

```
1import requests
2from bs4 import BeautifulSoup
3import pandas as pd
4
5def scrape_indeed_jobs(search_query, location):
6    """Scrapes job listings from Indeed based on
7
8    # Indeed search URL
9    URL = f"https://in.indeed.com/jobs?q={search_
10
11    # Set User-Agent to avoid blocking
12    headers = {
13        "User-Agent": "Mozilla/5.0 (Windows NT 10
14    }
15
16    # Send request
17    response = requests.get(URL, headers=headers)
18    soup = BeautifulSoup(response.text, "html.par
19
20    jobs = []
21
22    # Find all
23    for job_card in soup.find_all("div", class_="
24        # Extract Job Title
25        title = job_card.find("h2", {"data-testic
26        title = title.text.strip() if title else
27
28        # Extract Company Name
29        company = job_card.find("span", class_="j
30        company = company.text.strip() if company
31
32        # Extract Location
33        location = job_card.find("div", {"data-te
34        location = location.text.strip() if locat
35
36        jobs.append([title, company, location])
37
```

Run this cell to mount your Google Drive.  
Learn more

Dismiss

```

38     # Convert to DataFrame
39     df = pd.DataFrame(jobs, columns=["Job Title",
40
41     return df
42
43
44 # Scrape Data Analyst jobs in Mumbai
45 job_data = scrape_indeed_jobs("data analyst", "Mu
46
47 # Display first 5 rows
48 print(job_data.head())
49
50 # Save to CSV
51 job_data.to_csv("job_data.csv", index=False)
52 print("Data saved successfully!")
53

```

Empty DataFrame  
Columns: [Job Title, Company, Location]  
Index: []  
Data saved successfully!

```

1 # Install ChromeDriver and Chromium in Google Col
2 !apt-get update
3 !apt-get install chromium-driver chromium-chromedriver
4 !pip install selenium webdriver-manager
5
6 # Set ChromeDriver path for Colab
7 import os
8 os.environ["PATH"] += ":/usr/lib/chromium-browser
9

```

Run this cell to mount your Google Drive.  
Learn more

Get:1 <https://cloud.r-project.org/bin/linux/ubuntu> jammy-cran40/ InRelease [3,632 B]  
Get:2 <https://security.ubuntu.com/ubuntu> jammy-security InRelease [129 kB]  
Get:3 [https://developer.download.nvidia.com/compute/cuda/repos/ubuntu2204/x86\\_64](https://developer.download.nvidia.com/compute/cuda/repos/ubuntu2204/x86_64) InRelease [1,581 B]  
Get:4 <https://r2u.stat.illinois.edu/ubuntu> jammy InRelease [6,555 B]  
Hit:5 <http://archive.ubuntu.com/ubuntu> jammy InRelease  
Get:6 <http://archive.ubuntu.com/ubuntu> jammy-updates InRelease [128 kB]  
Get:7 [https://developer.download.nvidia.com/compute/cuda/repos/ubuntu2204/x86\\_64](https://developer.download.nvidia.com/compute/cuda/repos/ubuntu2204/x86_64) Packages [1,319 kB]  
Hit:8 <https://ppa.launchpadcontent.net/deadsnakes/ppa/ubuntu> jammy InRelease  
Hit:9 <https://ppa.launchpadcontent.net/graphics-drivers/ppa/ubuntu> jammy InRelease  
Hit:10 <https://ppa.launchpadcontent.net/ubuntugis/ppa/ubuntu> jammy InRelease  
Get:11 <https://r2u.stat.illinois.edu/ubuntu> jammy/main amd64 Packages [2,661 kB]  
Get:12 <http://archive.ubuntu.com/ubuntu> jammy-backports InRelease [127 kB]  
Get:13 <https://r2u.stat.illinois.edu/ubuntu> jammy/main all Packages [8,704 kB]  
Get:14 <http://security.ubuntu.com/ubuntu> jammy-security/restricted amd64 Packages [3,664 kB]  
Get:15 <http://archive.ubuntu.com/ubuntu> jammy-updates/universe amd64 Packages [1,532 kB]  
Get:16 <http://security.ubuntu.com/ubuntu> jammy-security/universe amd64 Packages [1,235 kB]  
Get:17 <http://security.ubuntu.com/ubuntu> jammy-security/main amd64 Packages [2,639 kB]  
Get:18 <http://archive.ubuntu.com/ubuntu> jammy-updates/restricted amd64 Packages [3,813 kB]  
Get:19 <http://archive.ubuntu.com/ubuntu> jammy-updates/main amd64 Packages [2,950 kB]  
Fetched 28.9 MB in 3s (10.1 MB/s)  
Reading package lists... Done  
W: Skipping acquire of configured file 'main/source/Sources' as repository '<https://r2u.stat.illinois.edu/ubuntu> jammy InRelease'  
Reading package lists... Done

```

Building dependency tree... Done
Reading state information... Done
The following additional packages will be installed:
  apparmor chromium-browser libfuse3-3 liblz2-2 snapd squashfs-tools systemd-hwe-hwdb udev
Suggested packages:
  apparmor-profiles-extra apparmor-utils fuse3 zenity | kdialog
The following NEW packages will be installed:
  apparmor chromium-browser chromium-chromedriver libfuse3-3 liblz2-2 snapd squashfs-tools
  systemd-hwe-hwdb udev
0 upgraded, 9 newly installed, 0 to remove and 32 not upgraded.
Need to get 30.1 MB of archives.
After this operation, 123 MB of additional disk space will be used.
Get:1 http://archive.ubuntu.com/ubuntu jammy-updates/main amd64 apparmor amd64 3.0.4-2ubuntu2.4 [598 kB]
Get:2 http://archive.ubuntu.com/ubuntu jammy/main amd64 liblz2-2 amd64 2.10-2build3 [53.7 kB]
Get:3 http://archive.ubuntu.com/ubuntu jammy/main amd64 squashfs-tools amd64 1:4.5-3build1 [159 kB]
Get:4 http://archive.ubuntu.com/ubuntu jammy-updates/main amd64 udev amd64 249.11-0ubuntu3.12 [1,557 kB]
Get:5 http://archive.ubuntu.com/ubuntu jammy/main amd64 libfuse3-3 amd64 3.10.5-1build1 [81.2 kB]
Get:6 http://archive.ubuntu.com/ubuntu jammy-updates/main amd64 snapd amd64 2.66.1+22.04 [27.6 MB]
Get:7 http://archive.ubuntu.com/ubuntu jammy-updates/universe amd64 chromium-browser amd64 1:85.0.4183.83-0ubuntu2.22.04.1 [49.2
Get:8 http://archive.ubuntu.com/ubuntu jammy-updates/universe amd64 chromium-chromedriver amd64 1:85.0.4183.83-0ubuntu2.22.04.1
Get:9 http://archive.ubuntu.com/ubuntu jammy-updates/main amd64 systemd-hwe-hwdb all 249.11.5 [3,228 B]
Fetched 30.1 MB in 2s (17.9 MB/s)
Preconfiguring packages ...
Selecting previously unselected package apparmor.
(Reading database ... 124947 files and directories currently installed.)
Preparing to unpack .../0-apparmor_3.0.4-2ubuntu2.4_amd64.deb ...
Unpacking apparmor (3.0.4-2ubuntu2.4) ...
Selecting previously unselected package liblz2-2:amd64.
Preparing to unpack .../1-liblz2-2_2.10-2build3_amd64.deb ...
Unpacking liblz2-2:amd64 (2.10-2build3) ...
Selecting previously unselected package squashfs-tools.
Preparing to unpack .../2-squashfs-tools_1%3a4.5-3build1_amd64.deb ...
Unpacking squashfs-tools (1:4.5-3build1) ...

