**Subject: Adv. ML / ML for DS (MS-DS)**          **Instructor: Dr. M. Ishtiaq**

# Assignment 1

The data has been loaded and stored in the drive and mounting drive.

```
[ ] #mounting drive for initializing dataframe
    from google.colab import drive
    drive.mount('/content/drive')

    Mounted at /content/drive

[ ] #Loading dataset
    data = pd.read_csv('/content/drive/MyDrive/Classroom/Adv. ML   ML for DS/Assignment01/data.csv')

    #printing head of dataset
    print(data.head())

    #printing description of dataset
    print(data.describe())

    #printing info of the dataset
    data.info()
```

Head:

```
        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]
```

Description:

Info:

[5 rows x 33 columns]

|  | id | radius_mean | texture_mean | perimeter_mean | area_mean \ |
|---|---|---|---|---|---|
| count | 5.690000e+02 | 569.000000 | 569.000000 | 569.000000 | 569.000000 |
| mean | 3.037183e+07 | 14.127292 | 19.289649 | 91.969033 | 654.889104 |
| std | 1.250206e+08 | 3.524049 | 4.301036 | 24.298981 | 351.914129 |
| min | 8.670000e+03 | 6.981000 | 9.710000 | 43.790000 | 143.500000 |
| 25% | 8.692180e+05 | 11.700000 | 16.170000 | 75.170000 | 420.300000 |
| 50% | 9.060240e+05 | 13.370000 | 18.840000 | 86.240000 | 551.100000 |
| 75% | 8.813129e+06 | 15.780000 | 21.800000 | 104.100000 | 782.700000 |
| max | 9.113205e+08 | 28.110000 | 39.280000 | 188.500000 | 2501.000000 |

|  | smoothness_mean | compactness_mean | concavity_mean | concave points_mean \ |
|---|---|---|---|---|
| count | 569.000000 | 569.000000 | 569.000000 | 569.000000 |
| mean | 0.096360 | 0.104341 | 0.088799 | 0.048919 |
| std | 0.014064 | 0.052813 | 0.079720 | 0.038803 |
| min | 0.052630 | 0.019380 | 0.000000 | 0.000000 |
| 25% | 0.086370 | 0.064920 | 0.029560 | 0.020310 |
| 50% | 0.095870 | 0.092630 | 0.061540 | 0.033500 |
| 75% | 0.105300 | 0.130400 | 0.130700 | 0.074000 |
| max | 0.163400 | 0.345400 | 0.426800 | 0.201200 |

|  | symmetry_mean | ... | texture_worst | perimeter_worst | area_worst \ |
|---|---|---|---|---|---|
| count | 569.000000 | ... | 569.000000 | 569.000000 | 569.000000 |
| mean | 0.181162 | ... | 25.677223 | 107.261213 | 880.583128 |
| std | 0.027414 | ... | 6.146258 | 33.602542 | 569.356993 |
| min | 0.106000 | ... | 12.020000 | 50.410000 | 185.200000 |
| 25% | 0.161900 | ... | 21.080000 | 84.110000 | 515.300000 |
| 50% | 0.179200 | ... | 25.410000 | 97.660000 | 686.500000 |
| 75% | 0.195700 | ... | 29.720000 | 125.400000 | 1084.000000 |
| max | 0.304000 | ... | 49.540000 | 251.200000 | 4254.000000 |

|  | smoothness_worst | compactness_worst | concavity_worst \ |
|---|---|---|---|
| count | 569.000000 | 569.000000 | 569.000000 |
| mean | 0.132369 | 0.254265 | 0.272188 |
| std | 0.022832 | 0.157336 | 0.208624 |
| min | 0.071170 | 0.027290 | 0.000000 |
| 25% | 0.116600 | 0.147200 | 0.114500 |
| 50% | 0.131300 | 0.211900 | 0.226700 |
| 75% | 0.146000 | 0.339100 | 0.382900 |
| max | 0.222600 | 1.058000 | 1.252000 |

|  | concave points_worst | symmetry_worst | fractal_dimension_worst \ |
|---|---|---|---|
| count | 569.000000 | 569.000000 | 569.000000 |
| mean | 0.114606 | 0.290076 | 0.083946 |
| std | 0.065732 | 0.061867 | 0.018061 |
| min | 0.000000 | 0.156500 | 0.055040 |
| 25% | 0.064930 | 0.250400 | 0.071460 |
| 50% | 0.099930 | 0.282200 | 0.080040 |
| 75% | 0.161400 | 0.317900 | 0.092080 |
| max | 0.291000 | 0.663800 | 0.207500 |

[8 rows x 32 columns]

