# Data and Artificial Intelligence
# Cyber Shujaa Program

## Week 8 Assignment
## Classification Models

**Student Name:** Shirleen Nanetia Simon

**Student ID:** CS-DA01-25077

## Introduction

The purpose of this assignment was to apply supervised machine learning classification techniques to the Wine dataset from scikit-learn. The goal was to build and evaluate several models, analyze their performance using standard metrics, and determine which model performs best under similar conditions. This hands-on exercise reinforced practical skills in data preprocessing, exploratory data analysis (EDA), model implementation, and evaluation.
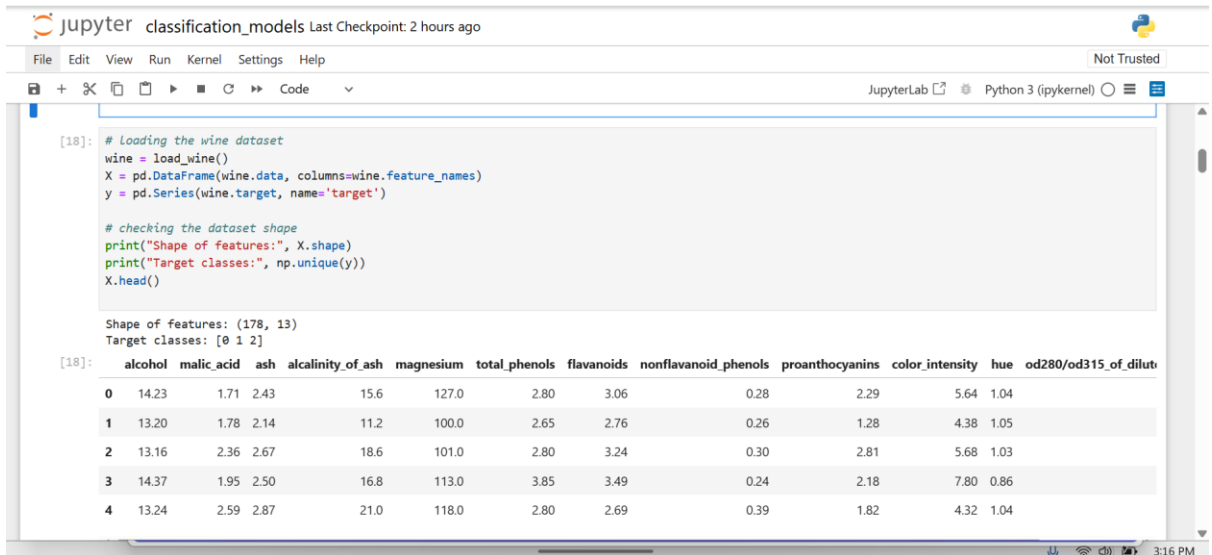
The following classification models were implemented:

- Logistic Regression
- Decision Tree
- Random Forest
- k-Nearest Neighbors (KNN)
- Naive Bayes
- Support Vector Machine (SVM)

