

1. Importing the data

```
import pandas as pd
import numpy as np
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.preprocessing import LabelEncoder
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score, classification_report

data_path = r"C:\Users\User\Desktop\job_descriptions.csv"

In [1]: data.head()

Out[2]:
```

	Job Id	Experience	Qualifications	Salary Range	Location	Country	latitude	longitude	Work Type	Company Size	...	Contact	Job Title	Role	Job Portal	Job Description	Benefits	skills
0	108984354011562	5 to 15 Years	M.Tech	58K - 65K	Douglas	ile of Man	54.2361	-4.5481	Intern	26801	...	001.381.930-7517177	Digital Marketing Specialist	Social Media Manager	Snagajob	Social Media Managers oversee an organization's...	[Flexible Spending Accounts (FSAs), Relocation...	Social media platforms (e.g., Facebook, Twitter...
1	395456664776	2 to 12 Years	BCA - 1st	56K - 61K	Ashgabat	Turkmenistan	38.9697	59.5563	Intern	100340	...	461-509-4216	Frontend Web Developer	Frontend Web Developer	Indeed	Frontend Web Developers design and implement...	[Health Insurance, Retirement Plans, Paid Time...	HTML, CSS, JavaScript, Frontend frameworks (e.g., React...
2	48164067296353	0 to 12 Years	PHD	61K - 101K	Macao	Macao SAR, China	22.1887	113.5439	Temporary	84525	...	9687618505	Operations Manager	Quality Control Manager	Jobs2Careers	Quality Control Managers establish and enforce...	[Legal Assistance, Bonuses and Incentive Progra...	Quality control processes and methodologies (e.g., Six Sigma, ISO 9001)
3	68819267473044	4 to 11 Years	PHD	65K - 81K	Ponto-Novo	Benin	9.2077	2.1518	Full-Time	129896	...	+1-800-643-5431/47576	Network Engineer	Wireless Network Engineer	FlexJobs	Wireless Network Engineers design, implement, and maintain...	[Transportation Benefits, Professional Development...	Wireless network design and architecture (e.g., 4G/LTE, 5G)
4	1170576136588	10 to 12 Years	MBA	64K - 87K	Santiago	Chile	-35.6751	-71.5429	Intern	53944	...	343.975.4702/63940	Event Manager	Conference Manager	Jobs2Careers	A Conference Manager coordinates and manages c...	[Flexible Spending Accounts (FSAs), Relocation...	Event planning, logistics, budget management

5 rows x 23 columns

```
In [3]: print(data.shape)
(161940, 23)

In [4]: # Reduce the number of rows to around 3000
data = data.sample(frac=0.02, random_state=42)

In [5]: print(data.columns)

Index(['Job Id', 'Experience', 'Qualifications', 'Salary Range', 'Location', 'Country', 'latitude', 'longitude', 'Work Type', 'Job Posting Date', 'Company Size', 'Contact', 'Job Title', 'Role', 'Job Portal', 'Job Description', 'Benefits', 'skills', 'responsibilities', 'Company', 'Company Profile'],
      dtype='object')
```

```
In [6]: data.info()

<class 'pandas.core.frame.DataFrame'>
Int64Index: 3000 entries, 161940 to 161889
Data columns (total 23 columns):
#   Column              Non-Null Count  Dtype
---  --
0   Job Id              3000 non-null   int64
1   Experience           3000 non-null   object
2   Qualifications       3000 non-null   object
3   Salary Range        3000 non-null   object
4   Location            3000 non-null   object
5   Country             3000 non-null   object
6   latitude            3000 non-null   float64
7   longitude           3000 non-null   float64
8   Work Type           3000 non-null   object
9   Company Size        3000 non-null   int64
10  Job Posting Date     3000 non-null   object
11  Preference           3000 non-null   object
12  Contact Person       3000 non-null   object
13  Contact             3000 non-null   object
14  Job Title           3000 non-null   object
15  Role                3000 non-null   object
16  Job Portal          3000 non-null   object
17  Job Description      3000 non-null   object
18  Benefits            3000 non-null   object
19  skills              3000 non-null   object
20  Responsibilities       3000 non-null   object
21  Company             3000 non-null   object
22  Company Profile     2991 non-null   object
dtypes: float64(2), int64(2), object(19)
memory usage: 562.3+ KB
```