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 569 entries, 0 to 568
Data columns (total 33 columns):
 #   Column                   Non-Null Count  Dtype
---  ------                   --------------  -----
 0   id                       569 non-null    int64
 1   diagnosis                569 non-null    object
 2   radius_mean              569 non-null    float64
 3   texture_mean             569 non-null    float64
 4   perimeter_mean           569 non-null    float64
 5   area_mean                569 non-null    float64
 6   smoothness_mean          569 non-null    float64
 7   compactness_mean         569 non-null    float64
 8   concavity_mean           569 non-null    float64
 9   concave points_mean      569 non-null    float64
 10  symmetry_mean            569 non-null    float64
 11  fractal_dimension_mean   569 non-null    float64
 12  radius_se                569 non-null    float64
 13  texture_se               569 non-null    float64
 14  perimeter_se             569 non-null    float64
 15  area_se                  569 non-null    float64
 16  smoothness_se            569 non-null    float64
 17  compactness_se           569 non-null    float64
 18  concavity_se             569 non-null    float64
 19  concave points_se        569 non-null    float64
 20  symmetry_se              569 non-null    float64
 21  fractal_dimension_se     569 non-null    float64
 22  radius_worst             569 non-null    float64
 23  texture_worst            569 non-null    float64
 24  perimeter_worst          569 non-null    float64
 25  area_worst               569 non-null    float64
 26  smoothness_worst         569 non-null    float64
 27  compactness_worst        569 non-null    float64
 28  concavity_worst          569 non-null    float64
 29  concave points_worst     569 non-null    float64
 30  symmetry_worst           569 non-null    float64
 31  fractal_dimension_worst  569 non-null    float64
 32  Unnamed: 32              0 non-null      float64
dtypes: float64(31), int64(1), object(1)
memory usage: 146.8+ KB
```

after that below phases has been followed.

## Phase 1: Exploratory Data Analysis (EDA) & Preprocessing

Exploratory Data Analysis (EDA) helps understand the dataset, detect patterns, and identify potential issues before model training. Below is a step-by-step guide to performing EDA on the selected.

➢ **Check missing values and handle them** (Found none, except an extra column, which was dropped)

```
#Checking null values
print(data.isnull().sum)
```

| <bound method DataFrame.sum of | | | id | diagnosis | radius_mean | texture_mean | perimeter_mean | area_mean \ |
|---|---|---|---|---|---|---|---|---|
| 0 | False | False | False | False | False | False | False |
| 1 | False | False | False | False | False | False | False |
| 2 | False | False | False | False | False | False | False |
| 3 | False | False | False | False | False | False | False |
| 4 | False | False | False | False | False | False | False |
| .. | ... | ... | ... | ... | ... | ... | ... |
| 564 | False | False | False | False | False | False | False |
| 565 | False | False | False | False | False | False | False |
| 566 | False | False | False | False | False | False | False |
| 567 | False | False | False | False | False | False | False |
| 568 | False | False | False | False | False | False | False |

|  | smoothness_mean | compactness_mean | concavity_mean | concave points_mean \ |
|---|---|---|---|---|
| 0 | False | False | False | False |
| 1 | False | False | False | False |
| 2 | False | False | False | False |
| 3 | False | False | False | False |
| 4 | False | False | False | False |
| .. | ... | ... | ... | ... |
| 564 | False | False | False | False |
| 565 | False | False | False | False |
| 566 | False | False | False | False |
| 567 | False | False | False | False |
| 568 | False | False | False | False |

|  | ... | radius_worst | texture_worst | perimeter_worst | area_worst \ |
|---|---|---|---|---|---|
| 0 | ... | False | False | False | False |
| 1 | ... | False | False | False | False |
| 2 | ... | False | False | False | False |
| 3 | ... | False | False | False | False |

## ➢ Visualize class distribution



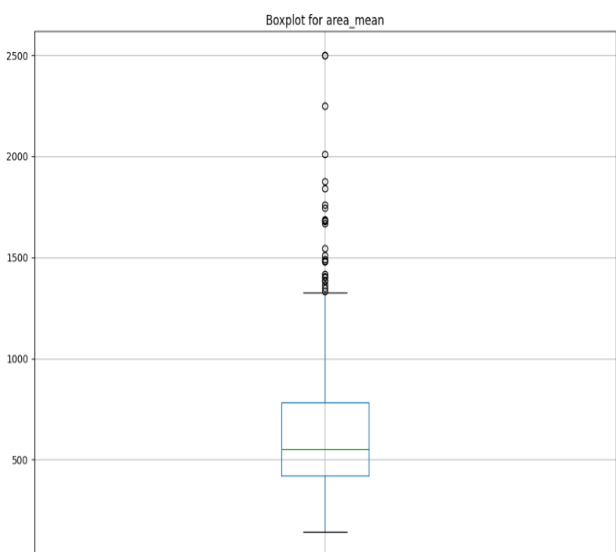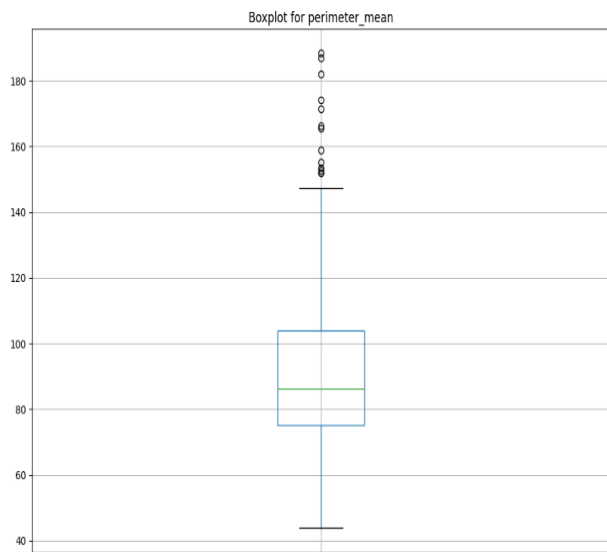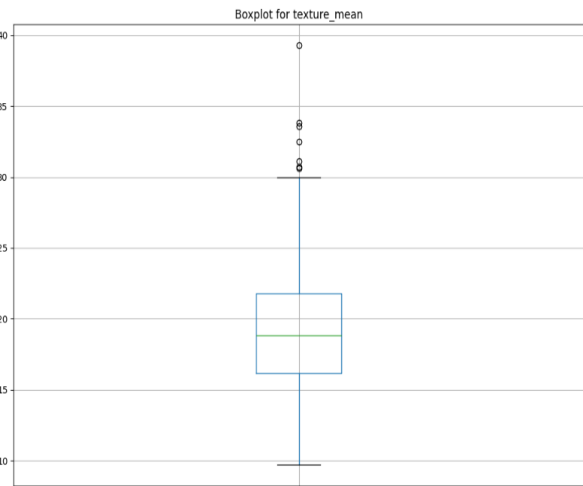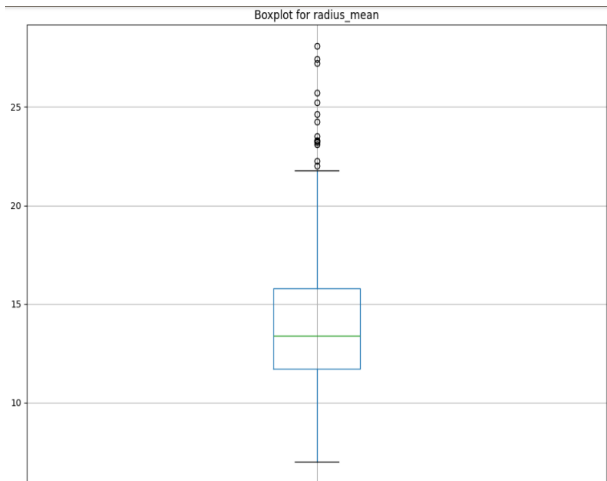## ➢ Perform correlation analysis (Correlation has been plotted, while fixing mix-up of values)

➢ **Identify outliers** (Outliers been identified via box-plot, some of the plots are pasted below, rest could be checked from the notebook submitted)