## Tasks Completed

### 1. Dataset Loading and Exploration

- The Wine dataset was loaded from scikit-learn.
- The dataset contains **178 samples** with **13 features** and a target variable representing three wine classes.
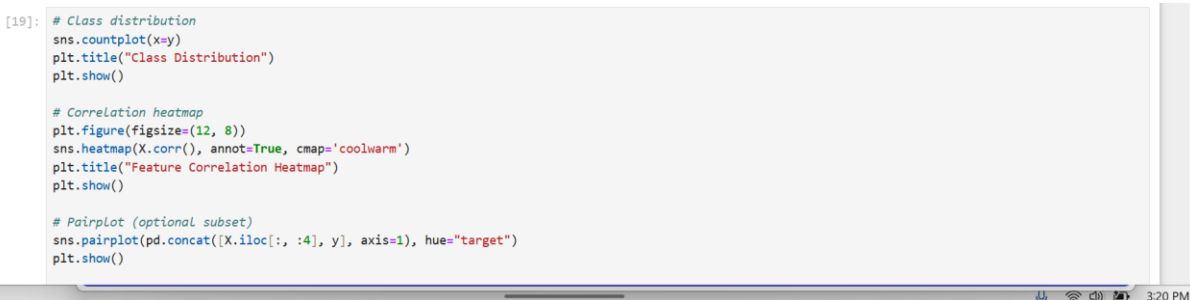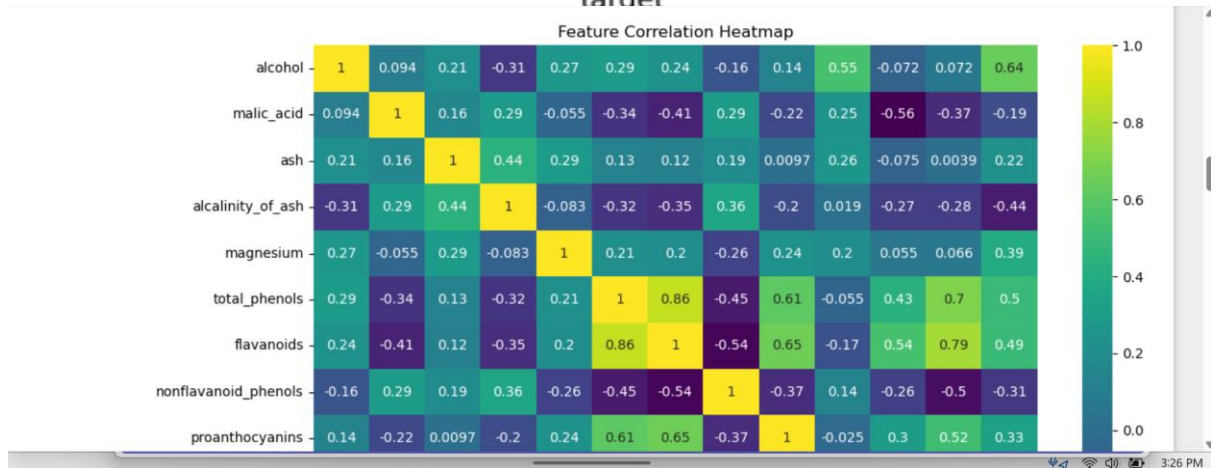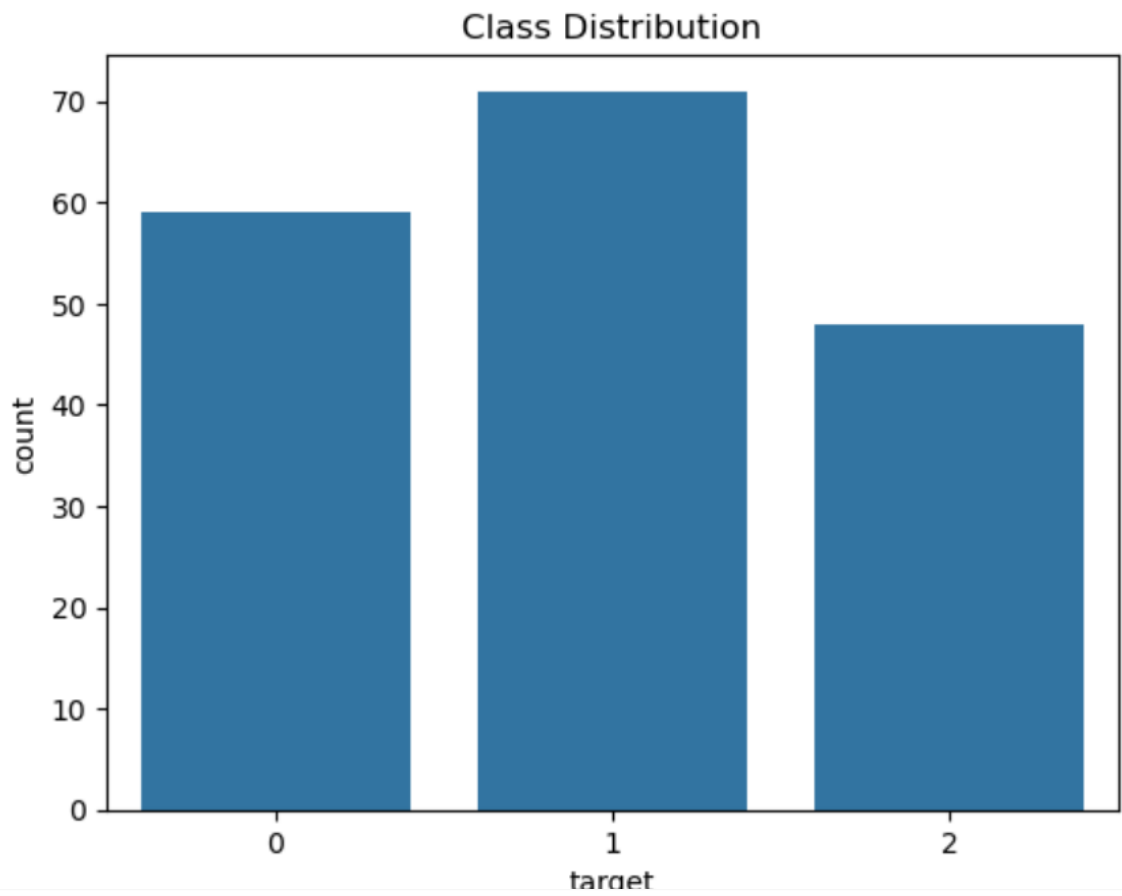