```

```

1 from selenium import webdriver
2 from selenium.webdriver.common.by import By
3 from selenium.webdriver.support.ui import WebDriverWait
4 from selenium.webdriver.support import expected_conditions as
5 import pandas as pd
6 import time
7
8 # Set up Selenium
9 options = webdriver.ChromeOptions()
10 options.add_argument("--headless") # Run without
11 options.add_argument("--no-sandbox")
12 options.add_argument("--disable-dev-shm-usage")
13 options.add_argument("user-agent=Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/91.0.4472.164 Safari/537.36")
14
15 # Use Chromium in Google Colab
16 driver = webdriver.Chrome(options=options)
17
18 def scrape_naukri_jobs(search_query, location):
19     """Scrapes job listings from Naukri.com using Selenium"""
20
21     URL = f"https://www.naukri.com/{search_query}"
22     driver.get(URL)
23

```

Run this cell to mount your Google Drive.  
Learn more

```
24 # Wait until job cards load
25 try:
26     WebDriverWait(driver, 15).until(
27         EC.presence_of_element_located((By.CL
28     )
29 except:
30     print("✖ Job listings did not load!")
31     return pd.DataFrame(columns=["Job Title",
32
33 jobs = []
34
35 # Extract job listings
36 job_cards = driver.find_elements(By.CLASS_NAM
37 for job in job_cards:
38     try:
39         title = job.find_element(By.CLASS_NAM
40     except:
41         title = "Not available"
42
43     try:
44         company = job.find_element(By.CLASS_N
45     except:
46         company = "Not available"
47
48     try:
49         location = job.find_element(By.CLASS_
50     except:
51         location = "Not available"
52
53     try:
54         salary = job.find_element(By.CLASS_NA
55     except:
56         salary = "Not available"
57
58     try:
59         skills = ", ".join([skill.text.strip(
60     except:
61         skills = "Not available"
```

Run this cell to mount your Google Drive.  
Learn more

```

62
63     jobs.append([title, company, location, sa
64
65     df = pd.DataFrame(jobs, columns=["Job Title",
66
67     return df
68
69
70# Scrape Data Analyst jobs in Mumbai from Naukri
71job_data = scrape_naukri_jobs("data-analyst", "mu
72
73# Close browser
74driver.quit()
75
76# Display first 5 rows
77print(job_data.head())
78
79# Save Data to CSV
80job_data.to_csv("naukri_jobs.csv", index=False)
81print("Data saved successfully!")
82

```



Run this cell to mount your Google Drive.  
Learn more

```

0      Data Analyst (World Panel)
1      GN- Public Service - Data & Analytics
2      Data Analyst - Monitoring & Reporting
3      Opening For Data Analyst -(DME) For Mumbai 100...
4      Data Analyst

      Company      Location      Salary \
0      Numerator      Hybrid - Mumbai      3-3.5 Lacs PA
1      Accenture      Mumbai, Gurugram, Bengaluru      Not disclosed
2      Aditya Birla Education Trust      Mumbai (All Areas)(Worli)      5-7 Lacs PA
3      Inland World Logistics      Mumbai      Not disclosed
4      HH Consultancy      Remote      7-9 Lacs PA

      Skills
0      Advanced Excel, R, Mac, Python, SQL, Excel, Pa...
1      data quality, data modeling, dashboards, artif...
2      Data Analysis Tools, Data Analysis, Programmin...
3      Data Analysis, Data Management, Data Reporting...
4      Data Analysis, Data Manipulation, Data Managem...
Data saved successfully!