```
✓  Identify outliers

    #Identifying outliers
    numerical_features = data.select_dtypes(include=np.number).drop(columns=['id'])

    #Creating boxplots for each feature individually
    for feature in numerical_features.columns:
        plt.figure(figsize=(10, 8))
        data.boxplot(column=feature)
        plt.title(f'Boxplot for {feature}')
        plt.xticks(rotation=45, ha='right')
        plt.tight_layout()
        plt.show()
```
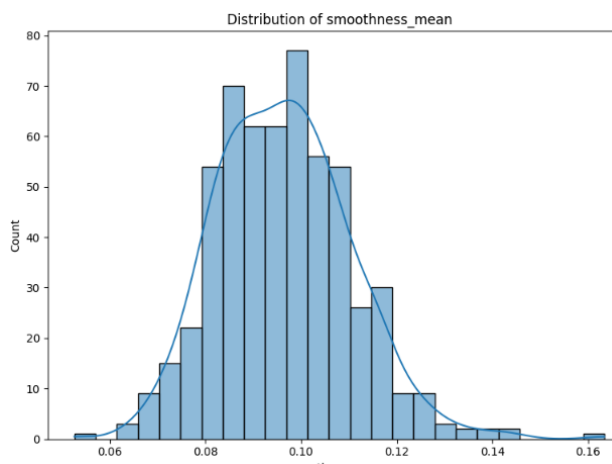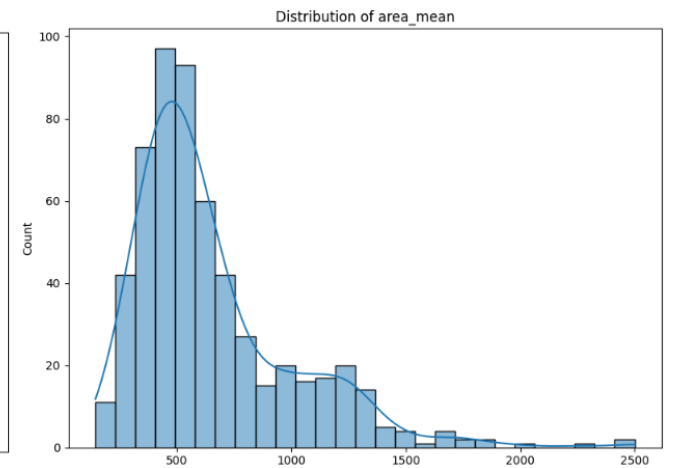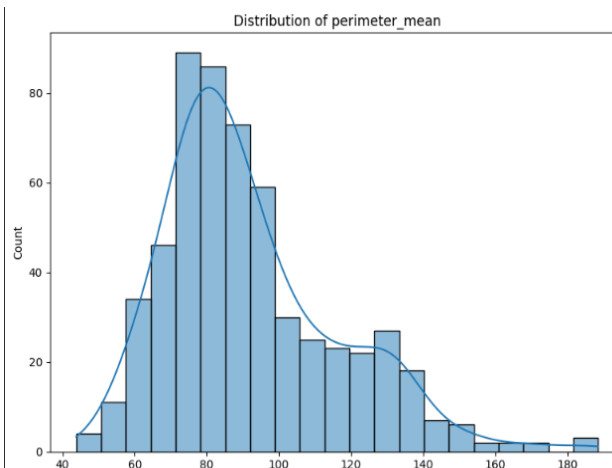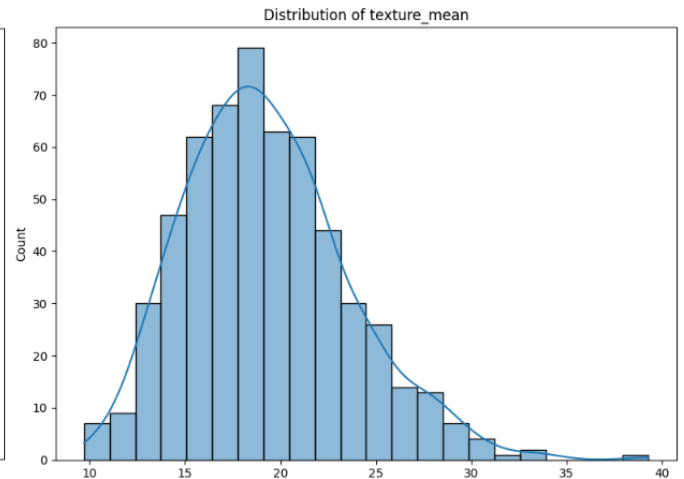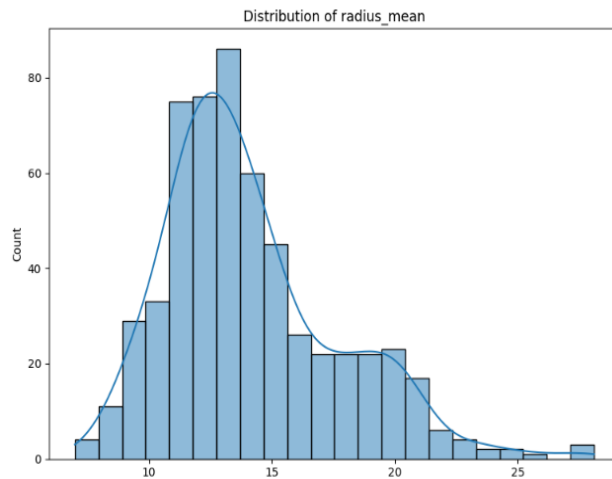


Boxplot for radius_mean



Boxplot for texture_mean



Boxplot for perimeter_mean



Boxplot for area_mean



Boxplot for smoothness_mean

➢ **Visualize feature distributions** (Feature distribution has been visualised; some of the distributions are here under; the rest could be analysed from the notebook)



