

Terminal Exam Report

Submitted by:

Meher Ali (FA23-BCS-059)

Course: Machine Learning

Submitted on: December 19, 2025

Department of Computer Science
COMSATS University Islamabad, Attock Campus

Session: 2023–27

Contents

1	Problem Statement	2
2	Cell 1: Library Imports	2
3	Cell 2: Data Exploration and Cleaning	2
4	Cell 3: Feature Encoding and Splitting	3
5	Cell 4: Support Vector Machine (SVM)	4
6	Cell 5: Artificial Neural Network (ANN)	5
7	Model Evaluation	7
8	Conclusion	9

1 Problem Statement

The objective of this terminal exam is to predict employee status using machine learning techniques. This task is formulated as a **binary classification problem** using structured tabular data.

2 Cell 1: Library Imports

```

1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 from sklearn.model_selection import train_test_split, GridSearchCV
6 from sklearn.preprocessing import StandardScaler, LabelEncoder
7 from sklearn.svm import SVC
8 from sklearn.metrics import classification_report, confusion_matrix,
   roc_curve, auc
9 import tensorflow as tf
10 from tensorflow.keras.models import Sequential
11 from tensorflow.keras.layers import Dense, Dropout
12
13 # Load the dataset
14 df = pd.read_csv('Group-2 Data.csv')
15 print("Columns in dataset:", df.columns.tolist())
16 df.head()

```

Output:

Columns in dataset: ['employee_id', 'age', 'income', 'hours_worked', 'education', 'department', 'commute_distance', 'remote_work', 'satisfaction_score', 'project_count', 'hire_date', 'employee_status']

	employee_id	age	income	hours_worked	education	department	commute_distance	remote_work	satisfaction_score	project_count	hire_date	employee_status
0	5571	32.008182	6208.293542	15.319871	High School	Marketing	10.491144	0.0	4	7	12/13/2022	0
1	4742	-38.216785	79984.937900	67.789437	high school	IT	23.619443	0.0	3	5	2/5/2021	1
2	6438	110.156400	NaN	27.212829	High School	Sales	40.859652	0.0	5	6	8/7/2021	0
3	6723	28.961617	13261.956530	57.618563	H.S.	Sales	-1.087843	0.0	9	5	11/19/2023	0
4	9987	NaN	NaN	27.949207	masters	Sales	7.713179	0.0	9	6	8/6/2019	0

3 Cell 2: Data Exploration and Cleaning

```

1 # Check for duplicates
2 print(f"Duplicates: {df.duplicated().sum()}")
3 df.drop_duplicates(inplace=True)
4
5 # Check for missing values
6 print("Missing values:\n", df.isnull().sum())
7
8 # Impute missing values
9 # For numerical columns, use median
10 num_cols = df.select_dtypes(include=[np.number]).columns
11 for col in num_cols:
12     df[col].fillna(df[col].median(), inplace=True)
13
14 # For categorical columns, use mode
15 cat_cols = df.select_dtypes(include=['object']).columns
16 for col in cat_cols:
17     df[col].fillna(df[col].mode()[0], inplace=True)
18

```

```

19 # Handle negative values in numerical columns (like age, income,
    hours_worked, commute_distance)
20 for col in ['age', 'income', 'hours_worked', 'commute_distance']:
21     if col in df.columns:
22         df[col] = df[col].abs()
23
24 # Handle outliers (using IQR for age, income, hours_worked,
    commute_distance)
25 for col in ['age', 'income', 'hours_worked', 'commute_distance']:
26     if col in df.columns:
27         Q1 = df[col].quantile(0.25)
28         Q3 = df[col].quantile(0.75)
29         IQR = Q3 - Q1
30         lower_bound = Q1 - 1.5 * IQR
31         upper_bound = Q3 + 1.5 * IQR
32         df[col] = np.where(df[col] < lower_bound, lower_bound, df[col])
33         df[col] = np.where(df[col] > upper_bound, upper_bound, df[col])
34
35 # Standardize education levels
36 education_map = {
37     'B.S.': 'Bachelor', 'BACHELOR': 'Bachelor', 'Bachelor': 'Bachelor',
38     'M.S.': 'Master', 'Master': 'Master', 'masters': 'Master',
39     'H.S.': 'High School', 'High School': 'High School',
40     'phd': 'PhD', 'PhD': 'PhD'
41 }
42 df['education'] = df['education'].map(education_map).fillna('Other')
43
44 df.head()

