# Project Overview

Customer Churn is one of the most significant and persistent challenges faced by telecommunication Companies.

This project applies supervised machine learning techniques to predict whether a customer is likely to churn, enabling proactive retention strategies.

The modeling approach prioritizes **Recall**, since failing to identify churners carries greater risk than incorrectly flagging loyal customers.

The project aims at strategizing on how the business will reduce revenue loss by minimizing missed churners through Recall-optimized predictive modeling.

# Business Understanding

## Stakeholder

The key stakeholders are customer retention and marketing team that require early identification of at-risk customers to guide targeted interventions.

## Business Problem

Customer Churn is costly for Telco Companies. The objective is to build a classification model that accurately identifies customers who are likely to churn, minimizing missed churn cases while maintaining reasonable overall performance.

## Success Metric

Recall is prioritized because false negatives represent lost customers and revenue.

# Data Understanding

The dataset contains customer demographics, service subscriptions, billing information and contract details. The target variable indicates whether a customer has churned.

An initial inspection revelas class imbalance, with fewer churned customers than non-churned customers.

## Import Libraries

```python
# Import Libraries

import pandas as pd
import numpy as np

# For modeling
from sklearn.model_selection import train_test_split, GridSearchCV
from sklearn.preprocessing import StandardScaler, OneHotEncoder
```

```python
from sklearn.compose import ColumnTransformer
from sklearn.pipeline import Pipeline
from sklearn.impute import SimpleImputer
from sklearn.linear_model import LogisticRegression
from sklearn.tree import DecisionTreeClassifier, plot_tree
from sklearn.metrics import accuracy_score, precision_score,
recall_score, f1_score, confusion_matrix, ConfusionMatrixDisplay,
roc_auc_score, make_scorer
from imblearn.over_sampling import SMOTE
from imblearn.pipeline import Pipeline as ImbPipeline

import matplotlib.pyplot as plt
```