```
∨  Visualize feature distributions

[ ]  #Visualizing feature distributions
     for feature in numerical_features.columns:
         plt.figure(figsize=(8, 6))
         sns.histplot(data[feature], kde=True)
         plt.title(f'Distribution of {feature}')
         plt.tight_layout();
```

> **Normalize the dataset** (The dataset has been normalized by via standard scaler using fit transform)

```
✓  Normalize the dataset

[ ]  #selectign numerical columns for scalling
     numerical_features = data.select_dtypes(include=np.number).drop(columns=['id'])

     #Normalizing the dataset
     scaler = StandardScaler()
     data[numerical_features.columns] = scaler.fit_transform(data[numerical_features.columns])
```

> **Split data into training & testing sets** (The data has been split into training and testing sets)

```
✓  Split data into training & testing sets

[ ]  #Splitting data into training and testing sets
     X = data.drop(columns=['id', 'diagnosis', 'diagnosis_encoded'])
     y = data['diagnosis_encoded']

[ ]  # Normalize features
     scaler = StandardScaler()
     X_scaled = scaler.fit_transform(X)

[ ]  # Split data into training and test sets
     X_train, X_test, y_train, y_test = train_test_split(X_scaled, y, test_size=0.2, random_state=42)
```

**Phase 2: PCA for dimensionality reduction**
> **Analyze the impact of reducing the number of features on classification performance.**
> **Compute the explained variance ratio for different numbers of principal components.**
> **Select the number of components that retain 95% of variance.**

The PCA has been applied for dimensionality reduction, the variance ratio has been depicted, and components containing 95pc of variance have been selected.

```
✓  Phase 2: PCA for dimensionality reduction

Analyze the impact of reducing the number of features on classification performance.

Compute the explained variance ratio for different number of principal components.

✓  Select the number of components that retain 95% of variance.