```

```

1import pandas as pd
2
3df = pd.read_csv("naukri_jobs.csv")
4
5# Top 5 highest-paying jobs
6print(df[df['Salary'] != "Not disclosed"].sort_va
7

```

```

8# Most common skills
9from collections import Counter
10skills = ", ".join(df["Skills"].dropna()).split("
11top_skills = Counter(skills).most_common(10)
12print("Top In-Demand Skills:", top_skills)
13

```

Job Title Company \

17	Data Analyst	WebMD
4	Data Analyst	HH Consultancy
6	Data Analyst	ETG
9	Data Analyst    German	mnc
8	Data Analyst	Good2Great Industries Pvt Ltd

Location Salary \

17	Mumbai (All Areas)	9.5-19.5 Lacs PA
4	Remote	7-9 Lacs PA
6	Navi Mumbai	7-9 Lacs PA
9	Mumbai (All Areas), Pune	7-14 Lacs PA
8	Mumbai (All Areas)	6-8 Lacs PA

Skills

```

17 SQL, Power BI, Data Analyst, Tableau, Bi, Data...
4 Data Analysis, Data Manipulation, Data Managem...
6 SQL, Python, Power BI, data modeling, dashboar...
9 German Language, Translation, German, Language...
8 Data Analysis, Advanced Excel, Statistical Dat...
Top In-Demand Skills: [('Data', 8), ('Data analysis', 7), ('Data Analysis', 7), ('Analysis', 6), ('Data Analyst', 5), ('Data Managem

```

```

1import pandas as pd
2import numpy as np
3import matplotlib.pyplot as plt
4import seaborn as sns
5inputfile= "/c
6df=pd.read_csv
7
8# Remove "Lacs PA" and clean the column
9df["Salary"] = df["Salary"].str.replace(" Lacs PA
10
11# Handle cases where salary is missing or "Not di
12df["Salary"] = df["Salary"].apply(lambda x: x.spl
13df["Salary"] = pd.to_numeric(df["Salary"], errors
14
15# Fill missing salaries with the mean salary
16df["Salary"].fillna(df["Salary"].mean(), inplace=
17
18print(df.head()) # Check the fixed Salary column
19

```

Job Title \

0	Data Analyst (World Panel by Kantar)
1	GN- Public Service - Data & AI - Analyst
2	Data Analyst - Monitoring & Evaluation
3	Opening For Data Analyst -(DME) For Mumbai loc...

```

4                                     Data Analyst

0                                     Company      Location      Salary \
1                                     Numerator      Hybrid - Mumbai 3.000000
2                                     Accenture  Mumbai, Gurugram, Bengaluru 5.833333
3   Aditya Birla Education Trust      Mumbai (All Areas)(Worli) 5.000000
4   Inland World Logistics      Mumbai 5.833333
5   HH Consultancy      Remote 7.000000

                                     Skills
0   Advanced Excel, R, Mac, Python, SQL, Excel, Pa...
1   data quality, data modeling, dashboards, artif...
2   Data Analysis Tools, Data Analysis, Programmin...
3   Data Analysis, Data Management, Data Reporting...
4   Data Analysis, Data Manipulation, Data Managem...
<ipython-input-6-752260c4ff9f>:16: FutureWarning: A value is trying to be set on a copy of a DataFrame or Series through chained ass
The behavior will change in pandas 3.0. This inplace method will never work because the intermediate object on which we are setting

For example, when doing 'df[col].method(value, inplace=True)', try using 'df.method({col: value}, inplace=True)' or df[col] = df[col]

df["Salary"].fillna(df["Salary"].mean(), inplace=True)

```


## 1 Top Hiring Companies

```

1import seaborn as sns
2import matplotlib.pyplot as plt
3
4# Count the top 10 hiring companies
5top_companies = df["Company"].value_counts().head(10)
6
7# Set figure size
8plt.figure(figsize=(10, 5))
9
10# Use Seaborn
11sns.barplot(x=range(1, 11), y=top_companies.values, yerr=top_companies.values)
12
13# Set labels and title
14plt.xlabel("Number of Jobs")
15plt.ylabel("Company")
16plt.title("Top 10 Hiring Companies")
17
18# Show the plot
19plt.show()
20

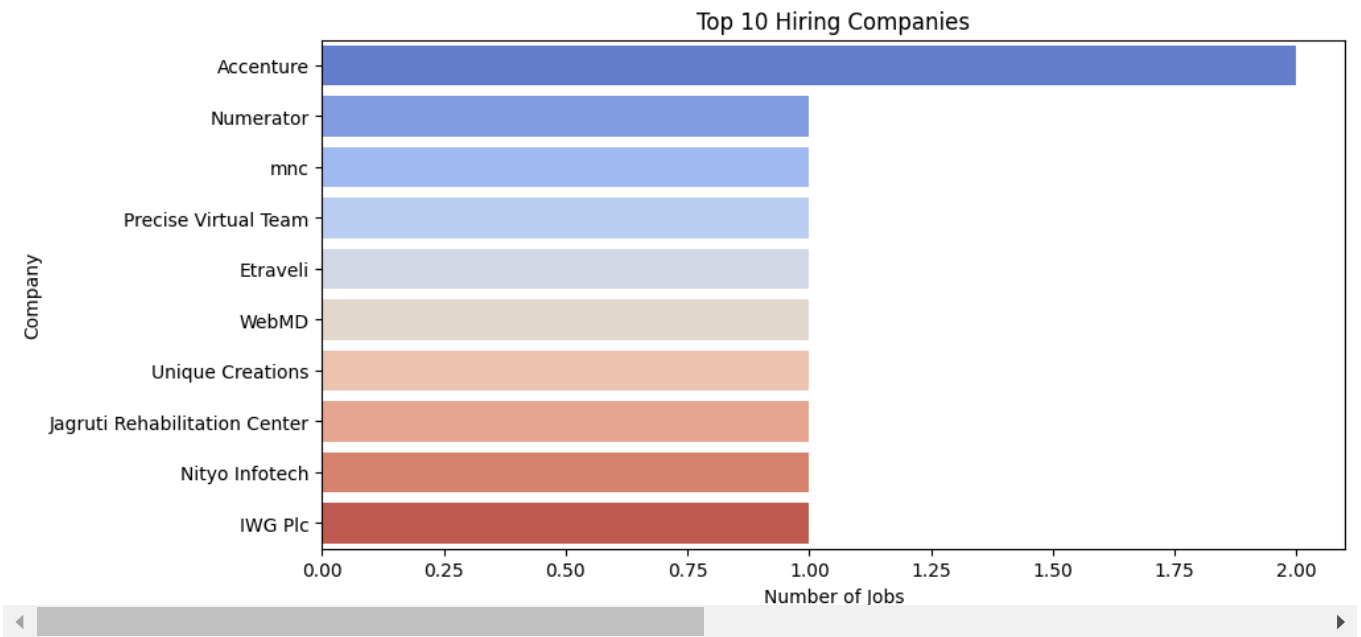
```

