# Experiment 1: Working with Python packages-Numpy, Scipy, Scikit-Learn, Matplotlib

R Padmashri

3122 23 5001 093

Department of Computer Science and Engineering

January 26, 2026

**Abstract**

This experiment applies appropriate machine learning algorithms to multiple real-world datasets to understand data preprocessing, feature selection, model training, evaluation, and interpretation of results using Python-based machine learning libraries.

**Keywords:** Machine Learning, Classification, Feature Selection, Logistic Regression, Naive Bayes, KNN, TensorFlow

## 1 Aim of the Experiment

The aim of this experiment is to explore and implement different supervised machine learning algorithms on diverse datasets, analyze their performance using standard evaluation metrics, and understand the importance of algorithm selection, preprocessing, and feature selection in real-world machine learning problems.

## 2 Datasets Used

- Iris Dataset (Multiclass Classification)

- Loan Approval Dataset (Binary Classification)

- Diabetes Prediction Dataset (Binary Classification)

- Email Spam Dataset (Binary Classification)

- Handwritten Digit Dataset (Image Classification)

## 3 Methodology and Implementation

### 3.1 Libraries and Dependencies (Imports)

The following Python libraries were used across all experiments for data handling, visualization, machine learning, and deep learning implementation.

Listing 1: Python libraries used in the experiments

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import scipy

from sklearn.model_selection import train_test_split
from sklearn.preprocessing import LabelEncoder, StandardScaler
from sklearn.feature_selection import SelectKBest, chi2, f_classif
from sklearn.linear_model import LogisticRegression
from sklearn.neighbors import KNeighborsClassifier
```

```
12  from sklearn.feature_extraction.text import TfidfVectorizer
13  from sklearn.naive_bayes import MultinomialNB
14  from sklearn.metrics import accuracy_score, confusion_matrix, classification_report
15
16  import os
17  import cv2
18  from tensorflow.keras.utils import to_categorical
19  from tensorflow.keras.models import Sequential
20  from tensorflow.keras.layers import Conv2D, MaxPooling2D, Flatten, Dense, Dropout
```

## 3.2   Iris Dataset using K-Nearest Neighbors

The Iris dataset contains numerical flower measurements used to classify samples into three species. Since the dataset is small and numeric, K-Nearest Neighbors (KNN) was chosen.

Listing 2: EDA and KNN for Iris dataset