[ ]  #Analyzing impact of reducing no. of features
     from sklearn.decomposition import PCA

     pca = PCA(n_components=0.95)  # Retain 95% of the variance
     X_train_pca = pca.fit_transform(X_train)
     X_test_pca = pca.transform(X_test)

     # Explained variance ratio
     print("Explained variance ratio:", pca.explained_variance_ratio_)
     print("Number of components selected:", pca.n_components_)

⊡▸  Explained variance ratio: [0.4325482  0.19732988 0.09865421 0.0625777  0.053537   0.04102624
      0.02274171 0.01636683 0.01404727 0.01198947]
     Number of components selected: 10
```
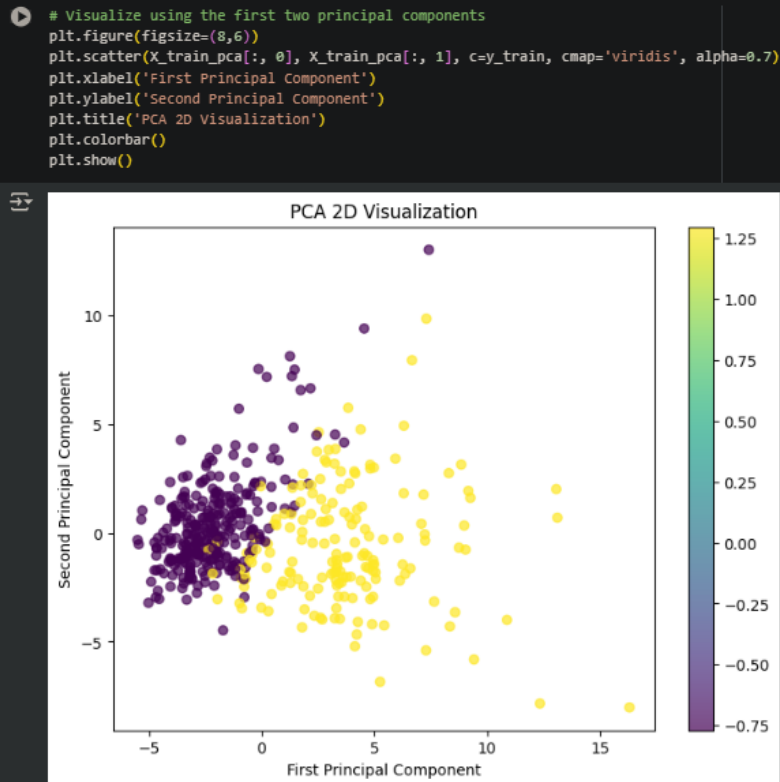
➢ **Visualize the data in 2D using the first two principal components.** (The visualization of the data has been done using first two principal components)



**Comparison:**

The comparison of the models has been carried out, and the results are as under:

➢ **SVM**

```
[ ] y_pred = svm_model.predict(X_test)
    print(classification_report(y_test, y_pred))
    print(confusion_matrix(y_test, y_pred))
```

```
              precision    recall  f1-score   support

           0       0.96      0.97      0.97        71
           1       0.95      0.93      0.94        43

    accuracy                           0.96       114
   macro avg       0.96      0.95      0.95       114
weighted avg       0.96      0.96      0.96       114

[[69  2]
 [ 3 40]]
```

## ➢ **Decision Tree**

```python
from sklearn.tree import DecisionTreeClassifier, export_graphviz
import graphviz

# Without pruning
dt_model = DecisionTreeClassifier(criterion='gini', random_state=42)
dt_model.fit(X_train, y_train)
print("Decision Tree (unpruned) accuracy:", dt_model.score(X_test, y_test))

# With pruning parameters (pre-pruning)
dt_pruned = DecisionTreeClassifier(criterion='gini', max_depth=5, min_samples_split=5, min_samples_leaf=2, random_state=42)
dt_pruned.fit(X_train, y_train)
print("Decision Tree (pruned) accuracy:", dt_pruned.score(X_test, y_test))

# Visualize the pruned tree
dot_data = export_graphviz(dt_pruned, out_file=None,
                           feature_names=X.columns,
                           class_names=True,
                           filled=True, rounded=True,
                           special_characters=True)
graph = graphviz.Source(dot_data)
graph.render("decision_tree")
```

```
Decision Tree (unpruned) accuracy: 0.9473684210526315
Decision Tree (pruned) accuracy: 0.956140350877193
'decision_tree.pdf'
```

## Evaluation

```
y_pred_dt = dt_pruned.predict(X_test)
print(classification_report(y_test, y_pred_dt))
print(confusion_matrix(y_test, y_pred_dt))
```

```
              precision    recall  f1-score   support

           0       0.96      0.97      0.97        71
           1       0.95      0.93      0.94        43

    accuracy                           0.96       114
   macro avg       0.96      0.95      0.95       114
weighted avg       0.96      0.96      0.96       114

[[69  2]
 [ 3 40]]
```

➢ **Neural Network**

## Neural Network

```python
import tensorflow as tf
from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import Dense, Dropout
from tensorflow.keras.optimizers import Adam

# Define the model architecture (using original data)
model = Sequential([
    Dense(64, activation='relu', input_shape=(X_train.shape[1],)),
    Dropout(0.2),
    Dense(32, activation='relu'),
    Dense(1, activation='sigmoid')  # For binary classification
])

model.compile(optimizer=Adam(learning_rate=0.001), loss='binary_crossentropy', metrics=['accuracy'])
history = model.fit(X_train, y_train, epochs=100, batch_size=32, validation_split=0.2, verbose=0)
loss, accuracy = model.evaluate(X_test, y_test)
print("Neural Network (original) accuracy:", accuracy)