Run this cell to mount your Google Drive.  
Learn more

 <ipython-input-10-8abe968e1c3c>:11: FutureWarning:

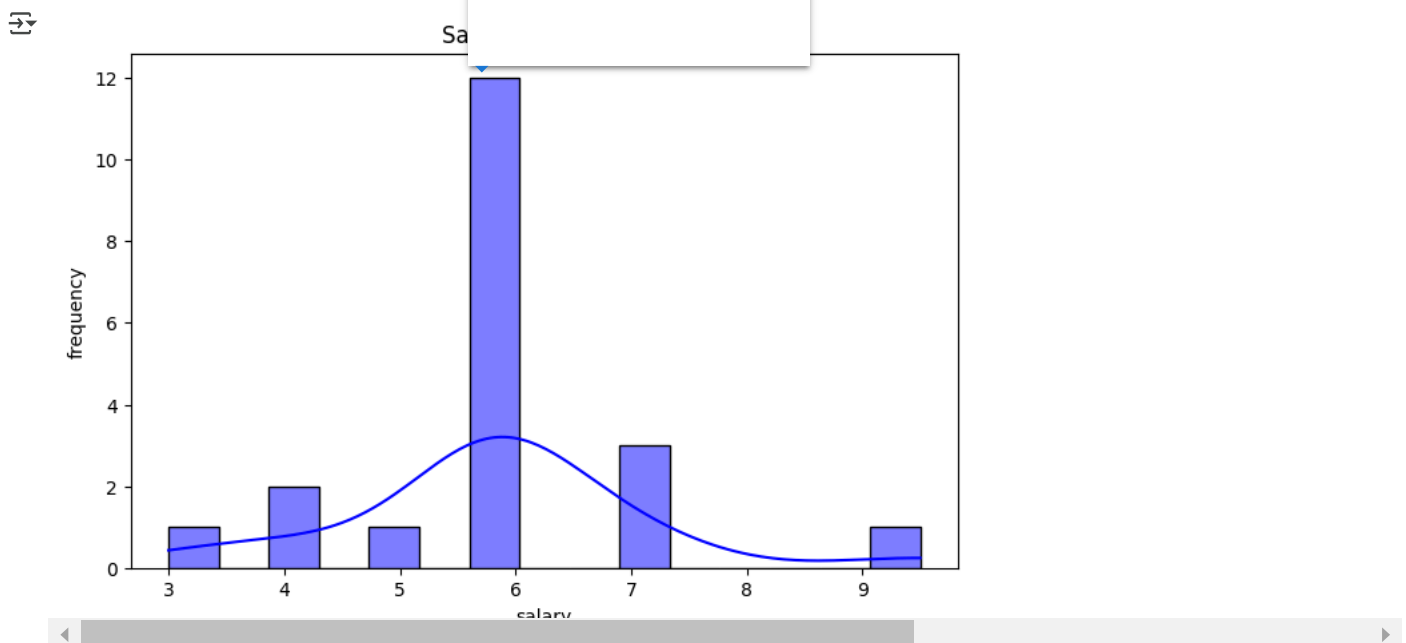
Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign the `y` variable to `hue` and set `le

```
sns.barplot(x=top_companies.values, y=top_companies.index, palette="coolwarm")
```



## 2 Salary Distribution Analysis

```
1 plt.figure(figsize=(8,5))
2 sns.histplot(df['Salary'],bins=15,kde=True,color=
3 plt.xlabel('salary')
4 plt.ylabel('frequency')
5 plt.title('Salary Distribution')
6 plt.show()
```



## 3 Most In-Demand Skills

```
1 from collections import Counter
2
3 all_skills = Counter(df['Skills'].dropna()).most_common(10)
```