Cleaning the data

```
In [8]: # checking for missing values
data.isna().sum()

Out[8]:
Job Id      0
Experience   0
Qualifications  0
Salary Range  0
Location     0
Country      0
latitude     0
longitude    0
Work Type   0
Company Size  0
Job Posting Date  0
Preference   0
Contact Person  0
Contact      0
Job Title    0
Role         0
Job Portal   0
Job Description  0
Benefits     0
skills       0
Responsibilities  0
Company      0
Company Profile  0
dtype: int64
```

```
In [9]: # Dropping unnecessary variables
# Drop unnecessary columns
unnecessary_columns = ['Job Id', 'Salary Range', 'Location', 'Country', 'latitude', 'longitude', 'Work Type', 'Job Posting Date', 'Preference', 'Contact Person', 'Contact', 'Role', 'Job Portal', 'Benefits', 'Company', 'Company Profile']
data.drop(unnecessary_columns, inplace=True)
```

```
Out[15]: data.head(1)

Out[15]:
```

	Experience	Qualifications	Company Size	Job Title	Job Description	skills	Responsibilities	Combined_Text
161878	1 to 9 Years	BA	95517	Human Resources Manager	Talent Acquisition Managers oversee the recruitment and hiring process for an organization.	Talent sourcing, interviewing, onboarding, recruiting.	Lead talent acquisition efforts, including recruiting, hiring, and onboarding new employees.	Talent acquisition manager oversee recruitment...

pre-processing the data

```
In [10]: import nltk
from nltk.corpus import stopwords
from nltk.tokenize import word_tokenize
from nltk.stem import WordNetLemmatizer
import string

# Download stopwords and initialize lemmatizer
nltk.download('stopwords')
nltk.download('punkt')
stop_words = set(stopwords.words('english'))
lemmatizer = WordNetLemmatizer()

# Text preprocessing
def preprocess(text):
    # Remove punctuation marks
    text = text.translate(str.maketrans('', '', string.punctuation))

    # Tokenization
    tokens = word_tokenize(text.lower())

    # Remove stopwords and lemmatization
    tokens = [lemmatizer.lemmatize(token) for token in tokens if token not in stop_words]

    return " ".join(tokens)

[nltk_data] Downloading package stopwords to
[nltk_data] C:\Users\User\AppData\Local\nltk_data...
[nltk_data] Package stopwords is already up-to-date!
[nltk_data] Downloading package punkt to
[nltk_data] C:\Users\User\AppData\Local\nltk_data...
[nltk_data] Package punkt is already up-to-date!
[nltk_data] Downloading package wordnet to
[nltk_data] C:\Users\User\AppData\Local\nltk_data...
[nltk_data] Package wordnet is already up-to-date!
```

Feature Engineering

```
In [11]: # Apply preprocessing to relevant columns and concatenate text
data['Combined_Text'] = data[['Job Description', 'skills', 'Responsibilities', 'Job Title']].apply(preprocess, axis=1)

# Feature extraction
# TF-IDF vectorization for combined text data
tfidf_vectorizer = TfidfVectorizer(max_features=1000)
X_text = tfidf_vectorizer.fit_transform(data['Combined_Text']).toarray()

# Encoding target variable
label_encoder = LabelEncoder()
y = label_encoder.fit_transform(data['Job Title'])
```

