# Contents

# 1  Introduction

Customer churn refers to the phenomenon in which customers discontinue a service. In the highly competitive telecom industry, predicting churn in advance enables service providers to take proactive retention actions. Machine learning provides an effective framework for learning churn patterns from historical customer behavior data.

This project focuses on the academic study, experimentation, and analysis of an existing telecom customer churn prediction system. The project was explored by modifying models, preprocessing pipelines, and evaluation strategies to observe performance changes and draw conclusions aligned with machine learning concepts covered in the course. The dataset and problem domain are credited to **Telecom Egypt (WE)**.

# 2  Learning Framework (Task–Experience–Performance)

The supervised learning formulation of the problem is defined as follows:

- **Task (T):** Binary classification to predict whether a customer will churn.

- **Experience (E):** Historical labeled telecom customer data.

- **Performance (P):** Accuracy, Precision, Recall, F1-score, and ROC-AUC on a held-out test set.

This formulation anchors the project in formal supervised learning theory.

# 3 Implementation Environment and Libraries

The project is implemented in Python using standard data science and machine learning libraries, primarily from the `scikit-learn` ecosystem.

```python
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
from IPython.display import display

from sklearn.model_selection import train_test_split, cross_val_score
from sklearn.compose import ColumnTransformer
from sklearn.preprocessing import OneHotEncoder, MinMaxScaler
from sklearn.pipeline import Pipeline
from sklearn.impute import SimpleImputer
from sklearn.metrics import (
    confusion_matrix,
    classification_report,
    ConfusionMatrixDisplay,
    accuracy_score,
    precision_score,
    recall_score,
    f1_score,
    roc_auc_score,
)
from sklearn.linear_model import LogisticRegression
from sklearn.tree import DecisionTreeClassifier
from sklearn.neighbors import KNeighborsClassifier
from sklearn.svm import SVC
from sklearn.neural_network import MLPClassifier
from sklearn.ensemble import RandomForestClassifier
from sklearn.cluster import KMeans

import joblib

plt.style.use("seaborn-v0_8")
pd.set_option("display.max_columns", None)
RANDOM_STATE = 42
```

Figure 1: Imported libraries and project configuration.

# 4 Data Loading and Inspection

The dataset is loaded from a CSV file, with date fields parsed explicitly to support time-aware features. Initial inspection verifies the dataset structure and feature types.

```python
Load and Inspect Data

# Load data
df = pd.read_csv("telecom_churn_full.csv", parse_dates=["signup_date","last_activity_ts"])

                                                                                        Python

# 1) Basic info
print("Shape:", df.shape)
print("Data Types:")
print(df.dtypes)

                                                                                        Python
```

Figure 2: Dataset loading and basic inspection.

# 5 Preprocessing Pipeline

Separate preprocessing pipelines are defined for numerical and categorical features using a `ColumnTransformer`.

## 5.1 Numerical Features

Numerical features are imputed using the median and normalized using Min–Max scaling.

```python
numeric_transformer = Pipeline(steps=[
    ("imputer", SimpleImputer(strategy="median")),
    ("scaler", MinMaxScaler())
])
```

```python
categorical_transformer = Pipeline(steps=[
    ("imputer", SimpleImputer(strategy="most_frequent")),
    ("onehot", OneHotEncoder(handle_unknown="ignore"))
])
```

```python
preprocess = ColumnTransformer(
    transformers=[
        ("num", numeric_transformer, numeric_features),
        ("cat", categorical_transformer, categorical_features)
    ]
)
```

Figure 3: Numerical preprocessing pipeline.

