This document provides an overview of the Python script developed to analyze and compare the performance of different machine learning models for predicting cancer types based on a given dataset. The main sections of the script include data import, manipulation, model training, and result storage.

**Data Description**

The dataset used in this analysis is the Breast Cancer Wisconsin (Diagnostic) dataset. It contains 569 instances of breast cancer cases, each described by 32 features. The target variable is diagnosis, which indicates whether the cancer is malignant (M) or benign (B).

**Columns:**

**id**: Identifier for each instance

**diagnosis**: Target variable (M for malignant, B for benign)

30 numerical features describing the characteristics of the cell nuclei present in the image

**Data Import and Manipulation**

The dataset is imported using pandas, and the necessary preprocessing steps are performed, including encoding the target variable and scaling the features.

```
1    import pandas as pd
2    from sklearn.model_selection import train_test_split
3    from sklearn.preprocessing import StandardScaler
4    from sklearn.svm import SVC
5    from sklearn.ensemble import RandomForestClassifier
6    from sklearn.neural_network import MLPClassifier
7    from sklearn.metrics import accuracy_score, classification_report
8    import matplotlib.pyplot as plt
9    import seaborn as sns
10
11   # Load the dataset
12   data = pd.read_csv('breastcancer.csv')
13
14   data = data.drop(columns=['Unnamed: 32'])
15
16   # Brief description of the dataset
17   print(data.head())
18   print(data.info())
19   print(data.describe())
20
21   # Target variable and feature columns
22   target = 'diagnosis'
23   features = data.columns.drop(['id', target])
24
25   #Data Manipulation
26   # Encode the target variable
27   data[target] = data[target].map({'M': 1, 'B': 0})
28
29   # Split the data into training and testing sets
30   X = data[features]
31   y = data[target]
32   X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
33
34   # Scale the features
35   scaler = StandardScaler()
36   X_train_scaled = scaler.fit_transform(X_train)
37   X_test_scaled = scaler.transform(X_test)
38
```

## Model Training
Three models were trained and evaluated: ***Support Vector Machines
(SVM), Random Forest,*** and ***Neural Network***.

## SVM Analysis
Support Vector Machines (SVM) are effective in <u>high-dimensional spaces</u>
and versatile with <u>different kernels</u>. We used three different kernels: *linear*,
*polynomial*, and *radial basis function (RBF)*.

```
41   #Model Building
42
43   # Train and evaluate SVM with different kernels
44   kernels = ['linear', 'poly', 'rbf']
45   svm_results = {}
46   for kernel in kernels:
47       svm = SVC(kernel=kernel)
48       svm.fit(X_train_scaled, y_train)
49       y_pred = svm.predict(X_test_scaled)
50       accuracy = accuracy_score(y_test, y_pred)
51       svm_results[kernel] = accuracy
52       print(f'SVM with {kernel} kernel accuracy: {accuracy}')
53       print(classification_report(y_test, y_pred))
54
```

## Neural Network Regression Analysis

Neural networks are capable of *capturing complex patterns* and have a flexible architecture. We used a *Multi-Layer Perceptron (MLP)* classifier for this analysis.

```
62
63    # Train and evaluate Neural Network
64    nn = MLPClassifier(hidden_layer_sizes=(100,), max_iter=1000, random_state=42)
65    nn.fit(X_train_scaled, y_train)
66    y_pred_nn = nn.predict(X_test_scaled)
67    nn_accuracy = accuracy_score(y_test, y_pred_nn)
68    print(f'Neural Network accuracy: {nn_accuracy}')
69    print(classification_report(y_test, y_pred_nn))
70
```

## Random Forest Analysis

Random Forest is robust to *overfitting* and handles *large datasets* well. It was trained and evaluated as follows:

```
54
55    # Train and evaluate Random Forest
56    rf = RandomForestClassifier(n_estimators=100, random_state=42)
57    rf.fit(X_train, y_train)
58    y_pred_rf = rf.predict(X_test)
59    rf_accuracy = accuracy_score(y_test, y_pred_rf)
60    print(f'Random Forest accuracy: {rf_accuracy}')
61    print(classification_report(y_test, y_pred_rf))
62
```