- Exploratory Data Analysis (EDA) was performed:
  - Checked for missing values and found none.
  - Visualized feature distributions using histograms and pair plots.
  - Analyzed the correlation matrix to understand relationships between variables.
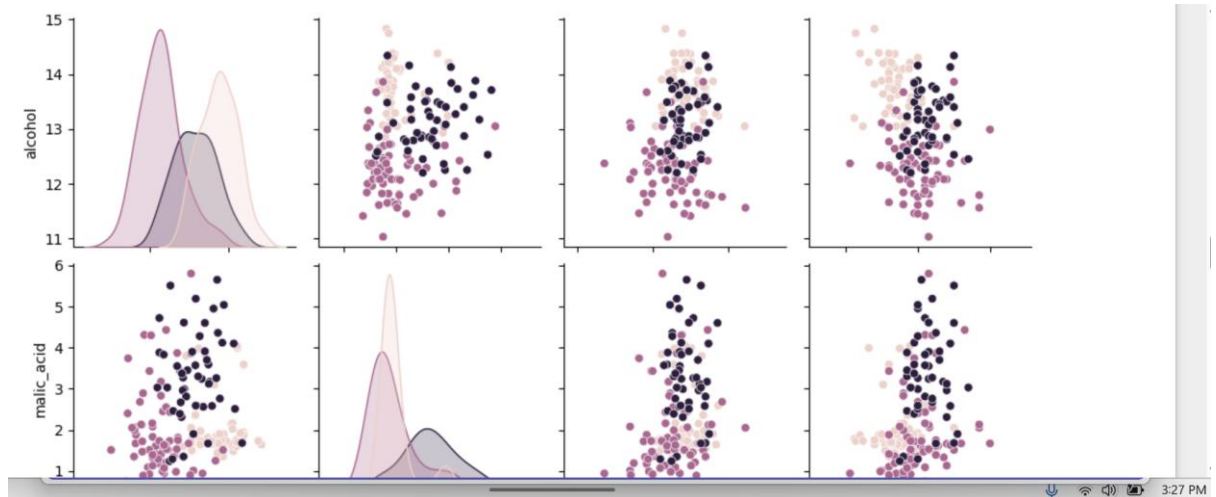  - Observed class distributions to confirm a balanced dataset

## Class Distribution



## Feature Correlation Heatmap

## 2. Data Preparation

- Features and target variables were separated.
- The data was split into **training and testing sets (70:30 split)** using train_test_split.
- Feature scaling was applied where necessary, especially for models sensitive to feature magnitudes (e.g., Logistic Regression, KNN, and SVM).

```python
# Standardize the features
scaler = StandardScaler()
X_scaled = scaler.fit_transform(X)

# Train/Test split
X_train, X_test, y_train, y_test = train_test_split(X_scaled, y, test_size=0.3, random_state=42)
```

## 3. Model Building , Training and Evaluation

Six classification algorithms were trained on the Wine dataset:

- **Logistic Regression:** A linear model suitable for multiclass problems.
- **Decision Tree Classifier:** A non-linear model capable of capturing complex patterns.
- **Random Forest Classifier:** An ensemble of decision trees to improve accuracy and reduce overfitting.
- **K-Nearest Neighbors (KNN):** A distance-based classifier relying on neighboring data points.
- **Naive Bayes Classifier:** A probabilistic model assuming feature independence.
- **Support Vector Machine (SVM):** A classifier maximizing margins between classes.

For each model, the following metrics were computed:

- **Accuracy Score:** Overall correctness of predictions.
- **Classification Report:** Precision, recall, and F1-score per class.
- **Confusion Matrix:** Detailed comparison of true vs. predicted values, visualized using heatmaps.

```python
def plot_conf_matrix(y_true, y_pred, title):
    cm = confusion_matrix(y_true, y_pred)
    plt.figure(figsize=(5, 4))
    sns.heatmap(cm, annot=True, fmt='d', cmap='Blues', xticklabels=wine.target_names, yticklabels=wine.target_names)
    plt.xlabel('Predicted')
    plt.ylabel('Actual')
    plt.title(f'Confusion Matrix: {title}')
    plt.show()
```

```python
results = pd.DataFrame(columns=['Model', 'Accuracy'])

def evaluate_model(name, model):
    model.fit(X_train, y_train)
    y_pred = model.predict(X_test)
    acc = accuracy_score(y_test, y_pred)
    print(f"🔍 {name}\n", classification_report(y_test, y_pred))
    plot_conf_matrix(y_test, y_pred, name)
    results.loc[len(results)] = [name, acc]
```