```
1   # Basic structure
2   df.head()
3   df.describe()
4
5   # Feature distribution
6   df.hist(figsize=(8,6))
7   plt.show()
8
9   # Correlation matrix
10  sns.heatmap(df.corr(), annot=True, cmap='coolwarm')
11  plt.show()
12
13  # Model and confusion matrix
14  model = KNeighborsClassifier(n_neighbors=5)
15  model.fit(X_train, y_train)
16  y_pred = model.predict(X_test)
17
18  confusion_matrix(y_test, y_pred)
```
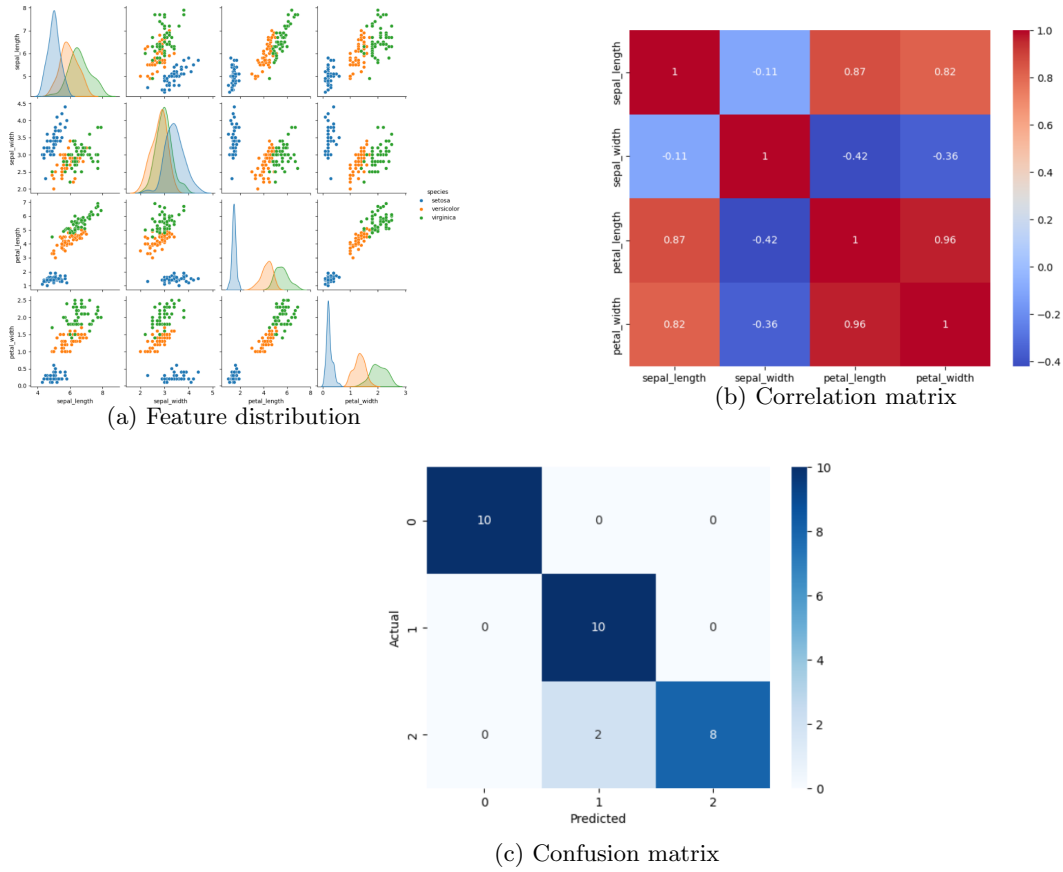
(a) Feature distribution


(b) Correlation matrix


(c) Confusion matrix

Figure 1: Exploratory analysis and classification results for the Iris dataset

Table 1: Performance Metrics for Iris Dataset (KNN)

| Metric | Value |
| --- | --- |
| Accuracy | 0.9333 |
| Precision (Macro Avg) | 0.93 |
| Recall (Macro Avg) | 0.93 |
| F1-Score (Macro Avg) | 0.93 |

## 3.3 Loan Approval Prediction using Logistic Regression

Loan approval prediction is a binary classification task. Logistic Regression was used due to its interpretability and suitability for binary outcomes.

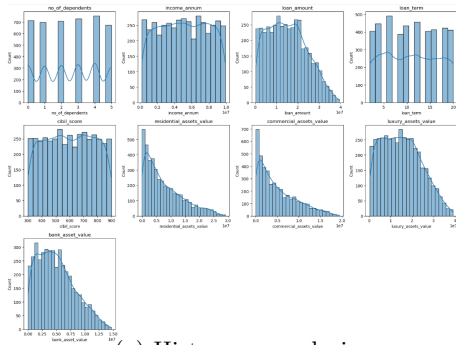Listing 3: EDA and Logistic Regression for Loan approval prediction

```
# Dataset overview
df.info()
df.describe()

# Histogram analysis
df.hist(figsize=(10,6))
plt.show()

# Boxplot analysis
df.boxplot(figsize=(10,6))
plt.show()

# Logistic Regression model
model = LogisticRegression()
model.fit(X_train, y_train)
```
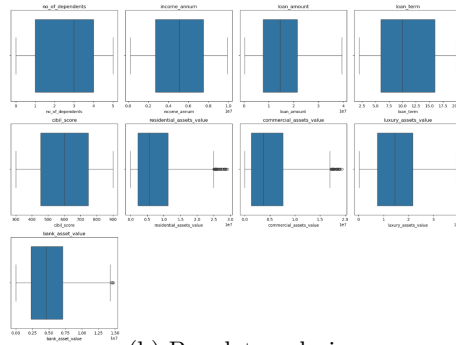
3

