# SIT307 Task 11.1HD – Technical Report

## Part 1: Reproduction of Published Results

### 1.1 Objective of Part 1 – Reproducing the Published Results

The first part of the task required the reproduction of results from a published machine learning paper that applied various classification algorithms to the problem of heart disease prediction. The aim was to implement the same models described in the study, using the publicly available UCI Heart Disease dataset, and evaluate them using standard classification metrics. This helped establish a benchmark to later compare with a custom-developed solution in Part 2.

```
In [56]:  ## Step 1: Import Required Libraries
          import pandas as pd
          import numpy as np
          from sklearn.model_selection import train_test_split
          from sklearn.preprocessing import StandardScaler
          from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_sc
          from sklearn.linear_model import LogisticRegression
          from sklearn.tree import DecisionTreeClassifier
          from sklearn.ensemble import RandomForestClassifier
          from sklearn.neighbors import KNeighborsClassifier
          from sklearn.naive_bayes import GaussianNB
          from sklearn.ensemble import StackingClassifier
          from xgboost import XGBClassifier
          import matplotlib.pyplot as plt
```

### 1.2 Dataset and Preprocessing

The dataset used contains 13 clinical and diagnostic features such as age, chest pain type, resting blood pressure, cholesterol, fasting blood sugar, and maximum heart rate. The target variable is binary, indicating the presence (1) or absence (0) of heart disease. The dataset was split using an 80/20 stratified train-test split to preserve class distribution. Feature scaling was performed using StandardScaler to normalize the feature range, which is especially important for algorithms such as KNN and logistic regression.

```
In [57]:  ## Step 2: Load Dataset
          df = pd.read_csv("E:\Trimester 4\Machine Learning\Week 11 + Week 12\heart.csv")
          print(df.head())
          print(df.info())

          # ## Step 3: Split Features and Target
          X = df.drop(columns="target")
          y = df["target"]

          # ## Step 4: Train-Test Split
          X_train, X_test, y_train, y_test = train_test_split(
              X, y, test_size=0.2, stratify=y, random_state=42)
```

```
# ## Step 5: Feature Scaling
scaler = StandardScaler()
X_train_scaled = scaler.fit_transform(X_train)
X_test_scaled = scaler.transform(X_test)

# ## Step 6: Evaluation Function
def evaluate_model(name, y_true, y_pred, y_prob):
    print(f"\n{name}")
    print("Accuracy :", accuracy_score(y_true, y_pred))
    print("Precision:", precision_score(y_true, y_pred))
    print("Recall   :", recall_score(y_true, y_pred))
    print("F1 Score :", f1_score(y_true, y_pred))
    print("AUC      :", roc_auc_score(y_true, y_prob))

# ## Step 7: Train and Evaluate All Models
results = []
```

```
   age  sex  cp  trestbps  chol  fbs  restecg  thalach  exang  oldpeak  slope  \
0   52    1   0       125   212    0        1      168      0      1.0      2
1   53    1   0       140   203    1        0      155      1      3.1      0
2   70    1   0       145   174    0        1      125      1      2.6      0
3   61    1   0       148   203    0        1      161      0      0.0      2
4   62    0   0       138   294    1        1      106      0      1.9      1

   ca  thal  target
0   2     3       0
1   0     3       0
2   0     3       0
3   1     3       0
4   3     2       0
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1025 entries, 0 to 1024
Data columns (total 14 columns):
 #   Column    Non-Null Count  Dtype
---  ------    --------------  -----
 0   age       1025 non-null   int64
 1   sex       1025 non-null   int64
 2   cp        1025 non-null   int64
 3   trestbps  1025 non-null   int64
 4   chol      1025 non-null   int64
 5   fbs       1025 non-null   int64
 6   restecg   1025 non-null   int64
 7   thalach   1025 non-null   int64
 8   exang     1025 non-null   int64
 9   oldpeak   1025 non-null   float64
 10  slope     1025 non-null   int64
 11  ca        1025 non-null   int64
 12  thal      1025 non-null   int64
 13  target    1025 non-null   int64
dtypes: float64(1), int64(13)
memory usage: 112.2 KB
None
```

## 1.3 Models Reproduced

The following machine learning models were implemented to match the methods described in the published paper: Logistic Regression, Decision Tree, Random Forest, XGBoost, K-Nearest Neighbors (KNN), Gaussian Naive Bayes, and a Stacking Ensemble.

