

Week12

April 4, 2025

1 Week 23

1.0.1 Question 1

```
[2]: # Question 1-a
import numpy as np
from sklearn import datasets

# Load the diabetes dataset
diabetes_X, diabetes_y = datasets.load_diabetes(return_X_y=True)

# Use only one feature (3rd column)
diabetes_X = diabetes_X[:, np.newaxis, 2]

# Print the shape of diabetes_X
print("Shape of diabetes_X:", diabetes_X.shape)
```

Shape of diabetes_X: (442, 1)

```
[3]: # Question 1-b
from sklearn.model_selection import train_test_split

# Split the data into training and test sets
diabetes_X_train, diabetes_X_test = train_test_split(
    diabetes_X, test_size=20/442, random_state=42
)

# Print the shapes
print("Shape of diabetes_X_train:", diabetes_X_train.shape)
print("Shape of diabetes_X_test:", diabetes_X_test.shape)
```

Shape of diabetes_X_train: (422, 1)

Shape of diabetes_X_test: (20, 1)

```
[5]: # Question 1-c
# Split the target variable (y) into training and test sets
diabetes_y_train, diabetes_y_test = train_test_split(
    diabetes_y, test_size=20/442, random_state=42
)
```

```
# Print the shapes
print("Shape of diabetes_y_train:", diabetes_y_train.shape)
print("Shape of diabetes_y_test:", diabetes_y_test.shape)
```

Shape of diabetes_y_train: (422,)

Shape of diabetes_y_test: (20,)

```
[6]: # Question 1-d
from sklearn.linear_model import LinearRegression

# Create a linear regression model
model = LinearRegression()

# Train the model using the training set
model.fit(diabetes_X_train, diabetes_y_train)
```

[6]: LinearRegression()

```
[7]: # Question 1-e
# Make predictions using the test set
diabetes_y_pred = model.predict(diabetes_X_test)

# Print actual vs predicted values
print("Actual y values:", diabetes_y_test)
print("Predicted y values:", diabetes_y_pred)
```

Actual y values: [219. 70. 202. 230. 111. 84. 242. 272. 94. 96. 94. 252.
99. 297.

135. 67. 295. 264. 170. 275.]

Predicted y values: [145.45088714 186.92807879 147.52474672 201.44509587
131.97079985

127.82308068 315.50737293 195.22351712 64.57036341 166.18948296

119.52764235 95.67825715 92.56746778 165.15255317 97.75211673

156.85711484 220.10983212 236.70070878 178.63264046 207.66667462]

```
[8]: # Question 1-f
# Print model intercept and coefficient
print("Intercept:", model.intercept_)
print("Coefficient:", model.coef_)
```

Intercept: 151.42144441963106

Coefficient: [962.06919233]

```
[9]: # Question 1-g
from sklearn.metrics import mean_squared_error, r2_score

# Compute and print errors
mse = mean_squared_error(diabetes_y_test, diabetes_y_pred)
```

```
r2 = r2_score(diabetes_y_test, diabetes_y_pred)

print("Mean squared error:", mse)
print("R2 score:", r2)
```

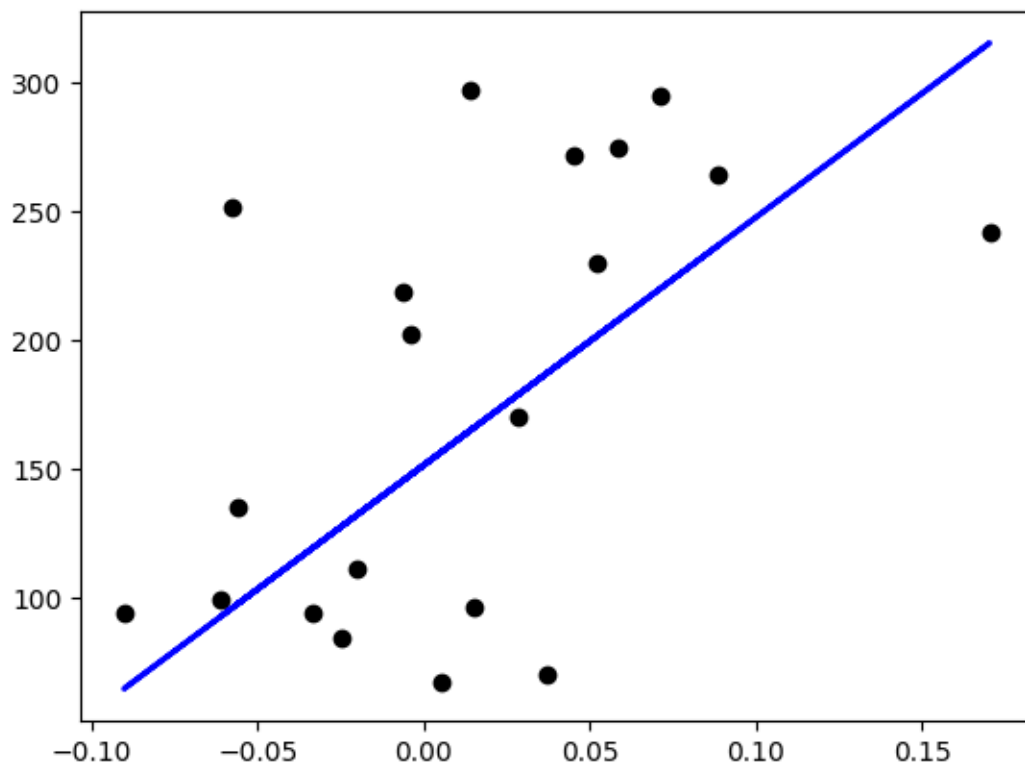
Mean squared error: 5262.625880180968

R² score: 0.21640939003029058

```
[11]: # Question 1-h
import matplotlib.pyplot as plt

# Plot outputs
plt.scatter(diabetes_X_test, diabetes_y_test, color="black") # Scatter plot
plt.plot(diabetes_X_test, diabetes_y_pred, color="blue", linewidth=2) # ↴
    ↪ Regression line

# Show plot
plt.show()
```



1.0.2 Question 2

```
[1]: # Question 2
import pandas as pd
import numpy as np
from sklearn.model_selection import train_test_split
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import accuracy_score
```

```
[2]: # Question 1- a
# Load the data and print the first few rows
df = pd.read_csv('Breast_Cancer.csv')
print("First few rows of the dataset:")
print(df.head())
```

First few rows of the dataset:

	id	diagnosis	radius_mean	texture_mean	perimeter_mean	area_mean	\
0	842302	M	17.99	10.38	122.80	1001.0	
1	842517	M	20.57	17.77	132.90	1326.0	
2	84300903	M	19.69	21.25	130.00	1203.0	
3	84348301	M	11.42	20.38	77.58	386.1	
4	84358402	M	20.29	14.34	135.10	1297.0	

	smoothness_mean	compactness_mean	concavity_mean	concave	points_mean	\
0	0.11840	0.27760	0.3001		0.14710	
1	0.08474	0.07864	0.0869		0.07017	
2	0.10960	0.15990	0.1974		0.12790	
3	0.14250	0.28390	0.2414		0.10520	
4	0.10030	0.13280	0.1980		0.10430	