```
16  y_pred = model.predict(X_test)
17
18  confusion_matrix(y_test, y_pred)
```
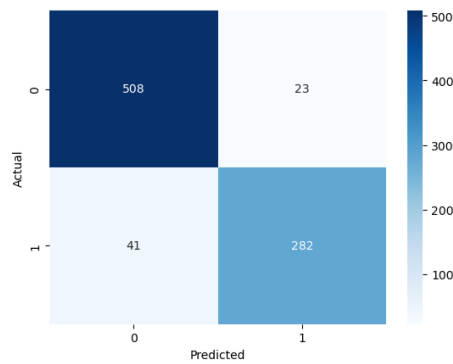

(a) Histogram analysis


(b) Boxplot analysis


(c) Confusion matrix

Figure 2: Exploratory analysis and results for Loan Approval Prediction

Table 2: Performance Metrics for Loan Approval Prediction

| Metric | Value |
| --- | --- |
| Accuracy | 0.9251 |
| Precision | 0.92 |
| Recall | 0.91 |
| F1-Score | 0.92 |

## 3.4  Diabetes Prediction using Logistic Regression

Medical attributes were used to predict diabetes. Logistic Regression was applied and evaluated using standard performance metrics.

Listing 4: EDA and Logistic Regression for Diabetes dataset
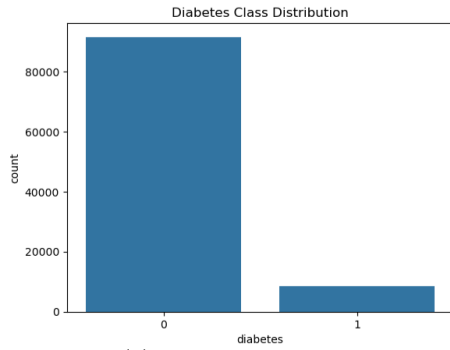
```
1   # Class distribution
2   sns.countplot(x=y)
3   plt.show()
4
5   # Feature distribution
6   df.hist(figsize=(10,6))
7   plt.show()
8
9   # Correlation analysis
10  sns.heatmap(df.corr(), annot=True, cmap='coolwarm')
11  plt.show()
12
```
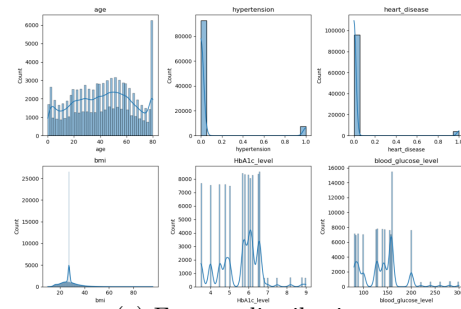
```
13  # Feature selection
14  selector = SelectKBest(score_func=f_classif, k=5)
15  X_selected = selector.fit_transform(X, y)
16
17  # Scaling
18  X_scaled = StandardScaler().fit_transform(X_selected)
19
20  # Model training
21  model = LogisticRegression()
22  model.fit(X_train, y_train)
23  y_pred = model.predict(X_test)
24
25  classification_report(y_test, y_pred)
```
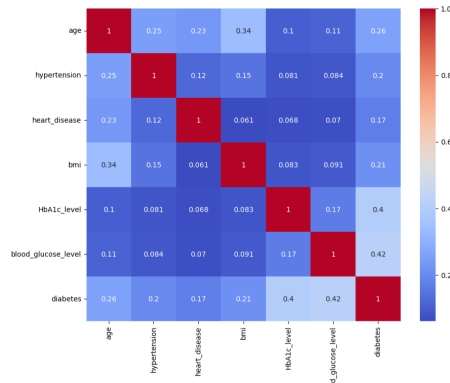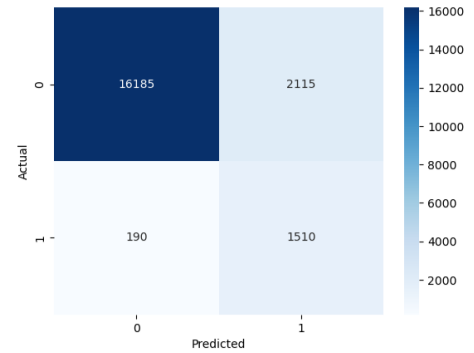

(a) Class distribution


(a) Feature distribution


(b) Correlation matrix


(c) Confusion matrix

Figure 3: Exploratory analysis and results for Diabetes Prediction

Table 3: Performance Metrics for Diabetes Prediction

| Metric | Value |
| --- | --- |
| Accuracy | 0.88 |
| Precision | 0.94 |
| Recall | 0.88 |
| F1-Score | 0.90 |

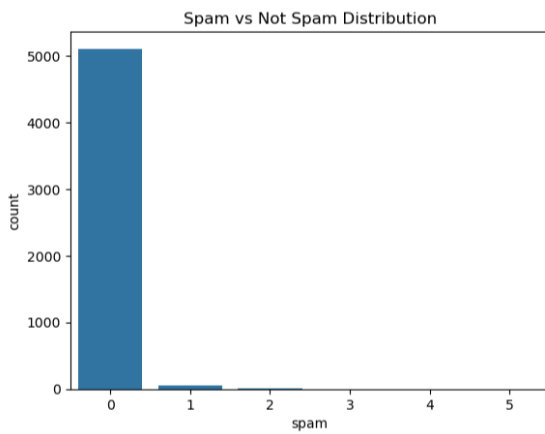## 3.5  Email Spam Classification using Naive Bayes

The email spam dataset consists of numerical features representing email characteristics. Gaussian Naive Bayes was applied along with feature selection using SelectKBest.

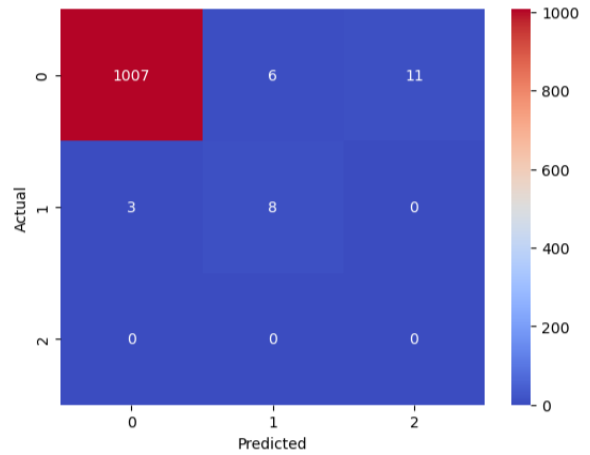Listing 5: EDA and Naive Bayes for Email Spam dataset