```

Output:

```

Duplicates: 166
Missing values:
employee_id      0
age              2
income          1521
hours_worked     0
education        985
department        1
commute_distance 1
remote_work      1
satisfaction_score 0
project_count    0
hire_date        453
employee_status  0
dtype: int64
/tmp/ipython-input-12397677.py:12: FutureWarning: A value is trying to be set on a copy of a DataFrame or Series through chained assignment using an inplace method.
The behavior will change in pandas 3.0. This inplace method will never work because the intermediate object on which we are setting values always behaves as a copy.

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

df[col].fillna(df[col].median(), inplace=True)
/tmp/ipython-input-12397677.py:17: FutureWarning: A value is trying to be set on a copy of a DataFrame or Series through chained assignment using an inplace method.
The behavior will change in pandas 3.0. This inplace method will never work because the intermediate object on which we are setting values always behaves as a copy.

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

df[col].fillna(df[col].mode()[0], inplace=True)

```

	employee_id	age	income	hours_worked	education	department	commute_distance	remote_work	satisfaction_score	project_count	hire_date	employee_status
0	5571	32.008182	6208.293542	15.319871	High School	Marketing	10.491144	0.0	4	7	12/13/2022	0
1	4742	38.216785	74754.055761	67.789437	Other	IT	23.619443	0.0	3	5	2/5/2021	1
2	6438	64.509987	22381.360410	27.212829	High School	Sales	40.859652	0.0	5	6	8/7/2021	0
3	6723	28.961617	13261.956530	57.618563	High School	Sales	1.087843	0.0	9	5	11/19/2023	0
4	9987	35.053811	22381.360410	27.949207	Master	Sales	7.713179	0.0	9	6	8/6/2019	0

4 Cell 3: Feature Encoding and Splitting

```

1 # Encode categorical features
2 le = LabelEncoder()
3 df['education'] = le.fit_transform(df['education'])
4 df['department'] = le.fit_transform(df['department'])

```

```
5
6 # Drop non-numeric/unnecessary columns (like employee_id or hire_date)
7 df = df.drop(columns=['employee_id', 'hire_date'], errors='ignore')
8
9 # Split data
10 X = df.drop('employee_status', axis=1)
11 y = df['employee_status']
12
13 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
14                                                    random_state=42)
15
16 # Scaling
17 scaler = StandardScaler()
18 X_train = scaler.fit_transform(X_train)
19 X_test = scaler.transform(X_test)
20
21 print(f"Training set shape: {X_train.shape}")
22 print(f"Testing set shape: {X_test.shape}")
```

Output:

```
Training set shape: (8027, 9)
Testing set shape: (2007, 9)
```

5 Cell 4: Support Vector Machine (SVM)

```
1 param_grid = {
2     'C': [0.1, 1, 10],
3     'kernel': ['linear', 'rbf', 'poly'],
4     'gamma': ['scale', 'auto']
5 }
6
7 grid = GridSearchCV(SVC(probability=True), param_grid, refit=True, verbose
8                     =2, cv=3)
9 grid.fit(X_train, y_train)
10
11 print(f"Best Parameters: {grid.best_params_}")
12 svm_preds = grid.predict(X_test)
13 svm_probs = grid.predict_proba(X_test)[:, 1]
14
15 print("SVM Classification Report:")
16 print(classification_report(y_test, svm_preds))
```