## 5.2 Categorical Features

Categorical features are imputed using the most frequent value and encoded using one-hot encoding.

```python
models = {
    "log_reg_baseline": LogisticRegression(max_iter=800, class_weight="balanced", n_jobs=-1, solver="lbfgs"),
    "decision_tree": DecisionTreeClassifier(
        max_depth=8, min_samples_leaf=10, class_weight="balanced", random_state=RANDOM_STATE
    ),
    "random_forest": RandomForestClassifier(
        n_estimators=350, random_state=RANDOM_STATE, class_weight="balanced_subsample"
    ),
    "knn": KNeighborsClassifier(n_neighbors=7, weights="distance"),
    "svm_rbf": SVC(
        probability=True,
        kernel="rbf",
        class_weight="balanced",
        C=1.5,
        gamma="scale",
        random_state=RANDOM_STATE,
    ),
    "mlp": MLPClassifier(
        hidden_layer_sizes=(128, 64),
        activation="relu",
        solver="adam",
        alpha=1e-4,
        learning_rate="adaptive",
        max_iter=400,
        early_stopping=True,
        random_state=RANDOM_STATE,
    ),
}

display_name_map = {
    "log_reg_baseline": "Logistic Regression (Baseline)",
    "decision_tree": "Decision Tree Classifier",
    "random_forest": "Random Forest Classifier",
    "knn": "K-Nearest Neighbors (KNN)",
    "svm_rbf": "Support Vector Machine (RBF Kernel)",
    "mlp": "Neural Network (Multilayer Perceptron)",
}
```

Figure 4: Categorical preprocessing pipeline.

## 5.3    Combined Preprocessing

```python
metrics_records = []
pipelines = {}

for name, model in models.items():
    clf = Pipeline(steps=[("preprocess", preprocess), ("model", model)])
    display_name = display_name_map.get(name, name)
    # 5-fold CV Balanced Accuracy
    bal_acc = cross_val_score(clf, X_train, y_train, cv=5, scoring="balanced_accuracy")
    print(f"{display_name} 5-fold CV Balanced Accuracy: {bal_acc.mean():.3f} +/- {bal_acc.std():.3f}")
    # Fit full train
    clf.fit(X_train, y_train)
    pipelines[name] = clf

    # Held-out test predictions
    y_pred = clf.predict(X_test)
    y_proba = clf.predict_proba(X_test)[:, 1]

    metrics_records.append({
        "model_key": name,
        "model": display_name,
        "accuracy": accuracy_score(y_test, y_pred),
        "precision": precision_score(y_test, y_pred, zero_division=0),
        "recall": recall_score(y_test, y_pred, zero_division=0),
        "f1": f1_score(y_test, y_pred, zero_division=0),
        "roc_auc": roc_auc_score(y_test, y_proba),
    })

    print(classification_report(y_test, y_pred, digits=3))

    cm = confusion_matrix(y_test, y_pred)
    fig, ax = plt.subplots()
    disp = ConfusionMatrixDisplay(confusion_matrix=cm, display_labels=clf.classes_)
    disp.plot(cmap="Blues", ax=ax)
    ax.set_title(f"Confusion Matrix - {display_name}")
    plt.show()

metrics_df = pd.DataFrame(metrics_records).set_index("model").sort_values("roc_auc", ascending=False)
display(metrics_df)
```

Figure 5: ColumnTransformer combining numerical and categorical pipelines.

Min–Max normalization is required for distance-based and margin-based models, while tree-based models remain scale-invariant.

# 6 Supervised Learning Models

Multiple supervised classifiers were studied to reflect course coverage.

```python
# Logistic Regression churn probabilities (sample)
log_reg_probs = pipelines["log_reg_baseline"].predict_proba(X_test)[:, 1]
prob_sample = pd.DataFrame({
    "true_churn": y_test.reset_index(drop=True),
    "log_reg_probability": log_reg_probs
})
display(prob_sample.head(10))
```

Figure 6: Defined supervised learning models and display names.

The models include Logistic Regression (baseline), Decision Tree, Random Forest, K-Nearest Neighbors, Support Vector Machine (RBF kernel), and a Neural Network (MLP).

# 7 Model Training and Evaluation

Each model is trained within a unified pipeline that includes preprocessing. Five-fold cross-validation using balanced accuracy is performed, followed by evaluation on a held-out test set.

```python
Model Evaluation & Selection

    sorted_metrics = metrics_df.sort_values("roc_auc", ascending=False)
    best_model_display = sorted_metrics.index[0]
    best_model_key = sorted_metrics["model_key"].iloc[0]
    best_auc = sorted_metrics.iloc[0]["roc_auc"]
    print(f"Best model by ROC-AUC: {best_model_display} ({best_auc:.3f})")

    best_model = pipelines[best_model_key]

                                                                        Python

Best model by ROC-AUC: log_reg_baseline (0.735)


    # K-Means clustering for customer segmentation
    kmeans_pipeline = Pipeline(steps=[
        ("preprocess", preprocess),
        ("kmeans", KMeans(n_clusters=3, n_init=20, random_state=RANDOM_STATE)),
    ])

    cluster_labels = kmeans_pipeline.fit_predict(X)

    cluster_df = df_new.copy()
    cluster_df["cluster"] = cluster_labels
    cluster_churn = cluster_df.groupby("cluster")["churn"].agg(["count", "mean"])
    cluster_churn = cluster_churn.rename(columns={"mean": "churn_rate"})

    display(cluster_churn)

                                                                        Python
```

Figure 7: Model training, cross-validation, and evaluation pipeline.