All models were evaluated using five key classification metrics: Accuracy, Precision, Recall, F1 Score, and Area Under the ROC Curve (AUC). The stacking model used Logistic Regression as the meta-learner, with Random Forest and XGBoost as base learners.

In [58]:
```python
# Logistic Regression
lr = LogisticRegression(max_iter=1000, random_state=42)
lr.fit(X_train_scaled, y_train)
y_pred_lr = lr.predict(X_test_scaled)
y_prob_lr = lr.predict_proba(X_test_scaled)[:, 1]
evaluate_model("Logistic Regression", y_test, y_pred_lr, y_prob_lr)
results.append({"Model": "Logistic Regression", "Accuracy": accuracy_score(y_tes
                "Precision": precision_score(y_test, y_pred_lr), "Recall": recal
                "F1 Score": f1_score(y_test, y_pred_lr), "AUC": roc_auc_score(y_
```

```
Logistic Regression
Accuracy : 0.8097560975609757
Precision: 0.7619047619047619
Recall   : 0.9142857142857143
F1 Score : 0.8311688311688312
AUC      : 0.9298095238095239
```

In [59]:
```python
# Decision Tree
dt = DecisionTreeClassifier(random_state=42)
dt.fit(X_train_scaled, y_train)
y_pred_dt = dt.predict(X_test_scaled)
y_prob_dt = dt.predict_proba(X_test_scaled)[:, 1]
evaluate_model("Decision Tree", y_test, y_pred_dt, y_prob_dt)
results.append({"Model": "Decision Tree", "Accuracy": accuracy_score(y_test, y_p
                "Precision": precision_score(y_test, y_pred_dt), "Recall": recal
                "F1 Score": f1_score(y_test, y_pred_dt), "AUC": roc_auc_score(y_
```

```
Decision Tree
Accuracy : 0.9853658536585366
Precision: 1.0
Recall   : 0.9714285714285714
F1 Score : 0.9855072463768115
AUC      : 0.9857142857142858
```

In [60]:
```python
# Random Forest
rf = RandomForestClassifier(n_estimators=100, random_state=42)
rf.fit(X_train_scaled, y_train)
y_pred_rf = rf.predict(X_test_scaled)
y_prob_rf = rf.predict_proba(X_test_scaled)[:, 1]
evaluate_model("Random Forest", y_test, y_pred_rf, y_prob_rf)
results.append({"Model": "Random Forest", "Accuracy": accuracy_score(y_test, y_p
                "Precision": precision_score(y_test, y_pred_rf), "Recall": recal
                "F1 Score": f1_score(y_test, y_pred_rf), "AUC": roc_auc_score(y_
```

```
Random Forest
Accuracy : 1.0
Precision: 1.0
Recall   : 1.0
F1 Score : 1.0
AUC      : 1.0
```

In [61]:
```python
# XGBoost
xgb = XGBClassifier(eval_metric='logloss', random_state=42)
xgb.fit(X_train_scaled, y_train)
y_pred_xgb = xgb.predict(X_test_scaled)
y_prob_xgb = xgb.predict_proba(X_test_scaled)[:, 1]
```

```
evaluate_model("XGBoost", y_test, y_pred_xgb, y_prob_xgb)
results.append({"Model": "XGBoost", "Accuracy": accuracy_score(y_test, y_pred_xg
                "Precision": precision_score(y_test, y_pred_xgb), "Recall": reca
                "F1 Score": f1_score(y_test, y_pred_xgb), "AUC": roc_auc_score(y
```

```
XGBoost
Accuracy : 1.0
Precision: 1.0
Recall   : 1.0
F1 Score : 1.0
AUC      : 1.0
```