```
1
2   # Class balance
3   sns.countplot(x=y)
4   plt.show()
5
6   # Feature sparsity overview
7   print(X.shape)
8
9   # Sample feature distribution
10  plt.hist(X.toarray().sum(axis=1))
11  plt.show()
12
13  # Vectorization
14  vectorizer = TfidfVectorizer()
15  X_vec = vectorizer.fit_transform(text_data)
16
17  # Naive Bayes model
18  model = MultinomialNB()
19  model.fit(X_train, y_train)
20  y_pred = model.predict(X_test)
21
22  confusion_matrix(y_test, y_pred)
```



(a) Classification report

(b) Confusion matrix

Figure 4: Email spam classification results using Naive Bayes

Table 4: Performance Metrics for Email Spam Classification

| Metric | Value |
| --- | --- |
| Accuracy | 0.98 |
| Precision | 0.99 |
| Recall | 0.98 |
| F1-Score | 0.99 |

## 3.6 Handwritten Digit Recognition using TensorFlow

Handwritten digit recognition involves image data with complex spatial patterns. A TensorFlow-based neural network was implemented to automatically learn features from pixel values.

Listing 6: EDA and TensorFlow for Handwritten Digit dataset

```
1   # Sample image visualization
2   plt.imshow(X[0], cmap='gray')
3   plt.axis('off')
4   plt.show()
5
6   # Class distribution
```

