```
from google.colab import drive
drive.mount('/content/drive')
%cd /content/drive/MyDrive/Colab Notebooks
```

Mounted at /content/drive
/content/drive/MyDrive/Colab Notebooks

```python
# === 1. Install necessary libraries ===
!pip install xgboost lightgbm scikit-learn matplotlib tensorflow

# === 2. Imports ===
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
from sklearn.ensemble import RandomForestClassifier
from xgboost import XGBClassifier
from sklearn.feature_selection import mutual_info_classif
from lightgbm import LGBMClassifier
from sklearn.metrics import accuracy_score
from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import Dense, Dropout

# Load your CSV from this path
df = pd.read_csv("CBioFiltered_Genes_Subset.csv")

# === 4. Preprocessing ===
df_cleaned = df.drop(columns=["Unnamed: 0", "batch"])
X = df_cleaned.drop(columns=["group"])
y = df_cleaned["group"].astype(int)

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

scaler = StandardScaler()
X_train_scaled = scaler.fit_transform(X_train)
X_test_scaled = scaler.transform(X_test)

# === 5. Feature selection ===
rfc = RandomForestClassifier(n_estimators=100, random_state=42)
rfc.fit(X_train_scaled, y_train)
rfc_importance = rfc.feature_importances_

xgb = XGBClassifier(use_label_encoder=False, eval_metric='logloss', random_state=42)
xgb.fit(X_train_scaled, y_train)
xgb_importance = xgb.feature_importances_

mi = mutual_info_classif(X_train_scaled, y_train, random_state=42)

# Combine using voting
rfc_norm = rfc_importance / np.max(rfc_importance)
xgb_norm = xgb_importance / np.max(xgb_importance)
mi_norm = mi / np.max(mi)
combined_score = (rfc_norm + xgb_norm + mi_norm) / 3

feature_importance_df = pd.DataFrame({
    'Feature': X.columns,
    'RFC': rfc_norm,
    'XGBoost': xgb_norm,
    'MI': mi_norm,
    'MeanScore': combined_score
}).sort_values(by="MeanScore", ascending=False)

top_features = feature_importance_df.head(30)["Feature"].values

# Plot top features
plt.figure(figsize=(12, 8))
plt.barh(feature_importance_df.head(30)['Feature'], feature_importance_df.head(30)['MeanScore'], color='skyblue')
plt.xlabel("Importance Score (Normalized Average)")
plt.title("Top 20 Important Features (RFC + XGBoost + MI)")
plt.gca().invert_yaxis()
plt.tight_layout()
plt.show()

# === 6. Prepare top features for modeling ===
X_train_top = X_train_scaled[:, [X.columns.get_loc(f) for f in top_features]]
X_test_top = X_test_scaled[:, [X.columns.get_loc(f) for f in top_features]]

# === 7. Evaluate models ===

# Random Forest
```

```python
rfc_model = RandomForestClassifier(random_state=42)
rfc_model.fit(X_train_top, y_train)
rfc_pred = rfc_model.predict(X_test_top)
print("Random Forest Accuracy:", accuracy_score(y_test, rfc_pred))

# XGBoost
xgb_model = XGBClassifier(use_label_encoder=False, eval_metric='logloss', random_state=42)
xgb_model.fit(X_train_top, y_train)
xgb_pred = xgb_model.predict(X_test_top)
print("XGBoost Accuracy:", accuracy_score(y_test, xgb_pred))

# LightGBM
lgbm_model = LGBMClassifier(random_state=42)
lgbm_model.fit(X_train_top, y_train)
lgbm_pred = lgbm_model.predict(X_test_top)
print("LightGBM Accuracy:", accuracy_score(y_test, lgbm_pred))

# ANN (Keras)
model = Sequential()
model.add(Dense(64, input_dim=X_train_top.shape[1], activation='relu'))
model.add(Dropout(0.3))
model.add(Dense(32, activation='relu'))
model.add(Dense(1, activation='sigmoid'))  # Use '1' output for binary

model.compile(optimizer='adam', loss='binary_crossentropy', metrics=['accuracy'])
model.fit(X_train_top, y_train, epochs=50, batch_size=16, verbose=0)
ann_loss, ann_acc = model.evaluate(X_test_top, y_test, verbose=0)
print("ANN Accuracy:", ann_acc)
```
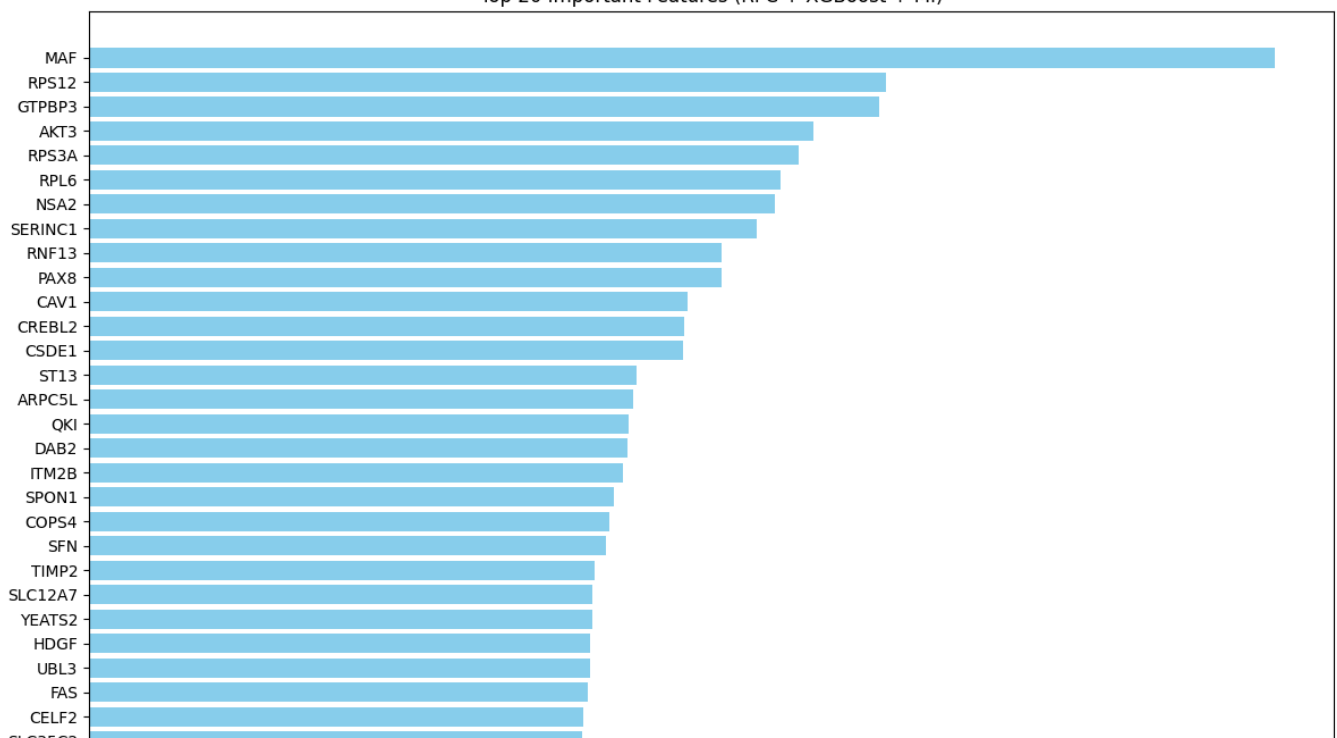
Requirement already satisfied: xgboost in /usr/local/lib/python3.11/dist-packages (2.1.4)
Requirement already satisfied: lightgbm in /usr/local/lib/python3.11/dist-packages (4.5.0)
Requirement already satisfied: scikit-learn in /usr/local/lib/python3.11/dist-packages (1.6.1)
Requirement already satisfied: matplotlib in /usr/local/lib/python3.11/dist-packages (3.10.0)
Requirement already satisfied: tensorflow in /usr/local/lib/python3.11/dist-packages (2.18.0)
Requirement already satisfied: numpy in /usr/local/lib/python3.11/dist-packages (from xgboost) (2.0.2)
Requirement already satisfied: nvidia-nccl-cu12 in /usr/local/lib/python3.11/dist-packages (from xgboost) (2.21.5)
Requirement already satisfied: scipy in /usr/local/lib/python3.11/dist-packages (from xgboost) (1.15.2)
Requirement already satisfied: joblib>=1.2.0 in /usr/local/lib/python3.11/dist-packages (from scikit-learn) (1.4.2)
Requirement already satisfied: threadpoolctl>=3.1.0 in /usr/local/lib/python3.11/dist-packages (from scikit-learn) (3.6.0)
Requirement already satisfied: contourpy>=1.0.1 in /usr/local/lib/python3.11/dist-packages (from matplotlib) (1.3.2)
Requirement already satisfied: cycler>=0.10 in /usr/local/lib/python3.11/dist-packages (from matplotlib) (0.12.1)
Requirement already satisfied: fonttools>=4.22.0 in /usr/local/lib/python3.11/dist-packages (from matplotlib) (4.57.0)