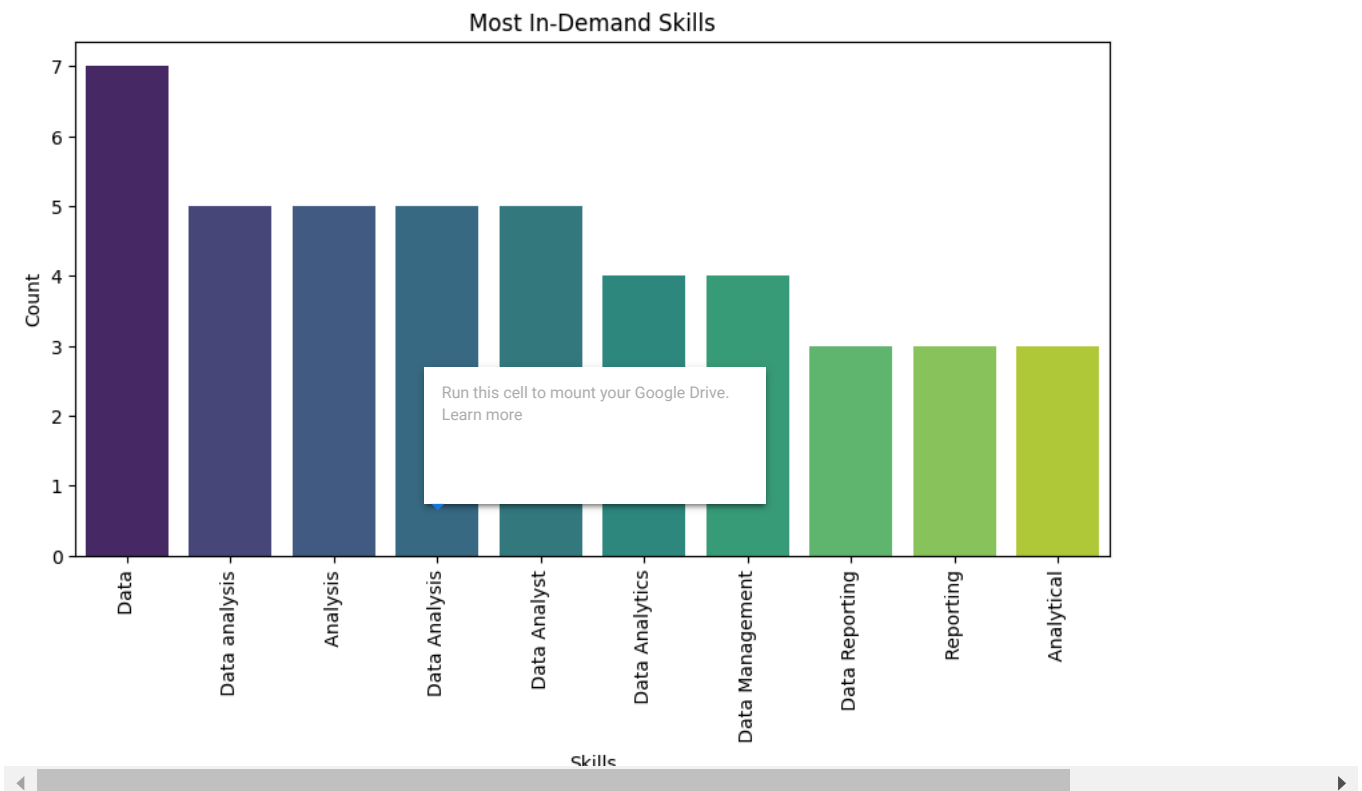
```

3 all_skills= , .join(all_skills).dropna()).split
4 top_skills=Counter(all_skills).most_common(10)
5
6 skills_df=pd.DataFrame(top_skills,columns=['Skill
7
8 plt.figure(figsize=(10,5))
9 sns.barplot(x=skills_df['Skills'],y=skills_df['Co
10 plt.xlabel('Skills')
11 plt.ylabel('Count')
12 plt.title('Most In-Demand Skills')
13 plt.xticks(rotation=90)
14 plt.show()

```

<ipython-input-15-49f1bb55ffbc>:9: FutureWarning:

Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign the `x` variable to `hue` and set `le`  
 sns.barplot(x=skills\_df['Skills'],y=skills\_df['Count'],palette="viridis")



```

1 from google.colab import drive
2 drive.mount('/content/drive')

```

#### 4 Job Location Trends

```


1 top_locations = df["Location"].value_counts().head
2
3 plt.figure(figsize=(10, 5))
4 sns.barplot(x=top_locations.values, y=top_locatic

```

```

5plt.xlabel("Number of Job Postings")
6plt.ylabel("Location")
7plt.title("Top 10 Job Locations for Data Analysts")
8plt.show()
9

```

 <ipython-input-16-5aff73691f4d>:4: FutureWarning:

Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign the `y` variable to `hue` and set `legend` to `True` to get a legend for your categories.  
 sns.barplot(x=top\_locations.values, y=top\_locations.index, palette="mako")




```

1# Save the cleaned dataset for visualization in P
2df.to_csv("cleaned_jobs.csv", index=False)
3
4print("Cleaned data saved successfully!")

