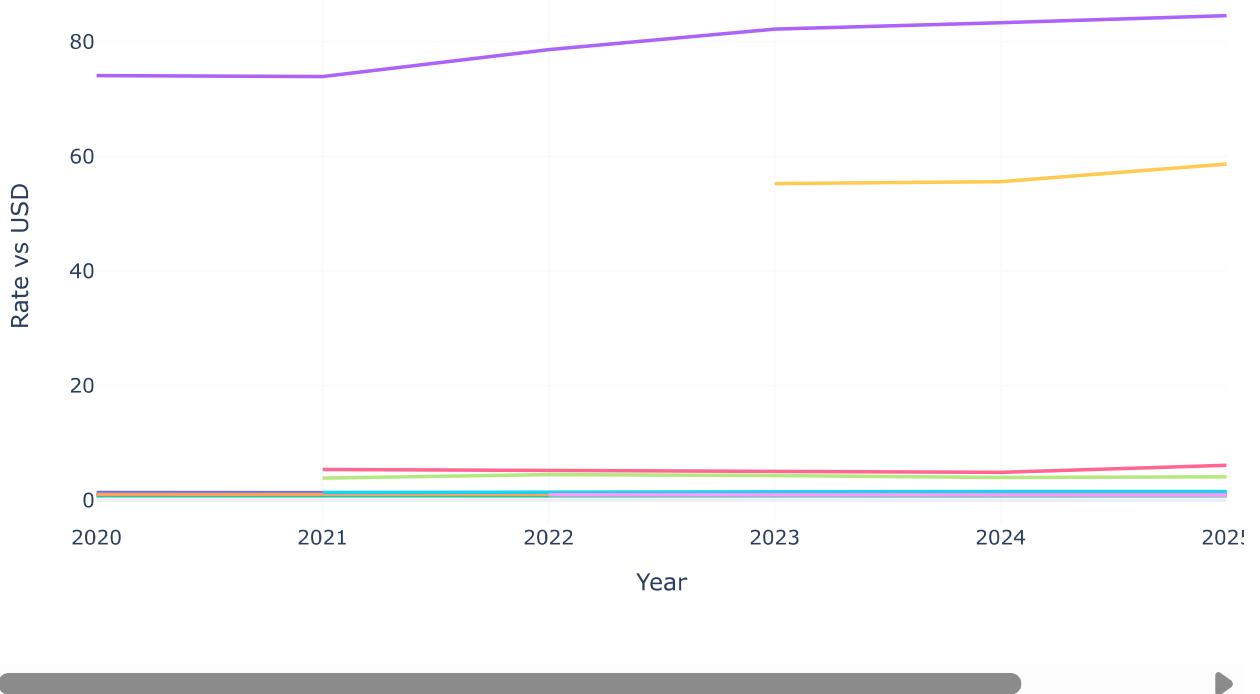
8. Currency Analysis

```
In [32]:  
1 # Top salary currencies  
2 currency_counts = df['salary_currency'].value_counts().head(10)  
3  
4 fig = px.pie(values=currency_counts.values, names=currency_counts.index,  
5                 title='Distribution of Salary Currencies',  
6                 template='plotly_white')  
7  
8 fig.update_traces(textposition='inside', textinfo='percent+label')  
9 fig.update_layout(width=700, height=500)  
10 fig.show()  
11  
12 # Exchange rate analysis (implied from salary and salary_in_usd)  
13 df['implied_exchange_rate'] = df['salary'] / df['salary_in_usd']  
14 top_currencies = df['salary_currency'].value_counts().head(10).index.tolist()  
15 exchange_rates = df[df['salary_currency'].isin(top_currencies)].groupby(['work_year', 'salary_c  
16  
17 fig = px.line(exchange_rates, x='work_year', y='implied_exchange_rate', color='salary_currency'  
18                 title='Implied Exchange Rate Trends (2020-2025)',  
19                 labels={'implied_exchange_rate': 'Rate vs USD', 'work_year': 'Year'},  
20                 template='plotly_white')  
21  
22 fig.update_layout(width=900, height=500, hovermode='x unified')  
23 fig.show()
```