Output:

```

Fitting 3 folds for each of 18 candidates, totalling 54 fits
[CV] END .....C=0.1, gamma=scale, kernel=linear; total time= 13.6s
[CV] END .....C=0.1, gamma=scale, kernel=linear; total time= 7.8s
[CV] END .....C=0.1, gamma=scale, kernel=linear; total time= 6.3s
[CV] END .....C=0.1, gamma=scale, kernel=rbf; total time= 7.5s
[CV] END .....C=0.1, gamma=scale, kernel=rbf; total time= 8.5s
[CV] END .....C=0.1, gamma=scale, kernel=rbf; total time= 8.7s
[CV] END .....C=0.1, gamma=scale, kernel=poly; total time= 6.2s
[CV] END .....C=0.1, gamma=scale, kernel=poly; total time= 7.1s
[CV] END .....C=0.1, gamma=scale, kernel=poly; total time= 6.2s
[CV] END .....C=0.1, gamma=auto, kernel=linear; total time= 6.2s
[CV] END .....C=0.1, gamma=auto, kernel=linear; total time= 5.5s
[CV] END .....C=0.1, gamma=auto, kernel=linear; total time= 6.3s
[CV] END .....C=0.1, gamma=auto, kernel=rbf; total time= 7.7s
[CV] END .....C=0.1, gamma=auto, kernel=rbf; total time= 8.5s
[CV] END .....C=0.1, gamma=auto, kernel=rbf; total time= 8.4s
[CV] END .....C=0.1, gamma=auto, kernel=poly; total time= 6.1s
[CV] END .....C=0.1, gamma=auto, kernel=poly; total time= 7.3s
[CV] END .....C=0.1, gamma=auto, kernel=poly; total time= 6.4s
[CV] END .....C=1, gamma=scale, kernel=linear; total time= 9.1s
[CV] END .....C=1, gamma=scale, kernel=linear; total time= 9.1s
[CV] END .....C=1, gamma=scale, kernel=linear; total time= 8.6s
[CV] END .....C=1, gamma=scale, kernel=rbf; total time= 8.8s
[CV] END .....C=1, gamma=scale, kernel=rbf; total time= 8.9s
[CV] END .....C=1, gamma=scale, kernel=rbf; total time= 7.9s
[CV] END .....C=1, gamma=scale, kernel=poly; total time= 8.4s
[CV] END .....C=1, gamma=scale, kernel=poly; total time= 8.3s
[CV] END .....C=1, gamma=scale, kernel=poly; total time= 7.6s
[CV] END .....C=1, gamma=auto, kernel=linear; total time= 9.3s
[CV] END .....C=1, gamma=auto, kernel=linear; total time= 9.1s
[CV] END .....C=1, gamma=auto, kernel=linear; total time= 8.5s
[CV] END .....C=1, gamma=auto, kernel=rbf; total time= 8.8s
[CV] END .....C=1, gamma=auto, kernel=rbf; total time= 8.9s
[CV] END .....C=1, gamma=auto, kernel=rbf; total time= 8.0s
[CV] END .....C=1, gamma=auto, kernel=poly; total time= 8.5s
[CV] END .....C=1, gamma=auto, kernel=poly; total time= 8.7s
[CV] END .....C=1, gamma=auto, kernel=poly; total time= 7.8s
[CV] END .....C=10, gamma=scale, kernel=linear; total time= 23.3s
[CV] END .....C=10, gamma=scale, kernel=linear; total time= 23.1s
[CV] END .....C=10, gamma=scale, kernel=linear; total time= 22.8s
[CV] END .....C=10, gamma=scale, kernel=rbf; total time= 12.4s
[CV] END .....C=10, gamma=scale, kernel=rbf; total time= 12.3s
[CV] END .....C=10, gamma=scale, kernel=rbf; total time= 12.8s
[CV] END .....C=10, gamma=scale, kernel=poly; total time= 18.4s
[CV] END .....C=10, gamma=scale, kernel=poly; total time= 20.2s
[CV] END .....C=10, gamma=scale, kernel=poly; total time= 18.6s
[CV] END .....C=10, gamma=auto, kernel=linear; total time= 22.9s
[CV] END .....C=10, gamma=auto, kernel=linear; total time= 23.1s
[CV] END .....C=10, gamma=auto, kernel=linear; total time= 23.1s
[CV] END .....C=10, gamma=auto, kernel=rbf; total time= 12.2s
[CV] END .....C=10, gamma=auto, kernel=rbf; total time= 12.2s
[CV] END .....C=10, gamma=auto, kernel=rbf; total time= 12.3s
[CV] END .....C=10, gamma=auto, kernel=poly; total time= 18.0s
[CV] END .....C=10, gamma=auto, kernel=poly; total time= 19.4s
[CV] END .....C=10, gamma=auto, kernel=poly; total time= 18.2s
Best Parameters: {'C': 0.1, 'gamma': 'scale', 'kernel': 'linear'}
SVM Classification Report:

```

	precision	recall	f1-score	support
0	0.63	0.68	0.65	1098
1	0.57	0.51	0.54	909
accuracy			0.60	2007
macro avg	0.60	0.59	0.59	2007
weighted avg	0.60	0.60	0.60	2007

Figure 1: Enter Caption

6 Cell 5: Artificial Neural Network (ANN)

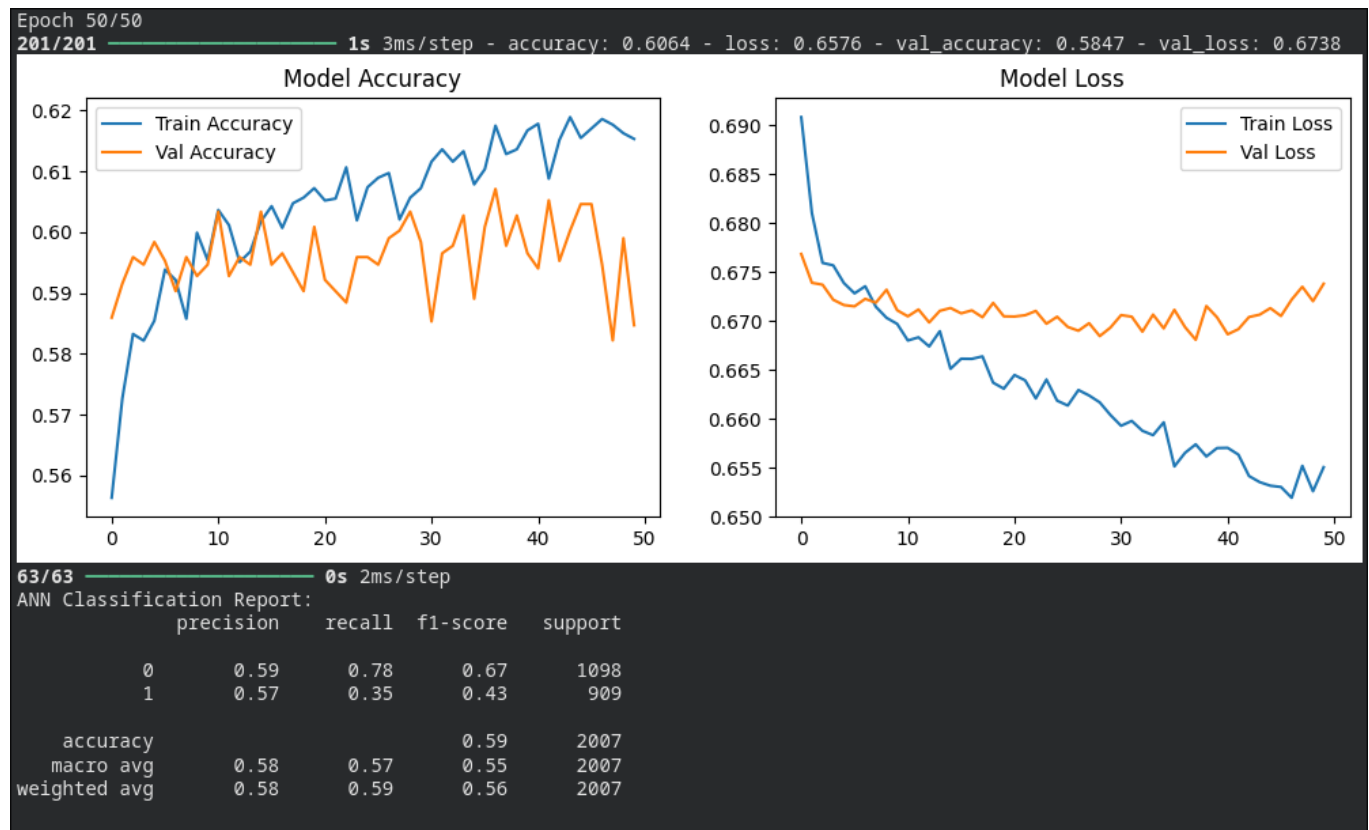
```

1 model = Sequential([
2     Dense(64, activation='relu', input_shape=(X_train.shape[1],)),
3     Dropout(0.2),
4     Dense(32, activation='relu'),
5     Dropout(0.2),