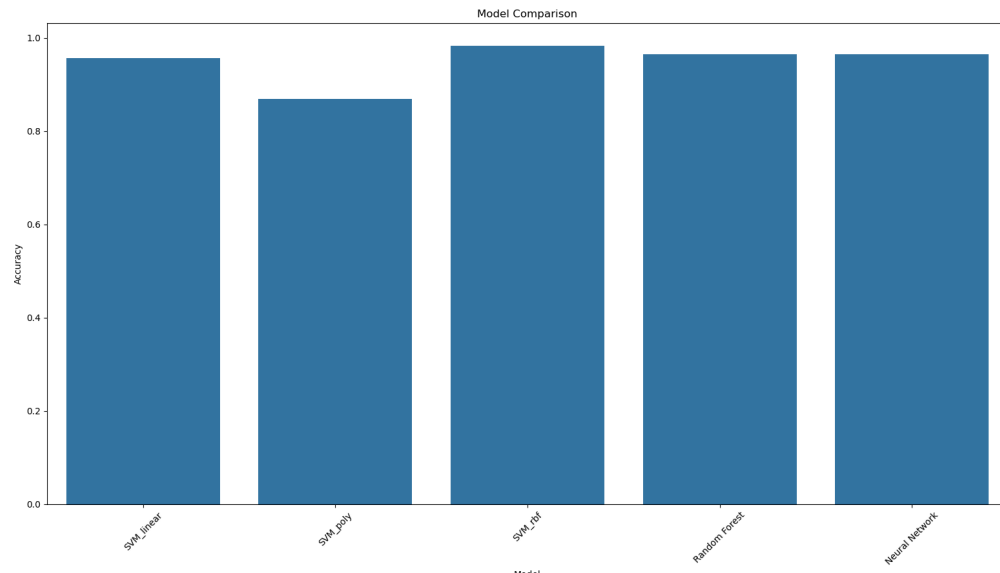## Grid Search for parameters

```
63
64    # Grid search for Neural Network
65    param_grid = {
66        'hidden_layer_sizes': [(50,), (100,), (150,)],
67        'activation': ['tanh', 'relu'],
68        'solver': ['sgd', 'adam'],
69        'alpha': [0.0001, 0.05],
70        'learning_rate': ['constant', 'adaptive'],
71    }
72
73    grid_search = GridSearchCV(MLPClassifier(max_iter=1000, random_state=42), param_grid, n_jobs=-1, cv=3)
74    grid_search.fit(X_train_scaled, y_train)
75
76    print(f'Best parameters found: {grid_search.best_params_}')
77    best_nn = grid_search.best_estimator_
78
```

## Comparison

The performance of the three models was compared based on their accuracy

```
86
87   # Store results in a dictionary
88   results = {
89       'SVM_linear': svm_results['linear'],
90       'SVM_poly': svm_results['poly'],
91       'SVM_rbf': svm_results['rbf'],
92       'Random Forest': rf_accuracy,
93       'Neural Network': nn_accuracy
94   }
95
96   # Compare the performance
97   print("Model Comparison:")
98   for model, result in results.items():
99       print(f"{model}: {result}")
100
101  # Plot the results
102  import seaborn as sns
103  import matplotlib.pyplot as plt
104
105  model_names = list(results.keys())
106  accuracies = list(results.values())
107
108  sns.barplot(x=model_names, y=accuracies)
109  plt.xlabel('Model')
110  plt.ylabel('Accuracy')
111  plt.title('Model Comparison')
112  plt.xticks(rotation=45)
113  plt.show()
114
```

## Accuracy plot:



## Accuracy of each models:

| | |
|---|---|
| SVM with Linear kernel function: | 0.956140350877193 |
| SVM with Polynomial kernel function | 0.86842105263157 |
| SVM with Radial Basis function | 0.9824561403508771 |
| Random Forest | 0.9649122807017544 |
| Neural Network | 0.9649122807017544 |

## Discussion

**Support Vector Machine (SVM)**

- **Strengths**: SVMs are effective in high-dimensional spaces, making them suitable for complex classification problems. They are versatile, with options to apply different kernel functions to separate data that is not linearly separable.
- **Weaknesses**: SVMs can be computationally expensive, particularly with large datasets, and are sensitive to parameter settings, which can impact performance if not optimized carefully.

**Random Forest**

- **Strengths**: Random Forests are robust to overfitting, especially with complex datasets, and they handle large datasets well due to their ensemble approach.
- **Weaknesses**: Random Forest models can be less interpretable than simpler models and require careful parameter tuning, especially when dealing with high-dimensional data.

**Neural Network**

- **Strengths**: Neural Networks are highly flexible and can capture complex, non-linear patterns in data, which makes them powerful for various tasks in machine learning.
- **Weaknesses**: They often require large datasets to generalize well and are computationally intensive, requiring significant resources for both training and fine-tuning.

## Real-World Applications

Accurate cancer type prediction has significant implications in the medical field. Early and accurate diagnosis can lead to better treatment plans and improved patient outcomes. Machine learning models, like the ones analyzed in this report, can assist doctors in making more informed decisions, potentially saving lives. These models can be integrated into diagnostic tools to provide real-time analysis and support to medical professionals.

References

[kaggle.com](kaggle.com)
[Paper on Benign and Malignant tumor cells](Paper on Benign and Malignant tumor cells)
**"Deep Learning"** by Ian Goodfellow, Yoshua Bengio, and Aaron Courville
**"Random Forests"** by Leo Breiman (2001, *Machine Learning*)

## Conclusion

The model with the highest accuracy is the **SVM with Radial Basis Function (RBF) kernel**, achieving an accuracy of 0.9825. This suggests that thoughtful model selection, combined with careful parameter tuning, is essential for achieving optimal performance in machine learning tasks.