# MACHINE LEARNING PROJECT DOCUMENTATION

## MODEL REFINEMENT

### 1. OVERVIEW

The model refinement phase is dedicated to improving the accuracy, efficiency, and reliability of the machine learning model after initial exploration. Refinement ensures that the model is well-generalized, avoids overfitting, and achieves optimal performance before testing and deployment.

### 2. MODEL EVALUATION

The initial evaluation indicated that the baseline model provided moderate accuracy but exhibited weaknesses such as high variance and limited ability to generalize. Key issues identified were:
- Overfitting on training data.
- Poor recall on minority classes.
- Sensitivity to noisy features.
- Metrics such as accuracy, F1-score, and ROC-AUC revealed clear opportunities for refinement.

### 3. REFINEMENT TECHNIQUES

- Algorithm comparison: Tested alternative models (Random Forest, Gradient Boosting, Logistic Regression).
- Ensemble methods: Combined predictions from multiple classifiers to balance bias and variance.
- Regularization: Applied L1/L2 regularization to improve generalization.
- Feature engineering refinement: Added interaction features and reduced dimensionality with PCA.

### 4. HYPERPARAMETER TUNING

Extensive tuning was carried out using Grid Search and Random Search. Examples:
- Random Forest $\rightarrow$ n_estimators=300, max_depth=15.
- Gradient Boosting $\rightarrow$ learning_rate=0.05, max_depth=8.
- Logistic Regression $\rightarrow$ C=0.5, penalty="l2".
- These adjustments led to significant improvements in recall and reduced error rates across classes.

### 5. CROSS-VALIDATION

Cross-validation was enhanced by moving from k-fold CV to stratified k-fold CV, ensuring balanced representation of all classes in each fold. This produced more reliable and stable performance estimates.

## 6. FEATURE SELECTION

Feature selection techniques, such as Recursive Feature Elimination (RFE) and importance ranking, were used to identify the most informative features. Irrelevant and redundant variables were dropped, reducing training complexity and boosting performance consistency.

# TEST SUBMISSION

## 1) OVERVIEW

The test submission phase assessed the refined model on a previously unseen dataset to verify its performance and readiness for deployment.

## 2) DATA PREPARATION FOR TESTING

The test dataset was cleaned, normalized, and transformed in the same manner as the training dataset. Care was taken to avoid data leakage and ensure consistent preprocessing.

## 3) MODEL APPLICATION

```
# Apply trained model to test dataset
y_test_pred = model.predict(X_test)
y_test_prob = model.predict_proba(X_test)[:,1]
```

## 4) TEST METRICS

- ➤ Accuracy = 0.87
- ➤ Precision = 0.81
- ➤ Recall = 0.84
- ➤ F1-score = 0.82
- ➤ ROC-AUC = 0.90

These results were consistent with training and validation metrics, indicating strong generalization.

## 5) MODEL DEPLOYMENT

The refined model was packaged and integrated into a simple web-based dashboard for demonstration. It was also exported using Pickle/Joblib for API deployment.

## 6) CODE IMPLEMENTATION

```
def apply_model(model, X):
    """
        Apply trained model to new input data
    """
    return model.predict(X), model.predict_proba(X)
```

## 7) CONCLUSION

The refinement phase enhanced the baseline model significantly by tuning hyperparameters, improving feature selection, and applying ensemble strategies. Testing confirmed the model's ability to generalize with balanced accuracy and recall. The final model is optimized for both deployment and real-world use cases.

## 8) REFERENCES

*Géron, A. (2019). Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow. O'Reilly Media.*
*Pedregosa, F., et al. (2011). Scikit-learn: Machine Learning in Python. Journal of Machine Learning Research.*
*Libraries: scikit-learn, xgboost, matplotlib, pandas.*