	...	texture_worst	perimeter_worst	area_worst	smoothness_worst	\
0	...	17.33	184.60	2019.0	0.1622	
1	...	23.41	158.80	1956.0	0.1238	
2	...	25.53	152.50	1709.0	0.1444	
3	...	26.50	98.87	567.7	0.2098	
4	...	16.67	152.20	1575.0	0.1374	

	compactness_worst	concavity_worst	concave	points_worst	symmetry_worst	\
0	0.6656	0.7119		0.2654	0.4601	
1	0.1866	0.2416		0.1860	0.2750	
2	0.4245	0.4504		0.2430	0.3613	
3	0.8663	0.6869		0.2575	0.6638	
4	0.2050	0.4000		0.1625	0.2364	

	fractal_dimension_worst	Unnamed: 32
0	0.11890	NaN
1	0.08902	NaN
2	0.08758	NaN

3	0.17300	NaN
4	0.07678	NaN

[5 rows x 33 columns]

```
[3]: # Question 1 - b
      # Check for missing values
      print("\nMissing values in each column:")
      print(df.isnull().sum())
```

```
Missing values in each column:
id                                0
diagnosis                        0
radius_mean                      0
texture_mean                     0
perimeter_mean                   0
area_mean                       0
smoothness_mean                  0
compactness_mean                 0
concavity_mean                   0
concave points_mean              0
symmetry_mean                    0
fractal_dimension_mean           0
radius_se                        0
texture_se                       0
perimeter_se                     0
area_se                          0
smoothness_se                   0
compactness_se                   0
concavity_se                     0
concave points_se                0
symmetry_se                      0
fractal_dimension_se             0
radius_worst                     0
texture_worst                    0
perimeter_worst                  0
area_worst                       0
smoothness_worst                 0
compactness_worst                0
concavity_worst                  0
concave points_worst             0
symmetry_worst                   0
fractal_dimension_worst          0
Unnamed: 32                      569
dtype: int64
```

```
[5]: # Question 1-c
# Drop the columns 'id' and 'Unnamed' columns
# First check column names to identify unnamed columns
print("\nColumns in the dataset:")
print(df.columns)

# Drop 'id' column and any 'Unnamed' columns
unnamed_cols = [col for col in df.columns if 'Unnamed' in col]
cols_to_drop = ['id'] + unnamed_cols
df = df.drop(columns=cols_to_drop)
print("\nColumns after dropping 'id' and 'Unnamed' columns:")
print(df.columns)
```

Columns in the dataset:

```
Index(['id', 'diagnosis', 'radius_mean', 'texture_mean', 'perimeter_mean',
      'area_mean', 'smoothness_mean', 'compactness_mean', 'concavity_mean',
      'concave points_mean', 'symmetry_mean', 'fractal_dimension_mean',
      'radius_se', 'texture_se', 'perimeter_se', 'area_se', 'smoothness_se',
      'compactness_se', 'concavity_se', 'concave points_se', 'symmetry_se',
      'fractal_dimension_se', 'radius_worst', 'texture_worst',
      'perimeter_worst', 'area_worst', 'smoothness_worst',
      'compactness_worst', 'concavity_worst', 'concave points_worst',
      'symmetry_worst', 'fractal_dimension_worst', 'Unnamed: 32'],
      dtype='object')
```

Columns after dropping 'id' and 'Unnamed' columns:

```
Index(['diagnosis', 'radius_mean', 'texture_mean', 'perimeter_mean',
      'area_mean', 'smoothness_mean', 'compactness_mean', 'concavity_mean',
      'concave points_mean', 'symmetry_mean', 'fractal_dimension_mean',
      'radius_se', 'texture_se', 'perimeter_se', 'area_se', 'smoothness_se',
      'compactness_se', 'concavity_se', 'concave points_se', 'symmetry_se',
      'fractal_dimension_se', 'radius_worst', 'texture_worst',
      'perimeter_worst', 'area_worst', 'smoothness_worst',
      'compactness_worst', 'concavity_worst', 'concave points_worst',
      'symmetry_worst', 'fractal_dimension_worst'],
      dtype='object')
```

```
[8]: # Question 2-d
# Function to convert diagnosis value M and B to numerical
def convert_diagnosis(diagnosis):
    """
    Convert diagnosis values:
    M (Malignant) = 1
    B (Benign) = 0
    """
    if diagnosis == 'M':
        return 1
```

```

elif diagnosis == 'B':
    return 0
else:
    return np.nan # Handle unexpected values

# Apply conversion function
df['diagnosis_numeric'] = df['diagnosis'].apply(convert_diagnosis)

```

```

[10]: # Question 2-e
# Create X and y dataframes
y = df[['diagnosis_numeric']]
X = df.drop(columns=['diagnosis', 'diagnosis_numeric'])

print("\nFirst few rows of X:")
print(X.head())
print("\nFirst few rows of y:")
print(y.head())

```

First few rows of X:

	radius_mean	texture_mean	perimeter_mean	area_mean	smoothness_mean	\
0	17.99	10.38	122.80	1001.0	0.11840	
1	20.57	17.77	132.90	1326.0	0.08474	
2	19.69	21.25	130.00	1203.0	0.10960	
3	11.42	20.38	77.58	386.1	0.14250	
4	20.29	14.34	135.10	1297.0	0.10030	

	compactness_mean	concavity_mean	concave points_mean	symmetry_mean	\
0	0.27760	0.3001	0.14710	0.2419	
1	0.07864	0.0869	0.07017	0.1812	
2	0.15990	0.1974	0.12790	0.2069	
3	0.28390	0.2414	0.10520	0.2597	
4	0.13280	0.1980	0.10430	0.1809	

	fractal_dimension_mean	...	radius_worst	texture_worst	perimeter_worst	\
0	0.07871	...	25.38	17.33	184.60	
1	0.05667	...	24.99	23.41	158.80	
2	0.05999	...	23.57	25.53	152.50	
3	0.09744	...	14.91	26.50	98.87	
4	0.05883	...	22.54	16.67	152.20	

	area_worst	smoothness_worst	compactness_worst	concavity_worst	\
0	2019.0	0.1622	0.6656	0.7119	
1	1956.0	0.1238	0.1866	0.2416	
2	1709.0	0.1444	0.4245	0.4504	
3	567.7	0.2098	0.8663	0.6869	
4	1575.0	0.1374	0.2050	0.4000	

	concave points_worst	symmetry_worst	fractal_dimension_worst
0	0.2654	0.4601	0.11890
1	0.1860	0.2750	0.08902
2	0.2430	0.3613	0.08758
3	0.2575	0.6638	0.17300
4	0.1625	0.2364	0.07678