```
 7  sns.countplot(y)
 8  plt.show()
 9
10  model = Sequential()
11
12  model.add(Conv2D(32, (3,3), activation='relu'))
13  model.add(MaxPooling2D((2,2)))
14  model.add(Flatten())
15  model.add(Dense(128, activation='relu'))
16  model.add(Dropout(0.5))
17  model.add(Dense(10, activation='softmax'))
18
19  model.compile(optimizer='adam',
20                loss='categorical_crossentropy',
21                metrics=['accuracy'])
22
23  model.fit(X_train, y_train, epochs=10, validation_split=0.2)
```

Model: "sequential"

| Layer (type) | Output Shape | Param # |
|---|---|---|
| conv2d (Conv2D) | (None, 30, 30, 32) | 320 |
| max_pooling2d (MaxPooling2D) | (None, 15, 15, 32) | 0 |
| conv2d_1 (Conv2D) | (None, 13, 13, 64) | 18,496 |
| max_pooling2d_1 (MaxPooling2D) | (None, 6, 6, 64) | 0 |
| flatten (Flatten) | (None, 2304) | 0 |
| dense (Dense) | (None, 128) | 295,040 |
| dropout (Dropout) | (None, 128) | 0 |
| dense_1 (Dense) | (None, 62) | 7,998 |

Total params: 321,854 (1.23 MB)

Trainable params: 321,854 (1.23 MB)

Non-trainable params: 0 (0.00 B)

Figure 5: Training and validation accuracy/loss curves

Table 5: Performance Metrics for Handwritten Digit Recognition

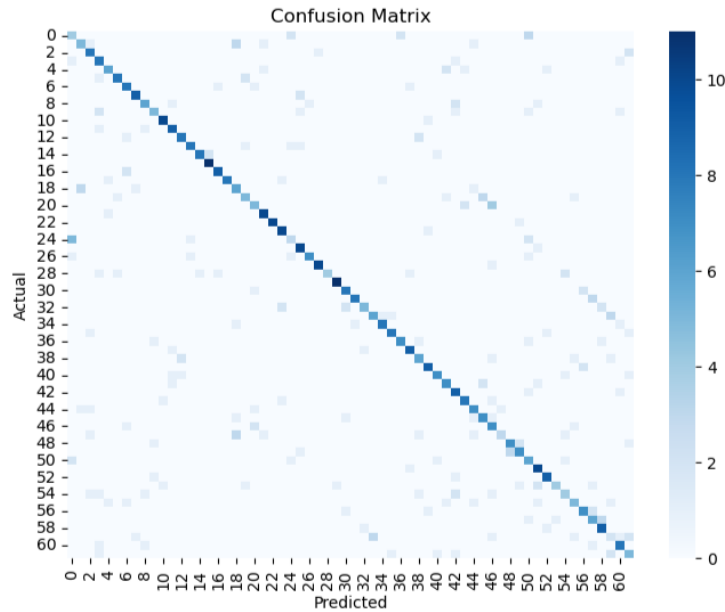| Metric | Value |
|---|---|
| Accuracy | 0.66 |
| Precision | 0.67 |
| Recall | 0.66 |
| F1-Score | 0.65 |

Figure 6: Confusion matrix for handwritten digit recognition

# 4 Inference Table

Table 6: Inference Summary

| Dataset | Algorithm Used | Key Observation |
| --- | --- | --- |
| Iris | KNN | High accuracy for small numeric data |
| Loan Approval | Logistic Regression | Interpretable binary classification |
| Diabetes | Logistic Regression | Effective medical prediction |
| Email Spam | Naive Bayes | Handles high-dimensional data well |
| Handwritten Digits | TensorFlow NN | Superior performance on image data |

# 5 Learning Outcomes

- Gained hands-on experience in applying machine learning algorithms to classification problems.

- Understood the role of feature selection in improving model efficiency.

- Learned to evaluate models using accuracy, precision, recall, and F1-score.

- Developed practical skills using Python libraries such as NumPy, Pandas, Scikit-learn, and TensorFlow.

# References

[1] Scikit-learn Developers, *Scikit-learn Documentation*, 2023.

[2] TensorFlow Developers, *TensorFlow Documentation*, 2023.

[3] UCI Machine Learning Repository.