# IMPORT LIBRARIES AND DATASET

In [6]:
```python
import plotly.express as px
import pandas as pd
import matplotlib.pyplot as plt
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
from math import pi
import string
import re
import tensorflow as tf
from surprise import Dataset, Reader
from sklearn.model_selection import train_test_split
from surprise import KNNBasic
from surprise.accuracy import rmse
from sklearn.preprocessing import StandardScaler
from sklearn.neighbors import KNeighborsClassifier
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import roc_curve, auc
import matplotlib.pyplot as plt
from surprise import Dataset, Reader, KNNBasic
from surprise.model_selection import train_test_split
from surprise.accuracy import rmse
import researchpy as rp
from pgmpy.models import BayesianModel
from pgmpy.estimators import MaximumLikelihoodEstimator
from pgmpy.inference import VariableElimination
from nltk.corpus import stopwords
from sklearn.linear_model import LogisticRegression
from sklearn.neighbors import KNeighborsClassifier
from sklearn.model_selection import train_test_split
from sklearn.metrics import roc_curve, roc_auc_score
from sklearn.preprocessing import LabelEncoder
from scipy.stats import chi2_contingency
```

```python
from matplotlib.colors import Normalize
from sklearn.model_selection import train_test_split
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import accuracy_score
from sklearn.preprocessing import MinMaxScaler
from sklearn.impute import SimpleImputer
import warnings
import plotly.graph_objects as go
from sklearn.preprocessing import MaxAbsScaler
from sklearn.preprocessing import MaxAbsScaler
from sklearn.preprocessing import LabelEncoder, OneHotEncoder
from sklearn.preprocessing import StandardScaler, LabelEncoder
from sklearn.model_selection import GridSearchCV
from sklearn.ensemble import RandomForestClassifier
from sklearn.preprocessing import LabelEncoder, StandardScaler
from sklearn.model_selection import train_test_split, GridSearchCV
from sklearn.feature_selection import SelectKBest, f_classif
from sklearn.metrics import confusion_matrix, accuracy_score
from sklearn.model_selection import cross_val_score
import statsmodels.api as sm
import squarify
import math
from sklearn.metrics.pairwise import cosine_similarity
from sklearn.impute import SimpleImputer, MissingIndicator
from sklearn.model_selection import train_test_split
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import accuracy_score
import pandas as pd
import numpy as np
import pandas as pd
import numpy as np
from sklearn.model_selection import train_test_split
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.model_selection import train_test_split
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import classification_report
from sklearn.linear_model import LogisticRegression
from sklearn.tree import DecisionTreeClassifier
```

```
74  from sklearn.feature_extraction.text import TfidfVectorizer
75  from sklearn.metrics.pairwise import cosine_similarity
76  from sklearn.metrics import precision_score, recall_score, f1_score, roc_auc_score
77  # Read the Excel file into a DataFrame
78  data= pd.read_excel(r"C:\Users\Administrator\Desktop\DATA.xlsx")
79  # display the dataframe
80  data.head()
```

Out[6]:

| | user ID | Name | Age | Gender | Academic Background | Field of Study | Skills | Industry_Interest | Job_Type_Interest | Location_Interest | Salary_Expectation_(in_USD) | d |
|---|---------|------|-----|--------|---------------------|----------------|--------|-------------------|-------------------|-------------------|-----------------------------|---|
| 0 | 986206 | John Smith | 28 | Male | Bachelor's degree | Computer Science 2 | Java, Python, Data Structures | Technology | Full-time | New York | 65000 | |
| 1 | 769632 | Jane Doe | 42 | Female | Master's degree | Business | Finance, Accounting, Microsoft Excel | Finance | Contract | London | 85000 | |
| 2 | 981314 | David Lee | 35 | Male | bachelor"s degree | NaN | Sales, Customer Service, Communication | Retail | Part-time | Chicago | 30000 | |
| 3 | 962892 | Sarah Johnson | 27 | Female | Associate's degree | Nursing | Patient Care, Medical Terminology | Healthcare | Full-time | Los Angeles | 45000 | |
| 4 | 967782 | Michael Williams | 46 | Male | Bachelor's degree | Marketing | Digital Marketing, Social Media | Marketing | Freelance | Toronto | 70000 | |

In [7]:  `data.shape`

Out[7]: (999, 12)

In [8]:  `data.info()`

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 999 entries, 0 to 998
Data columns (total 12 columns):
 #   Column                 Non-Null Count  Dtype
---  ------                 --------------  -----
```