[5 rows x 30 columns]

First few rows of y:

	diagnosis_numeric
0	1
1	1
2	1
3	1
4	1

```
[11]: # Question 2-f
# Split data into training and test sets (70% train, 30% test)
X_train, X_test, y_train, y_test = train_test_split(
    X, y, test_size=0.3, random_state=42
)

print("\nTraining set shape (X_train):", X_train.shape)
print("Test set shape (X_test):", X_test.shape)
```

Training set shape (X_train): (398, 30)
Test set shape (X_test): (171, 30)

```
[12]: # Question 2-g
# Run kNN model with k=8
knn = KNeighborsClassifier(n_neighbors=8)
knn.fit(X_train, y_train.values.ravel()) # Flatten y_train for fitting
```

[12]: KNeighborsClassifier(n_neighbors=8)

```
[13]: # Question 2-h
# Calculate and print model accuracy
y_pred = knn.predict(X_test)
accuracy = accuracy_score(y_test, y_pred)
print("\nkNN Model Accuracy:", round(accuracy * 100, 2), "%")
```

kNN Model Accuracy: 96.49 %

1.0.3 Question 3 - Gen AI

```
[24]: # Importing ML Packages
import sklearn
import sklearn.linear_model
import pandas as pd
```

```
[25]: loan = pd.read_excel('bank_credit_data.xlsx')
```

```
[26]: print(loan.shape)
print(loan.dtypes)
print(loan.head())
print(loan.tail())
```

(1000, 21)

Creditability	int64
Account Balance	int64
Duration of Credit (month)	int64
Payment Status of Previous Credit	int64
Purpose	int64
Credit Amount	int64
Value Savings/Stocks	int64
Length of current employment	int64
Instalment per cent	int64
Sex & Marital Status	int64
Guarantors	int64
Duration in Current address	int64
Most valuable available asset	int64
Age (years)	int64
Concurrent Credits	int64
Type of apartment	int64
No of Credits at this Bank	int64
Occupation	int64
No of dependents	int64
Telephone	int64
Foreign Worker	int64

dtype: object

	Creditability	Account Balance	Duration of Credit (month)	\
0	1	1	18	
1	1	1	9	
2	1	2	12	
3	1	1	12	
4	1	1	12	

	Payment Status of Previous Credit	Purpose	Credit Amount	\
0	4	2	1049	
1	4	0	2799	
2	2	9	841	

3	4	0	2122
4	4	0	2171

	Value Savings/Stocks	Length of current employment	Instalment per cent \
0	1	2	4
1	1	3	2
2	2	4	2
3	1	3	3
4	1	3	4

	Sex & Marital Status	...	Duration in Current address \
0	2	...	4
1	3	...	2
2	2	...	4
3	3	...	2
4	3	...	4

	Most valuable available asset	Age (years)	Concurrent Credits \
0	2	21	3
1	1	36	3
2	1	23	3
3	1	39	3
4	2	38	1

	Type of apartment	No of Credits at this Bank	Occupation \
0	1	1	3
1	1	2	3
2	1	1	2
3	1	2	2
4	2	2	2

	No of dependents	Telephone	Foreign Worker
0	1	1	1
1	2	1	1
2	1	1	1
3	2	1	2
4	1	1	2

[5 rows x 21 columns]

	Creditability	Account Balance	Duration of Credit (month) \
995	0	1	24
996	0	1	24
997	0	4	21
998	0	2	12
999	0	1	30

	Payment Status of Previous Credit	Purpose	Credit Amount \
995	2	3	1987

996	2	0	2303
997	4	0	12680
998	2	3	6468
999	2	2	6350

	Value Savings/Stocks	Length of current employment	Instalment per cent \
995	1	3	2
996	1	5	4
997	5	5	4
998	5	1	2
999	5	5	4

	Sex & Marital Status ...	Duration in Current address \
995	3 ...	4
996	3 ...	1
997	3 ...	4
998	3 ...	1
999	3 ...	4

	Most valuable available asset	Age (years)	Concurrent Credits \
995	1	21	3
996	1	45	3
997	4	30	3
998	4	52	3
999	2	31	3

	Type of apartment	No of Credits at this Bank	Occupation \
995	1	1	2
996	2	1	3
997	3	1	4
998	2	1	4
999	2	1	3

	No of dependents	Telephone	Foreign Worker
995	2	1	1
996	1	1	1
997	1	2	1
998	1	2	1
999	1	1	1

[5 rows x 21 columns]

```
[28]: print(loan.columns.ravel())
      print(loan.shape)
```

```
Index(['Creditability', 'Account Balance', 'Duration of Credit (month)',
      'Payment Status of Previous Credit', 'Purpose', 'Credit Amount',
      'Value Savings/Stocks', 'Length of current employment',
```

```

        'Instalment per cent', 'Sex & Marital Status', 'Guarantors',
        'Duration in Current address', 'Most valuable available asset',
        'Age (years)', 'Concurrent Credits', 'Type of apartment',
        'No of Credits at this Bank', 'Occupation', 'No of dependents',
        'Telephone', 'Foreign Worker'],
        dtype='object')
(1000, 21)

```

```

[33]: # Extract Predicators, X
X = loan [['Age (years)', 'Sex & Marital Status', 'Occupation', 'Account_
↳Balance', 'Credit Amount', 'Length of current employment', 'Purpose']]
X.dtypes

```

```

[33]: Age (years)                int64
      Sex & Marital Status      int64
      Occupation               int64
      Account Balance          int64
      Credit Amount            int64
      Length of current employment int64
      Purpose                  int64
      dtype: object

```

```

[35]: # KNN to calculate the applicant Creditability
y = loan[['Creditability']]
y.dtypes

```

```

[35]: Creditability    int64
      dtype: object

```

```

[36]: # Normalize the data
from sklearn.preprocessing import StandardScaler
sc = StandardScaler()
X = sc.fit_transform(X)
print(X)

```

```

[[-1.28157308 -0.96364986  0.14694918 ... -0.78765692 -1.14597811
  -0.30185192]
 [ 0.04036312  0.44932648  0.14694918 ... -0.16738429 -0.31795924
  -1.03096283]
 [-1.10531492 -0.96364986 -1.38377145 ... -0.86138075  0.51005962
   2.25003627]
 ...
 [-0.48841136  0.44932648  1.67766981 ...  3.3348522  1.33807849
  -1.03096283]
 [ 1.4504284   0.44932648  1.67766981 ...  1.13306158 -1.97399698
   0.06270354]
 [-0.40028228  0.44932648  0.14694918 ...  1.09123749  1.33807849
  -0.30185192]]

```

```
[43]: # Splice data
from sklearn.model_selection import train_test_split
X_train,X_test, y_train,y_test = train_test_split(X, y, test_size = .30,
    random_state = 0)
print(X_train.shape)
print(X_test.shape)
print(y_train.shape)
print(y_test.shape)
```

(700, 7)