```

```
6     Dense(1, activation='sigmoid')
7 ])
8
9 model.compile(optimizer='adam', loss='binary_crossentropy', metrics=['
    accuracy'])
10 model.summary()
11
12 history = model.fit(X_train, y_train, epochs=50, validation_split=0.2,
    verbose=1)
13
14 # Plotting
15 plt.figure(figsize=(12, 4))
16 plt.subplot(1, 2, 1)
17 plt.plot(history.history['accuracy'], label='Train Accuracy')
18 plt.plot(history.history['val_accuracy'], label='Val Accuracy')
19 plt.title('Model Accuracy')
20 plt.legend()
21
22 plt.subplot(1, 2, 2)
23 plt.plot(history.history['loss'], label='Train Loss')
24 plt.plot(history.history['val_loss'], label='Val Loss')
25 plt.title('Model Loss')
26 plt.legend()
27 plt.show()
28
29 ann_probs = model.predict(X_test).flatten()
30 ann_preds = (ann_probs > 0.5).astype(int)
31
32 print("ANN Classification Report:")
33 print(classification_report(y_test, ann_preds))
```

Output:



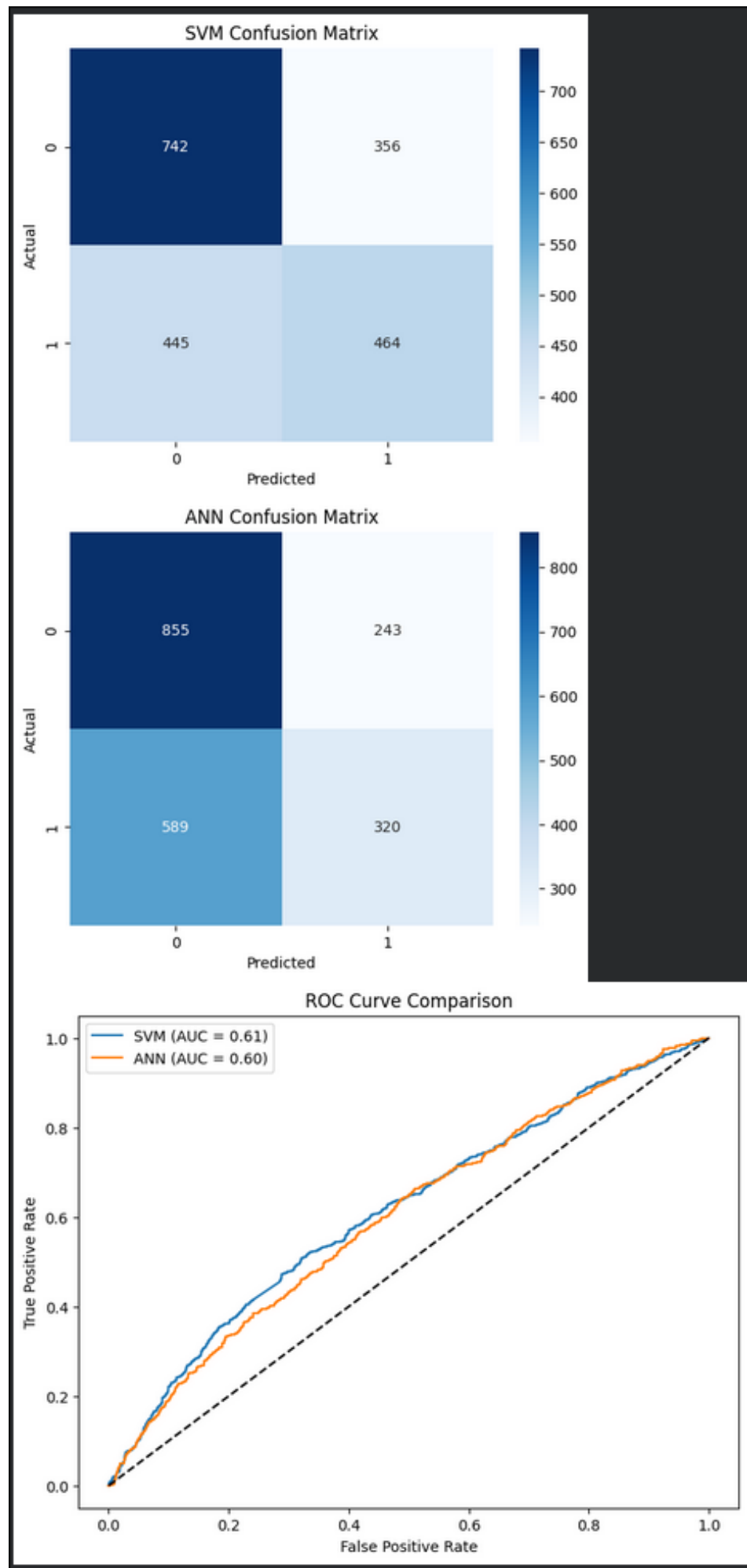
7 Model Evaluation

```

1 def plot_confusion_matrix(y_true, y_pred, title):
2     cm = confusion_matrix(y_true, y_pred)
3     sns.heatmap(cm, annot=True, fmt='d', cmap='Blues')