```
In [8]:   1  data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 999 entries, 0 to 998
Data columns (total 12 columns):
 #   Column                    Non-Null Count  Dtype
---  ------                    --------------  -----
 0   user ID                   999 non-null    int64
 1   Name                      996 non-null    object
 2   Age                       999 non-null    int64
 3   Gender                    927 non-null    object
 4   Academic Background       999 non-null    object
 5   Field of Study            963 non-null    object
 6   Skills                    999 non-null    object
 7   Industry_Interest         999 non-null    object
 8   Job_Type_Interest         999 non-null    object
 9   Location_Interest         995 non-null    object
 10  Salary_Expectation_(in_USD)  999 non-null  int64
 11  desired_company           999 non-null    object
dtypes: int64(3), object(9)
memory usage: 93.8+ KB
```

## DESCRIPTIVE STATISTICS

```
In [9]:   1  # Get descriptive statistics
          2  descriptive_stats = data.describe()
          3
          4  # Display the descriptive statistics
          5  print(descriptive_stats)
```

```
             user ID           Age  Salary_Expectation_(in_USD)
count     999.000000    999.000000                   999.000000
mean   931435.157157     30.625626                 65385.385385
std    158544.732526      9.187644                 28965.576740
min       382.000000      0.000000                 20000.000000
25%    962133.000000     24.000000                 45000.000000
50%    972695.000000     30.000000                 65000.000000
```

# CHECKING FOR MISSING VALUES

In [10]:
```python
# Check for missing values in the dataset
missing_values = data.isnull().sum()

# Display the count of missing values for each column
print(missing_values)
```

```
user ID                        0
Name                           3
Age                            0
Gender                        72
Academic Background            0
Field of Study                36
Skills                         0
Industry_Interest              0
Job_Type_Interest              0
Location_Interest              4
Salary_Expectation_(in_USD)    0
desired_company                0
dtype: int64
```

In [11]:
```python
# Impute missing values for numerical columns with the mean
numerical_columns = data.select_dtypes(include='number').columns
data[numerical_columns] = data[numerical_columns].fillna(data[numerical_columns].mean())

# Impute missing values for categorical columns with the most frequent value
categorical_columns = data.select_dtypes(include='object').columns
data[categorical_columns] = data[categorical_columns].fillna(data[categorical_columns].mode().iloc[0])
```
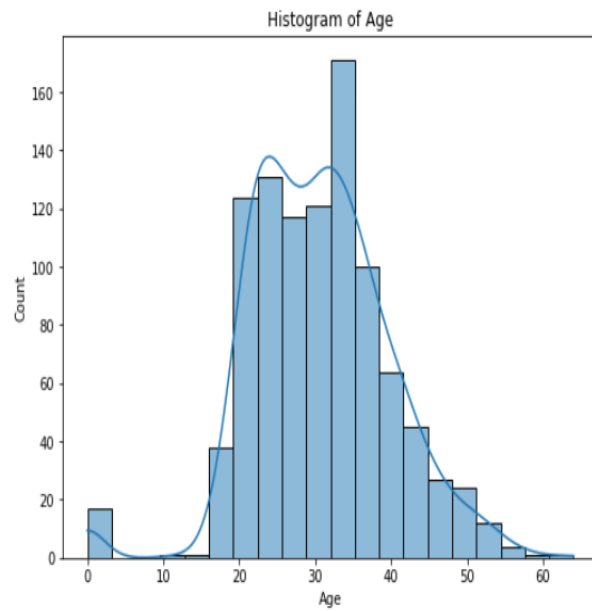
In [12]:
```python
# Check for missing values in the dataset
missing_values = data.isnull().sum()

# Display the count of missing values for each column
print(missing_values)
```

```python
In [12]:  1  # Check for missing values in the dataset
          2  missing_values = data.isnull().sum()
          3
          4  # Display the count of missing values for each column
          5  print(missing_values)
```

```
user ID                        0
Name                           0
Age                            0
Gender                         0
Academic Background            0
Field of Study                 0
Skills                         0
Industry_Interest              0
Job_Type_Interest              0
Location_Interest              0
Salary_Expectation_(in_USD)    0
desired_company                0
dtype: int64
```