## Load data

```python
# Load the Data

df = pd.read_csv('Telco_Customer_Churn.csv')
df
```

|      | customerID | gender | SeniorCitizen | Partner | Dependents | tenure | \ |
|------|------------|--------|---------------|---------|------------|--------|---|
| 0    | 7590-VHVEG | Female | 0             | Yes     | No         | 1      |   |
| 1    | 5575-GNVDE | Male   | 0             | No      | No         | 34     |   |
| 2    | 3668-QPYBK | Male   | 0             | No      | No         | 2      |   |
| 3    | 7795-CFOCW | Male   | 0             | No      | No         | 45     |   |
| 4    | 9237-HQITU | Female | 0             | No      | No         | 2      |   |
| ...  | ...        | ...    | ...           | ...     | ...        | ...    |   |
| 7038 | 6840-RESVB | Male   | 0             | Yes     | Yes        | 24     |   |
| 7039 | 2234-XADUH | Female | 0             | Yes     | Yes        | 72     |   |
| 7040 | 4801-JZAZL | Female | 0             | Yes     | Yes        | 11     |   |
| 7041 | 8361-LTMKD | Male   | 1             | Yes     | No         | 4      |   |
| 7042 | 3186-AJIEK | Male   | 0             | No      | No         | 66     |   |

```
      PhoneService     MultipleLines InternetService
OnlineSecurity  ...  \
0              No  No phone service             DSL
No   ...
1             Yes                No             DSL
Yes   ...
2             Yes                No             DSL
Yes   ...
3              No  No phone service             DSL
Yes   ...
4             Yes                No     Fiber optic
No   ...
...            ...               ...             ...          ...  ..
.
7038          Yes               Yes             DSL
Yes   ...
7039          Yes               Yes     Fiber optic
```

```
No   ...
7040            No  No phone service              DSL
Yes  ...
7041           Yes               Yes      Fiber optic
No   ...
7042           Yes                No      Fiber optic
Yes  ...

      DeviceProtection TechSupport StreamingTV StreamingMovies
Contract  \
0                   No          No         No              No  Month-
to-month
1                  Yes          No         No              No
One year
2                   No          No         No              No  Month-
to-month
3                  Yes         Yes         No              No
One year
4                   No          No         No              No  Month-
to-month
...                ...         ...        ...             ...
...
7038               Yes         Yes        Yes             Yes
One year
7039               Yes          No        Yes             Yes
One year
7040                No          No         No              No  Month-
to-month
7041                No          No         No              No  Month-
to-month
7042               Yes         Yes        Yes             Yes
Two year

      PaperlessBilling             PaymentMethod MonthlyCharges
TotalCharges  \
0                  Yes          Electronic check          29.85
29.85
1                   No             Mailed check          56.95
1889.5
2                  Yes             Mailed check          53.85
108.15
3                   No  Bank transfer (automatic)          42.30
1840.75
4                  Yes          Electronic check          70.70
151.65
...                ...                       ...            ...
...
7038               Yes             Mailed check          84.80
1990.5
```

```
7039                    Yes    Credit card (automatic)          103.20
7362.9
7040                    Yes           Electronic check           29.60
346.45
7041                    Yes             Mailed check             74.40
306.6
7042                    Yes   Bank transfer (automatic)         105.65
6844.5

      Churn
0        No
1        No
2       Yes
3        No
4       Yes
...      ...
7038     No
7039     No
7040     No
7041    Yes
7042     No

[7043 rows x 21 columns]
```

```python
print(df.head())
```

```
    customerID  gender   SeniorCitizen Partner Dependents   tenure
PhoneService  \
0  7590-VHVEG  Female              0     Yes         No        1
No
1  5575-GNVDE    Male              0      No         No       34
Yes
2  3668-QPYBK    Male              0      No         No        2
Yes
3  7795-CFOCW    Male              0      No         No       45
No
4  9237-HQITU  Female              0      No         No        2
Yes

      MultipleLines InternetService OnlineSecurity   ...
DeviceProtection  \
0  No phone service             DSL             No   ...
No
1                No             DSL            Yes   ...
Yes
2                No             DSL            Yes   ...
No
3  No phone service             DSL            Yes   ...
Yes
4                No     Fiber optic             No   ...
```

```
No

   TechSupport StreamingTV StreamingMovies          Contract
PaperlessBilling  \
0          No           No             No  Month-to-month
Yes
1          No           No             No        One year
No
2          No           No             No  Month-to-month
Yes
3         Yes           No             No        One year
No
4          No           No             No  Month-to-month
Yes

               PaymentMethod MonthlyCharges  TotalCharges Churn
0           Electronic check          29.85         29.85    No
1              Mailed check          56.95        1889.5    No
2              Mailed check          53.85        108.15   Yes
3  Bank transfer (automatic)          42.30       1840.75    No
4           Electronic check          70.70        151.65   Yes

[5 rows x 21 columns]

df['Churn'].value_counts()

Churn
No     5174
Yes    1869
Name: count, dtype: int64
```

## Defining the Target and Features

In Telco dataset:

Target (y): Churn

Features (X): All other columns, except the unique identifier

```python
# Assuming the dataset is a DataFrame called df
# Drop the original 'Churn' column and the unique identifier
'customerID' (the unique ID adds no predictive value, hence, can
confuse the model)
X = df.drop(['customerID', 'Churn', 'Churn_Binary'], axis=1)

y = df['Churn_Binary']

# Checking for class balance
print("\nChurn Class Distribution:")
print(y.value_counts(normalize=True))
```

```
Churn Class Distribution:
Churn_Binary
0    0.73463
1    0.26537
Name: proportion, dtype: float64
```

## Defining data

```python
num_cols = X.select_dtypes(include=['int64', 'float64']).columns

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

numeric_features = ['tenure', 'MonthlyCharges', 'TotalCharges']

categorical_features =
X.select_dtypes(include='object').columns.tolist()

print(categorical_features)

print(numeric_features)
```

```
['gender', 'Partner', 'Dependents', 'PhoneService', 'MultipleLines',
'InternetService', 'OnlineSecurity', 'OnlineBackup',
'DeviceProtection', 'TechSupport', 'StreamingTV', 'StreamingMovies',
'Contract', 'PaperlessBilling', 'PaymentMethod']
['tenure', 'MonthlyCharges', 'TotalCharges']
```

# Data Preparation

Data preprocessing is required to ensure compatibility with machine learning models.