Requirement already satisfied: kiwisolver>=1.3.1 in /usr/local/lib/python3.11/dist-packages (from matplotlib) (1.4.8)
Requirement already satisfied: packaging>=20.0 in /usr/local/lib/python3.11/dist-packages (from matplotlib) (24.2)
Requirement already satisfied: pillow>=8 in /usr/local/lib/python3.11/dist-packages (from matplotlib) (11.2.1)
Requirement already satisfied: pyparsing>=2.3.1 in /usr/local/lib/python3.11/dist-packages (from matplotlib) (3.2.3)
Requirement already satisfied: python-dateutil>=2.7 in /usr/local/lib/python3.11/dist-packages (from matplotlib) (2.9.0.post0)
Requirement already satisfied: absl-py>=1.0.0 in /usr/local/lib/python3.11/dist-packages (from tensorflow) (1.4.0)
Requirement already satisfied: astunparse>=1.6.0 in /usr/local/lib/python3.11/dist-packages (from tensorflow) (1.6.3)
Requirement already satisfied: flatbuffers>=24.3.25 in /usr/local/lib/python3.11/dist-packages (from tensorflow) (25.2.10)
Requirement already satisfied: gast!=0.5.0,!=0.5.1,!=0.5.2,>=0.2.1 in /usr/local/lib/python3.11/dist-packages (from tensorflow) (0.6
Requirement already satisfied: google-pasta>=0.1.1 in /usr/local/lib/python3.11/dist-packages (from tensorflow) (0.2.0)
Requirement already satisfied: libclang>=13.0.0 in /usr/local/lib/python3.11/dist-packages (from tensorflow) (18.1.1)
Requirement already satisfied: opt-einsum>=2.3.2 in /usr/local/lib/python3.11/dist-packages (from tensorflow) (3.4.0)
Requirement already satisfied: protobuf!=4.21.0,!=4.21.1,!=4.21.2,!=4.21.3,!=4.21.4,!=4.21.5,<6.0.0dev,>=3.20.3 in /usr/local/lib/py
Requirement already satisfied: requests<3,>=2.21.0 in /usr/local/lib/python3.11/dist-packages (from tensorflow) (2.32.3)
Requirement already satisfied: setuptools in /usr/local/lib/python3.11/dist-packages (from tensorflow) (75.2.0)
Requirement already satisfied: six>=1.12.0 in /usr/local/lib/python3.11/dist-packages (from tensorflow) (1.17.0)
Requirement already satisfied: termcolor>=1.1.0 in /usr/local/lib/python3.11/dist-packages (from tensorflow) (3.1.0)
Requirement already satisfied: typing-extensions>=3.6.6 in /usr/local/lib/python3.11/dist-packages (from tensorflow) (4.13.2)
Requirement already satisfied: wrapt>=1.11.0 in /usr/local/lib/python3.11/dist-packages (from tensorflow) (1.17.2)
Requirement already satisfied: grpcio<2.0,>=1.24.3 in /usr/local/lib/python3.11/dist-packages (from tensorflow) (1.71.0)
Requirement already satisfied: tensorboard<2.19,>=2.18 in /usr/local/lib/python3.11/dist-packages (from tensorflow) (2.18.0)
Requirement already satisfied: keras>=3.5.0 in /usr/local/lib/python3.11/dist-packages (from tensorflow) (3.8.0)
Requirement already satisfied: h5py>=3.11.0 in /usr/local/lib/python3.11/dist-packages (from tensorflow) (3.13.0)
Requirement already satisfied: ml-dtypes<0.5.0,>=0.4.0 in /usr/local/lib/python3.11/dist-packages (from tensorflow) (0.4.1)
Requirement already satisfied: tensorflow-io-gcs-filesystem>=0.23.1 in /usr/local/lib/python3.11/dist-packages (from tensorflow) (0
Requirement already satisfied: wheel<1.0,>=0.23.0 in /usr/local/lib/python3.11/dist-packages (from astunparse>=1.6.0->tensorflow) (6
Requirement already satisfied: rich in /usr/local/lib/python3.11/dist-packages (from keras>=3.5.0->tensorflow) (13.9.4)
Requirement already satisfied: namex in /usr/local/lib/python3.11/dist-packages (from keras>=3.5.0->tensorflow) (0.0.9)
Requirement already satisfied: optree in /usr/local/lib/python3.11/dist-packages (from keras>=3.5.0->tensorflow) (0.15.0)
Requirement already satisfied: charset-normalizer<4,>=2 in /usr/local/lib/python3.11/dist-packages (from requests<3,>=2.21.0->tensor
Requirement already satisfied: idna<4,>=2.5 in /usr/local/lib/python3.11/dist-packages (from requests<3,>=2.21.0->tensorflow) (3.10
Requirement already satisfied: urllib3<3,>=1.21.1 in /usr/local/lib/python3.11/dist-packages (from requests<3,>=2.21.0->tensorflow)
Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.11/dist-packages (from requests<3,>=2.21.0->tensorflow)
Requirement already satisfied: markdown>=2.6.8 in /usr/local/lib/python3.11/dist-packages (from tensorboard<2.19,>=2.18->tensorflow
Requirement already satisfied: tensorboard-data-server<0.8.0,>=0.7.0 in /usr/local/lib/python3.11/dist-packages (from tensorboard<2
Requirement already satisfied: werkzeug>=1.0.1 in /usr/local/lib/python3.11/dist-packages (from tensorboard<2.19,>=2.18->tensorflow
Requirement already satisfied: MarkupSafe>=2.1.1 in /usr/local/lib/python3.11/dist-packages (from werkzeug>=1.0.1->tensorboard<2.19,
Requirement already satisfied: markdown-it-py>=2.2.0 in /usr/local/lib/python3.11/dist-packages (from rich->keras>=3.5.0->tensorflow
Requirement already satisfied: pygments<3.0.0,>=2.13.0 in /usr/local/lib/python3.11/dist-packages (from rich->keras>=3.5.0->tensorfl
Requirement already satisfied: mdurl~=0.1 in /usr/local/lib/python3.11/dist-packages (from markdown-it-py>=2.2.0->rich->keras>=3.5.6
/usr/local/lib/python3.11/dist-packages/xgboost/core.py:158: UserWarning: [04:41:28] WARNING: /workspace/src/learner.cc:740:
Parameters: { "use_label_encoder" } are not used.