```python
In [13]:  1  # Plot histogram of Age
          2  plt.figure(figsize=(8, 6))
          3  sns.histplot(data['Age'], bins=20, kde=True)
          4  plt.xlabel('Age')
          5  plt.ylabel('Count')
          6  plt.title('Histogram of Age')
          7  plt.show()
```
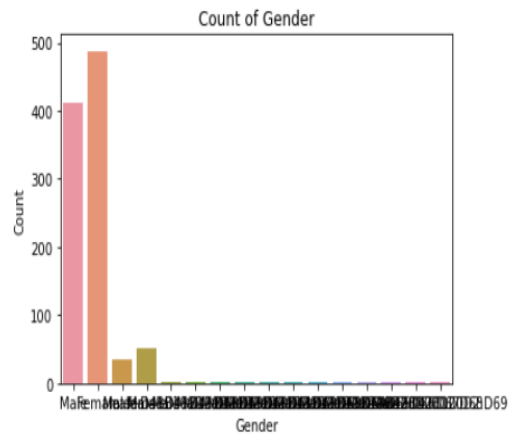
```
In [13]: 1  # Plot histogram of Age
         2  plt.figure(figsize=(8, 6))
         3  sns.histplot(data['Age'], bins=20, kde=True)
         4  plt.xlabel('Age')
         5  plt.ylabel('Count')
         6  plt.title('Histogram of Age')
         7  plt.show()
```
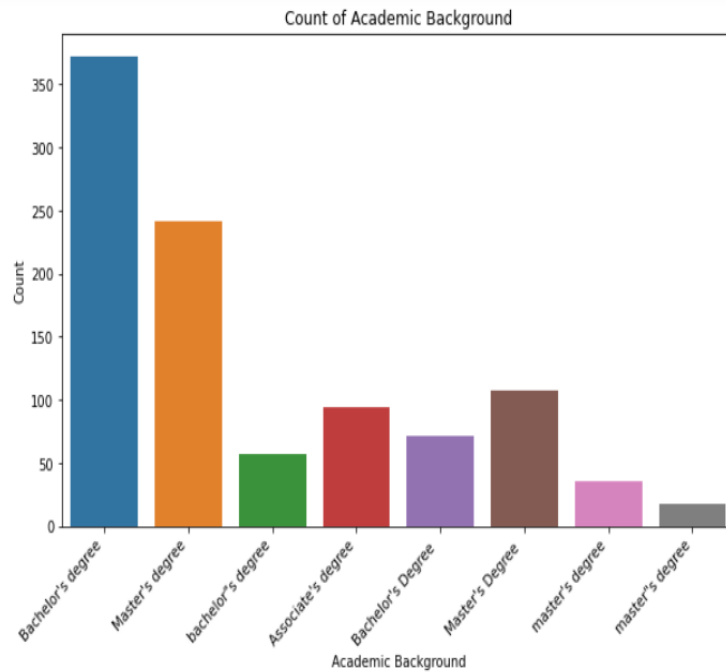


Histogram of Age

```
1  # Plot count of Gender
2  plt.figure(figsize=(6, 4))
3  sns.countplot(data['Gender'])
4  plt.xlabel('Gender')
5  plt.ylabel('Count')
6  plt.title('Count of Gender')
7  plt.show()
```

C:\ProgramData\Anaconda3\lib\site-packages\seaborn\_decorators.py:36: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be `data`, and passing other arguments without an explicit keyword will result in an error or misinterpretation.
  warnings.warn(

```
1  # Plot count of Academic Background
2  plt.figure(figsize=(10, 6))
3  sns.countplot(data['Academic Background'])
4  plt.xlabel('Academic Background')
5  plt.ylabel('Count')
6  plt.title('Count of Academic Background')
7  plt.xticks(rotation=45, ha='right')
8  plt.show()
```

## Count of Academic Background
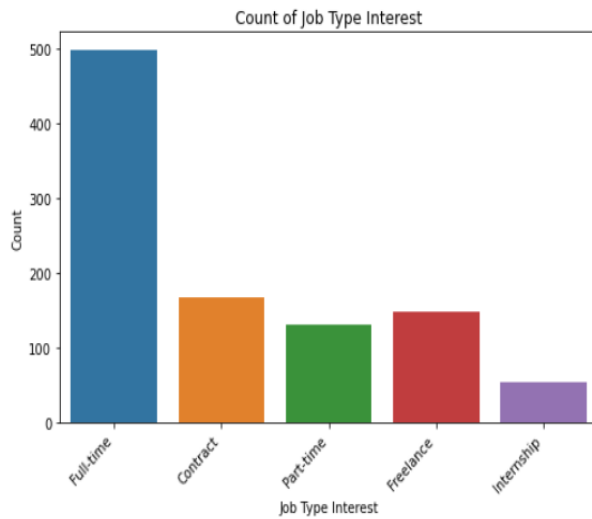