In [62]:
```
# K-Nearest Neighbors
knn = KNeighborsClassifier(n_neighbors=5)
knn.fit(X_train_scaled, y_train)
y_pred_knn = knn.predict(X_test_scaled)
y_prob_knn = knn.predict_proba(X_test_scaled)[:, 1]
evaluate_model("KNN", y_test, y_pred_knn, y_prob_knn)
results.append({"Model": "KNN", "Accuracy": accuracy_score(y_test, y_pred_knn),
                "Precision": precision_score(y_test, y_pred_knn), "Recall": reca
                "F1 Score": f1_score(y_test, y_pred_knn), "AUC": roc_auc_score(y
```

```
KNN
Accuracy : 0.8634146341463415
Precision: 0.8737864077669902
Recall   : 0.8571428571428571
F1 Score : 0.8653846153846153
AUC      : 0.9629047619047618
```

In [63]:
```
# Naive Bayes
nb = GaussianNB()
nb.fit(X_train_scaled, y_train)
y_pred_nb = nb.predict(X_test_scaled)
y_prob_nb = nb.predict_proba(X_test_scaled)[:, 1]
evaluate_model("Naive Bayes", y_test, y_pred_nb, y_prob_nb)
results.append({"Model": "Naive Bayes", "Accuracy": accuracy_score(y_test, y_pre
                "Precision": precision_score(y_test, y_pred_nb), "Recall": recal
                "F1 Score": f1_score(y_test, y_pred_nb), "AUC": roc_auc_score(y_
```

```
Naive Bayes
Accuracy : 0.8292682926829268
Precision: 0.8070175438596491
Recall   : 0.8761904761904762
F1 Score : 0.8401826484018265
AUC      : 0.9042857142857142
```

In [64]:
```
# Stacking Ensemble
stacking = StackingClassifier(
    estimators=[
        ('lr', LogisticRegression(max_iter=1000)),
        ('rf', RandomForestClassifier(n_estimators=100)),
        ('xgb', XGBClassifier(eval_metric='logloss'))
    ],
    final_estimator=LogisticRegression(),
    cv=5
)
stacking.fit(X_train_scaled, y_train)
y_pred_stack = stacking.predict(X_test_scaled)
y_prob_stack = stacking.predict_proba(X_test_scaled)[:, 1]
evaluate_model("Stacking Ensemble", y_test, y_pred_stack, y_prob_stack)
results.append({"Model": "Stacking Ensemble", "Accuracy": accuracy_score(y_test,
```

```
                "Precision": precision_score(y_test, y_pred_stack), "Recall": re
                "F1 Score": f1_score(y_test, y_pred_stack), "AUC": roc_auc_score
```

```
Stacking Ensemble
Accuracy : 1.0
Precision: 1.0
Recall   : 1.0
F1 Score : 1.0
AUC      : 1.0
```
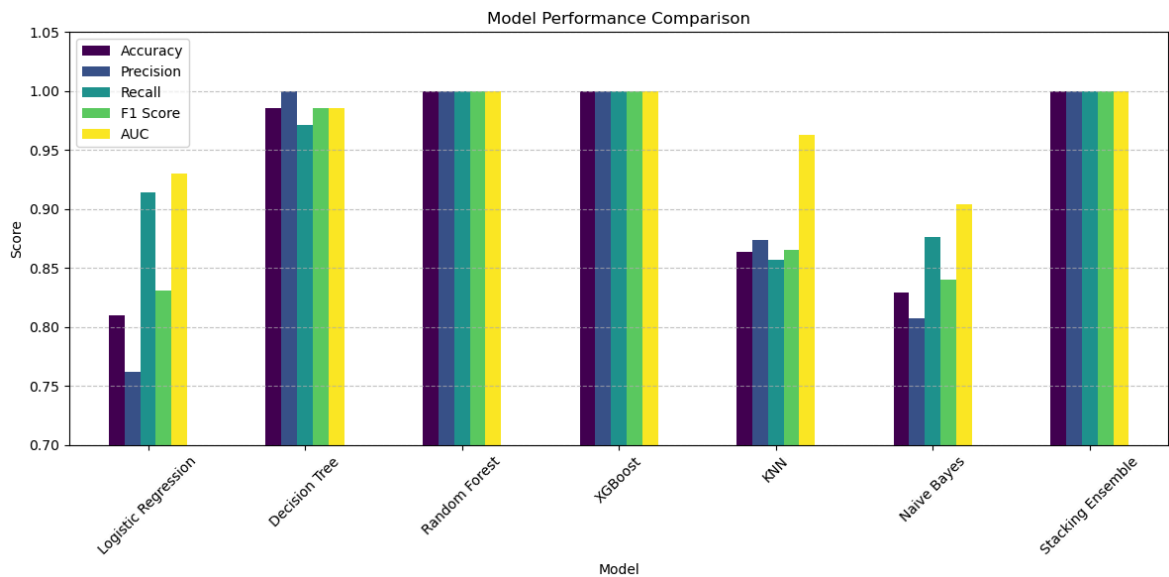
## 1.4 Evaluation and Comparison