```

Run this cell to mount your Google Drive.  
 Learn more

 Cleaned data saved successfully!

```

1# Import necessary libraries
2import pandas as pd
3import numpy as np
4import seaborn as sns
5import matplotlib.pyplot as plt
6from sklearn.model_selection import train_test_split
7from sklearn.linear_model import LinearRegressor
8from sklearn.ensemble import RandomForestRegressor
9from sklearn.metrics import mean_absolute_error
10from sklearn.preprocessing import MultiLabelBinarizer
11from sklearn.decomposition import PCA
12

```

```
13# Load dataset
14file_path = "/content/cleaned_naukri_jobs.csv"
15df = pd.read_csv(file_path)
16
17# One-hot encoding for categorical variables
18df = pd.get_dummies(df, columns=["Company", "Location"])
19
20# Ensure Salary is numeric
21df["Salary"] = pd.to_numeric(df["Salary"], errors="coerce")
22
23# Drop rows with missing Salary
24df.dropna(subset=["Salary"], inplace=True)
25
26# Convert Skills into Numeric Format
27df["Skills"] = df["Skills"].fillna("").apply(lambda x: x.split(", "))
28mlb = MultiLabelBinarizer()
29skills_encoded = pd.DataFrame(mlb.fit_transform(df["Skills"], df.index).toarray(),
30                               columns=mlb.get_feature_names_out(), index=df.index)
31df = pd.concat([df, skills_encoded], axis=1)
32df.drop(columns=["Skills"], inplace=True)
33
34# Remove constant columns (if any)
35constant_columns = df.columns[df.nunique() == 1]
36df.drop(columns=constant_columns, inplace=True)
37
38# Identify and remove low-variance features
39threshold = 0.95
40low_variance_cols = [col for col in df.columns if df[col].nunique() < threshold]
41df.drop(columns=low_variance_cols, inplace=True)
42print(f"Dropped {len(low_variance_cols)} low-variance columns")
43
44# Prepare data for modeling
45X = df.drop(columns=["Salary"])
46y = df["Salary"]
47X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
48
49# Train Linear Regression Model
50lin_model = LinearRegression()
51lin_model.fit(X_train, y_train)
```

Run this cell to mount your Google Drive.  
Learn more

```
51 y_pred_lin = lin_model.predict(X_test)
52 mae_lin = mean_absolute_error(y_test, y_pred_lin)
53 print(f"Linear Regression MAE: {mae_lin:.2f}")
54
55 # Train Random Forest Model
56 rf_model = RandomForestRegressor(n_estimators=100)
57 rf_model.fit(X_train, y_train)
58 y_pred_rf = rf_model.predict(X_test)
59 mae_rf = mean_absolute_error(y_test, y_pred_rf)
60 print(f"Random Forest MAE: {mae_rf:.2f}")
61
62 # Feature Importance (Random Forest)
63 feature_importance_rf = pd.Series(rf_model.feature_
64
65 # Plot Feature Importance
66 plt.figure(figsize=(12, 6))
67 sns.barplot(x=feature_importance_rf[:20].values,
68 plt.xlabel("Feature Importance Score")
69 plt.ylabel("Feature")
70 plt.title("Top 20 Most Important Features (Random
71 plt.show()
72
73 # Apply PCA for Dimensionality Reduction
74 scaler = StandardScaler()
75 X_scaled = scaler.fit_transform(X)
76
77 pca = PCA(n_components=0.95) # Keep 95% variance
78 X_pca = pca.fit_transform(X_scaled)
79
80 print(f"Original feature count: {X.shape[1]}")
81 print(f"Reduced feature count after PCA: {X_pca.s
82
83 # Predict Salary for a new job with Python, SQL,
84 new_job_skills = pd.DataFrame([[0] * len(X.columns)
85 for skill in ["Python", "SQL", "Power BI"]]:
86     if skill in new_job_skills.columns:
87         new_job_skills[skill] = 1
88
```

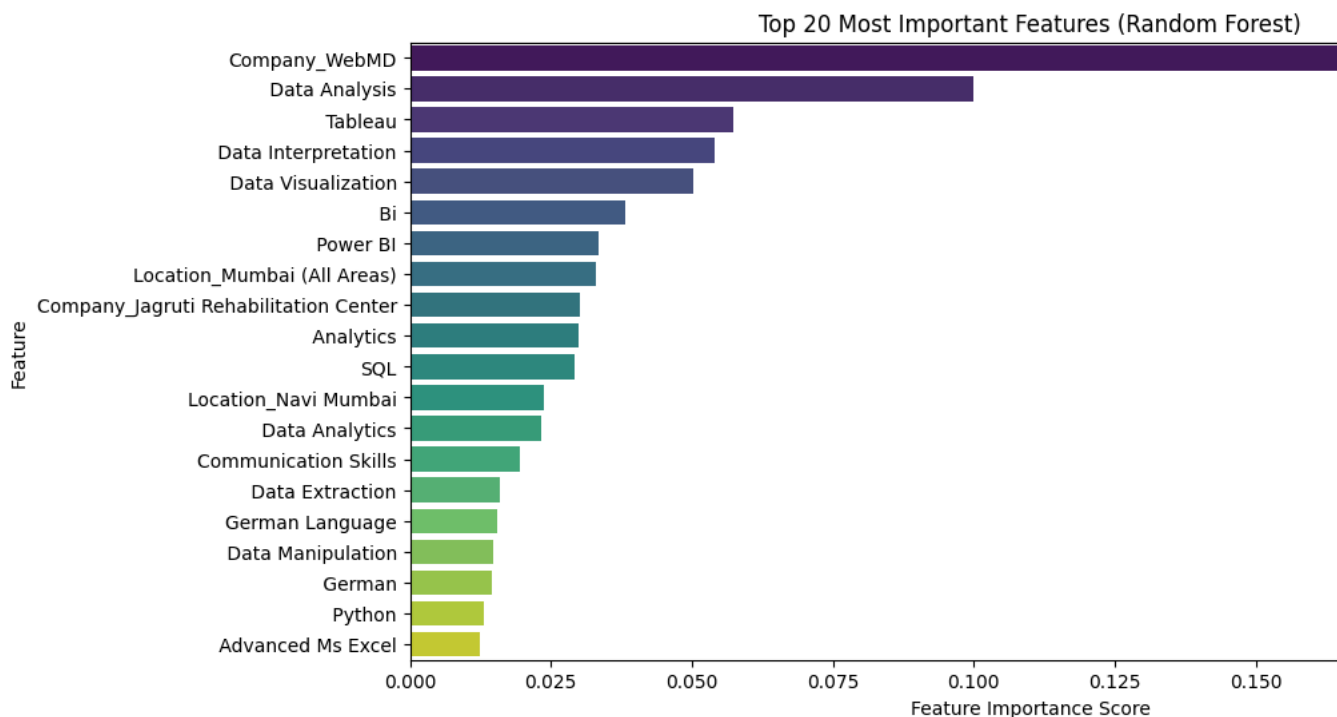
Run this cell to mount your Google Drive.  
[Learn more](#)

```
89 predicted_salary = rf_model.predict(new_job_skill
90 print(f"💰 Predicted Salary for Python, SQL, Pow
91
```