Model training

```
In [13]: # Train-test split
X_train, X_test, y_train, y_test = train_test_split(X_text, y, test_size=0.2, random_state=42)

# Model training
model = LogisticRegression(max_iter=1000)
model.fit(X_train, y_train)

# Model evaluation
y_pred = model.predict(X_test)
accuracy = accuracy_score(y_test, y_pred)
print("Logistic Regression Accuracy:", accuracy)
print(classification_report(y_test, y_pred))

from sklearn.metrics import roc_auc_score

# Train and evaluate Support Vector Machine (SVM) model
svm_model = SVC()
svm_model.fit(X_train, y_train)
svm_y_pred = svm_model.predict(X_test)
svm_accuracy = accuracy_score(y_test, svm_y_pred)
print("SVM Accuracy:", svm_accuracy)
print(classification_report(y_test, svm_y_pred))

# Train and evaluate Random Forest model
rf_model = RandomForestClassifier()
rf_model.fit(X_train, y_train)
rf_y_pred = rf_model.predict(X_test)
rf_accuracy = accuracy_score(y_test, rf_y_pred)
print("Random Forest Accuracy:", rf_accuracy)
print(classification_report(y_test, rf_y_pred))

logistic_regression_accuracy = 0.9083333333333333
svm_accuracy = 0.9946666666666667
```

	precision	recall	f1-score	support
0	1.00	0.50	0.67	6
1	0.60	1.00	0.75	3
2	1.00	1.00	1.00	9
3	1.00	0.33	0.50	3
4	1.00	1.00	1.00	11
5	1.00	1.00	1.00	2
6	1.00	1.00	1.00	7
7	1.00	1.00	1.00	2
8	1.00	1.00	1.00	2
9	1.00	1.00	1.00	6
10	1.00	1.00	1.00	3
11	0.67	1.00	0.80	5
12	1.00	1.00	1.00	5
13	1.00	1.00	1.00	2
14	0.50	1.00	0.67	3
15	1.00	1.00	1.00	5
16	0.50	1.00	0.67	2
17	1.00	1.00	1.00	5
18	1.00	1.00	1.00	6
19	1.00	1.00	1.00	1
20	1.00	1.00	1.00	8
21	1.00	1.00	1.00	2
22	1.00	1.00	1.00	6
23	1.00	1.00	1.00	7
24	1.00	1.00	1.00	9
25	1.00	1.00	1.00	6
26	1.00	1.00	1.00	3
27	1.00	1.00	1.00	3
28	1.00	1.00	1.00	5
29	1.00	1.00	1.00	7
30	1.00	1.00	1.00	5
31	1.00	1.00	1.00	10
32	1.00	1.00	1.00	2
33	1.00	1.00	1.00	4
34	1.00	1.00	1.00	2
35	1.00	1.00	1.00	7
36	1.00	1.00	1.00	1
37	1.00	1.00	1.00	7
38	1.00	1.00	1.00	1
39	1.00	1.00	1.00	5
40	1.00	1.00	1.00	7
41	1.00	1.00	1.00	1
42	1.00	1.00	1.00	8
43	1.00	1.00	1.00	5
44	1.00	1.00	1.00	6
45	1.00	1.00	1.00	1
46	1.00	1.00	1.00	5
47	1.00	1.00	1.00	1
48	1.00	1.00	1.00	5
49	1.00	1.00	1.00	1
50	0.89	1.00	0.94	8
51	0.89	1.00	0.94	8
52	0.00	0.00	0.00	1
53	1.00	1.00	1.00	1
54	1.00	1.00	1.00	4
55	1.00	1.00	1.00	4
56	1.00	1.00	1.00	4
57	1.00	1.00	1.00	1
58	1.00	1.00	1.00	6
59	1.00	1.00	1.00	1
60	1.00	1.00	1.00	9
61	0.53	1.00	0.69	10
62	1.00	1.00	1.00	5
63	1.00	1.00	1.00	5
64	1.00	1.00	1.00	2
65	0.75	1.00	0.83	3
66	0.00	0.00	0.00	1
67	1.00	0.71	0.83	7
68	0.86	1.00	0.92	6
69	1.00	1.00	1.00	4
70	1.00	1.00	1.00	4
71	0.82	1.00	0.90	9
72	1.00	1.00	1.00	3
73	1.00	0.67	0.80	3
74	1.00	1.00	1.00	1
75	0.55	1.00	0.71	6
76	0.33	1.00	0.52	4
77	1.00	0.67	0.80	3
78	1.00	1.00	1.00	1
79	1.00	1.00	1.00	5
80	1.00	1.00	1.00	3
81	1.00	1.00	1.00	1
82	0.00	0.00	0.00	1
83	0.00	0.00	0.00	1
84	0.83	1.00	0.91	5
85	1.00	0.50	0.67	6
86	1.00	1.00	1.00	3
87	1.00	1.00	1.00	5
88	1.00	1.00	1.00	5
89	1.00	1.00	1.00	5
90	1.00	1.00	1.00	7
91	1.00	0.80	0.89	5
92	0.73	1.00	0.84	6
93	0.00	0.00	0.00	4
94	1.00	1.00	1.00	4
95	1.00	1.00	1.00	1
96	1.00	1.00	1.00	1
97	1.00	1.00	1.00	4
98	1.00	1.00	1.00	1
99	1.00	1.00	1.00	4
100	1.00	1.00	1.00	9
101	1.00	1.00	1.00	2
102	1.00	1.00	1.00	2
103	1.00	1.00	1.00	6
104	1.00	1.00	1.00	4
105	1.00	1.00	1.00	2
106	1.00	1.00	1.00	2
107	0.69	1.00	0.84	3
108	1.00	0.33	0.50	3
109	0.00	0.00	0.00	2
110	0.50	0.50	0.50	4
111	1.00	1.00	1.00	4
112	1.00	1.00	1.00	2
113	1.00	1.00	1.00	4
114	1.00	1.00	1.00	2
115	1.00	1.00	1.00	5
116	1.00	1.00	1.00	11
117	1.00	1.00	1.00	1
118	1.00	1.00	1.00	2
119	1.00	1.00	1.00	1
120	1.00	1.00	1.00	3
121	1.00	1.00	1.00	3
122	1.00	1.00	1.00	3
123	1.00	1.00	1.00	4
124	1.00	1.00	1.00	8
125	1.00	1.00	1.00	2
126	1.00	1.00	1.00	2
127	1.00	1.00	1.00	2
128	1.00	1.00	1.00	2
129	1.00	1.00	1.00	1
130	1.00	1.00	1.00	2
131	1.00	1.00	1.00	2
132	1.00	1.00	1.00	8
133	1.00	1.00	1.00	8
134	1.00	1.00	1.00	2
135	1.00	1.00	1.00	2
136	1.00	1.00	1.00	2
137	1.00	1.00	1.00	2
138	1.00	1.00	1.00	3
139	1.00	1.00	1.00	9
140	1.00	1.00	1.00	4
141	1.00	1.00	1.00	19
142	1.00	1.00	1.00	2
143	1.00	1.00	1.00	2
144	1.00	1.00	1.00	1
145	1.00	1.00	1.00	1
146	0.00	0.00	0.00	1
accuracy				
macro avg	0.86	0.86	0.84	600
weighted avg	0.91	0.91	0.89	600