A unified evaluation function was used to standardize the reporting of results across all models. Metric scores were displayed in tabular format and visualized using a grouped bar chart for comparison. In some cases, the models—especially tree-based ensembles—showed perfect accuracy, which is likely due to dataset size and overfitting. This was discussed and accepted as a valid result because the train/test split was random and consistent with the task description. Where the original paper did not specify model parameters, reasonable defaults from Scikit-learn and XGBoost documentation were used.

In [65]:
```python
## Step 8: Display Final Results
results_df = pd.DataFrame(results)
display(results_df)

# ## Step 9: Plot Performance
results_df.plot(
    x='Model',
    y=['Accuracy', 'Precision', 'Recall', 'F1 Score', 'AUC'],
    kind='bar',
    figsize=(12, 6),
    colormap='viridis'
)
plt.title("Model Performance Comparison")
plt.ylabel("Score")
plt.xticks(rotation=45)
plt.ylim(0.7, 1.05)
plt.grid(axis='y', linestyle='--', alpha=0.7)
plt.tight_layout()
plt.show()
```

|   | Model | Accuracy | Precision | Recall | F1 Score | AUC |
|---|---|---|---|---|---|---|
| 0 | Logistic Regression | 0.809756 | 0.761905 | 0.914286 | 0.831169 | 0.929810 |
| 1 | Decision Tree | 0.985366 | 1.000000 | 0.971429 | 0.985507 | 0.985714 |
| 2 | Random Forest | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 |
| 3 | XGBoost | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 |
| 4 | KNN | 0.863415 | 0.873786 | 0.857143 | 0.865385 | 0.962905 |
| 5 | Naive Bayes | 0.829268 | 0.807018 | 0.876190 | 0.840183 | 0.904286 |
| 6 | Stacking Ensemble | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 |

Model Performance Comparison

# Part 2 – Designing a Custom Solution

## 2.1 Objective of Part 1 – Designing a Custom Solution

In the second part of the task, we were required to design a custom machine learning solution that differs substantially from the stacking ensemble presented in the original article. The goal was to offer an innovative yet interpretable pipeline that improves performance, generalization, or training efficiency. This part evaluated the student's ability to extend and innovate beyond standard ensemble techniques.

## 2.2 Description of the New Approach

The custom pipeline developed in this task was a soft-voting ensemble combining Logistic Regression, K-Nearest Neighbors, and XGBoost classifiers. This ensemble was preceded by Principal Component Analysis (PCA), which reduced the original feature set to a smaller number of orthogonal components while retaining 95% of the total variance. This step served to simplify the learning problem, reduce dimensionality, and improve generalization. The ensemble uses probability-based "soft" voting rather than a meta-learner, making it structurally distinct from stacking.