  warnings.warn(smsg, UserWarning)

Top 20 Important Features (RFC + XGBoost + MI)

Importance Score (Normalized Average)

```
Random Forest Accuracy: 0.9327731092436975
XGBoost Accuracy: 0.9327731092436975
[LightGBM] [Info] Number of positive: 422, number of negative: 53
[LightGBM] [Info] Auto-choosing col-wise multi-threading, the overhead of testing was 0.000204 seconds.
You can set `force_col_wise=true` to remove the overhead.
[LightGBM] [Info] Total Bins 4758
[LightGBM] [Info] Number of data points in the train set: 475, number of used features: 30
[LightGBM] [Info] [binary:BoostFromScore]: pavg=0.888421 -> initscore=2.074713
[LightGBM] [Info] Start training from score 2.074713
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
```

```
/usr/local/lib/python3.11/dist-packages/xgboost/core.py:158: UserWarning: [04:41:52] WARNING: /workspace/src/learner.cc:740:
Parameters: { "use_label_encoder" } are not used.

  warnings.warn(smsg, UserWarning)
/usr/local/lib/python3.11/dist-packages/sklearn/utils/deprecation.py:151: FutureWarning: 'force_all_finite' was renamed to 'ensure_a
  warnings.warn(
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
LightGBM Accuracy: 0.9411764705882353
/usr/local/lib/python3.11/dist-packages/sklearn/utils/deprecation.py:151: FutureWarning: 'force_all_finite' was renamed to 'ensure_a
  warnings.warn(
/usr/local/lib/python3.11/dist-packages/keras/src/layers/core/dense.py:87: UserWarning: Do not pass an `input_shape`/`input_dim` arg
  super().__init__(activity_regularizer=activity_regularizer, **kwargs)
ANN Accuracy: 0.9747899174690247
```

```python
from sklearn.metrics import confusion_matrix, roc_curve, auc, ConfusionMatrixDisplay
import seaborn as sns

# Plotting function for ROC curves
def plot_roc_curves(models, model_names, X_test, y_test):
    plt.figure(figsize=(10, 7))
    for model, name in zip(models, model_names):
        if hasattr(model, "predict_proba"):
            y_proba = model.predict_proba(X_test)[:, 1]
        else:  # ANN model returns probability directly
            y_proba = model.predict(X_test).ravel()
        fpr, tpr, _ = roc_curve(y_test, y_proba)
        roc_auc = auc(fpr, tpr)
        plt.plot(fpr, tpr, label=f'{name} (AUC = {roc_auc:.2f})')

    plt.plot([0, 1], [0, 1], 'k--')
    plt.xlabel("False Positive Rate")
    plt.ylabel("True Positive Rate")
    plt.title("ROC Curves")
    plt.legend(loc="lower right")
    plt.grid(True)
    plt.show()

# Plotting function for confusion matrices
def plot_confusion_matrices(models, model_names, X_test, y_test):
    for model, name in zip(models, model_names):
        if name == "ANN":
            y_pred = (model.predict(X_test).ravel() > 0.5).astype(int)
        else:
            y_pred = model.predict(X_test)
        cm = confusion_matrix(y_test, y_pred)
        disp = ConfusionMatrixDisplay(confusion_matrix=cm)
        disp.plot(cmap='Blues')
        plt.title(f"Confusion Matrix - {name}")
        plt.show()

# List of models and their names
models = [rfc_model, xgb_model, lgbm_model, model]
model_names = ["Random Forest", "XGBoost", "LightGBM", "ANN"]

# Plot Confusion Matrices
plot_confusion_matrices(models, model_names, X_test_top, y_test)

# Plot ROC Curves
plot_roc_curves(models, model_names, X_test_top, y_test)
```
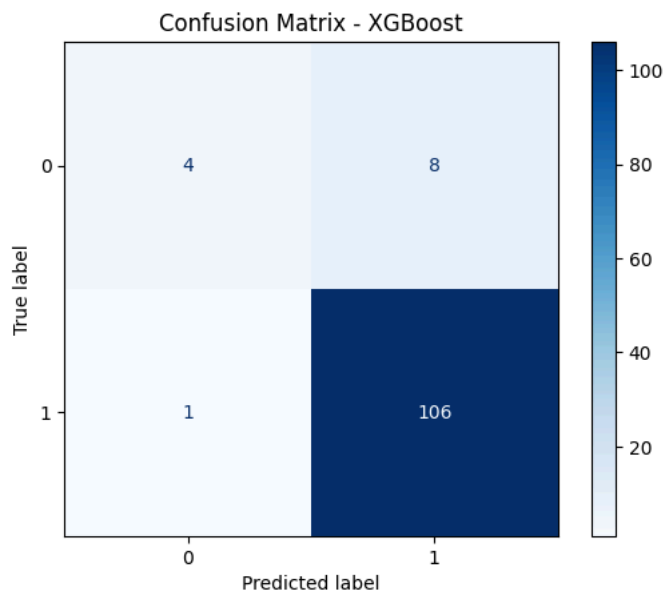
## Confusion Matrix - Random Forest