C:\Users\User\anaconda3\Lib\site-packages\sklearn\metrics_classification.py:1469: UndefinedMetricWarning: Precision and F-score are ill-defined and being set to 0.0 in labels with no predicted samples. Use 'zero_division' parameter to control this behavior.

Warning: UndefinedMetricWarning: Precision and F-score are ill-defined and being set to 0.0 in labels with no predicted samples. Use 'zero_division' parameter to control this behavior.

C:\Users\User\anaconda3\Lib\site-packages\sklearn\metrics_classification.py:1469: UndefinedMetricWarning: Precision and F-score are ill-defined and being set to 0.0 in labels with no predicted samples. Use 'zero_division' parameter to control this behavior.

C:\Users\User\anaconda3\Lib\site-packages\sklearn\metrics_classification.py:1469: UndefinedMetricWarning: Precision and F-score are ill-defined and being set to 0.0 in labels with no predicted samples. Use 'zero_division' parameter to control this behavior.

Warning: UndefinedMetricWarning: Precision and F-score are ill-defined and being set to 0.0 in labels with no predicted samples. Use 'zero_division' parameter to control this behavior.

SVM Accuracy: 0.9946666666666667

66	0.00	0.00	0.00	1
67	1.00	0.71	0.83	7
68	0.86	1.00	0.92	6
69	1.00	1.00	1.00	4
70	1.00	1.00	1.00	2
71	0.82	1.00	0.90	9
73	1.00	1.00	1.00	5