## Logistics regression

```python
[13]: evaluate_model("Logistic Regression", LogisticRegression(max_iter=1000))
```

```
🔍 Logistic Regression
              precision    recall  f1-score   support

           0       1.00      1.00      1.00        19
           1       1.00      0.95      0.98        21
           2       0.93      1.00      0.97        14

    accuracy                           0.98        54
   macro avg       0.98      0.98      0.98        54
weighted avg       0.98      0.98      0.98        54
```



Confusion Matrix: Logistic Regression

Confusion Matrix: Logistic Regression

## Desicion trees

```
[14]: evaluate_model("Decision Tree", DecisionTreeClassifier())
```

🔍 Decision Tree

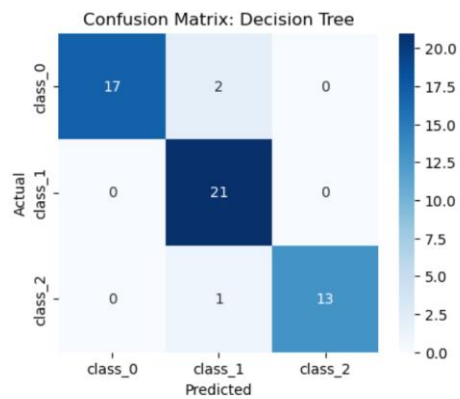|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0 | 1.00 | 0.89 | 0.94 | 19 |
| 1 | 0.88 | 1.00 | 0.93 | 21 |
| 2 | 1.00 | 0.93 | 0.96 | 14 |
|  |  |  |  |  |
| accuracy |  |  | 0.94 | 54 |
| macro avg | 0.96 | 0.94 | 0.95 | 54 |
| weighted avg | 0.95 | 0.94 | 0.94 | 54 |

### Confusion Matrix: Decision Tree

| | class_0 | class_1 | class_2 |
|---|---|---|---|
| class_0 | 17 | 2 | 0 |
| class_1 | 0 | 21 | 0 |
| class_2 | 0 | 1 | 13 |

Actual / Predicted

## randomforest

```
[15]: evaluate_model("Random Forest", RandomForestClassifier(n_estimators=100))
```

🔍 Random Forest

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0 | 1.00 | 1.00 | 1.00 | 19 |
| 1 | 1.00 | 1.00 | 1.00 | 21 |
| 2 | 1.00 | 1.00 | 1.00 | 14 |
|  |  |  |  |  |
| accuracy |  |  | 1.00 | 54 |
| macro avg | 1.00 | 1.00 | 1.00 | 54 |
| weighted avg | 1.00 | 1.00 | 1.00 | 54 |

### Confusion Matrix: Random Forest

| | class_0 | class_1 | class_2 |
|---|---|---|---|
| class_0 | 19 | 0 | 0 |
| class_1 | 0 | 21 | 0 |
| class_2 | 0 | 0 | 14 |

Actual / Predicted

## k-Nearest Neighbors (KNN) ¶

```
[16]: evaluate_model("KNN", KNeighborsClassifier(n_neighbors=5))
```

KNN

```
              precision    recall  f1-score   support

           0       0.95      1.00      0.97        19
           1       1.00      0.90      0.95        21
           2       0.93      1.00      0.97        14

    accuracy                           0.96        54
   macro avg       0.96      0.97      0.96        54
weighted avg       0.97      0.96      0.96        54
```