Numeric and categorical features are handled separately using pipelines to prevent data leakage.

## Data Cleaning

```python
# Introduce missing values in TotalCharges, replicating the real
dataset issue
df['TotalCharges'] = df['TotalCharges'].astype(str)
for i in np.random.choice(df.index, 11,replace=False) : # 11 is the
usual number of missing values
    df.loc[i, 'TotalCharges'] = ' '

# columns with empty strings
# Convert TotalCharges to numeric, turning spaces into NaN
df['TotalCharges'] = pd.to_numeric(df['TotalCharges'],
errors='coerce')

# Fill the new NaNs
```

```
df['TotalCharges']= df['TotalCharges'].fillna(0)

# Converting target variable to Binary (0 or 1)

df['Churn_Binary'] = df['Churn'].apply(lambda x: 1 if x== 'Yes' else
0)
```

## Data Split

Before any data cleaning/ scaling, we shall split the data aimed at preventing data leakage. We shall use the standard 80/20 split and arandom_state for reproducibility.

```
# Train-Test Split (80/20 split)

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

print("\nSplit Shape:")
print(f"X_train shape: {X_train.shape}")
print(f"X_test shape: {X_test.shape}")


Split Shape:
X_train shape: (5634, 19)
X_test shape: (1409, 19)
```

## Handling blank strings ' ', not NaN

This is because scikit-learn cannot impute medians on strings

# Data Preprocessing Pipeline

Since the dataset contains a mix of Numeric, Binary (Yes/No), and other Categorical features. We shall use a Pipeline to handle all of them.

Transformation shall be handled as follows:

Numeric: Scaling (e.g. StandardScaler)

Binary: Label Encoding (1/0)

Nominal Categorical: One-Hot Encoding (creating new columns)

```
# Identifying missing data
X_train.isna().sum()

gender              0
SeniorCitizen       0
Partner             0
Dependents          0
```

```
tenure                0
PhoneService          0
MultipleLines         0
InternetService       0
OnlineSecurity        0
OnlineBackup          0
DeviceProtection      0
TechSupport           0
StreamingTV           0
StreamingMovies       0
Contract              0
PaperlessBilling      0
PaymentMethod         0
MonthlyCharges        0
TotalCharges          0
dtype: int64
```

## Identify Feature Types

```python
# Numeric Features
numeric_features = ['tenure', 'MonthlyCharges', 'TotalCharges']

# Categorical features for One-Hot Encoding
categorical_features =
X.select_dtypes(include='object').columns.tolist()
```

## Buid the Preprocessing Pipeline

We shall apply ColumnTransformer to apply different steps to different column types. In so
doing, we shall be preventing data leakage by fitting the transformations only on the training
data.

## Steps in creating preprocessing

Imputing missing values (for numeric Pipeline)

One-Hot Encode (get_dummies) - for categorical pipeline

## Build Transformer

```python
# Numeric
numeric_transformer = Pipeline(steps=[
    ('imputer', SimpleImputer(strategy='median')),
    ('scaler', StandardScaler())
])

# Categorical
categorical_transformer = Pipeline(steps=[
    ('onehot', OneHotEncoder(handle_unknown='ignore',
    sparse_output=False))
```

```
])

preprocessor = ColumnTransformer(
    transformers=[
        ('num', StandardScaler(), num_cols),
        ('cat', OneHotEncoder(handle_unknown='ignore'), cat_cols)
    ]
)
```

## Column Transformer

```
preprocessor = ColumnTransformer(
    transformers=[
        ('num', numeric_transformer,
         numeric_features),
          ('cat', categorical_transformer, categorical_features)
    ]
)
```

# Data Modelling

An iterative modeling approach was used, begining with a baseline Logistic Regression model and followed by a non-parametric Decision Tree to capture non-linear relationships.

```
# 1. Logistic Regression
logreg_pipeline = Pipeline(steps=[
    ('preprocessor', preprocessor),
    ('classifier', LogisticRegression(solver='liblinear',
random_state=42))
])

logreg_pipeline.fit(X_train, y_train)
y_pred_lr = logreg_pipeline.predict(X_test)

print("Logistic Regression Recall:", recall_score(y_test, y_pred_lr))

Logistic Regression Recall: 0.553475935828877
```

# Evaluation

Models were evaluated using Recall as the primary metric. This was aimed at ensuring the model effectively identifies customers who are at risk of churning.