Distribution of Salary Currencies



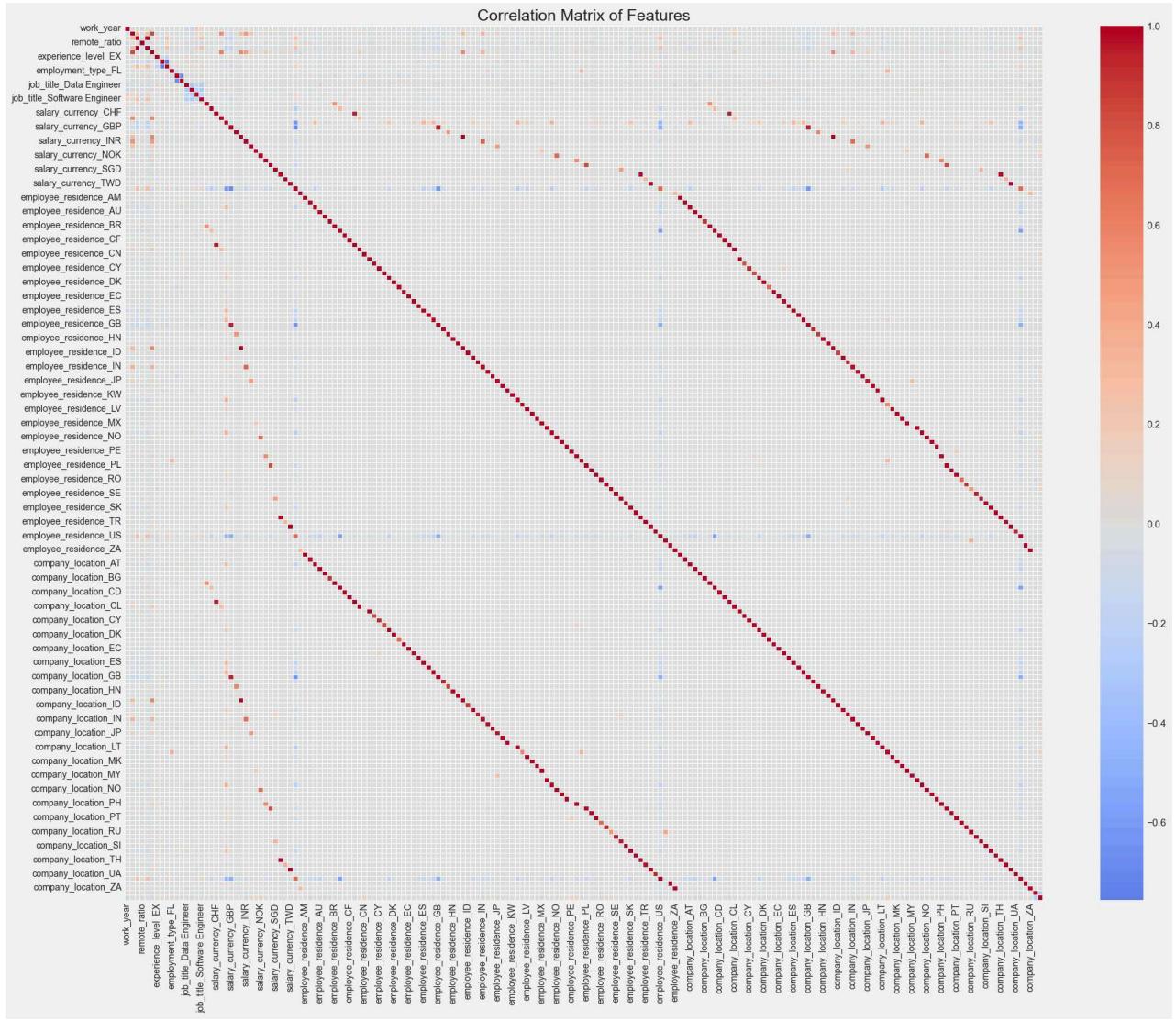
Implied Exchange Rate Trends (2020-2025)



USD dominates as the primary currency for data science salaries globally. The exchange rate analysis shows relative currency strength over time, with some currencies showing depreciation against the USD.