# Using PCA-transformed data
model_pca = Sequential([
    Dense(64, activation='relu', input_shape=(X_train_pca.shape[1],)),
    Dropout(0.2),
    Dense(32, activation='relu'),
    Dense(1, activation='sigmoid')
])
model_pca.compile(optimizer=Adam(learning_rate=0.001), loss='binary_crossentropy', metrics=['accuracy'])
history_pca = model_pca.fit(X_train_pca, y_train, epochs=100, batch_size=32, validation_split=0.2, verbose=0)
loss_pca, accuracy_pca = model_pca.evaluate(X_test_pca, y_test)
print("Neural Network (PCA-transformed) accuracy:", accuracy_pca)
```

```
/usr/local/lib/python3.11/dist-packages/keras/src/layers/core/dense.py:87: UserWarning: Do not pass an `input_shape`/`input_dim` argument to a lay
  super().__init__(activity_regularizer=activity_regularizer, **kwargs)
4/4 ━━━━━━━━━━━━━━━━━━━━ 0s 15ms/step - accuracy: 0.9651 - loss: 0.1999
Neural Network (original) accuracy: 0.9649122953414917
4/4 ━━━━━━━━━━━━━━━━━━━━ 0s 9ms/step - accuracy: 0.9558 - loss: 0.1994
Neural Network (PCA-transformed) accuracy: 0.9649122953414917
```

## Phase 3: Model Comparison & Final Report
➢ **Model Comparison**

## Model Comparison

```python
# Model Comparison
models = {
    "SVM": svm_model,
    "SVM (PCA)": svm_pca,
    "Decision Tree": dt_pruned,
    "Neural Network": model,
    "Neural Network (PCA)": model_pca
}

results = []
for name, model in models.items():
    if name.endswith("(PCA)"):
        y_pred = model.predict(X_test_pca)
    else:
        y_pred = model.predict(X_test)

    if name in ["Neural Network", "Neural Network (PCA)"]:
        y_pred_binary = (y_pred > 0.5).astype(int)  # Convert probabilities to binary predictions
    else:
        y_pred_binary = y_pred

    report = classification_report(y_test, y_pred_binary, output_dict=True)
    results.append({
        "Model": name,
        "Accuracy": report['accuracy'],
        "Precision": report['weighted avg']['precision'],
        "Recall": report['weighted avg']['recall'],
        "F1-score": report['weighted avg']['f1-score'],
        "ROC-AUC": roc_auc_score(y_test, y_pred) if name in ["SVM", "SVM (PCA)", "Neural Network", "Neural Network (PCA)"] else roc_auc_score(y_test, y_pred_binary) # Calculate ROC-AUC if applicable
    })

results_df = pd.DataFrame(results)
print(results_df)
```

The SVM and Neural Network has been the top performers among all.

```
4/4 ──────────── 0s 26ms/step
4/4 ──────────── 0s 29ms/step
                  Model  Accuracy  Precision    Recall  F1-score   ROC-AUC
0                   SVM  0.956140   0.956088  0.956140  0.956036  0.951032
1             SVM (PCA)  0.982456   0.982456  0.982456  0.982456  0.981330
2         Decision Tree  0.956140   0.956088  0.956140  0.956036  0.951032
3        Neural Network  0.964912   0.964912  0.964912  0.964912  0.992794
4  Neural Network (PCA)  0.964912   0.965828  0.964912  0.965060  0.992466
```

➢ **Confusion Metrices**

```python
# Confusion Matrices
fig, axes = plt.subplots(1, len(models), figsize=(20, 5))
fig.suptitle("Confusion Matrices")

for i, (name, model) in enumerate(models.items()):
    if name.endswith("(PCA)"):
        y_pred = model.predict(X_test_pca)
    else:
        y_pred = model.predict(X_test)

    if name in ["Neural Network", "Neural Network (PCA)"]:
        y_pred_binary = (y_pred > 0.5).astype(int)  # Convert probabilities to binary predictions
    else:
        y_pred_binary = y_pred

    cm = confusion_matrix(y_test, y_pred_binary)
    sns.heatmap(cm, annot=True, fmt="d", cmap="Blues", ax=axes[i])
    axes[i].set_title(name)
    axes[i].set_xlabel("Predicted Label")
    axes[i].set_ylabel("True Label")