|  | Predicted 0 | Predicted 1 |
|---|---|---|
| True 0 | 4 | 8 |
| True 1 | 0 | 107 |

## Confusion Matrix - XGBoost

|  | Predicted 0 | Predicted 1 |
|---|---|---|
| True 0 | 4 | 8 |
| True 1 | 1 | 106 |

/usr/local/lib/python3.11/dist-packages/sklearn/utils/deprecation.py:151: FutureWarning: 'force_all_finite' was renamed to 'ensur
  warnings.warn(

## Confusion Matrix - LightGBM

|  | Predicted 0 | Predicted 1 |
|---|---|---|
| True 0 | 4 | 8 |
| True 1 | 0 | 107 |

4/4 ━━━━━━━━━━━━━━━━━━━━ 0s 29ms/step

## Confusion Matrix - ANN

|  | Predicted 0 | Predicted 1 |
|---|---|---|
| True 0 | 8 | 4 |

True label

Predicted label

/usr/local/lib/python3.11/dist-packages/sklearn/utils/deprecation.py:151: FutureWarning: 'force_all_finite' was renamed to 'ensur
  warnings.warn(
**4/4** ──────────────── **0s** 24ms/step



ROC Curves

True Positive Rate

False Positive Rate

Random Forest (AUC = 0.90)
XGBoost (AUC = 0.92)
LightGBM (AUC = 0.94)
ANN (AUC = 0.88)

```python
top_features = feature_importance_df.head(30)["Feature"].values


print("Top 20 features based on ensemble importance voting:")
for i, feature in enumerate(top_features, 1):
    print(f"{i}. {feature}")
```

Top 20 features based on ensemble importance voting:
  1. MAF
  2. RPS12
  3. GTPBP3
  4. AKT3
  5. RPS3A
  6. RPL6
  7. NSA2
  8. SERINC1
  9. RNF13
  10. PAX8
  11. CAV1
  12. CREBL2
  13. CSDE1
  14. ST13
  15. ARPC5L
  16. QKI
  17. DAB2
  18. ITM2B
  19. SPON1
  20. COPS4
  21. SFN
  22. TIMP2
  23. SLC12A7
  24. YEATS2
  25. HDGF
  26. UBL3
  27. FAS
  28. CELF2
  29. SLC35C2
  30. BNC2


```python
import seaborn as sns

# Select top 20 features from the original (unscaled) data for interpretability
X_top_20 = X[top_features]
X_top_20["group"] = y.values  # Add target label back

# Set up the plot
plt.figure(figsize=(16, 10))
sns.heatmap(X_top_20.groupby("group").mean(), cmap="viridis", annot=True, fmt=".2f", cbar=True)

plt.title("Mean Feature Expression Heatmap by Group (Top 20 Features)")
plt.xlabel("Feature")
plt.ylabel("Group")
plt.tight_layout()
plt.show()
```

Mean Feature Expression Heatmap by Group (Top 20 Features)

```
import seaborn as sns

# Select and prepare top features
X_top_20 = X[top_features]
X_top_20["group"] = y.values

# Transpose the grouped mean for rotated heatmap
heatmap_data = X_top_20.groupby("group").mean().T

# Plot
plt.figure(figsize=(12, 14))
sns.heatmap(heatmap_data, cmap="viridis", annot=True, fmt=".2f", cbar=True)

plt.title("Top 20 Feature Mean Expression by Group (Rotated Heatmap)")
plt.xlabel("Group")
plt.ylabel("Feature")
plt.tight_layout()
plt.show()
```

