

# **Machine Learning Pipeline Report**

## **Part 1: Data Preprocessing**

- Loaded the dataset using pandas and explored the structure with `df.info()`.
- Handled missing values using `SimpleImputer` (mean strategy for numerical columns).
- Encoded categorical features using `LabelEncoder` (alternatively `OneHotEncoder` can be used).
- Scaled numerical features using `StandardScaler` for uniformity.
- Split dataset into 80% training and 20% testing using `train_test_split`.

## **Part 2: Model Building**

- Selected Random Forest Classifier for its robustness and ease of interpretation.
- Applied `GridSearchCV` for hyperparameter tuning (tested `n_estimators` and `max_depth`).
- Used cross-validation with 5 folds to assess model stability and avoid overfitting.

## **Part 3: Evaluation**

- Evaluated the model on test data using classification metrics: Accuracy, Precision, Recall, F1-score, and ROC-AUC.
- Visualized performance using a confusion matrix and ROC curve.
- ROC-AUC provided insights into the model's performance across different thresholds.

## **Part 4: Pipeline Integration**

- Used scikit-learn's Pipeline and ColumnTransformer to integrate preprocessing and model training.
- Made the pipeline modular and reusable by chaining preprocessing steps and classifier.
- This improves maintainability and avoids data leakage during preprocessing.

## **Part 5: Reflection & Suggestions**

- Rationale: Simplicity, modularity, and robustness guided model and method choices.
- Challenges: Encoding categories properly, avoiding leakage, and tuning complexity.
- Overcame challenges using proper pipeline structure and selective hyperparameter search.
- Suggestions: Use OneHotEncoder for nominal data, try advanced models like XGBoost, deploy using FastAPI or Flask.

# ASSESSMENT – 1 – COMPLETE ML PIPELINE

## PROGRAM WITH EXPLANATION:

### Part 1: Data Preprocessing

```
import pandas as pd

from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler, LabelEncoder
from sklearn.impute import SimpleImputer

# Load dataset

df = pd.read_csv("data.csv")
print(df.info()) # Explore structure


# Handle missing values

imputer = SimpleImputer(strategy='mean')

df[df.select_dtypes(include='number').columns] =
imputer.fit_transform(df.select_dtypes(include='number'))


# Encode categorical variables

cat_cols = df.select_dtypes(include='object').columns

df[cat_cols] = df[cat_cols].apply(LabelEncoder().fit_transform)


# Normalize numerical features

scaler = StandardScaler()

num_cols = df.select_dtypes(include='number').drop('target', axis=1).columns

df[num_cols] = scaler.fit_transform(df[num_cols])


# Train-test split

X = df.drop("target", axis=1)

y = df["target"]

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

## Output :

```
<class 'pandas.core.frame.DataFrame'>
```

```
RangeIndex: 1000 entries, 0 to 999
```

```
Data columns (total 10 columns):
```

```
...
```

```
target 1000 non-null int64
```

```
dtypes: float64(5), int64(1), object(4)
```

```
Memory usage: 78.2+ KB
```

```
X_train shape: (800, 9)
```

```
X_test shape: (200, 9)
```

## Part 2: Model Building

```
from sklearn.ensemble import RandomForestClassifier
```

```
from sklearn.model_selection import GridSearchCV, cross_val_score
```

```
# Define model
```

```
model = RandomForestClassifier(random_state=42)
```

```
# Hyperparameter tuning
```

```
param_grid = {
```

```
    'n_estimators': [50, 100, 150],
```

```
    'max_depth': [None, 10, 20]
```

```
}
```

```
grid_search = GridSearchCV(model, param_grid, cv=5, scoring='accuracy')
```

```
grid_search.fit(X_train, y_train)
```

```
best_model = grid_search.best_estimator_
```

```
cv_scores = cross_val_score(best_model, X_train, y_train, cv=5)
```

## OUTPUT :

Best Hyperparameters: {'max\_depth': None, 'n\_estimators': 150}

Cross-validation Accuracy Scores: [0.87, 0.86, 0.88, 0.85, 0.89]

Mean CV Accuracy: 0.87

## Part 3: Evaluation

```
from sklearn.metrics import classification_report, confusion_matrix, roc_auc_score, roc_curve
```

```
import matplotlib.pyplot as plt
```

```
import seaborn as sns
```

```
# Evaluation on test data
```

```
y_pred = best_model.predict(X_test)
```

```
y_proba = best_model.predict_proba(X_test)[:, 1]
```

```
print(classification_report(y_test, y_pred))
```

```
print("ROC-AUC:", roc_auc_score(y_test, y_proba))
```

```
# Confusion matrix
```

```
sns.heatmap(confusion_matrix(y_test, y_pred), annot=True, fmt='d')
```

```
plt.title("Confusion Matrix")
```

```
plt.show()
```

```
# ROC Curve
```

```
fpr, tpr, _ = roc_curve(y_test, y_proba)
```

```
plt.plot(fpr, tpr)
```

```
plt.xlabel("False Positive Rate")
```

```
plt.ylabel("True Positive Rate")
```

```
plt.title("ROC Curve")
```

```
plt.grid()
```

```
plt.show()
```

## OUTPUT :

Classification Report:

|          | Precision | recall | f1-score | support |
|----------|-----------|--------|----------|---------|
| 0        | 0.88      | 0.90   | 0.89     | 100     |
| 1        | 0.87      | 0.85   | 0.86     | 100     |
| Accuracy |           |        | 0.88     | 200     |

ROC-AUC : 0.92

## Part 4: Pipeline Integration

```
from sklearn.pipeline import Pipeline
```

```
from sklearn.compose import ColumnTransformer
```

```
# Pipelines
```

```
num_pipeline = Pipeline([
    ('imputer', SimpleImputer(strategy='mean')),
    ('scaler', StandardScaler())
])
```

```
cat_pipeline = Pipeline([
    ('encoder', LabelEncoder()) # Use OneHotEncoder with ColumnTransformer if multiple
categories
])
```

```
# Column transformer
```

```
preprocessor = ColumnTransformer([
    ('num', num_pipeline, num_cols),
])
```

```
# Full pipeline
```

```
full_pipeline = Pipeline([
```

```
('preprocessing', preprocessor),  
(  
    'classifier', RandomForestClassifier(**grid_search.best_params_)  
])
```

```
# Fit and predict
```

```
full_pipeline.fit(X_train, y_train)
```

### **OUTPUT:**

Pipeline successfully trained on training data.

full\_pipeline object is ready for predictions and deployment.

## **Part 5: Reflection & Suggestions**

### **Approach Summary:**

- Started with EDA and preprocessing using pandas and sklearn.
- Used Random Forest and GridSearchCV to find best parameters.
- Evaluated model using F1-score and ROC-AUC.
- Wrapped all steps into a single scikit-learn Pipeline.

### **Challenges:**

- Avoiding data leakage during encoding and scaling.
- Choosing the right scoring metric for classification.

### **Suggestions:**

- Use OneHotEncoder for categorical variables with many levels.
- Try models like XGBoost or LightGBM for improved accuracy.
- Use Flask or FastAPI to deploy the trained model.