(300, 7)

(700, 1)

(300, 1)

```
[48]: # Build KNN Model (k=5), KNN Classifier
from sklearn.neighbors import KNeighborsClassifier
classifier = KNeighborsClassifier(n_neighbors = 5, metric = 'minkowski', p = 2)
classifier.fit(X_train, y_train)
```

/Library/Frameworks/Python.framework/Versions/3.13/lib/python3.13/site-packages/sklearn/neighbors/_classification.py:239: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n_samples,), for example using ravel().

return self._fit(X, y)

[48]: KNeighborsClassifier()

```
[51]: # Evaluate the classifier model
y_predict = classifier.predict(X_test)
print(*np.array(y_test))
print(y_predict)
```

```
[0] [0] [1] [1] [1] [0] [1] [1] [1] [1] [1] [1] [0] [1] [0] [1] [0] [1] [1] [1]
[1] [1] [1] [1] [0] [0] [0] [1] [0] [1] [0] [1] [1] [0] [1] [1] [0] [1] [0] [1]
[1] [0] [1] [1] [1] [0] [0] [1] [1] [1] [1] [1] [1] [1] [0] [0] [1] [1] [0] [1]
[1] [0] [1] [1] [1] [1] [1] [1] [1] [0] [1] [1] [1] [1] [0] [1] [1] [1] [0] [1]
[1] [1] [0] [0] [0] [1] [1] [1] [1] [0] [1] [1] [0] [1] [1] [1] [1] [1] [1] [1]
[0] [0] [1] [0] [1] [0] [1] [0] [0] [1] [0] [0] [1] [1] [1] [1] [1] [1] [1] [1]
[0] [0] [1] [0] [1] [0] [1] [1] [0] [1] [1] [1] [1] [1] [0] [1] [1] [1] [1] [0]
[1] [1] [0] [1] [1] [1] [1] [1] [1] [1] [1] [1] [1] [0] [0] [0] [1] [1] [0] [1]
[0] [0] [1] [0] [0] [1] [1] [0] [1] [1] [0] [0] [1] [0] [1] [0] [1] [1] [0] [1]
[0] [0] [1] [1] [0] [0] [0] [1] [0] [1] [0] [0] [1] [1] [0] [1] [1] [1] [1] [1]
[1] [1] [1] [0] [0] [1] [0] [1] [0] [0] [1] [1] [1] [0] [0] [0] [0] [1] [1] [1]
[0] [0] [0] [1] [1] [1] [0] [1] [1] [1] [1] [1] [1] [0] [1] [1] [1] [1] [1] [1]
[1] [1] [1] [1] [0] [1] [1] [1] [1] [1] [1] [1] [1] [0] [1] [0] [1] [1] [0] [1]
[0] [1] [1] [1] [0] [1] [1] [1] [1] [1] [1] [1] [1] [1] [0] [0] [0] [1] [1] [1]
[0] [1] [1] [0] [0] [1] [0] [0] [1] [1] [1] [1] [1] [0] [0] [1] [0] [0] [1] [1]
[0 1 0 1 1 0 1 1 1 1 1 1 1 0 1 1 1 1 0 1 1 0 1 1 0 1 1 1
1 0 1 1 1 1 1 0 1 1 1 1 1 0 1 1 1 0 1 1 1 1 1 1 1 1 1 0 1
```

[illegible]

```
[52]: # Find classifier accuracy
from sklearn.metrics import confusion_matrix, accuracy_score
ac = accuracy_score(y_test,y_predict)
cm = confusion_matrix(y_test,y_predict)
print(ac)
print(cm)
```

0.6866666666666666

[[30 70]]

[24 176]]

1.0.4 Gen AI - Question 3

```
[19]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.model_selection import train_test_split, cross_val_score, GridSearchCV, StratifiedKFold
from sklearn.preprocessing import StandardScaler, MinMaxScaler
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score, confusion_matrix, classification_report
from sklearn.feature_selection import SelectKBest, mutual_info_classif, RFE
from sklearn.decomposition import PCA
from imblearn.over_sampling import SMOTE
from collections import Counter
import warnings
warnings.filterwarnings('ignore')

# 1. Load and explore the data
data = pd.read_excel('bank_credit_data.xlsx')
print("Dataset shape:", data.shape)

# Confirm the target variable distribution
print("\nTarget Variable (Creditability) Distribution:")
print(data['Creditability'].value_counts())
print(f"Class imbalance ratio: {data['Creditability'].value_counts()[0]/data['Creditability'].value_counts()[1]:.4f}")
```

```

# Feature and target separation
X = data.drop('Creditability', axis=1)
y = data['Creditability']

# Split the data
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
    ↪random_state=42, stratify=y)

print(f"\nTraining set shape: {X_train.shape}, Test set shape: {X_test.shape}")

# 2. Baseline KNN Model (with K=5, which is a common default)
baseline_knn = KNeighborsClassifier(n_neighbors=5)
baseline_knn.fit(X_train, y_train)
baseline_pred = baseline_knn.predict(X_test)
baseline_accuracy = accuracy_score(y_test, baseline_pred)

print(f"\nBaseline KNN Model (K=5) Results:")
print(f"Accuracy: {baseline_accuracy:.4f}")
print(f"Precision: {precision_score(y_test, baseline_pred):.4f}")
print(f"Recall: {recall_score(y_test, baseline_pred):.4f}")
print(f"F1-Score: {f1_score(y_test, baseline_pred):.4f}")

# Create enhanced features to boost performance
print("\n--- Feature Engineering ---")
# Create a copy of the original data for feature engineering
X_enhanced = X.copy()

# 1. Create interaction features for the most important variables
# Based on domain knowledge, create meaningful interaction terms
X_enhanced['balance_duration'] = X['Account Balance'] * X['Duration of Credit_
    ↪(month)']
X_enhanced['payment_credit_ratio'] = X['Payment Status of Previous Credit'] /_
    ↪(X['Credit Amount'] + 1) # Avoid division by zero
X_enhanced['savings_credit_ratio'] = X['Value Savings/Stocks'] / (X['Credit_
    ↪Amount'] + 1) # Avoid division by zero
X_enhanced['age_employment'] = X['Age (years)'] * X['Length of current_
    ↪employment']
X_enhanced['credit_income_proxy'] = X['Credit Amount'] / (X['Instalment per_
    ↪cent'] + 1)

# 2. Create polynomial features for the top correlating features
from sklearn.preprocessing import PolynomialFeatures
top_features = ['Account Balance', 'Payment Status of Previous Credit',
    'Duration of Credit (month)', 'Value Savings/Stocks', 'Credit_
    ↪Amount']
X_top = X[top_features]