## Top 20 Feature Mean Expression by Group (Rotated Heatmap)

| Feature | Group 0 | Group 1 |
|---------|---------|---------|
| MAF | 191.29 | 137.74 |
| RPS12 | 334.72 | 243.73 |
| GTPBP3 | 66.92 | 77.02 |
| AKT3 | 255.56 | 177.32 |
| RPS3A | 744.59 | 583.47 |
| RPL6 | 251.74 | 183.31 |
| NSA2 | 169.41 | 102.96 |
| SERINC1 | 99.14 | 70.27 |
| RNF13 | 128.24 | 98.18 |
| PAX8 | 304.14 | 385.94 |
| CAV1 | 229.67 | 125.11 |
| CREBL2 | 204.07 | 162.17 |
| CSDE1 | 278.27 | 219.73 |
| ST13 | 418.90 | 251.45 |
| ARPC5L | 138.38 | 167.02 |
| QKI | 454.17 | 390.90 |
| DAB2 | 278.02 | 228.05 |
| ITM2B | 411.59 | 296.49 |
| SPON1 | 162.08 | 345.55 |
| COPS4 | 77.97 | 64.00 |

```python
# === 2. Imports ===
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
from sklearn.ensemble import RandomForestClassifier
from xgboost import XGBClassifier
from sklearn.feature_selection import mutual_info_classif
from lightgbm import LGBMClassifier
from sklearn.metrics import accuracy_score
from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import Dense, Dropout


# === 3. Load Dataset ===
df = pd.read_csv("Predicted_OC_subset_with_labels.csv")

# === 4. Preprocessing ===
df_cleaned = df.drop(columns=["Unnamed: 0"], errors='ignore')  # Drop index col if exists
X = df_cleaned.drop(columns=["predicted_group"])  # Features
y = df_cleaned["predicted_group"].astype(int)      # Labels

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)


scaler = StandardScaler()
X_train_scaled = scaler.fit_transform(X_train)
X_test_scaled = scaler.transform(X_test)

# === 5. Feature selection via ensemble ===
rfc = RandomForestClassifier(n_estimators=100, random_state=42)
rfc.fit(X_train_scaled, y_train)
rfc_importance = rfc.feature_importances_

xgb = XGBClassifier(use_label_encoder=False, eval_metric='logloss', random_state=42)
xgb.fit(X_train_scaled, y_train)
xgb_importance = xgb.feature_importances_

mi = mutual_info_classif(X_train_scaled, y_train, random_state=42)

# Normalize and vote
rfc_norm = rfc_importance / np.max(rfc_importance)
xgb_norm = xgb_importance / np.max(xgb_importance)
mi_norm = mi / np.max(mi)
combined_score = (rfc_norm + xgb_norm + mi_norm) / 3

feature_importance_df = pd.DataFrame({
    'Feature': X.columns,
    'RFC': rfc_norm,
    'XGBoost': xgb_norm,
    'MI': mi_norm,
    'MeanScore': combined_score
}).sort_values(by="MeanScore", ascending=False)

top_features = feature_importance_df.head(30)["Feature"].values

# === 6. Visualize top features ===
plt.figure(figsize=(12, 8))
plt.barh(feature_importance_df.head(30)['Feature'], feature_importance_df.head(20)['MeanScore'], color='skyblue')
plt.xlabel("Importance Score (Normalized Average)")
plt.title("Top 20 Important Features (RFC + XGBoost + MI)")
plt.gca().invert_yaxis()
plt.tight_layout()
plt.show()

# === 7. Prepare top features for modeling ===
X_train_top = X_train_scaled[:, [X.columns.get_loc(f) for f in top_features]]
X_test_top = X_test_scaled[:, [X.columns.get_loc(f) for f in top_features]]

# === 8. Model Evaluations ===

# Random Forest
rfc_model = RandomForestClassifier(random_state=42)
rfc_model.fit(X_train_top, y_train)
rfc_pred = rfc_model.predict(X_test_top)
print("Random Forest Accuracy:", accuracy_score(y_test, rfc_pred))

# XGBoost
xgb_model = XGBClassifier(use_label_encoder=False, eval_metric='logloss', random_state=42)
xgb_model.fit(X_train_top, y_train)
xgb_pred = xgb_model.predict(X_test_top)
```

```python
print("XGBoost Accuracy:", accuracy_score(y_test, xgb_pred))


# LightGBM
lgbm_model = LGBMClassifier(random_state=42)
lgbm_model.fit(X_train_top, y_train)
lgbm_pred = lgbm_model.predict(X_test_top)
print("LightGBM Accuracy:", accuracy_score(y_test, lgbm_pred))


# ANN using Keras
model = Sequential()
model.add(Dense(64, input_dim=X_train_top.shape[1], activation='relu'))
model.add(Dropout(0.3))
model.add(Dense(32, activation='relu'))
model.add(Dense(1, activation='sigmoid'))  # Binary classification

model.compile(optimizer='adam', loss='binary_crossentropy', metrics=['accuracy'])
model.fit(X_train_top, y_train, epochs=50, batch_size=16, verbose=0)
ann_loss, ann_acc = model.evaluate(X_test_top, y_test, verbose=0)
print("ANN Accuracy:", ann_acc)
```

Top 20 Important Features (RFC + XGBoost + MI)