plt.tight_layout()
plt.show()
```
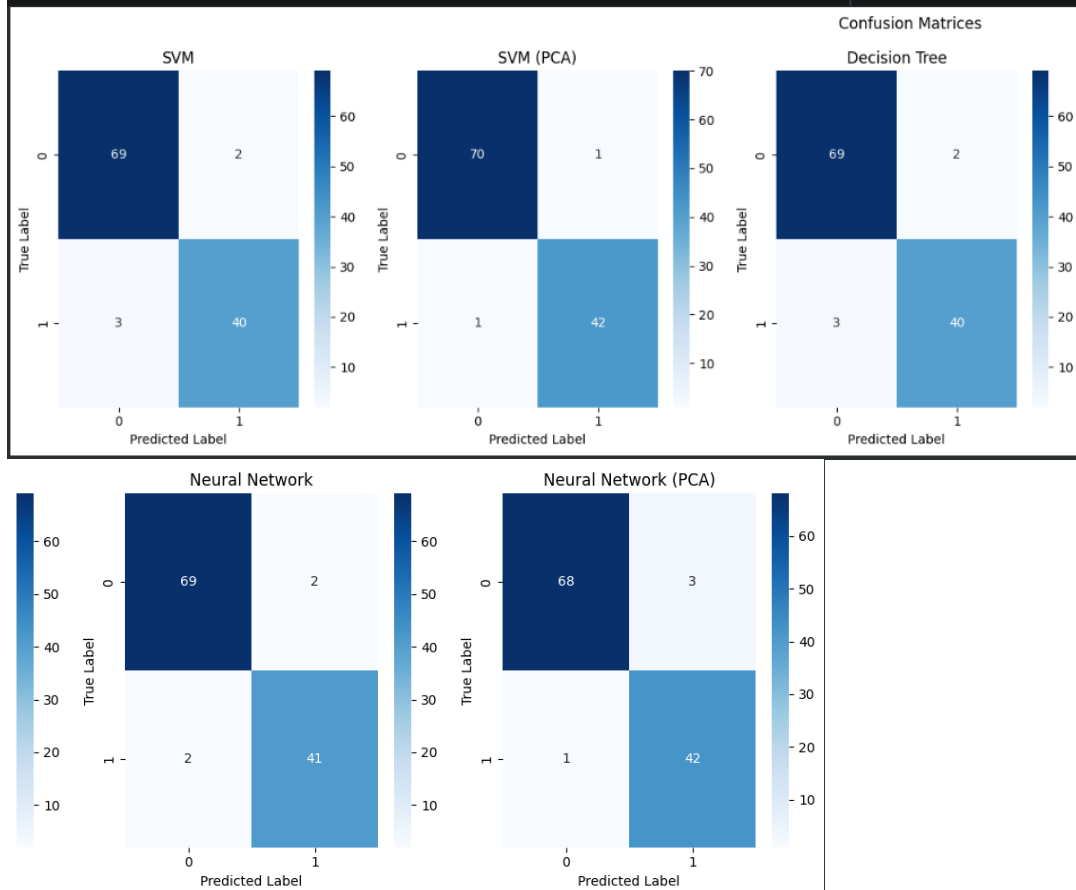
Confusion Matrices

SVM: True Label 0 / Predicted 0 = 69, Predicted 1 = 2; True Label 1 / Predicted 0 = 3, Predicted 1 = 40

SVM (PCA): True Label 0 / Predicted 0 = 70, Predicted 1 = 1; True Label 1 / Predicted 0 = 1, Predicted 1 = 42

Decision Tree: True Label 0 / Predicted 0 = 69, Predicted 1 = 2; True Label 1 / Predicted 0 = 3, Predicted 1 = 40

Neural Network: True Label 0 / Predicted 0 = 69, Predicted 1 = 2; True Label 1 / Predicted 0 = 2, Predicted 1 = 41

Neural Network (PCA): True Label 0 / Predicted 0 = 68, Predicted 1 = 3; True Label 1 / Predicted 0 = 1, Predicted 1 = 42

➢ **PCA Variance Plot**

After analyzing the plot it can be clearly seen that first few principal components explain significant portion of variance suggesting that PCA is likely to be effective without losing much information.

```
# PCA Variance Plot
plt.figure(figsize=(8, 6))
plt.plot(np.cumsum(pca.explained_variance_ratio_))
plt.xlabel("Number of Components")
plt.ylabel("Cumulative Explained Variance")
plt.title("PCA Variance Explained")
plt.grid(True)
plt.show()
```



> ## ROC Curves

## Roc Curves

```
# ROC Curves
fig, ax = plt.subplots(figsize=(8, 6))
for name, model in models.items():
    if name in ["SVM", "SVM (PCA)", "Neural Network", "Neural Network (PCA)"]:  # Models with predict_proba
        if name.endswith("(PCA)"):
            y_score = model.predict(X_test_pca)
            # For Keras models, predict already returns probabilities for binary classification
        else:
            y_score = model.predict(X_test)
            # For Keras models, predict already returns probabilities for binary classification

        # For binary classification, take the probability of class 1
        y_score = y_score[:, 0] if y_score.ndim > 1 else y_score

        fpr, tpr, _ = roc_curve(y_test, y_score)
        roc_auc = auc(fpr, tpr)
        ax.plot(fpr, tpr, label=f"{name} (AUC = {roc_auc:.2f})")