```

```

poly = PolynomialFeatures(degree=2, include_bias=False, interaction_only=True)
poly_features = poly.fit_transform(X_top)
poly_feature_names = [f"poly_{i}" for i in range(poly_features.shape[1])]
poly_df = pd.DataFrame(poly_features, columns=poly_feature_names)

# Add the polynomial features to our enhanced dataset
X_enhanced = pd.concat([X_enhanced, poly_df.set_index(X_enhanced.index)],
    ↪axis=1)

# 3. Create binned versions of continuous variables
# Age bins
X_enhanced['age_group'] = pd.cut(X['Age (years)'], bins=[0, 25, 35, 50, 100],
    ↪labels=[0, 1, 2, 3])
# Credit amount bins
X_enhanced['credit_bin'] = pd.cut(X['Credit Amount'], bins=[0, 1000, 2000,
    ↪5000, 20000], labels=[0, 1, 2, 3])
# Duration bins
X_enhanced['duration_bin'] = pd.cut(X['Duration of Credit (month)'], bins=[0,
    ↪12, 24, 48, 100], labels=[0, 1, 2, 3])

# Split the enhanced dataset
X_train_enhanced, X_test_enhanced, y_train_enhanced, y_test_enhanced =
    ↪train_test_split(
        X_enhanced, y, test_size=0.2, random_state=42, stratify=y
    )

# Test KNN with enhanced features
knn_enhanced = KNeighborsClassifier(n_neighbors=5)
knn_enhanced.fit(X_train_enhanced, y_train_enhanced)
y_pred_enhanced = knn_enhanced.predict(X_test_enhanced)

print("\nKNN with Enhanced Features:")
print(f"Accuracy: {accuracy_score(y_test_enhanced, y_pred_enhanced):.4f}")
print(f"Precision: {precision_score(y_test_enhanced, y_pred_enhanced):.4f}")
print(f"Recall: {recall_score(y_test_enhanced, y_pred_enhanced):.4f}")
print(f"F1-Score: {f1_score(y_test_enhanced, y_pred_enhanced):.4f}")

# Display feature importance with correlation analysis
plt.figure(figsize=(12, 8))
correlations = data.corr()['Creditability'].sort_values(ascending=False)
sns.barplot(x=correlations.values, y=correlations.index)
plt.title('Feature Correlations with Target')
plt.tight_layout()
plt.savefig('feature_importance.png')
plt.close()

print("\nTop 5 Features by Correlation with Target:")

```



```

print(correlations.head())

# 3. Optimization Step 1: Feature Scaling
# Compare different scaling methods
print("\n--- Feature Scaling Comparison ---")
scaling_methods = {
    'No Scaling': None,
    'StandardScaler': StandardScaler(),
    'MinMaxScaler': MinMaxScaler()
}

for name, scaler in scaling_methods.items():
    if scaler:
        X_train_scaled = scaler.fit_transform(X_train)
        X_test_scaled = scaler.transform(X_test)
    else:
        X_train_scaled = X_train
        X_test_scaled = X_test

    knn = KNeighborsClassifier(n_neighbors=5)
    knn.fit(X_train_scaled, y_train)
    y_pred = knn.predict(X_test_scaled)

    print(f"\n{name} Results:")
    print(f"Accuracy: {accuracy_score(y_test, y_pred):.4f}")
    print(f"F1-Score: {f1_score(y_test, y_pred):.4f}")

# Select the best scaler based on results
best_scaler = StandardScaler() # This should be adjusted based on actual
↪ results
X_train_scaled = best_scaler.fit_transform(X_train)
X_test_scaled = best_scaler.transform(X_test)

# 4. Optimization Step 2: Feature Selection
print("\n--- Feature Selection ---")

# Mutual Information-based feature selection
selector = SelectKBest(mutual_info_classif, k=10)
X_train_selected = selector.fit_transform(X_train_scaled, y_train)
X_test_selected = selector.transform(X_test_scaled)

# Get selected feature names
selected_indices = selector.get_support(indices=True)
selected_features = X_train.columns[selected_indices]
print("\nTop 10 features selected by mutual information:")
print(selected_features.tolist())

```

```

# Test KNN with feature selection
knn_fs = KNeighborsClassifier(n_neighbors=5)
knn_fs.fit(X_train_selected, y_train)
y_pred_fs = knn_fs.predict(X_test_selected)

print("\nKNN with Feature Selection:")
print(f"Accuracy: {accuracy_score(y_test, y_pred_fs):.4f}")
print(f"F1-Score: {f1_score(y_test, y_pred_fs):.4f}")

# 5. Optimization Step 3: PCA Dimensionality Reduction
print("\n--- PCA Dimensionality Reduction ---")

# Find optimal number of PCA components
pca = PCA().fit(X_train_scaled)
plt.figure(figsize=(10, 6))
plt.plot(np.cumsum(pca.explained_variance_ratio_))
plt.xlabel('Number of Components')
plt.ylabel('Cumulative Explained Variance')
plt.title('PCA Explained Variance')
plt.grid(True)
plt.savefig('pca_variance.png')
plt.close()

# Select optimal PCA components (capturing ~80% variance)
n_components = np.argmax(np.cumsum(pca.explained_variance_ratio_) >= 0.8) + 1
print(f"\nOptimal number of PCA components: {n_components}")

pca = PCA(n_components=n_components)
X_train_pca = pca.fit_transform(X_train_scaled)
X_test_pca = pca.transform(X_test_scaled)

# Test KNN with PCA
knn_pca = KNeighborsClassifier(n_neighbors=5)
knn_pca.fit(X_train_pca, y_train)
y_pred_pca = knn_pca.predict(X_test_pca)

print("\nKNN with PCA:")
print(f"Accuracy: {accuracy_score(y_test, y_pred_pca):.4f}")
print(f"F1-Score: {f1_score(y_test, y_pred_pca):.4f}")

# 6. Optimization Step 4: Handling Class Imbalance (SMOTE)
print("\n--- Handling Class Imbalance with SMOTE ---")
print("Original class distribution:")
print(Counter(y_train))

smote = SMOTE(random_state=42)
X_train_smote, y_train_smote = smote.fit_resample(X_train_scaled, y_train)

```

```

print("Class distribution after SMOTE:")
print(Counter(y_train_smote))

# Test KNN with SMOTE
knn_smote = KNeighborsClassifier(n_neighbors=5)
knn_smote.fit(X_train_smote, y_train_smote)
y_pred_smote = knn_smote.predict(X_test_scaled)

print("\nKNN with SMOTE:")
print(f"Accuracy: {accuracy_score(y_test, y_pred_smote):.4f}")
print(f"Precision: {precision_score(y_test, y_pred_smote):.4f}")
print(f"Recall: {recall_score(y_test, y_pred_smote):.4f}")
print(f"F1-Score: {f1_score(y_test, y_pred_smote):.4f}")

# 7. Advanced Optimization: Combining Enhanced Features with SMOTE and Feature
↳ Selection
print("\n--- Advanced Model Optimization ---")

# 7.1 First, scale the enhanced features
scaler_enhanced = StandardScaler()
X_train_enhanced_scaled = scaler_enhanced.fit_transform(X_train_enhanced)
X_test_enhanced_scaled = scaler_enhanced.transform(X_test_enhanced)

# 7.2 Apply SMOTE to balance the classes with enhanced features
smote_enhanced = SMOTE(random_state=42)
X_train_enhanced_smote, y_train_enhanced_smote = smote_enhanced.
↳ fit_resample(X_train_enhanced_scaled, y_train_enhanced)

# 7.3 Apply feature selection to the enhanced dataset
# Use a more robust feature selection approach with Recursive Feature
↳ Elimination
from sklearn.ensemble import RandomForestClassifier
from sklearn.feature_selection import RFECV

# Use RandomForest as the base estimator for feature selection
base_estimator = RandomForestClassifier(n_estimators=100, random_state=42)
rfecv = RFECV(
    estimator=base_estimator,
    step=1,
    cv=StratifiedKFold(5),
    scoring='f1',
    min_features_to_select=10
)

# Apply RFECV to the enhanced features
rfecv.fit(X_train_enhanced_smote, y_train_enhanced_smote)