Random Forest Accuracy: 0.4864864864864865
XGBoost Accuracy: 0.4594594594594595
[LightGBM] [Info] Number of positive: 72, number of negative: 74
[LightGBM] [Info] Auto-choosing col-wise multi-threading, the overhead of testing was 0.000137 seconds.
You can set `force_col_wise=true` to remove the overhead.
[LightGBM] [Info] Total Bins 1508
[LightGBM] [Info] Number of data points in the train set: 146, number of used features: 30
[LightGBM] [Info] [binary:BoostFromScore]: pavg=0.493151 -> initscore=-0.027399
[LightGBM] [Info] Start training from score -0.027399
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf

```
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
LightGBM Accuracy: 0.5675675675675675
/usr/local/lib/python3.11/dist-packages/xgboost/core.py:158: UserWarning: [04:29:10] WARNING: /workspace/src/learner.cc:740:
Parameters: { "use_label_encoder" } are not used.

  warnings.warn(smsg, UserWarning)
/usr/local/lib/python3.11/dist-packages/sklearn/utils/deprecation.py:151: FutureWarning: 'force_all_finite' was renamed to 'ensur
  warnings.warn(
/usr/local/lib/python3.11/dist-packages/sklearn/utils/deprecation.py:151: FutureWarning: 'force_all_finite' was renamed to 'ensur
  warnings.warn(
/usr/local/lib/python3.11/dist-packages/keras/src/layers/core/dense.py:87: UserWarning: Do not pass an `input_shape`/`input_dim`
  super().__init__(activity_regularizer=activity_regularizer, **kwargs)
ANN Accuracy: 0.45945945382118225
```

```python
# === 2. Imports ===
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
from sklearn.model_selection import train_test_split, GridSearchCV
from sklearn.preprocessing import StandardScaler
from sklearn.ensemble import RandomForestClassifier, VotingClassifier
from xgboost import XGBClassifier
from sklearn.feature_selection import mutual_info_classif
from lightgbm import LGBMClassifier
from sklearn.metrics import accuracy_score
from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import Dense, Dropout
from tensorflow.keras.callbacks import EarlyStopping


# === 3. Load Dataset ===
df = pd.read_csv("Predicted_OC_subset_with_labels.csv")


# === 4. Preprocessing ===
df_cleaned = df.drop(columns=["Unnamed: 0"], errors='ignore')  # Drop index col if exists
X = df_cleaned.drop(columns=["predicted_group"])  # Features
y = df_cleaned["predicted_group"].astype(int)      # Labels

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

scaler = StandardScaler()
X_train_scaled = scaler.fit_transform(X_train)
X_test_scaled = scaler.transform(X_test)


# === 5. Feature selection via ensemble voting ===
rfc = RandomForestClassifier(n_estimators=100, random_state=42)
rfc.fit(X_train_scaled, y_train)
rfc_importance = rfc.feature_importances_

xgb = XGBClassifier(use_label_encoder=False, eval_metric='logloss', random_state=42)
xgb.fit(X_train_scaled, y_train)
xgb_importance = xgb.feature_importances_

mi = mutual_info_classif(X_train_scaled, y_train, random_state=42)


# Normalize and combine
rfc_norm = rfc_importance / np.max(rfc_importance)
xgb_norm = xgb_importance / np.max(xgb_importance)
mi_norm = mi / np.max(mi)
combined_score = (rfc_norm + xgb_norm + mi_norm) / 3

feature_importance_df = pd.DataFrame({
    'Feature': X.columns,
    'RFC': rfc_norm,
    'XGBoost': xgb_norm,
    'MI': mi_norm,
    'MeanScore': combined_score
}).sort_values(by="MeanScore", ascending=False)

top_features = feature_importance_df.head(30)["Feature"].values

# === 6. Visualize top features ===
plt.figure(figsize=(12, 8))
plt.barh(feature_importance_df.head(30)['Feature'], feature_importance_df.head(30)['MeanScore'], color='skyblue')
plt.xlabel("Importance Score (Normalized Average)")
plt.title("Top 30 Important Features (RFC + XGBoost + MI)")
plt.gca().invert_yaxis()
plt.tight_layout()
plt.show()

# === 7. Prepare data with selected features ===
X_train_top = X_train_scaled[:, [X.columns.get_loc(f) for f in top_features]]
X_test_top = X_test_scaled[:, [X.columns.get_loc(f) for f in top_features]]

# === 8. Hyperparameter Tuning for RF ===
param_grid = {
    'n_estimators': [100, 200],
    'max_depth': [None, 10, 20],
    'min_samples_split': [2, 5]
}

grid_search = GridSearchCV(RandomForestClassifier(random_state=42), param_grid, cv=3, scoring='accuracy', n_jobs=-1)
grid_search.fit(X_train_top, y_train)
best_rf = grid_search.best_estimator_
```

```python
# === 9. Train XGBoost and LightGBM ===
xgb_model = XGBClassifier(use_label_encoder=False, eval_metric='logloss', random_state=42)
xgb_model.fit(X_train_top, y_train)

lgbm_model = LGBMClassifier(random_state=42)
lgbm_model.fit(X_train_top, y_train)

# === 10. Voting Classifier Ensemble ===
voting_model = VotingClassifier(estimators=[
    ('rf', best_rf),
    ('xgb', xgb_model),
    ('lgbm', lgbm_model)
], voting='soft')

voting_model.fit(X_train_top, y_train)
voting_pred = voting_model.predict(X_test_top)
print("Voting Classifier Accuracy:", accuracy_score(y_test, voting_pred))

# === 11. Improved ANN ===
model = Sequential()
model.add(Dense(128, input_dim=X_train_top.shape[1], activation='relu'))
model.add(Dropout(0.4))
model.add(Dense(64, activation='relu'))
model.add(Dropout(0.3))
model.add(Dense(32, activation='relu'))
model.add(Dense(1, activation='sigmoid'))  # Binary output

model.compile(optimizer='adam', loss='binary_crossentropy', metrics=['accuracy'])

early_stop = EarlyStopping(monitor='val_loss', patience=5, restore_best_weights=True)
model.fit(X_train_top, y_train, validation_split=0.2, epochs=100, batch_size=16, verbose=0, callbacks=[early_stop])

ann_loss, ann_acc = model.evaluate(X_test_top, y_test, verbose=0)
print("Improved ANN Accuracy:", ann_acc)
```

Top 30 Important Features (RFC + XGBoost + MI)

```
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
/usr/local/lib/python3.11/dist-packages/xgboost/core.py:158: UserWarning: [03:36:10] WARNING: /workspace/src/learner.cc:740:
Parameters: { "use_label_encoder" } are not used.

  warnings.warn(smsg, UserWarning)
/usr/local/lib/python3.11/dist-packages/sklearn/utils/deprecation.py:151: FutureWarning: 'force_all_finite' was renamed to 'ensur
  warnings.warn(
/usr/local/lib/python3.11/dist-packages/sklearn/utils/deprecation.py:151: FutureWarning: 'force_all_finite' was renamed to 'ensur
  warnings.warn(
/usr/local/lib/python3.11/dist-packages/keras/src/layers/core/dense.py:87: UserWarning: Do not pass an `input_shape`/`input_dim`
  super().__init__(activity_regularizer=activity_regularizer, **kwargs)
[LightGBM] [Info] Number of positive: 72, number of negative: 74
[LightGBM] [Info] Auto-choosing col-wise multi-threading, the overhead of testing was 0.000122 seconds.
You can set `force_col_wise=true` to remove the overhead.
[LightGBM] [Info] Total Bins 1508
[LightGBM] [Info] Number of data points in the train set: 146, number of used features: 30
[LightGBM] [Info] [binary:BoostFromScore]: pavg=0.493151 -> initscore=-0.027399
[LightGBM] [Info] Start training from score -0.027399
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
```

```
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
```

```
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
Voting Classifier Accuracy: 0.4594594594594595
Improved ANN Accuracy: 0.5945945978164673
```

```
# === Print Top 30 Features with Importances ===
print("\n=== Top 30 Features by Ensemble Importance ===")
print(feature_importance_df.head(30).to_string(index=False))
```