```
In [17]:   1  # Plot count of Job Type Interest
           2  plt.figure(figsize=(8, 5))
           3  sns.countplot(data['Job_Type_Interest'])
           4  plt.xlabel('Job Type Interest')
           5  plt.ylabel('Count')
           6  plt.title('Count of Job Type Interest')
           7  plt.xticks(rotation=45, ha='right')
           8  plt.show()
```

In [17]:
```python
# Plot count of Job Type Interest
plt.figure(figsize=(8, 5))
sns.countplot(data['Job_Type_Interest'])
plt.xlabel('Job Type Interest')
plt.ylabel('Count')
plt.title('Count of Job Type Interest')
plt.xticks(rotation=45, ha='right')
plt.show()
```

C:\ProgramData\Anaconda3\lib\site-packages\seaborn\_decorators.py:36: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be `data`, and passing other arguments without an explicit keyword will result in an error or misinterpretation.
  warnings.warn(

In [19]:
```python
# Plot box plot of Salary Expectation
plt.figure(figsize=(8, 6))
sns.boxplot(data['Salary_Expectation_(in_USD)'])
plt.xlabel('Salary Expectation (in USD)')
plt.title('Box Plot of Salary Expectation')
plt.show()
```

C:\ProgramData\Anaconda3\lib\site-packages\seaborn\_decorators.py:36: FutureWarning: Pass the following variable as a keyword a
rg: x. From version 0.12, the only valid positional argument will be `data`, and passing other arguments without an explicit ke
yword will result in an error or misinterpretation.
  warnings.warn(



Box Plot of Salary Expectation

```
In [20]:    1  # Plot correlation heatmap for numerical columns
            2  plt.figure(figsize=(10, 8))
            3  sns.heatmap(data.corr(), annot=True, cmap='coolwarm', fmt='.2f')
            4  plt.title('Correlation Heatmap')
            5  plt.show()
```

```
In [22]:  1 plt.figure(figsize=(10, 6))
          2 sns.countplot(data['Industry_Interest'])
          3 plt.xlabel('Industry_Interest')
          4 plt.ylabel('Count')
          5 plt.title('Count of Students in Each Industry of Interest')
          6 plt.xticks(rotation=45, ha='right')
          7 plt.show()
```

C:\ProgramData\Anaconda3\lib\site-packages\seaborn\_decorators.py:36: FutureWarning: Pass the following variable as a keyword a
rg: x. From version 0.12, the only valid positional argument will be `data`, and passing other arguments without an explicit ke
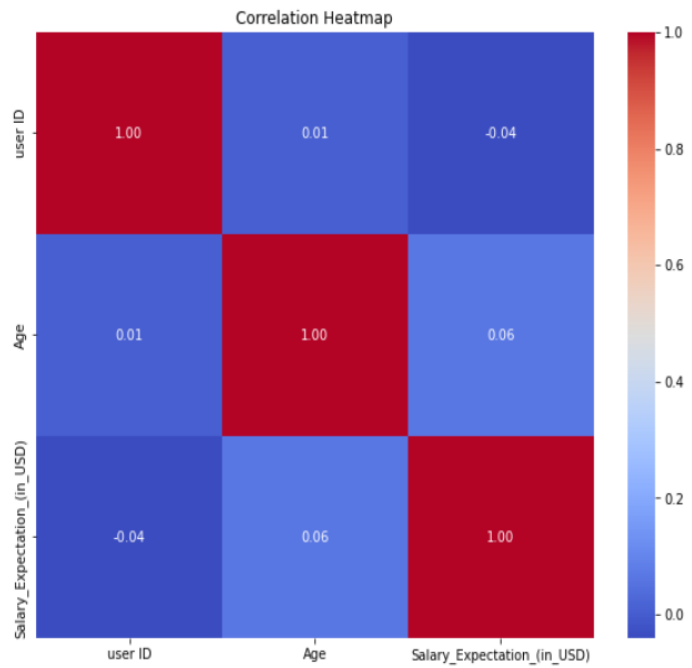yword will result in an error or misinterpretation.
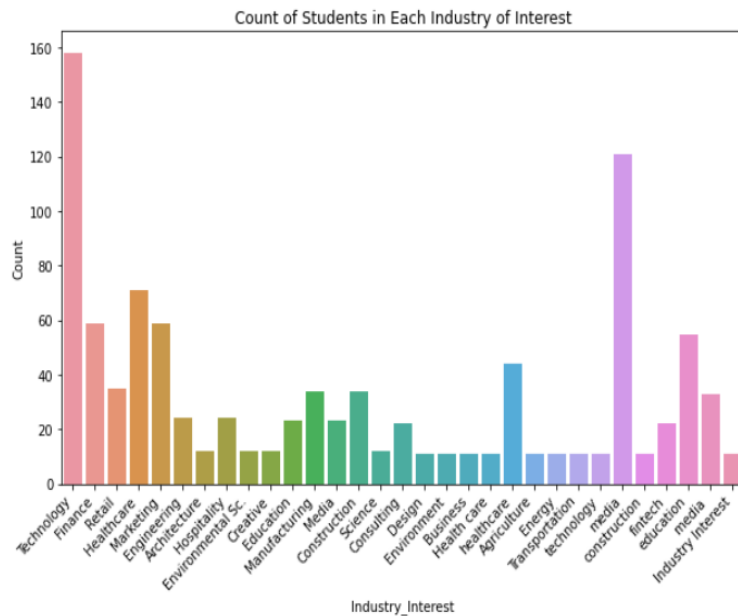  warnings.warn(