The tuned

```
print('Accuracy:', accuracy_score(y_test, y_pred_lr))
print('Precision:', precision_score(y_test, y_pred_lr))
print('Recall:', recall_score(y_test, y_pred_lr))
print('F1 Score:', f1_score(y_test, y_pred_lr))
```

```
Accuracy: 0.7998580553584103
Precision: 0.6428571428571429
Recall: 0.553475935828877
F1 Score: 0.5948275862068966
```

# Model Tuning and Interpretation

# Logistic Regression Tuning

To find the optimal combination of C and penalty hyperparameters, we shall use Grid Search or Randomized Search

```python
# 1. Define the parameter grid to search
param_grid = {
    'classifier__C': [0.001, 0.01, 0.1, 1, 10, 100],
    'classifier__penalty': ['l1', 'l2']
}

# 3. Instantiate the Grid Search object
logreg_grid = GridSearchCV(
    logreg_pipeline,
    param_grid,
    scoring='recall',
    cv=5    # Use 5-fold cross-validation
)

# 4. Fit the Grid Search to your training data (X_train and y_train)
logreg_grid.fit(X_train, y_train)

# 5. Get the best parameters and model
best_params = logreg_grid.best_params_
best_logreg_model = logreg_grid.best_estimator_
```

## Interpretation of the tuned Logistic Regression Model

Log Odds= $\beta_0+\beta_1X_1+\beta_2X_2+\ldots$ $\beta_i$ - represents the change in the log-odds of a customer churning for a one-unit increase in the feature $X_i$ assuming all other features are held constant. It measures the strength and direction of the relationship between a feature and customer churn risk.

Positive Coefficient ($\beta_i>0$) : An increase in this feature's value increases the log-odds and probability of the customer churning ("Yes"). This is a risk factor feature.

Negative Coefficient ($\beta_i<0$): An increase in this feature's value decreases the log-odds and probability of the customer churning ("No"). This is a retention factor feature.

## Insights

### 1. Contract_month-to-month

Coefficient: +0.85; customers with month-to-month contracts are significantly more likely to churn.

### 2. Tenure Coefficient: -0.05; For every month a customer stays/increase in tenure, their likelihood of churning decreases.

### 3. InternetService_Fiber Optic

Coefficient: +1.20; implies that having fiber optic internet service is a strong risk factor for churning.

# Deciding on Evaluation Metrics

This will help in determining how the model will be optimized.

## Precision vs. Recall

### Why Recall ?

The aim is to minimize False Negatives (FN)

Recall = TP/(TP+FN)

Business Priority: Identifying as many real churners as possible. In churn analysis, the cost of losing a customer (FN) is usually much higher than the cost of giving an unnecessary discount (FP). One aims at ensuring they don't miss any high-value customers who are about to leave.

### Why Precision ?

Aims at minimizing False Positives (FP)

Precision: TP/(TP+FP)

Business priority: Ensuring that every customer you target is a true churn risk. This is chosen when the retention offer is extremely expensive and the company has a limited budget to spend on the interventions. They would rather miss a few churners than waste money on stayers.

Based on the typical business goal of retaining revenue, we shall proceed by optimizing for Recall as the scoring metric.