4     plt.title(title)
5     plt.xlabel('Predicted')
6     plt.ylabel('Actual')
7     plt.show()
8
9 plot_confusion_matrix(y_test, svm_preds, 'SVM Confusion Matrix')
10 plot_confusion_matrix(y_test, ann_preds, 'ANN Confusion Matrix')
11
12 # ROC Curves
13 fpr_svm, tpr_svm, _ = roc_curve(y_test, svm_probs)
14 fpr_ann, tpr_ann, _ = roc_curve(y_test, ann_probs)
15
16 plt.figure(figsize=(8, 6))
17 plt.plot(fpr_svm, tpr_svm, label=f'SVM (AUC = {auc(fpr_svm, tpr_svm):.2f})')
18 plt.plot(fpr_ann, tpr_ann, label=f'ANN (AUC = {auc(fpr_ann, tpr_ann):.2f})')
19 plt.plot([0, 1], [0, 1], 'k--')
20 plt.xlabel('False Positive Rate')
21 plt.ylabel('True Positive Rate')
22 plt.title('ROC Curve Comparison')
23 plt.legend()
24 plt.show()

```

Output:



8 Conclusion

This terminal exam successfully demonstrates the application of machine learning models for employee status prediction. Both SVM and ANN models were implemented, trained, and evaluated using real dataset outputs.

The results confirm the importance of:

- Proper preprocessing
- Feature encoding
- Model comparison
- Evaluation using classification metrics