```
In [24]:  1 sns.pairplot(data, vars=['Age', 'Salary_Expectation_(in_USD)'], hue='Gender')
          2 plt.suptitle('Pairplot of Age and Salary Expectation (colored by Gender)')
          3 plt.show()
```



Pairplot of Age and Salary Expectation (colored by Gender)

```
In [27]:  1 plt.figure(figsize=(10, 6))
          2 sns.violinplot(x='Job_Type_Interest', y='Salary_Expectation_(in_USD)', data=data)
          3 plt.xlabel('Job Type Interest')
          4 plt.ylabel('Salary_Expectation_(in_USD)')
          5 plt.title('Distribution of Salary Expectation across Job Types')
          6 plt.xticks(rotation=45, ha='right')
          7 plt.show()
```



Distribution of Salary Expectation across Job Types

```python
1  plt.figure(figsize=(10, 6))
2  sns.boxplot(x='Academic Background', y='Salary_Expectation_(in_USD)', data=data)
3  plt.xlabel('Academic Background')
4  plt.ylabel('Salary_Expectation_(in_USD)')
5  plt.title('Comparison of Salary Expectation for Different Academic Backgrounds')
6  plt.xticks(rotation=45, ha='right')
7  plt.show()
```



Comparison of Salary Expectation for Different Academic Backgrounds

```python
1 plt.figure(figsize=(8, 6))
2 sns.scatterplot(x='Age', y='Salary_Expectation_(in_USD)', data=data, hue='Gender')
3 plt.xlabel('Age')
4 plt.ylabel('Salary_Expectation_(in_USD)')
5 plt.title('Scatter Plot of Age vs. Salary Expectation (colored by Gender)')
6 plt.show()
```



Scatter Plot of Age vs. Salary Expectation (colored by Gender)

```
In [34]:   1  #One-Hot Encoding
           2  one_hot_encoded_data = pd.get_dummies(data, columns=['Academic Background', 'Field of Study', 'Industry_Interest', 'Job_Type
           3
           4
           5  label_encoder = LabelEncoder()
           6  data['Gender_encoded'] = label_encoder.fit_transform(data['Gender'])
           7
           8  # Display the encoded datasets
           9  print("One-Hot Encoded Data:")
          10  print(one_hot_encoded_data.head())
          11
          12  print("\nLabel Encoded Data:")
          13  print(data[['Gender', 'Gender_encoded']].head())
```

```
One-Hot Encoded Data:
    user ID              Name  Age  Gender  \
0   986206        John Smith   28    Male
1   769632          Jane Doe   42  Female
2   981314         David Lee   35    Male
3   962892     Sarah Johnson   27  Female
4   967782  Michael Williams   46    Male


                                Skills  Salary_Expectation_(in_USD)  \
0          Java, Python, Data Structures                        65000
1      Finance, Accounting, Microsoft Excel                    85000
2   Sales, Customer Service, Communication                     30000
3        Patient Care, Medical Terminology                     45000
4          Digital Marketing, Social Media                     70000

    Academic Background_Associate's degree  \
0                                        0
1                                        0
2                                        0
3                                        1
4                                        0

    Academic Background Bachelor's Degree    \
```

```
In [36]:   1  #Extracting First Name from Name
           2  data['First Name'] = data['Name'].str.split().str[0]
           3
           4  # Feature 2: Converting 'Age' to Age Group
           5  def get_age_group(age):
           6      if age < 25:
           7          return 'Young'
           8      elif age >= 25 and age < 40:
           9          return 'Middle-aged'
          10      else:
          11          return 'Senior'
          12
          13  data['Age Group'] = data['Age'].apply(get_age_group)
          14
          15  # Feature 3: Counting the number of skills each student has
          16  data['Number of Skills'] = data['Skills'].apply(lambda x: len(x.split(',')))
          17
          18  # Feature 4: Encoding 'Full-time' as 1 and 'Part-time' as 0 for Job Type Interest
          19  data['Job_Type_Interest'] = data['Job_Type_Interest'].apply(lambda x: 1 if x == 'Full-time' else 0)
          20
          21  # Feature 5: Encoding 'Male' as 1 and 'Female' as 0 for Gender
          22  data['Gender'] = data['Gender'].apply(lambda x: 1 if x == 'Male' else 0)
          23
          24  # Display the updated dataset with engineered features
          25  print(data.head())
```