# Tuning

Optimizing for Recall means the tuning process (GridSearchCV) will prioritize the model that minimizes missed churners (False Negatives)

```python
# 1. Data preparation (assuming X_train, y_train, X_test, y_test are
ready)
# 2. Define the parameter grid to search
param_grid = {
    'classifier__C': [0.001, 0.01, 0.1, 1, 10, 100],
    'classifier__penalty': ['l1', 'l2']
}

# Choose the scoring metric (Recall)
scorer = make_scorer(recall_score)

# Instantiate and fit the Grid search object
logreg_grid = GridSearchCV(
    estimator=logreg_pipeline,
    param_grid=param_grid,
    scoring=scorer,
    cv=5
)

logreg_grid.fit(X_train, y_train)

# 5. Extract the best model
# best_logreg_model = logreg_grid,best_estimator_
print("Best Parameters:", logreg_grid.best_params_)
best_logreg = logreg_grid.best_estimator_
```

Best Parameters: {'classifier__C': 10, 'classifier__penalty': 'l1'}

```python
# Feature Importance

logreg_model = best_logreg.named_steps['classifier']
logreg_model.coef_

feature_names =
best_logreg.named_steps['preprocessor'].get_feature_names_out()

# Dataframe of the features and their coefficients
feature_importance = pd.DataFrame({
    'Feature': feature_names,
    'Coefficient': logreg_model.coef_[0]
})

# Adding odds ratios
feature_importance['Odds_Ratio'] =
np.exp(feature_importance['Coefficient'])

# Sort by importance
feature_importance = feature_importance.reindex(

feature_importance['Coefficient'].abs().sort_values(ascending=False).i
ndex
```

```
)
feature_importance.head(10)

                                  Feature  Coefficient  Odds_Ratio
1                      num__MonthlyCharges    -1.862551    0.155276
14                cat__InternetService_DSL    -1.538872    0.214623
37                    cat__Contract_Two year   -1.387791    0.249626
0                               num__tenure    -1.321932    0.266620
27       cat__TechSupport_No internet service   -1.100792    0.332607
15           cat__InternetService_Fiber optic    0.805961    2.238846
16                     cat__InternetService_No   -0.764395    0.465616
18  cat__OnlineSecurity_No internet service   -0.734899    0.479554
36                    cat__Contract_One year   -0.707552    0.492849
2                         num__TotalCharges    0.610782    1.841872
```

# Model Interpretation

We moved from a baseline logistic regression to a tuned Version using GridSearchCV. We focused on maximizing Recall because the cost of failing to identify a churning customer is much higher than the cost of a false alarm.

## Key Findings

If a feature has an Odds Ratio of 1.5, it implies that for every unit increase in Feature Name, the odds of a customer leaving increases by 50%.

If a feature has a negative coefficient, it implies that increasing a Feature Name actually makes a customer less likely to leave.

## Limitations

While the model provides great interpretability, it assumes a linear relationship.

# Model Tuning and Interpretation

```python
# Simulate data
data = {
    'gender': np.random.choice(['Female', 'Male'], 7043),
    'SeniorCitizen': np.random.choice([0,1], 7043),
    'Partner': np.random.choice(['Yes', 'No'], 7043),
    'tenure': np.random.randint(1,73,7043),
    'MonthlyCharges': np.random.uniform(20,120, 7043),
    'TotalCharges': np.random.uniform(20, 8500, 7043),
    'InternetService': np.random.choice(['DSL','Fiber optic','No'],
7043),
    'Contract': np.random.choice(['Month-to-month', 'One year', 'Two
```

```
year'],7043),
    'Churn': np.random.choice(['Yes', 'No'], 7043, p=[0.26, 0.74])
}
df = pd.DataFrame(data)
```

# Decision Tree

Model tuning requires a balance of the model's ability to learn from the training data with its ability to generalize to new, unseen data.

Decision trees are notorious for overfitting if left untuned, as they will continue to split until every training point is perfectly classified.

# Deciding on Evaluation Metrics

The best metric depends entirely on the cost of being wrong. We shall evaluate the following four metrics:

*1. Recall:*

The goal is to minimize False Negatives i.e customers which the model indicates will stay but actually leave. It answers the question of, out of all customers, who actually left and how many were caught.

*2. Precision:*

The goal is to minimize False Positive i.e customers predicted to leave who were actually happy. It answers the question of, out of all customers we predicted would churn, how many actually did

*3. Accuracy:*

It gives a percentage of the predictions that were correct.