# 8 Model Selection

Models are ranked based on ROC-AUC, which serves as the primary selection metric due to its robustness for probabilistic classification and imbalanced datasets.

```python
# Persist pipelines and evaluation splits for deployment
joblib.dump(pipelines, "all_churn_pipelines.pkl")
X_test.to_csv("X_test_churn.csv", index=False)
y_test.to_csv("y_test_churn.csv", index=False)
```

Figure 8: Model ranking and best model selection based on ROC-AUC.

# 9 Sample Predictions and Model Output Interpretation

Sample predictions from the best-performing model are displayed to illustrate how customer attributes are translated into churn probabilities and final predictions.



Figure 9: Sample predictions from the best-performing model (part 1).



Figure 10: Sample predictions from the best-performing model (part 2).

# 10 Unsupervised Learning: K-Means Clustering

K-Means clustering is applied for customer segmentation using the same preprocessing pipeline. Customers are grouped into three clusters, and churn rates are analyzed per cluster.

```python
models = {
    "log_reg_baseline": LogisticRegression(max_iter=800, class_weight="balanced", n_jobs=-1, solver="lbfgs"),
    "decision_tree": DecisionTreeClassifier(
        max_depth=8, min_samples_leaf=10, class_weight="balanced", random_state=RANDOM_STATE
    ),
    "random_forest": RandomForestClassifier(
        n_estimators=350, random_state=RANDOM_STATE, class_weight="balanced_subsample"
    ),
    "knn": KNeighborsClassifier(n_neighbors=7, weights="distance"),
    "svm_rbf": SVC(
        probability=True,
        kernel="rbf",
        class_weight="balanced",
        C=1.5,
        gamma="scale",
        random_state=RANDOM_STATE,
    ),
    "mlp": MLPClassifier(
        hidden_layer_sizes=(128, 64),
        activation="relu",
        solver="adam",
        alpha=1e-4,
        learning_rate="adaptive",
        max_iter=400,
        early_stopping=True,
        random_state=RANDOM_STATE,
    ),
}

display_name_map = {
    "log_reg_baseline": "Logistic Regression (Baseline)",
    "decision_tree": "Decision Tree Classifier",
    "random_forest": "Random Forest Classifier",
    "knn": "K-Nearest Neighbors (KNN)",
    "svm_rbf": "Support Vector Machine (RBF Kernel)",
    "mlp": "Neural Network (Multilayer Perceptron)",
}
```

Figure 11: K-Means clustering and churn rate analysis.

# 11 Deployment Preparation

All trained pipelines and evaluation splits are persisted using `joblib` to support deployment in a Flask web application.

```python
# Persist pipelines and evaluation splits for deployment
joblib.dump(pipelines, "all_churn_pipelines.pkl")
X_test.to_csv("X_test_churn.csv", index=False)
y_test.to_csv("y_test_churn.csv", index=False)
```

Figure 12: Persisting trained pipelines and evaluation data for deployment.

# 12 Academic Integrity and Project Scope

This project does not claim original authorship of the initial system. It represents an academic study of an existing machine learning project, where experimentation and analysis were performed to deepen understanding of course concepts.

# 13    Limitations and Future Improvements

Potential improvements aligned with course content include:

- Introducing a majority-class baseline.

- Using cross-validation for model selection.

- Performing hyperparameter sensitivity analysis.

- Tuning decision thresholds using ROC or precision–recall trade-offs.

- Justifying the number of clusters using silhouette or elbow methods.

# 14    Conclusion

This project demonstrates the application of supervised and unsupervised machine learning techniques to a real-world telecom churn problem. By studying and modifying an existing system, the project reinforces theoretical concepts while delivering a deployable machine learning solution.