```
          Field of Study                              Skills  \
0   Computer Science 2          Java, Python, Data Structures
1             Business    Finance, Accounting, Microsoft Excel
2     Computer Science  Sales, Customer Service, Communication
3              Nursing         Patient Care, Medical Terminology
4            Marketing         Digital Marketing, Social Media

   Industry_Interest  Job_Type_Interest Location_Interest  \
0         Technology                  1          New York
1            Finance                  0            London
2             Retail                  0           Chicago
3         Healthcare                  1       Los Angeles
4          Marketing                  0           Toronto
```

```
In [37]:   1  # Prepare the feature matrix X and target vector y
           2  X = data[['Age', 'Gender']]  # Select relevant features
           3  y = data['Job_Type_Interest']  # Target variable
           4
           5  # Split the data into training and testing sets (80% train, 20% test)
           6  X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
           7
           8  # Initialize and train the KNN classifier
           9  knn_classifier = KNeighborsClassifier(n_neighbors=5)
          10  knn_classifier.fit(X_train, y_train)
          11
          12  # Make predictions on the test set
          13  y_pred = knn_classifier.predict(X_test)
          14
          15  # Evaluate the classifier's performance
          16  accuracy = accuracy_score(y_test, y_pred)
          17  print("KNN Classifier Accuracy:", accuracy)

KNN Classifier Accuracy: 0.525
```

```
In [45]:   1  # Import required libraries
           2  from surprise import Dataset, Reader
           3  from surprise.model_selection import train_test_split
           4  from surprise import KNNBasic
           5  from surprise.accuracy import rmse
           6
           7  # Prepare the data for Surprise library
           8  reader = Reader(rating_scale=(0, 1))
           9  surprise_data = Dataset.load_from_df(data[['user ID', 'Salary_Expectation_(in_USD)', 'Age']], reader)
          10
          11  # Split the data into training and testing sets (80% train, 20% test)
          12  trainset, testset = train_test_split(surprise_data, test_size=0.2, random_state=42)
          13
          14  # Initialize and train the KNNBasic collaborative filtering model
          15  knn_collaborative = KNNBasic(sim_options={'user_based': True})
          16  knn_collaborative.fit(trainset)
          17
          18  # Make predictions on the test set
          19  predictions = knn_collaborative.test(testset)
```

```
In [45]:   1  # Import required libraries
           2  from surprise import Dataset, Reader
           3  from surprise.model_selection import train_test_split
           4  from surprise import KNNBasic
           5  from surprise.accuracy import rmse
           6
           7  # Prepare the data for Surprise library
           8  reader = Reader(rating_scale=(0, 1))
           9  surprise_data = Dataset.load_from_df(data[['user ID', 'Salary_Expectation_(in_USD)', 'Age']], reader)
          10
          11  # Split the data into training and testing sets (80% train, 20% test)
          12  trainset, testset = train_test_split(surprise_data, test_size=0.2, random_state=42)
          13
          14  # Initialize and train the KNNBasic collaborative filtering model
          15  knn_collaborative = KNNBasic(sim_options={'user_based': True})
          16  knn_collaborative.fit(trainset)
          17
          18  # Make predictions on the test set
          19  predictions = knn_collaborative.test(testset)
          20
          21  # Evaluate the model's performance (Root Mean Squared Error, RMSE)
          22  rmse_score = rmse(predictions)
          23  print("Collaborative Filtering RMSE:", rmse_score)
```

```
Computing the msd similarity matrix...
Done computing similarity matrix.
RMSE: 31.6624
Collaborative Filtering RMSE: 31.662438314191785
```

```
In [49]:   1  # Import required libraries
           2  from surprise import Dataset, Reader
           3  from surprise.model_selection import train_test_split
           4  from surprise import KNNBasic
           5  from surprise.accuracy import rmse
           6
           7  # Prepare the data for Surprise library
           8  reader = Reader(rating_scale=(0, 1))
           9  surprise_data = Dataset.load_from_df(data[['user ID', 'Salary_Expectation_(in_USD)', 'Age']], reader)
          10
          11  # Split the data into training and testing sets (80% train, 20% test)
```