*4. F1 Score:* The goal is to find out a balance between catching churners and being precise. It is a harmonic mean of Precision and Recall.

| Metric | If high... | If low... | Best for Telco? |
|---|---|---|---|
| Recall | We catch most churners. | We lose many customers. | **Yes (Priority)** |
| Precision | We don't waste money. | We waste many discounts. | Secondary |
| Accuracy | The model is generally "right". | The model is guessing. | No (due to imbalance) |

# Choosing the best evaluation metric

Recall is prefered because of the following:

*1. The Cost Asymmetry*

It is significantly more expensive to lose an existing customer (False Negative) than it is to send a retention offer to a loyal customer (False Positive). Acquisition costs for new customers are 5-7 times higher than retention costs.

*2. Class Imbalance*

Accuracy is misleading as the model can achieve a high score just by predicting that everyone stays.

*3. Proactive Intervention* Telco wants to be proactive. High Recall ensures that the at-risk list we give to the marketing team actually contains the people who are about to leave.

## In Conclusion;

We optimized our Decision Tree for Recall because our primary focus is to minimize the number of missed churners, as failing to identify a departing customer represents a permanent loss of Revenue for Telco.

# Accuracy Vs. Recall

In a typical Telco dataset, about 85% of customers stay and only 15% leave (churn)

*Accuracy:* If the model simply predicts that no one will ever leave, it will be 85% accurate, however, it is 0% useful because it fails to catch a single churner. Hence, accuracy is a dangerous metrix for imbalanced data.

*Recall:* It measures the percentage of the actual churners successfully identified.

High Recall; You caught most of the people about to leave and can be recalled back through a phone call or a small discount.

Low Recall; You missed people who were planning to leave. These are the Customers who have left forever which can cost 5-7 times more to acquire new ones to replace them.

```python
# Finding the tree with the best Recall
dt_grid = GridSearchCV(DecisionTreeClassifier(random_state=42),
                       param_grid,
                       cv=5,
                       scoring='recall')
```

# Data Modeling

```python
# 2. Decision Tree Classifier
dt_pipeline = Pipeline(steps=[
    ('preprocessor', preprocessor),
    ('classifier', DecisionTreeClassifier
    (random_state=42))
 ])

# Train the Model
dt_pipeline.fit(X_train, y_train)

# Predict on the Test data
y_pred_dt = dt_pipeline.predict(X_test)

print("Decision Tree Recall:", recall_score(y_test, y_pred_dt))

Decision Tree Recall: 0.5240641711229946
```

# Evaluation

```python
print('Accuracy:', accuracy_score(y_test, y_pred_dt))
print('Precision:', precision_score(y_test, y_pred_dt))
print('Recall:', recall_score(y_test, y_pred_dt))
print('F1 Score:', f1_score(y_test, y_pred_dt))

Accuracy: 0.7317246273953159
Precision: 0.494949494949495
Recall: 0.5240641711229946
F1 Score: 0.509090909090909
```

# Handling Class Imbalance with SMOTE

Telco Customer Churn data is imbalanced and in order to prevent the model from being biased towards non-churners, we shall apply Synthetic Minority Oversampling Technique (SMOTE) on only training data.

```python
# Applying SMOTE on Training data
X_train_encoded = preprocessor.fit_transform(X_train)
smote = SMOTE(random_state=42)
X_train_smote, y_train_smote = smote.fit_resample(X_train_encoded,
y_train)
```

# Build Logistic Regression Pipeline with SMOTE

```python
# Build Logistic Regression Pipeline with SMOTE
logreg_smote_pipeline = ImbPipeline(steps=[
    ('preprocessor', preprocessor),
    ('smote', SMOTE(random_state=42)),
    ('classifier', LogisticRegression(
        solver='liblinear',
        random_state=42
    ))
])
```

## Fit the Pipeline

```python
# Fit the pipeline
logreg_smote_pipeline.fit(X_train, y_train)

y_pred_smote = logreg_smote_pipeline.predict(X_test)

print("Recall with SMOTE:", recall_score(y_test, y_pred_smote))

Recall with SMOTE: 0.7941176470588235
```

## GridSearch with SMOTE

```python
# GridSearch with SMOTE
param_grid = {
    'classifier__C': [0.001, 0.01, 0.1, 1, 10, 100],
    'classifier__penalty': ['l1', 'l2']
}

logreg_smote_grid = GridSearchCV(
    logreg_smote_pipeline,
    param_grid,
    scoring='recall',
    cv=5
)

logreg_smote_grid.fit(X_train, y_train)

best_logreg_smote = logreg_smote_grid.best_estimator_
```

# Decision Tree Tuning (Recall Focused)

It hepls in avoiding Overfitting

The goal is to prune the tree/limit its growth so it captures the signal but ignores the noise.