↳ Dropped 0 low-variance columns.  
 Linear Regression MAE: 0.90  
 Random Forest MAE: 0.84  
 <ipython-input-22-ec954b8f84f4>:67: FutureWarning:

Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign the `y` variable to `hue` and set `le

```
sns.barplot(x=feature_importance_rf[:20].values, y=feature_importance_rf[:20].index, palette="viridis")
```



Original feature count: 148

Reduced feature count after PCA: 17

Run this cell to mount your Google Drive.  
[Learn more](#)

```
1!pip install xgboost scikit-learn pandas numpy ma
2
```

↳ Requirement already satisfied: xgboost in /usr/local/lib/python3.11/dist-packages (2.1.4)  
 Requirement already satisfied: scikit-learn in /usr/local/lib/python3.11/dist-packages (1.6.1)  
 Requirement already satisfied: pandas in /usr/local/lib/python3.11/dist-packages (2.2.2)  
 Requirement already satisfied: numpy in /usr/local/lib/python3.11/dist-packages (1.26.4)  
 Requirement already satisfied: matplotlib in /usr/local/lib/python3.11/dist-packages (3.10.0)  
 Requirement already satisfied: seaborn in /usr/local/lib/python3.11/dist-packages (0.13.2)  
 Requirement already satisfied: nvidia-nccl-cu12 in /usr/local/lib/python3.11/dist-packages (from xgboost) (2.21.5)  
 Requirement already satisfied: scipy in /usr/local/lib/python3.11/dist-packages (from xgboost) (1.13.1)  
 Requirement already satisfied: joblib>=1.2.0 in /usr/local/lib/python3.11/dist-packages (from scikit-learn) (1.4.2)  
 Requirement already satisfied: threadpoolctl>=3.1.0 in /usr/local/lib/python3.11/dist-packages (from scikit-learn) (3.5.0)  
 Requirement already satisfied: python-dateutil>=2.8.2 in /usr/local/lib/python3.11/dist-packages (from pandas) (2.8.2)  
 Requirement already satisfied: pytz>=2020.1 in /usr/local/lib/python3.11/dist-packages (from pandas) (2025.1)  
 Requirement already satisfied: tzdata>=2022.7 in /usr/local/lib/python3.11/dist-packages (from pandas) (2025.1)  
 Requirement already satisfied: contourpy>=1.0.1 in /usr/local/lib/python3.11/dist-packages (from matplotlib) (1.3.1)  
 Requirement already satisfied: cycler>=0.10 in /usr/local/lib/python3.11/dist-packages (from matplotlib) (0.12.1)  
 Requirement already satisfied: fonttools>=4.22.0 in /usr/local/lib/python3.11/dist-packages (from matplotlib) (4.56.0)  
 Requirement already satisfied: kiwisolver>=1.3.1 in /usr/local/lib/python3.11/dist-packages (from matplotlib) (1.4.8)  
 Requirement already satisfied: packaging>=20.0 in /usr/local/lib/python3.11/dist-packages (from matplotlib) (24.2)  
 Requirement already satisfied: pillow>=8 in /usr/local/lib/python3.11/dist-packages (from matplotlib) (11.1.0)  
 Requirement already satisfied: pyparsing>=2.3.1 in /usr/local/lib/python3.11/dist-packages (from matplotlib) (3.2.1)  
 Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.11/dist-packages (from python-dateutil>=2.8.2->pandas) (1.17.0)

```
1# Import necessary libraries
2import pandas as pd
3import numpy as np
4import seaborn as sns
```

```
5 import matplotlib.pyplot as plt
6 from sklearn.model_selection import train_test_sp
7 from sklearn.metrics import mean_absolute_error
8 from sklearn.preprocessing import MultiLabelBinar
9 import xgboost as xgb
10
11 # Load dataset
12 file_path = "/content/cleaned_naukri_jobs.csv"
13 df = pd.read_csv(file_path)
14
15 # One-hot encode categorical columns
16 df = pd.get_dummies(df, columns=["Company", "Loca
17
18 # Convert Salary column to numeric
19 df["Salary"] = pd.to_numeric(df["Salary"], errors
20
21 # Drop missing salary values
22 df.dropna(subset=["Salary"], inplace=True)
23
24 # Convert Skills into numeric format
25 df["Skills"] = df["Skills"].fillna("").apply(lamb
26
27 # MultiLabelB skills encoding
28 mlb = MultiLabelBinarizer()
29 skills_encoded = pd.DataFrame(mlb.fit_transform(d
30 df = pd.concat([df, skills_encoded], axis=1)
31
32 # Drop original Skills column
33 df.drop(columns=["Skills"], inplace=True)
34
35 # Define features (X) and target variable (y)
36 X = df.drop(columns=["Salary"])
37 y = df["Salary"]
38
39 # Train-test split (80-20)
40 X_train, X_test, y_train, y_test = train_test_spl
41
42
```