ax.plot([0, 1], [0, 1], "k--", label="Random")
ax.set_xlabel("False Positive Rate")
ax.set_ylabel("True Positive Rate")
ax.set_title("ROC Curves")
ax.legend(loc="lower right")
plt.show()
```
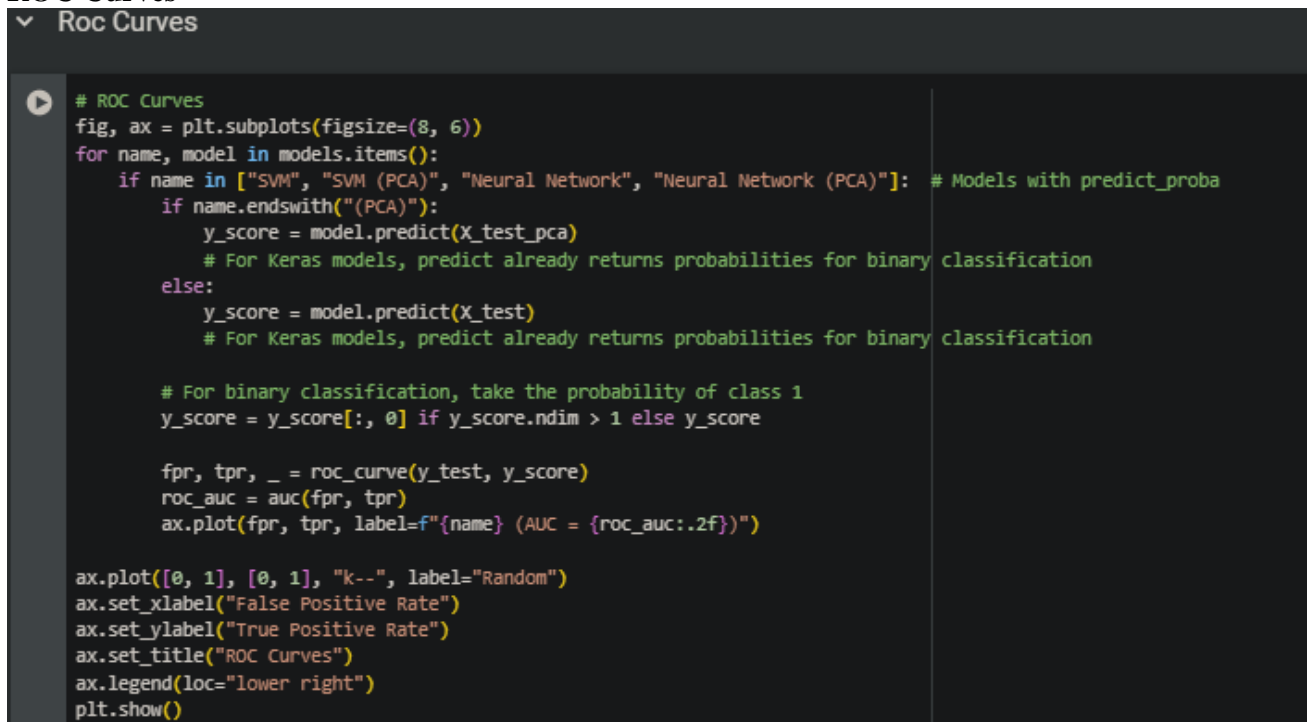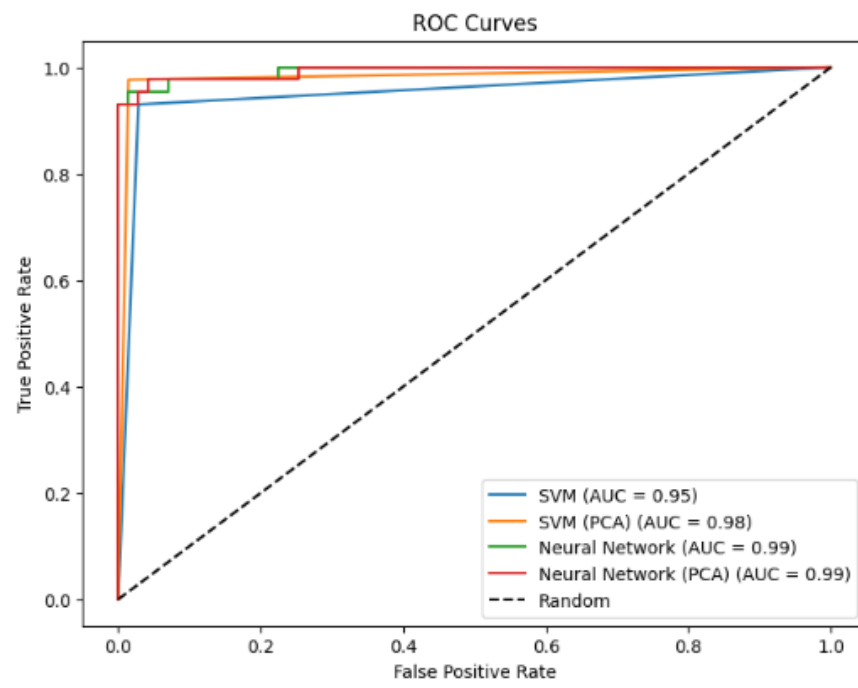
The ROC curves has been plotted with their corresponding Area Under the Curve (AUC) values, The SVM and Neural Network models appear to be the top performers, which has already been conveyed while model comparisons.

ROC Curves

# Conclusion:

**1. Which model performed best?**
The SVM and Neural Network models (both with and without PCA) achieved the best results in terms of accuracy, precision, recall, and F1-score. SVM with PCA showed similar performance to SVM without PCA, though its scores were slightly lower. The Neural Network with PCA did not significantly impact performance and performed comparably to the original model without PCA.

**2. How did PCA impact classification performance?**
In this case, PCA had minimal effect on classification performance and did not lead to significant improvements. While PCA can enhance model performance in certain situations, it seems that the original features contained crucial information that was not compromised during dimensionality reduction with PCA.

**3. Which approach is most suitable for medical diagnosis?**
Although SVM and Neural Networks (both with and without PCA) performed similarly, Neural Networks might be slightly more suitable for this specific case due to their higher accuracy and other metrics, especially when compared to SVM without PCA. However, all the models performed comparably, allowing us flexibility in our choice.