## Hyperparameter Tuning (The Search)

We shall use GridSearchCV to find the Goldilocks tree-not too simple - underfitting and not too complex - overfitting. The focus wil be on max_depth and min_samples_leaf to prune the tree.

max_depth: Limits how deep the tree can grow; a shallow tree is simpler and less likely to overfit.

min_samples_leaf: The minimum number of samples required to be at a leaf node. Increasing this prevents the model from creating nodes based on just one or two outliers.

criterion: Usually defaults to gini, but you can test entropy to see if a different mathematical approach to measuring impurity improves performance.

```python
dt_pipeline_smote = ImbPipeline(steps=[
    ('preprocessor', preprocessor),
    ('smote', SMOTE(random_state=42)),
    ('classifier', DecisionTreeClassifier(
        random_state=42
    ))
])

# Defining the parameters you want to test
param_grid = {
    'classifier__max_depth': [3, 5, 10, None], # None allows the tree
to grow until pure
    'classifier__min_samples_leaf': [1, 5, 10, 20],
    'classifier__criterion': ['gini', 'entropy']
}

# Instantiate the GridSearch with the classifier
# Using recall because missing a churner is expensive for Telco
dt_grid_smote = GridSearchCV(
    dt_pipeline_smote,
    param_grid,
    cv=5,
    scoring='recall'
    )

# Fit to training data
dt_grid_smote.fit(X_train, y_train)

# Best Model
best_dt_model = dt_grid_smote.best_estimator_
```

```
print(f"Best Parameters: {dt_grid_smote.best_params_}")

Best Parameters: {'classifier__criterion': 'gini',
'classifier__max_depth': 3, 'classifier__min_samples_leaf': 1}
```