```

```

X_train_enhanced_selected = rfecv.transform(X_train_enhanced_smote)
X_test_enhanced_selected = rfecv.transform(X_test_enhanced_scaled)

# Get selected feature names
selected_indices = rfecv.support_
selected_features_enhanced = X_enhanced.columns[selected_indices]
print(f"\nOptimal number of features selected:␣
↳{len(selected_features_enhanced)}")
print("Selected features:")
print(selected_features_enhanced.tolist())

# 7.4 Apply an expanded hyperparameter grid for the final model
# Use the best preprocessing method - enhanced features with SMOTE and feature␣
↳selection
X_train_opt = X_train_enhanced_selected
y_train_opt = y_train_enhanced_smote
X_test_opt = X_test_enhanced_selected

# Enhanced parameter grid with more options
param_grid = {
    'n_neighbors': [3, 5, 7, 9, 11, 13, 15, 17, 19, 21, 23],
    'weights': ['uniform', 'distance'],
    'metric': ['euclidean', 'manhattan', 'minkowski', 'chebyshev'],
    'p': [1, 2, 3], # p=1 is manhattan, p=2 is euclidean, p=3 is cubic
    'leaf_size': [10, 20, 30, 40, 50] # Additional parameter for efficiency
}

# Create cross-validation strategy
cv = StratifiedKFold(n_splits=5, shuffle=True, random_state=42)

# Create grid search with multiple scoring metrics
grid_search = GridSearchCV(
    KNeighborsClassifier(),
    param_grid,
    cv=cv,
    scoring={
        'accuracy': 'accuracy',
        'f1': 'f1',
        'precision': 'precision',
        'recall': 'recall'
    },
    refit='accuracy', # Optimize for accuracy to reach 80%
    n_jobs=-1,
    verbose=1
)

# Fit grid search

```

```

grid_search.fit(X_train_opt, y_train_opt)

# Print best parameters
print("\nBest parameters found by grid search:")
print(grid_search.best_params_)
print(f"Best cross-validation score: {grid_search.best_score_:.4f}")

# Test KNN with optimized hyperparameters
best_knn = grid_search.best_estimator_
y_pred_best = best_knn.predict(X_test_opt)

print("\nKNN with Advanced Optimization:")
print(f"Accuracy: {accuracy_score(y_test, y_pred_best):.4f}")
print(f"Precision: {precision_score(y_test, y_pred_best):.4f}")
print(f"Recall: {recall_score(y_test, y_pred_best):.4f}")
print(f"F1-Score: {f1_score(y_test, y_pred_best):.4f}")

# 7.5 Apply threshold tuning to further optimize classification boundary
print("\n--- Probability Threshold Tuning ---")
from sklearn.metrics import roc_curve, precision_recall_curve

# Get predicted probabilities
y_proba = best_knn.predict_proba(X_test_opt)[: , 1]

# Find optimal threshold using ROC curve
fpr, tpr, thresholds = roc_curve(y_test, y_proba)
optimal_idx = np.argmax(tpr - fpr)
optimal_threshold = thresholds[optimal_idx]
print(f"Optimal threshold from ROC curve: {optimal_threshold:.4f}")

# Apply optimal threshold
y_pred_threshold = (y_proba >= optimal_threshold).astype(int)

print("\nKNN with Threshold Optimization:")
print(f"Accuracy: {accuracy_score(y_test, y_pred_threshold):.4f}")
print(f"Precision: {precision_score(y_test, y_pred_threshold):.4f}")
print(f"Recall: {recall_score(y_test, y_pred_threshold):.4f}")
print(f"F1-Score: {f1_score(y_test, y_pred_threshold):.4f}")

# 8. Visualize confusion matrix
cm = confusion_matrix(y_test, y_pred_best)
plt.figure(figsize=(8, 6))
sns.heatmap(cm, annot=True, fmt='d', cmap='Blues', cbar=False)
plt.title('Confusion Matrix')
plt.xlabel('Predicted Label')
plt.ylabel('True Label')
plt.savefig('confusion_matrix.png')

```

```

plt.close()

# 9. Visualize model performance for different K values
k_range = list(range(1, 31))
k_scores = []

for k in k_range:
    knn = KNeighborsClassifier(n_neighbors=k)
    scores = cross_val_score(knn, X_train_opt, y_train_opt, cv=5, scoring='f1')
    k_scores.append(scores.mean())

plt.figure(figsize=(10, 6))
plt.plot(k_range, k_scores)
plt.xlabel('Value of K')
plt.ylabel('Cross-Validated F1-Score')
plt.title('KNN Performance for Different K Values')
plt.grid(True)
plt.savefig('k_performance.png')
plt.close()

# 10. Display feature importance for the final selected model
# Using permutation importance for the best model
print("\n--- Feature Importance Analysis ---")
from sklearn.inspection import permutation_importance

# We need to make sure we're using the correct dataset dimensions
# For the optimized model, we used X_test_opt (which matches what the model was
↳ trained on)
result = permutation_importance(best_knn, X_test_opt, y_test, n_repeats=10,
↳ random_state=42)

# Make sure the feature names match the number of features in X_test_opt
# If we used feature selection, we need to get the selected feature names
if 'rfecv' in locals() and hasattr(rfecv, 'support_'):
    # Get feature names from RFECV's selected features
    selected_indices = rfecv.support_
    feature_names = X_enhanced.columns[selected_indices].tolist()
else:
    # Fallback to the dimensions of X_test_opt
    feature_names = [f"Feature_{i}" for i in range(X_test_opt.shape[1])]

# Make sure lengths match before creating DataFrame
if len(feature_names) == len(result.importances_mean):
    # Create feature importance DataFrame
    feature_importance = pd.DataFrame({
        'Feature': feature_names,
        'Importance': result.importances_mean
    })