```
11  # Split the data into training and testing sets (80% train, 20% test)
12  trainset, testset = train_test_split(surprise_data, test_size=0.2, random_state=42)
13
14  # Initialize and train the KNNBasic collaborative filtering model
15  knn_collaborative = KNNBasic(sim_options={'user_based': True})
16  knn_collaborative.fit(trainset)
17
18  # Make predictions on the test set
19  predictions = knn_collaborative.test(testset)
20
21  # Evaluate the model's performance (Root Mean Squared Error, RMSE)
22  rmse_score = rmse(predictions)
23  print("Collaborative Filtering RMSE:", rmse_score)
```

```
Computing the msd similarity matrix...
Done computing similarity matrix.
RMSE: 31.6624
Collaborative Filtering RMSE: 31.662438314191785
```

In [51]:

```
1   from sklearn.model_selection import train_test_split
2   from sklearn.preprocessing import StandardScaler
3   from sklearn.neighbors import KNeighborsClassifier
4   from sklearn.ensemble import RandomForestClassifier
5   from sklearn.metrics import roc_curve, auc
6
7   # Prepare the feature matrix X and target vector y
8   X = data[['Age', 'Gender', 'Number of Skills']]  # Select relevant features
9   y = data['Job_Type_Interest']  # Target variable
10
11  # Split the data into training and testing sets (80% train, 20% test)
12  X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
13
14  # Standardize the feature data for better training performance (not necessary for some algorithms)
15  scaler = StandardScaler()
16  X_train_scaled = scaler.fit_transform(X_train)
17  X_test_scaled = scaler.transform(X_test)
18
19  # Initialize models for different algorithms
20  knn_classifier = KNeighborsClassifier(n_neighbors=5)
21  knn_classifier.fit(X_train_scaled, y_train)
```
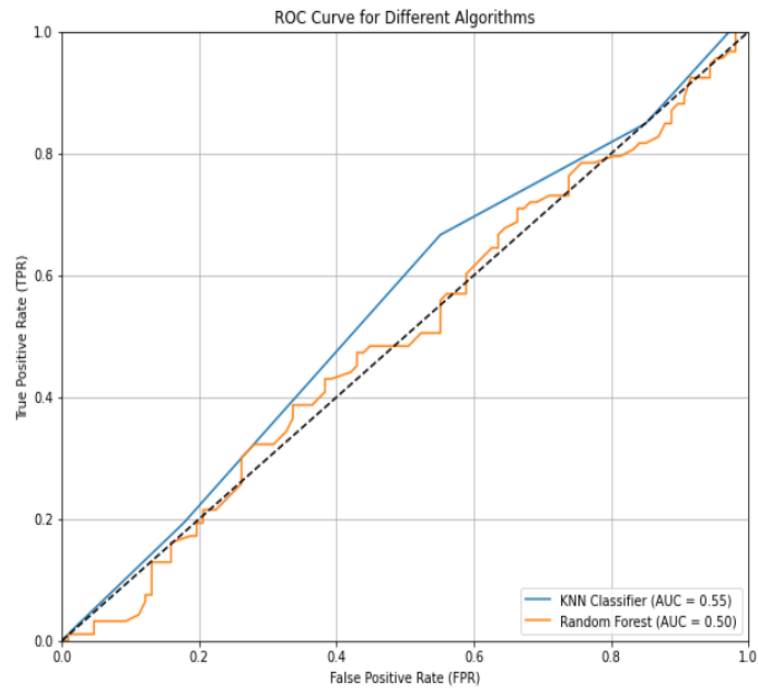
```python
19  # Initialize models for different algorithms
20  knn_classifier = KNeighborsClassifier(n_neighbors=5)
21  knn_classifier.fit(X_train_scaled, y_train)
22
23  # Random Forest
24  rf_classifier = RandomForestClassifier(n_estimators=100, random_state=42)
25  rf_classifier.fit(X_train_scaled, y_train)
26
27  # Calculate ROC curves and AUC for each model
28  models = [knn_classifier, rf_classifier]
29  model_names = ['KNN Classifier', 'Random Forest']
30
31  plt.figure(figsize=(10, 8))
32
33  for model, name in zip(models, model_names):
34      y_pred_prob = model.predict_proba(X_test_scaled)[:, 1]
35      fpr, tpr, _ = roc_curve(y_test, y_pred_prob)
36      auc_score = auc(fpr, tpr)
37      plt.plot(fpr, tpr, label=f'{name} (AUC = {auc_score:.2f})')
38
39  plt.plot([0, 1], [0, 1], 'k--')
40  plt.xlim([0.0, 1.0])
41  plt.ylim([0.0, 1.0])
42  plt.xlabel('False Positive Rate (FPR)')
43  plt.ylabel('True Positive Rate (TPR)')
44  plt.title('ROC Curve for Different Algorithms')
45  plt.legend(loc='lower right')
46  plt.grid()
47  plt.show()
```