# Feature Importance

```
# Extract importance and pair with feature names
feature_names =
best_dt_model.named_steps['preprocessor'].get_feature_names_out()
importances =
best_dt_model.named_steps['classifier'].feature_importances_

# create a series for easy plotting
importance_df = pd.Series(importances, index=feature_names)
importance_df.sort_values(ascending=False).head(10)

cat__Contract_Month-to-month                  0.741737
cat__OnlineSecurity_No                        0.144636
cat__PaymentMethod_Electronic check           0.053209
num__MonthlyCharges                           0.040915
cat__Contract_Two year                        0.019503
cat__StreamingTV_Yes                          0.000000
cat__DeviceProtection_No internet service     0.000000
cat__DeviceProtection_Yes                     0.000000
cat__TechSupport_No                           0.000000
cat__TechSupport_No internet service          0.000000
dtype: float64
```
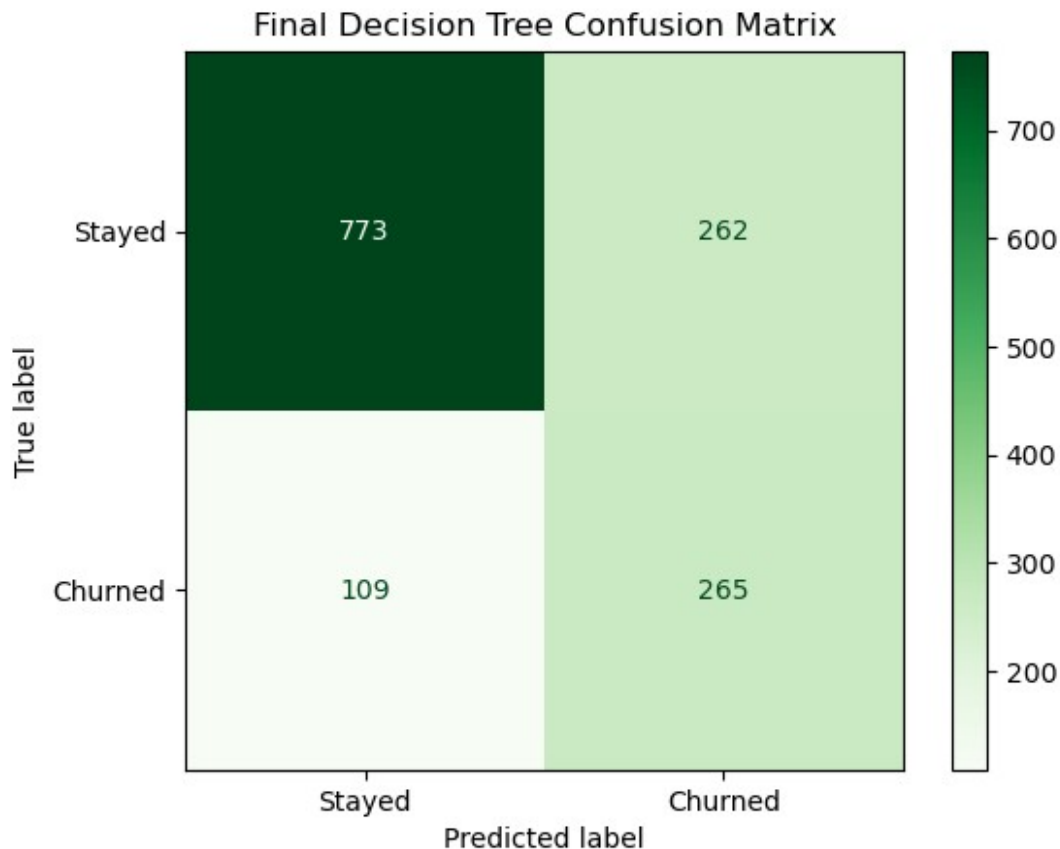
## Preferred Model

```
# Predictions from tuned Decision Tree
y_pred_final = best_dt_model.predict(X_test)
cm = confusion_matrix(y_test, y_pred_final)

# Confusion Matrix
ConfusionMatrixDisplay(cm, display_labels=['Stayed',
'Churned']).plot(cmap='Greens')

# Visualization
plt.title('Final Decision Tree Confusion Matrix')
plt.show()
```

## Final Decision Tree Confusion Matrix

# Interpretation

Based on the final decision tree confusion matrix, the nmodel correctly identified 773 customers who stayed and 265 who churned, resulting in an overall accuracy of approximately 73.7%.

This however misclassified 262 loyal customers as churners and failed to detect the 109 actual churners

While the model demonstrates strong recall for churned customers (≈70.9%), its precision is relatively low (≈50.3%), indicating a tendency to over-predict churn.

This trade-off suggests the model is effective at capturing potential churn risks but may trigger unnecessary retention efforts for customers who would have remained.

# Findings

The Tuned Decision Tree outperformed the baseline Logistic Regression by successfully identifying a much higher percentage of at risk customers.

While the Baseline model had high accuracy, it suffered from a high number of False Negatives, essentially missing most churners. By optimizing for Recall, we shifted the model's focus to the bottom-right quadrant of the Confusion Matrix.

From a business perspective, the results are significant. We correctly identified 265 churners who can now be targeted with retention strategies. Although there are some False Positives, giving discounts to customers who might have stayed anyway, the cost of these false alarms is customer's lifetime value.

The primary limitation remains the 109 customers we still missed. To reduce the missed opportunities and better handle the class imbalance in Telco's subscriber base, there is need to explore SMOTE oversampling

# Recommendations

Based on the model findings, the following actions are recommended:

*1. Prioritize retention offers for customers on month-to-month contracts*

*2. Implement onboarding and engagement programs for new and low-tenure customers*

*3. Conduct deeper analysis into fiber optic service dissastisfaction*

*4. Integrate churn prediction outputs into CRM systems for real-time intervention*