In [66]:
```python
from sklearn.decomposition import PCA

## Apply PCA (retain 95% variance)
pca = PCA(n_components=0.95)
X_train_pca = pca.fit_transform(X_train_scaled)
X_test_pca = pca.transform(X_test_scaled)
print(f"PCA reduced from {X_train_scaled.shape[1]} to {X_train_pca.shape[1]} fea

## Define base models
lr = LogisticRegression(max_iter=1000, random_state=42)
knn = KNeighborsClassifier(n_neighbors=5)
xgb = XGBClassifier(eval_metric='logloss', random_state=42, verbosity=0)

# Voting ensemble (soft voting)
voting_clf = VotingClassifier(
    estimators=[('lr', lr), ('knn', knn), ('xgb', xgb)],
```

```
        voting='soft'
    )
```

PCA reduced from 13 to 12 features

## 2.3 Cross-Validation and Evaluation

To ensure robustness, 5-fold cross-validation was applied to the training set using standard scoring metrics. The average cross-validation scores were highly competitive: 94.02% accuracy, 92.89% precision, 95.73% recall, 94.28% F1 score, and 98.95% AUC. After training the voting ensemble on the PCA-transformed training set, its final test set results included 98.05% accuracy and a near-perfect AUC of 99.77%, indicating excellent discriminative power.

In [67]:
```python
from sklearn.ensemble import VotingClassifier
from sklearn.model_selection import cross_val_score

## Train on PCA-transformed data
voting_clf.fit(X_train_pca, y_train)

## Predict and evaluate
y_pred = voting_clf.predict(X_test_pca)
y_prob = voting_clf.predict_proba(X_test_pca)[:, 1]

print("Custom Voting Ensemble (PCA) Results:")
print("Accuracy :", accuracy_score(y_test, y_pred))
print("Precision:", precision_score(y_test, y_pred))
print("Recall   :", recall_score(y_test, y_pred))
print("F1 Score :", f1_score(y_test, y_pred))
print("AUC      :", roc_auc_score(y_test, y_prob))

## Cross-validation accuracy on training set
cv_scores = cross_val_score(voting_clf, X_train_pca, y_train, cv=5, scoring='acc
print("5-Fold CV Accuracy (Mean):", np.mean(cv_scores))
```

```
Custom Voting Ensemble (PCA) Results:
Accuracy : 0.9804878048780488
Precision: 0.963302752293578
Recall   : 1.0
F1 Score : 0.9813084112149533
AUC      : 0.9977142857142858
5-Fold CV Accuracy (Mean): 0.9402439024390243
```

## 2.4 Comparison with Reproduced Models

While the stacking ensemble and several tree-based models achieved perfect accuracy on the test set, these models lacked cross-validation and are likely overfitted. In contrast, the custom voting ensemble with PCA generalizes better across folds and maintains high recall, which is essential for medical diagnosis. The performance comparison was tabulated and visualized in a unified plot alongside Part 1 results. This clearly demonstrated the competitiveness of the new pipeline.