Confusion Matrix: KNN

## Naive Bayes

```
[17]: evaluate_model("Naive Bayes", GaussianNB())
```

Naive Bayes

```
              precision    recall  f1-score   support

           0       1.00      1.00      1.00        19
           1       1.00      1.00      1.00        21
           2       1.00      1.00      1.00        14

    accuracy                           1.00        54
   macro avg       1.00      1.00      1.00        54
weighted avg       1.00      1.00      1.00        54
```
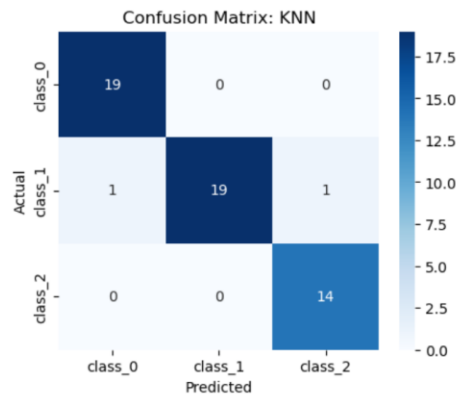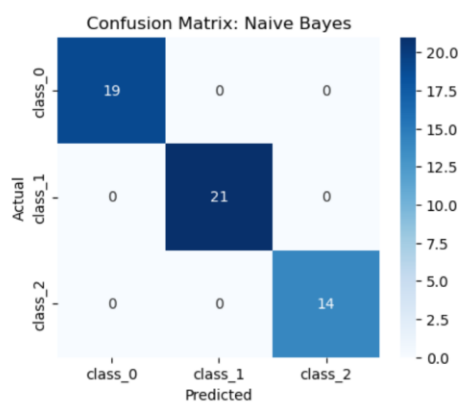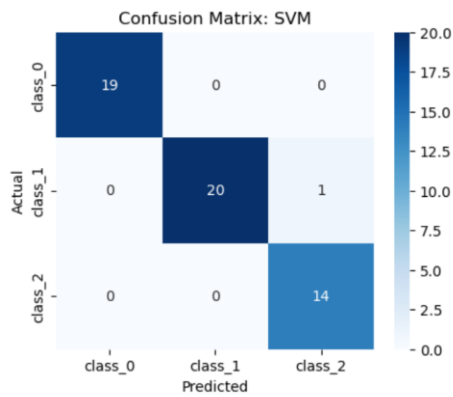
Confusion Matrix: Naive Bayes

```
[16]: evaluate_model("KNN", KNeighborsClassifier(n_neighbors=5))
```

## Support Vector Machine (SVM)

```
[8]: evaluate_model("SVM", SVC(kernel='linear'))
```

```
SVM
              precision    recall  f1-score   support

           0       1.00      1.00      1.00        19
           1       1.00      0.95      0.98        21
           2       0.93      1.00      0.97        14

    accuracy                           0.98        54
   macro avg       0.98      0.98      0.98        54
weighted avg       0.98      0.98      0.98        54
```
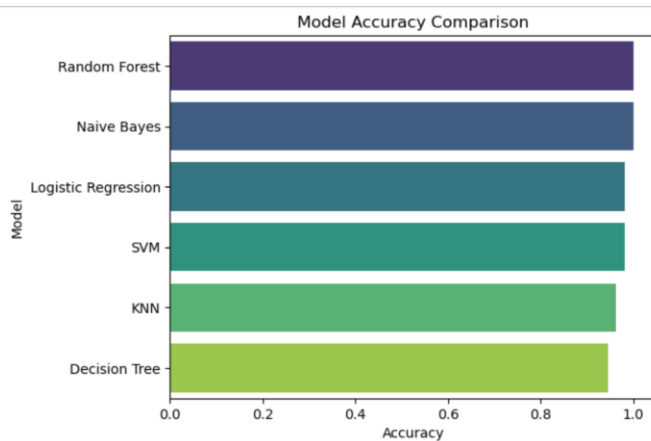
4:04 PM


Confusion Matrix: SVM

4:05 PM

## Summary of Model Performance (Accuracy):

```python
# Compare Model Performance
results.sort_values(by='Accuracy', ascending=False, inplace=True)
sns.barplot(data=results, x='Accuracy', y='Model', palette='viridis')
plt.title("Model Accuracy Comparison")
plt.show()

results
```

4:32 PM


Model Accuracy Comparison

4:32 PM

| | Model | Accuracy |
|---|---|---|
| 2 | Random Forest | 1.000000 |
| 4 | Naive Bayes | 1.000000 |
| 0 | Logistic Regression | 0.981481 |
| 5 | SVM | 0.981481 |
| 3 | KNN | 0.962963 |
| 1 | Decision Tree | 0.944444 |

## 5. Model Comparison and Insights

- After evaluating all models, Random Forest achieved the highest accuracy with strong precision and recall across all classes. It performed better due to its ensemble nature, reducing overfitting and capturing complex patterns.
- Logistic Regression and SVM also performed well.
- Naive Bayes, although simple, was surprisingly effective.
- Decision Tree and KNN had decent performance but showed more variance.

Based on these results, Random Forest is the most reliable choice for this dataset.

## Link to Code:

https://github.com/shirleensimon/classification-models/blob/main/classification_models.ipynb

## Conclusion

This assignment successfully demonstrated the process of applying multiple supervised classification algorithms to a real-world dataset. Naive Bayes and Random Forest emerged as the top performers in terms of accuracy and evaluation metrics. However, each model has its strengths and use cases.