9. Correlation Analysis

```
In [34]:  
1 # Create dummy variables for categorical features  
2 categorical_cols = ['experience_level', 'employment_type', 'job_title', 'salary_currency',  
3                     'employee_residence', 'company_location', 'company_size']  
4 numerical_cols = ['work_year', 'salary', 'salary_in_usd', 'remote_ratio']  
5  
6 # Select top job titles for simplified correlation analysis  
7 top_job_titles = df['job_title'].value_counts().head(5).index.tolist()  
8 df_corr = df[df['job_title'].isin(top_job_titles)].copy()  
9  
10 # Create dummy variables  
11 df_dummies = pd.get_dummies(df_corr, columns=categorical_cols, drop_first=True)  
12  
13 # Calculate and plot correlation matrix  
14 corr_matrix = df_dummies.corr()  
15 plt.figure(figsize=(20, 16))  
16 sns.heatmap(corr_matrix, annot=False, cmap='coolwarm', center=0, linewidths=0.5)  
17 plt.title('Correlation Matrix of Features', fontsize=18)  
18 plt.xticks(fontsize=10, rotation=90)  
19 plt.yticks(fontsize=10)  
20 plt.tight_layout()  
21 plt.show()  
22  
23 # Key salary correlations  
24 salary_corr = corr_matrix['salary_in_usd'].sort_values(ascending=False)  
25 print("Top 10 features positively correlated with salary:")  
26 print(salary_corr.head(11)) # Including salary_in_usd itself  
27 print("\nTop 10 features negatively correlated with salary:")  
28 print(salary_corr.tail(10))
```



Top 10 features positively correlated with salary:

salary_in_usd	1.000000
adjusted_salary	1.000000
salary	0.307103
experience_level_SE	0.283058
salary_currency_USD	0.272702
job_title_Software_Engineer	0.263706
employee_residence_US	0.259746
company_location_US	0.258906
experience_level_EX	0.101781
job_title_Engineer	0.091453

Name: salary_in_usd, dtype: float64

Top 10 features negatively correlated with salary:

employee_residence_ES	-0.069109
company_location_ES	-0.069288
employee_residence_FR	-0.069365
company_location_LT	-0.077144
employee_residence_LT	-0.077144
experience_level_MI	-0.140575
employee_residence_GB	-0.154869
company_location_GB	-0.155245
salary_currency_GBP	-0.164023
salary_currency_EUR	-0.195138

Name: salary_in_usd, dtype: float64

The correlation analysis reveals strong relationships between salary and factors like experience level, company size, and certain job titles. We can use these insights for feature selection in our models.

Conclusion and Key Insights

Our comprehensive analysis of the data science salary dataset has yielded several valuable insights:

1. Salary Determinants: Experience level, company location, and job category are the strongest predictors of salary in the data science field. Executive positions command significantly higher salaries than other levels.

2. Geographic Impact: The US, Switzerland, and Israel consistently offer the highest salaries. There's a substantial gap between top-paying countries and emerging tech markets.

3. Remote Work Premium: Fully remote positions tend to offer higher salaries on average, potentially reflecting the global competition for talent regardless of location.

4. Career Growth: The transition from Senior to Executive level shows the largest percentage increase in salary, highlighting the significant premium placed on leadership skills.

5. Job Title Differentiation: Machine Learning Engineers and Research Scientists command higher salaries than Data Analysts and general Software Engineers, reflecting the specialized skills required.

6. Company Size Effect: Larger companies generally offer higher compensation across all experience levels, with the gap most pronounced at the executive level.

Practical Applications

This analysis and the resulting model can be used for:

1. Career Planning: Understanding the financial impact of different career paths and transitions

2. Negotiation Support: Data-backed reference points for salary negotiations

3. Location Strategy: Evaluating remote work opportunities across different geographies

4. Recruitment Planning: Setting competitive compensation packages for data science roles

Future Work

To further enhance this analysis, we could:

1. Incorporate education level and specific skills data
2. Add cost of living adjustments to provide normalized comparisons
3. Include industry vertical information to differentiate between sectors
4. Develop a time-series model to forecast future salary trends

This notebook demonstrates a comprehensive approach to analyzing and modeling the data science salary dataset. The techniques used here can be applied to other salary datasets or extended with additional features for more nuanced insights.

If you found this analysis helpful, please upvote and share your thoughts in the comments!