Run this cell to mount your Google Drive.  
Learn more

```
43# Define XGBoost model
44xgb_model = xgb.XGBRegressor(objective="reg:squarederror")
45
46# Hyperparameter grid
47param_grid = {
48    "n_estimators": [100, 200, 300], # Number of
49    "learning_rate": [0.01, 0.1, 0.2], # Step size
50    "max_depth": [3, 5, 7] # Tree depth
51}
52
53# GridSearchCV to find the best hyperparameters
54grid_search = GridSearchCV(xgb_model, param_grid,
55                             cv=5, scoring='neg_mean_squared_error')
56grid_search.fit(X_train, y_train)
57
58# Best model after tuning
59best_xgb_model = grid_search.best_estimator_
60
61# Predict on test set
62y_pred_xgb = best_xgb_model.predict(X_test)
63
64# Calculate Mean Absolute Error (MAE)
65mae_xgb = mean_absolute_error(y_test, y_pred_xgb)
66print(f" ♦ XGBoost MAE: {mae_xgb:.2f}")
67
68# Get feature importance
69feature_importance_xgb = pd.Series(best_xgb_model.feature_importances_)
70feature_importance_xgb = feature_importance_xgb.sort_values(ascending=False)
71
72# Plot Top 20 Features
73plt.figure(figsize=(10, 6))
74sns.barplot(x=feature_importance_xgb[:20].values,
75            y=feature_importance_xgb.index[:20])
76plt.title("Top 20 Most Important Features (XGBoost)")
77plt.xlabel("Feature Importance Score")
78plt.ylabel("Feature")
79plt.show()
80
81# Create a new job sample with Python, SQL, Power BI, and Data Science skills
82new_job_skills = pd.DataFrame([['01' * len(X.columns)]])
```

Run this cell to mount your Google Drive.  
Learn more

```

81 new_job_skills = pd.DataFrame([[0]] * len(X.columns))
82
83 # Set required skills to 1
84 for skill in ["Python", "SQL", "Power BI"]:
85     if skill in new_job_skills.columns:
86         new_job_skills[skill] = 1
87
88 # Predict salary
89 predicted_salary_xgb = best_xgb_model.predict(new
90 print(f"$ Predicted Salary for Python, SQL, Powe
91
92
93
94

```

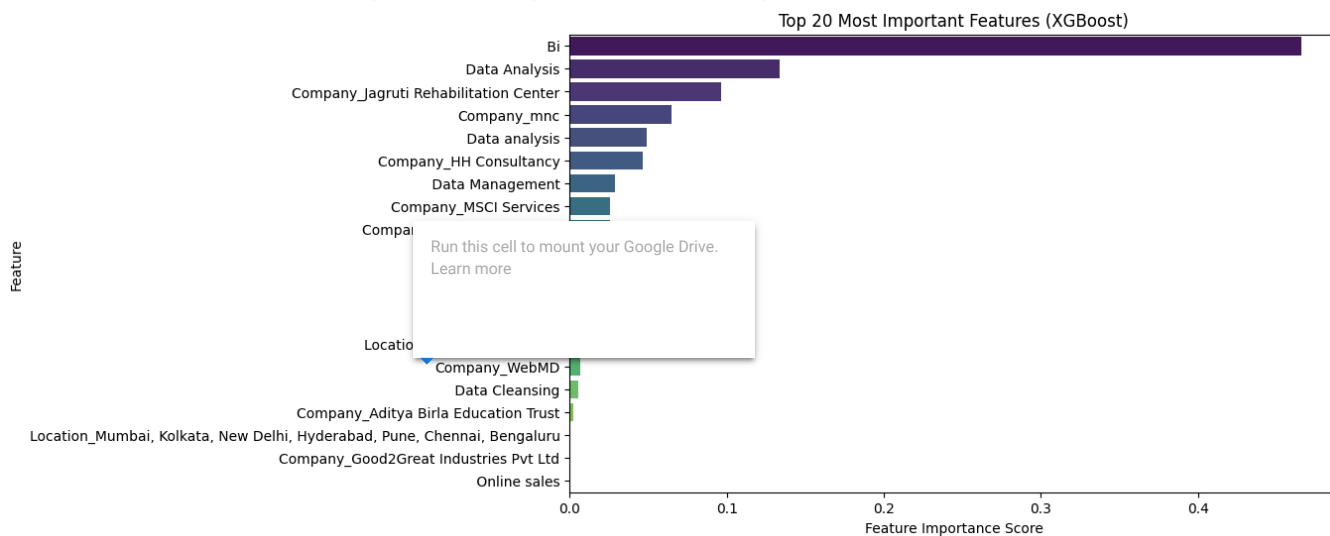
↗ Fitting 3 folds for each of 27 candidates, totalling 81 fits

◆ XGBoost Model MAE: 0.76

<ipython-input-24-b72690ff7ed9>:74: FutureWarning:

Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign the `y` variable to `hue` and set `le

sns.barplot(x=feature\_importance\_xgb[:20].values, y=feature\_importance\_xgb[:20].index, palette="viridis")



\$ Predicted Salary for Python, SQL, Power BI: 5.85 Lacs PA