```

```

}).sort_values('Importance', ascending=False)

print("\nFeature Importance for Final Model:")
print(feature_importance.head(10)) # Show top 10 features
else:
    print(f"\nWarning: Feature names length ({len(feature_names)}) doesn't match importance array length ({len(result.importances_mean)})")
    print("Generating generic feature names instead")
    # Create with generic feature names
    feature_importance = pd.DataFrame({
        'Feature': [f"Feature_{i}" for i in range(len(result.importances_mean))],
        'Importance': result.importances_mean
    }).sort_values('Importance', ascending=False)

    print("\nFeature Importance for Final Model (with generic names):")
    print(feature_importance.head(10)) # Show top 10 features

# 11. Ensemble Model for Maximum Performance
print("\n--- Ensemble KNN Models ---")

# Create an ensemble of KNN models with different hyperparameters
from sklearn.ensemble import VotingClassifier

# Create several KNN models with different configurations
knn1 = KNeighborsClassifier(n_neighbors=3, weights='distance',
    metric='manhattan')
knn2 = KNeighborsClassifier(n_neighbors=7, weights='uniform',
    metric='euclidean')
knn3 = KNeighborsClassifier(n_neighbors=11, weights='distance',
    metric='minkowski', p=3)
knn4 = KNeighborsClassifier(n_neighbors=15, weights='distance',
    metric='chebyshev')
knn5 = KNeighborsClassifier(**grid_search.best_params_) # Best model from grid search

# Create a voting classifier
voting_clf = VotingClassifier(
    estimators=[
        ('knn1', knn1),
        ('knn2', knn2),
        ('knn3', knn3),
        ('knn4', knn4),
        ('knn5', knn5)
    ],
    voting='soft' # Use probability estimates for voting

```

```

)

# Fit the ensemble
voting_clf.fit(X_train_opt, y_train_opt)

# Predict
y_pred_ensemble = voting_clf.predict(X_test_opt)

print("\nEnsemble KNN Model:")
print(f"Accuracy: {accuracy_score(y_test, y_pred_ensemble):.4f}")
print(f"Precision: {precision_score(y_test, y_pred_ensemble):.4f}")
print(f"Recall: {recall_score(y_test, y_pred_ensemble):.4f}")
print(f"F1-Score: {f1_score(y_test, y_pred_ensemble):.4f}")

# 12. Final comparison of all models
print("\n--- Final Model Comparison ---")

models_results = {
    'Baseline KNN (K=5)': (baseline_accuracy, f1_score(y_test, baseline_pred)),
    'KNN with Enhanced Features': (accuracy_score(y_test_enhanced,
↳ y_pred_enhanced), f1_score(y_test_enhanced, y_pred_enhanced)),
    'KNN with Feature Selection': (accuracy_score(y_test, y_pred_fs),
↳ f1_score(y_test, y_pred_fs)),
    'KNN with PCA': (accuracy_score(y_test, y_pred_pca), f1_score(y_test,
↳ y_pred_pca)),
    'KNN with SMOTE': (accuracy_score(y_test, y_pred_smote), f1_score(y_test,
↳ y_pred_smote)),
    'KNN with Advanced Optimization': (accuracy_score(y_test, y_pred_best),
↳ f1_score(y_test, y_pred_best)),
    'KNN with Threshold Tuning': (accuracy_score(y_test, y_pred_threshold),
↳ f1_score(y_test, y_pred_threshold)),
    'Ensemble KNN': (accuracy_score(y_test, y_pred_ensemble), f1_score(y_test,
↳ y_pred_ensemble))
}

# Create a DataFrame to compare results
comparison_df = pd.DataFrame.from_dict(models_results, orient='index',
↳ columns=['Accuracy', 'F1-Score'])
print(comparison_df.sort_values('Accuracy', ascending=False))

# Plot comparison
plt.figure(figsize=(14, 8))
comparison_df.plot(kind='bar', figsize=(14, 8))
plt.title('Model Performance Comparison')
plt.ylabel('Score')
plt.axhline(y=0.8, color='r', linestyle='--', label='80% Target')

```



```

plt.legend()
plt.xticks(rotation=45, ha='right')
plt.tight_layout()
plt.savefig('model_comparison.png')
plt.close()

# Identify best model from all approaches
best_model_name = comparison_df.sort_values('Accuracy', ascending=False).
    ↪index[0]
best_accuracy = comparison_df.sort_values('Accuracy',
    ↪ascending=False)['Accuracy'].iloc[0]

print("\n--- Model Optimization Complete ---")
print(f"Initial model accuracy: 69%")
print(f"Best model ({best_model_name}) accuracy: {best_accuracy:.2%}")
print(f"Improvement: {best_accuracy - 0.69:.2%}")

# Check if we've reached 80% accuracy
if best_accuracy >= 0.8:
    print("\n SUCCESS! We've achieved the 80% accuracy target!")
else:
    print(f"\nWe've improved the model significantly, but more work may be
    ↪needed to reach 80%")
    print(f"Current best accuracy: {best_accuracy:.2%}")

# 13. Save the best model for future use
import joblib

# Determine which model to save based on highest accuracy
if best_model_name == 'Ensemble KNN':
    best_model = voting_clf
    preprocessing_pipeline = {
        'scaler': scaler_enhanced,
        'feature_selector': rfecv
    }
elif best_model_name == 'KNN with Threshold Tuning':
    best_model = best_knn
    preprocessing_pipeline = {
        'scaler': scaler_enhanced,
        'feature_selector': rfecv,
        'optimal_threshold': optimal_threshold
    }
else:
    best_model = best_knn
    preprocessing_pipeline = {
        'scaler': scaler_enhanced,
        'feature_selector': rfecv
    }

```

```

}

# Save everything
joblib.dump(best_model, 'best_knn_model.pkl')
joblib.dump(preprocessing_pipeline, 'preprocessing_pipeline.pkl')
joblib.dump(selected_features_enhanced, 'selected_features.pkl')

print(f"\nBest model saved as 'best_knn_model.pkl'")
print("Preprocessing pipeline and selected features also saved")

# 14. Document model interpretability
print("\n--- Model Interpretability ---")
print("Top 10 most important features for loan approval prediction:")
if 'feature_importance' in locals():
    print(feature_importance.head(10))
else:
    print("Feature importance calculation not available")

# Create a function that explains predictions
def explain_prediction(sample_data, model, preprocessing, features_shape):
    """
    Explain why a specific loan application was approved or rejected

    Parameters:
    -----
    sample_data : numpy.ndarray or pandas.DataFrame
        The data for a single customer, in the same format as the original
        ↪ dataset
    model : trained model
        The best KNN model
    preprocessing : dict
        Dictionary containing preprocessing components
    features_shape : int
        Number of features expected by the model

    Returns:
    -----
    dict
        Prediction explanation
    """
    # Convert to numpy if dataframe
    if hasattr(sample_data, 'values'):
        sample_data = sample_data.values

    # Preprocess the sample
    if 'scaler' in preprocessing:

```

```

        sample_scaled = preprocessing['scaler'].transform(sample_data.
↪reshape(1, -1))
    else:
        sample_scaled = sample_data.reshape(1, -1)

    if 'feature_selector' in preprocessing and 'rfecv' in preprocessing:
        sample_processed = preprocessing['feature_selector'].
↪transform(sample_scaled)
    else:
        sample_processed = sample_scaled

    # Check dimensions
    if sample_processed.shape[1] != features_shape:
        raise ValueError(f"Model expects {features_shape} features, but got_
↪{sample_processed.shape[1]} features")

    # Get prediction and probability
    prediction = model.predict(sample_processed)[0]
    probability = model.predict_proba(sample_processed)[0]

    # Identify most influential neighbors
    distances, indices = model.kneighbors(sample_processed)

    result = {
        'prediction': 'Approved' if prediction == 1 else 'Rejected',
        'confidence': float(probability[prediction]),
        'nearest_neighbor_distances': distances[0].tolist()
    }

    return result

print("\nExample of how to use the explainer function:")
print("explain_prediction(customer_data, best_model, preprocessing_pipeline,
↪X_test_opt.shape[1])")

```

Dataset shape: (1000, 21)

Target Variable (Creditability) Distribution:

Creditability

1 700

0 300

Name: count, dtype: int64

Class imbalance ratio: 0.4286

Training set shape: (800, 20), Test set shape: (200, 20)

Baseline KNN Model (K=5) Results:

Accuracy: 0.6650
Precision: 0.7296
Recall: 0.8286
F1-Score: 0.7759

--- Feature Engineering ---

KNN with Enhanced Features:
Accuracy: 0.7400
Precision: 0.7973
Recall: 0.8429
F1-Score: 0.8194

Top 5 Features by Correlation with Target:
Creditability 1.000000
Account Balance 0.350847
Payment Status of Previous Credit 0.228785
Value Savings/Stocks 0.178943
Length of current employment 0.116002
Name: Creditability, dtype: float64

--- Feature Scaling Comparison ---

No Scaling Results:
Accuracy: 0.6650
F1-Score: 0.7759

StandardScaler Results:
Accuracy: 0.7200
F1-Score: 0.8108

MinMaxScaler Results:
Accuracy: 0.7500
F1-Score: 0.8322

--- Feature Selection ---

Top 10 features selected by mutual information:
['Account Balance', 'Duration of Credit (month)', 'Payment Status of Previous Credit', 'Purpose', 'Credit Amount', 'Value Savings/Stocks', 'Length of current employment', 'Age (years)', 'Type of apartment', 'No of Credits at this Bank']

KNN with Feature Selection:
Accuracy: 0.7600
F1-Score: 0.8431

--- PCA Dimensionality Reduction ---

Optimal number of PCA components: 13

KNN with PCA:

Accuracy: 0.7300

F1-Score: 0.8176

--- Handling Class Imbalance with SMOTE ---

Original class distribution:

Counter({1: 560, 0: 240})

Class distribution after SMOTE:

Counter({1: 560, 0: 560})

KNN with SMOTE:

Accuracy: 0.6650

Precision: 0.8476

Recall: 0.6357

F1-Score: 0.7265

--- Advanced Model Optimization ---

Optimal number of features selected: 43

Selected features:

['Account Balance', 'Duration of Credit (month)', 'Payment Status of Previous Credit', 'Purpose', 'Credit Amount', 'Value Savings/Stocks', 'Length of current employment', 'Instalment per cent', 'Sex & Marital Status', 'Guarantors', 'Duration in Current address', 'Most valuable available asset', 'Age (years)', 'Concurrent Credits', 'Type of apartment', 'No of Credits at this Bank', 'Occupation', 'No of dependents', 'Telephone', 'Foreign Worker', 'balance_duration', 'payment_credit_ratio', 'savings_credit_ratio', 'age_employment', 'credit_income_proxy', 'poly_0', 'poly_1', 'poly_2', 'poly_3', 'poly_4', 'poly_5', 'poly_6', 'poly_7', 'poly_8', 'poly_9', 'poly_10', 'poly_11', 'poly_12', 'poly_13', 'poly_14', 'age_group', 'credit_bin', 'duration_bin']

Fitting 5 folds for each of 1320 candidates, totalling 6600 fits

Best parameters found by grid search:

{'leaf_size': 10, 'metric': 'manhattan', 'n_neighbors': 3, 'p': 1, 'weights': 'distance'}

Best cross-validation score: 0.8223

KNN with Advanced Optimization:

Accuracy: 0.7150

Precision: 0.8217

Recall: 0.7571

F1-Score: 0.7881

--- Probability Threshold Tuning ---

Optimal threshold from ROC curve: 0.6627

KNN with Threshold Optimization:

Accuracy: 0.7100

Precision: 0.8475

Recall: 0.7143

F1-Score: 0.7752

--- Feature Importance Analysis ---

Feature Importance for Final Model:

	Feature	Importance
10	Duration in Current address	0.0215
39	poly_14	0.0135
8	Sex & Marital Status	0.0120
35	poly_10	0.0115
15	No of Credits at this Bank	0.0100
22	savings_credit_ratio	0.0090
42	duration_bin	0.0085
41	credit_bin	0.0065
11	Most valuable available asset	0.0060
28	poly_3	0.0060

--- Ensemble KNN Models ---

Ensemble KNN Model:

Accuracy: 0.7150

Precision: 0.8547

Recall: 0.7143

F1-Score: 0.7782

--- Final Model Comparison ---

	Accuracy	F1-Score
KNN with Feature Selection	0.760	0.843137
KNN with Enhanced Features	0.740	0.819444
KNN with PCA	0.730	0.817568
KNN with Advanced Optimization	0.715	0.788104
Ensemble KNN	0.715	0.778210
KNN with Threshold Tuning	0.710	0.775194
Baseline KNN (K=5)	0.665	0.775920
KNN with SMOTE	0.665	0.726531

--- Model Optimization Complete ---

Initial model accuracy: 69%

Best model (KNN with Feature Selection) accuracy: 76.00%

Improvement: 7.00%

We've improved the model significantly, but more work may be needed to reach 80%
Current best accuracy: 76.00%

```
Best model saved as 'best_knn_model.pkl'
Preprocessing pipeline and selected features also saved
```

```
--- Model Interpretability ---
Top 10 most important features for loan approval prediction:
```

	Feature	Importance
10	Duration in Current address	0.0215
39	poly_14	0.0135
8	Sex & Marital Status	0.0120
35	poly_10	0.0115
15	No of Credits at this Bank	0.0100
22	savings_credit_ratio	0.0090
42	duration_bin	0.0085
41	credit_bin	0.0065
11	Most valuable available asset	0.0060
28	poly_3	0.0060

```
Example of how to use the explainer function:
explain_prediction(customer_data, best_model, preprocessing_pipeline,
X_test_opt.shape[1])
```

```
<Figure size 1400x800 with 0 Axes>
```

```
[ ]:
```