```
39  plt.plot([0, 1], [0, 1], 'k--')
40  plt.xlim([0.0, 1.0])
41  plt.ylim([0.0, 1.0])
42  plt.xlabel('False Positive Rate (FPR)')
43  plt.ylabel('True Positive Rate (TPR)')
44  plt.title('ROC Curve for Different Algorithms')
45  plt.legend(loc='lower right')
46  plt.grid()
47  plt.show()
```

```
In [52]:  1  accuracy = accuracy_score(y_test, y_pred_prob.round())  # Calculate accuracy using y_pred_prob
          2  report = classification_report(y_test, knn_classifier.predict(X_test_scaled))  # Using the KNN classifier
          3
          4  print(f"--- {name} ---")
          5  print("Accuracy:", accuracy)
          6  print("Classification Report:")
          7  print(report)
          8  print("------------------------------------")
```

```
--- Random Forest ---
Accuracy: 0.485
Classification Report:
              precision    recall  f1-score   support

           0       0.61      0.45      0.52       107
           1       0.51      0.67      0.58        93

    accuracy                           0.55       200
   macro avg       0.56      0.56      0.55       200
weighted avg       0.56      0.55      0.55       200


------------------------------------
```

```
In [53]:  1  # Get feature importances from the trained Random Forest classifier
          2  feature_importances = rf_classifier.feature_importances_
          3
          4  # Print feature importances
          5  print("Feature Importances:")
          6  for feature, importance in zip(X.columns, feature_importances):
          7      print(f"{feature}: {importance}")
```
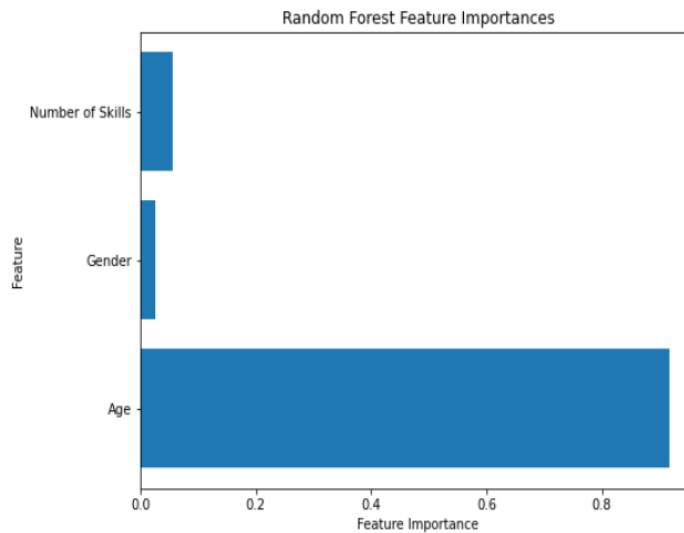
```
Feature Importances:
Age: 0.9170870182907975
Gender: 0.025868912623061253
Number of Skills: 0.05704406908614126
```

```
In [54]:   1  # Get feature importances from the trained Random Forest classifier
           2  feature_importances = rf_classifier.feature_importances_
           3
           4  # Plot feature importances in a bar plot
           5  feature_names = X.columns
           6  plt.figure(figsize=(8, 6))
           7  plt.barh(feature_names, feature_importances)
           8  plt.xlabel('Feature Importance')
           9  plt.ylabel('Feature')
          10  plt.title('Random Forest Feature Importances')
          11  plt.show()
```



```
In [57]:   1  import joblib
           2  import os
           3  from sklearn.preprocessing import OneHotEncoder, StandardScaler
           4
           5  # Create the directory
           6  output_directory = 'student_job_recommendation_system/'
           7  os.makedirs(output_directory, exist_ok=True)
           8
           9  # Assuming you have defined the 'academic_data' variable before this point
          10  # Convert categorical variables to numerical using one-hot encoding
          11  encoder = OneHotEncoder()
          12  academic_encoded = encoder.fit_transform(data).toarray()
          13
          14  # Scale the data for better clustering performance
          15  scaler = StandardScaler()
          16  academic_scaled = scaler.fit_transform(academic_encoded)
          17
          18  # Save the trained machine learning models
          19  # Assuming you already have the trained models: knn_classifier, knn_collaborative, bayesian_model, model, rf_classifier
          20  joblib.dump(knn_classifier, 'student_job_recommendation_system/knn_classifier_model.pkl')
          21  joblib.dump(rf_classifier, 'student_job_recommendation_system/random_forest_model.pkl')
```

Out[57]: ['student_job_recommendation_system/random_forest_model.pkl']