⮞
```
    === Top 30 Features by Ensemble Importance ===
      Feature      RFC  XGBoost       MI  MeanScore
        RFTN1 0.583471 0.417042 0.675741   0.558752
        CNPY3 0.951109 0.117479 0.560896   0.543161
         MCAM 0.428714 1.000000 0.180788   0.536501
       FKBP11 0.505867 0.763407 0.299585   0.522953
         GDI2 1.000000 0.082486 0.372769   0.485085
        BLVRA 0.646419 0.467693 0.267867   0.460660
          FTL 0.464914 0.000000 0.825350   0.430088
      PPP2R3C 0.585084 0.000000 0.662829   0.415971
      TMEM135 0.515351 0.000000 0.723869   0.413073
         DHX9 0.000000 0.657232 0.533063   0.396765
       PRPF19 0.000000 0.676384 0.457683   0.378022
        UBE2H 0.248739 0.000000 0.872064   0.373601
        UBE4B 0.319527 0.158539 0.592147   0.356738
         CLK4 0.561675 0.021279 0.479266   0.354074
        PPARD 0.375995 0.000000 0.685518   0.353837
        WDR37 0.516124 0.000000 0.529543   0.348556
       LARP4B 0.202867 0.209571 0.624410   0.345616
        VPS53 0.619001 0.013702 0.382622   0.338442
        TOP2B 0.000000 0.000000 1.000000   0.333333
        BBIP1 0.840277 0.158752 0.000000   0.333009
     RAB3GAP1 0.163612 0.000000 0.825869   0.329827
       DYNLT3 0.223822 0.077134 0.682958   0.327971
        RPL31 0.000000 0.000000 0.971875   0.323958
        CAND2 0.348786 0.229128 0.389635   0.322516
        SNX27 0.252671 0.000000 0.704502   0.319058
        ITGB4 0.430054 0.000000 0.524527   0.318194
        XRCC6 0.000000 0.000000 0.950620   0.316873
         PIGG 0.307097 0.000000 0.639732   0.315610
       NDUFB4 0.000000 0.586018 0.351224   0.312414
        NCOR1 0.374605 0.043673 0.517302   0.311860
```

```
from sklearn.metrics import classification_report, confusion_matrix

# === Helper function to evaluate a model ===
def evaluate_model(name, model, X_test, y_test):
    preds = model.predict(X_test)
    print(f"\n=== {name} ===")
    print("Accuracy:", accuracy_score(y_test, preds))
    print("Classification Report:")
    print(classification_report(y_test, preds))
    print("Confusion Matrix:")
    print(confusion_matrix(y_test, preds))

# Evaluate Random Forest (best_rf from GridSearch)
evaluate_model("Tuned Random Forest", best_rf, X_test_top, y_test)

# Evaluate XGBoost
evaluate_model("XGBoost", xgb_model, X_test_top, y_test)

# Evaluate LightGBM
evaluate_model("LightGBM", lgbm_model, X_test_top, y_test)

# Evaluate Voting Classifier
evaluate_model("Voting Classifier (Ensemble)", voting_model, X_test_top, y_test)

# Evaluate ANN separately (sigmoid output thresholded)
ann_pred_prob = model.predict(X_test_top).flatten()
ann_pred = (ann_pred_prob > 0.5).astype(int)
print("\n=== ANN (Keras) ===")
print("Accuracy:", accuracy_score(y_test, ann_pred))
print("Classification Report:")
print(classification_report(y_test, ann_pred))
print("Confusion Matrix:")
print(confusion_matrix(y_test, ann_pred))
```