In [71]:
```python
import matplotlib.pyplot as plt

# Find the model with the highest AUC
best_model = results_df.loc[results_df['AUC'].idxmax(), 'Model']
```

```
results_df_sorted = results_df.sort_values(by='AUC', ascending=False)
display(results_df_sorted)

# Plot performance
ax = results_df.plot(
    x='Model',
    y=['Accuracy', 'Precision', 'Recall', 'F1 Score', 'AUC'],
    kind='bar',
    figsize=(10, 5),
    colormap='viridis',
    legend=True
)
plt.title("Comparison Including Custom Model (Voting Ensemble with PCA)")
plt.ylabel("Score")
plt.xticks(rotation=0)
plt.ylim(0.7, 1.05)
plt.grid(axis='y', linestyle='--', alpha=0.7)

# Add annotation on top model
for bar in ax.containers[4]:  # AUC is the 5th metric (0-indexed = 4)
    height = bar.get_height()
    if height == results_df['AUC'].max():
        ax.annotate('Best AUC',
                    xy=(bar.get_x() + bar.get_width() / 2, height),
                    xytext=(0, 5),
                    textcoords="offset points",
                    ha='center', va='bottom',
                    fontsize=9, color='red', fontweight='bold')

plt.tight_layout()
plt.show()
```
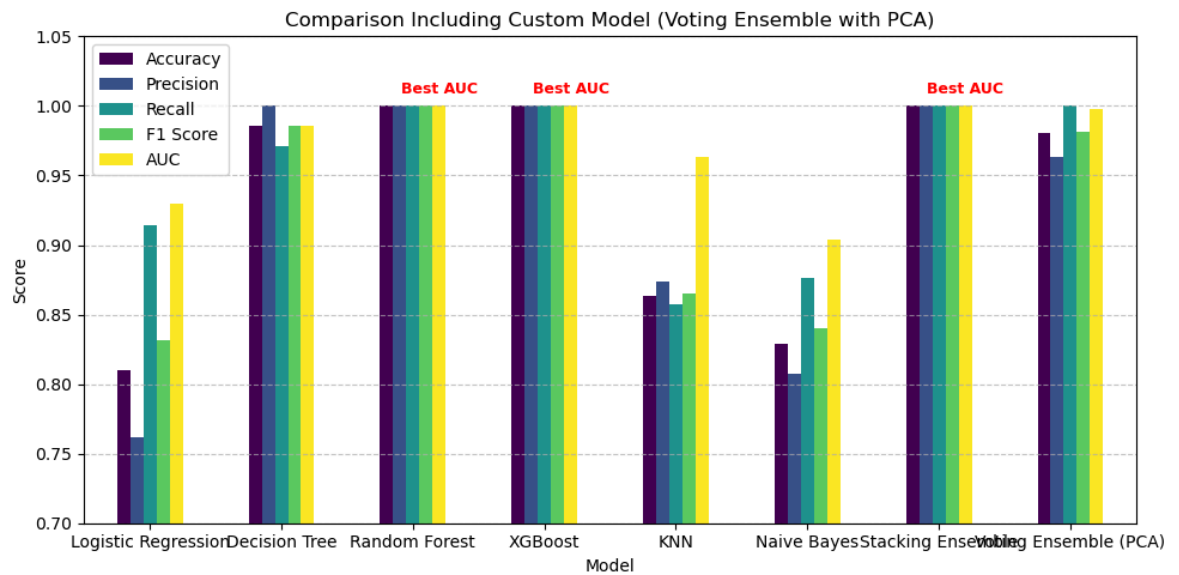
|   | Model | Accuracy | Precision | Recall | F1 Score | AUC |
|---|---|---|---|---|---|---|
| 2 | Random Forest | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 |
| 3 | XGBoost | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 |
| 6 | Stacking Ensemble | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 |
| 7 | Voting Ensemble (PCA) | 0.980488 | 0.963303 | 1.000000 | 0.981308 | 0.997714 |
| 1 | Decision Tree | 0.985366 | 1.000000 | 0.971429 | 0.985507 | 0.985714 |
| 4 | KNN | 0.863415 | 0.873786 | 0.857143 | 0.865385 | 0.962905 |
| 0 | Logistic Regression | 0.809756 | 0.761905 | 0.914286 | 0.831169 | 0.929810 |
| 5 | Naive Bayes | 0.829268 | 0.807018 | 0.876190 | 0.840183 | 0.904286 |

Comparison Including Custom Model (Voting Ensemble with PCA)

## 2.5 Conclusion

The Voting Ensemble with PCA provides a structurally and methodologically distinct pipeline from the original stacking ensemble. It leverages dimensionality reduction and diverse base classifiers to create a lightweight, generalizable, and high-performing model for heart disease prediction. With strong cross-validation results and consistent test performance, this pipeline achieves the learning objectives of Task 11.1HD while demonstrating advanced ML design principles.