⮞
```
    === Tuned Random Forest ===
    Accuracy: 0.4864864864864865
    Classification Report:
                  precision    recall  f1-score   support

               0       0.33      0.36      0.34        14
               1       0.59      0.57      0.58        23

        accuracy                           0.49        37
       macro avg       0.46      0.46      0.46        37
```

```
     weighted avg        0.49      0.49       0.49         37

   Confusion Matrix:
   [[ 5  9]
    [10 13]]

   === XGBoost ===
   Accuracy: 0.4594594594594595
   Classification Report:
               precision    recall  f1-score   support

            0        0.35      0.50       0.41         14
            1        0.59      0.43       0.50         23

     accuracy                            0.46         37
    macro avg        0.47      0.47       0.46         37
   weighted avg      0.50      0.46       0.47         37

   Confusion Matrix:
   [[ 7  7]
    [13 10]]

   === LightGBM ===
   Accuracy: 0.5675675675675675
   Classification Report:
               precision    recall  f1-score   support

            0        0.42      0.36       0.38         14
            1        0.64      0.70       0.67         23

     accuracy                            0.57         37
    macro avg        0.53      0.53       0.53         37
   weighted avg      0.56      0.57       0.56         37

   Confusion Matrix:
   [[ 5  9]
    [ 7 16]]

   === Voting Classifier (Ensemble) ===
   Accuracy: 0.4594594594594595
   Classification Report:
               precision    recall  f1-score   support

            0        0.31      0.36       0.33         14
            1        0.57      0.52       0.55         23
```

```python
from sklearn.metrics import roc_curve, auc
import matplotlib.pyplot as plt

# === Function to plot ROC ===
def plot_roc_curves(models, model_names, X_test, y_test):
    plt.figure(figsize=(10, 8))

    for model, name in zip(models, model_names):
        if name == "ANN":
            y_scores = model.predict(X_test).flatten()
        else:
            try:
                y_scores = model.predict_proba(X_test)[:, 1]
            except:
                y_scores = model.decision_function(X_test)

        fpr, tpr, _ = roc_curve(y_test, y_scores)
        roc_auc = auc(fpr, tpr)
        plt.plot(fpr, tpr, label=f'{name} (AUC = {roc_auc:.2f})')

    # Plot base line
    plt.plot([0, 1], [0, 1], 'k--', label='Random Classifier')
    plt.xlabel('False Positive Rate')
    plt.ylabel('True Positive Rate')
    plt.title('ROC Curves for All Models')
    plt.legend(loc='lower right')
    plt.grid(True)
    plt.tight_layout()
    plt.show()

# === Plot all ROC curves ===
plot_roc_curves(
    models=[best_rf, xgb_model, lgbm_model, voting_model, model],
    model_names=["Tuned RF", "XGBoost", "LightGBM", "Voting", "ANN"],
    X_test=X_test_top,
    y_test=y_test
)
```
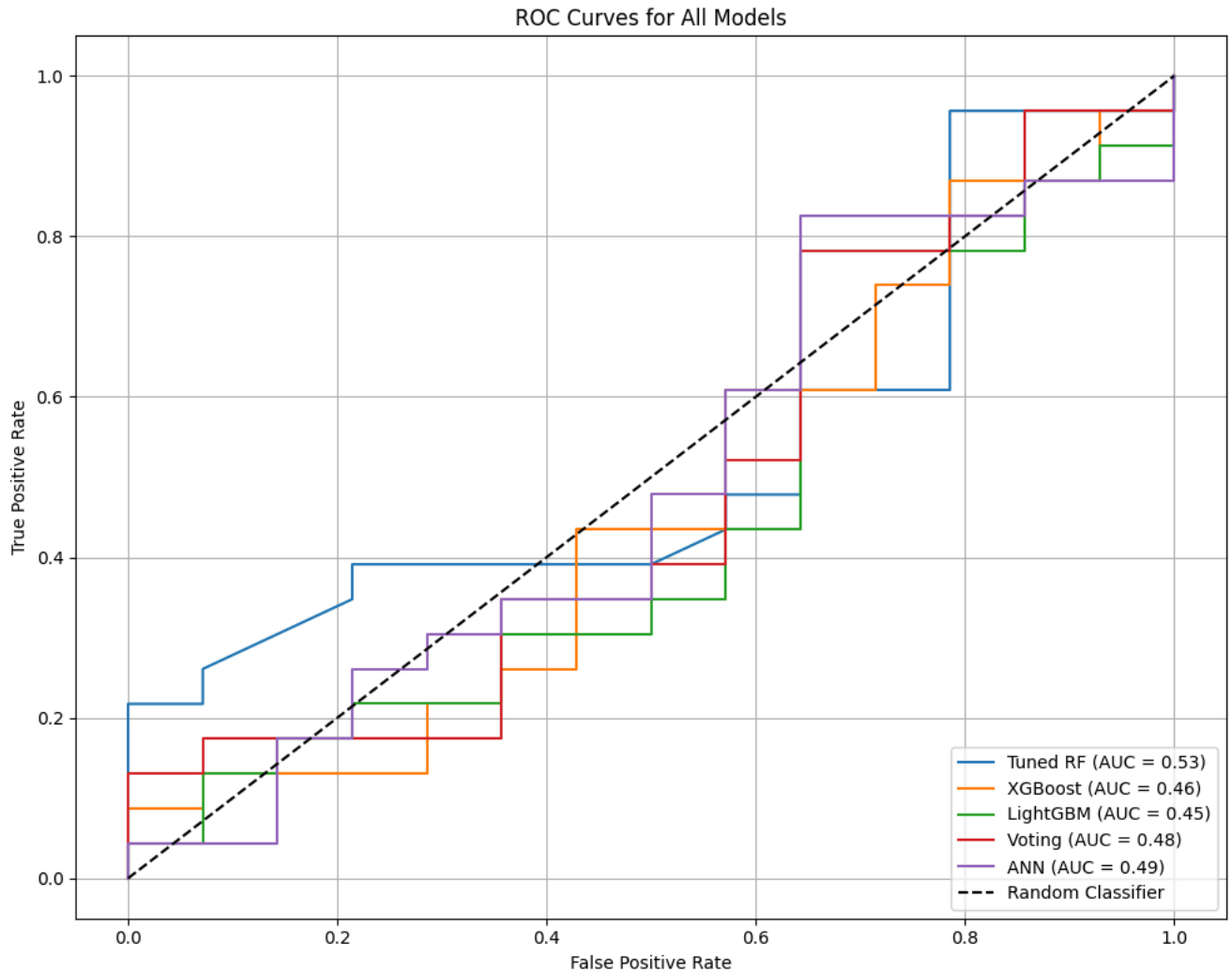
ROC Curves for All Models

```python
from sklearn.metrics import roc_curve, auc
import matplotlib.pyplot as plt

# === Function to plot ROC ===
def plot_roc_curves(models, model_names, X_test, y_test):
    plt.figure(figsize=(10, 8))

    for model, name in zip(models, model_names):
        if name == "ANN":
            y_scores = model.predict(X_test).flatten()
        else:
            try:
                y_scores = model.predict_proba(X_test)[:, 1]
            except:
                y_scores = model.decision_function(X_test)

        fpr, tpr, _ = roc_curve(y_test, y_scores)
        roc_auc = auc(fpr, tpr)
        plt.plot(fpr, tpr, label=f'{name} (AUC = {roc_auc:.3f})')

    plt.plot([0, 1], [0, 1], 'k--', label='Random Classifier')
    plt.xlabel('False Positive Rate')
    plt.ylabel('True Positive Rate')
    plt.title('ROC Curves for All Models (No Voting)')
    plt.legend(loc='lower right')
    plt.grid(True)
    plt.tight_layout()
    plt.show()

# === Use the function for your trained models ===
plot_roc_curves(
    models=[best_rf, xgb_model, lgbm_model, model],
    model_names=["Tuned RF", "XGBoost", "LightGBM", "ANN"],
    X_test=X_test_top,
```