# Sub task1

## EDA Overview

### 1. Dataset Information

- **Number of Features:** 15 (12 independent, 1 target: 'churn')
- **Total Records:** 10,000 (before splitting).

### 2. Data Types and Characteristics

- **Categorical Features (e.g., gender, region, plan_type):** 4
- **Numerical Features (e.g., monthly_charges, tenure):** 8
- **Target Variable:** Binary (1 = Churned, 0 = Not Churned).

### 3. Data Quality Check

- **Missing Values:** Handled during data cleaning. Minor missing entries in `age` and `plan_type` were imputed appropriately.
- **Duplicates:** None found.
- **Outliers:** Detected in `monthly_charges` and `tenure`, verified as valid values based on domain understanding.

## Statistical Summary

| Feature | Mean | Median | Std Dev | Min | Max |
|---|---|---|---|---|---|
| Monthly Charges | 65.87 | 70.00 | 24.23 | 10.00 | 125.00 |
| Tenure (months) | 32.12 | 29.00 | 22.14 | 1.00 | 72.00 |

- Customers with higher monthly charges seem more likely to churn.
- Churners typically have shorter tenure.

## Visualizations and Insights

### 1. Target Variable Distribution

- **Churn Rate:** 22.3% (imbalanced dataset). SMOTE applied for balancing.
- Imbalance handled effectively during modeling.

- Strong positive correlation: `tenure` and `contract_type`.
- `monthly_charges` is positively associated with `churn`.

### *3. Categorical Distribution*

- Customers on month-to-month contracts churn more frequently than those with annual or bi-annual plans.

### *4. Price Sensitivity Analysis*

- Hypothesis validated: Higher price sensitivity correlates with increased churn likelihood. This was computed using metrics like price elasticity and variance in spending habits.

## Action Points Based on EDA

1. Addressed class imbalance using SMOTE.
2. Engineered features to enhance signal strength:
   a. Derived `price_sensitivity` score.
   b. Aggregated usage trends for better granularity.
3. Removed features with minimal predictive power, such as `customer_id`.

## Conclusion

The EDA highlighted key drivers of churn, identified data quality issues, and shaped the approach for feature engineering and modeling. We are now well-positioned to build a predictive model with high interpretability and actionable insights for business decisions.

# Sub task 2

- o **Verifying the Hypothesis of Price Sensitivity and Churn**
- o **Objective**:
- o This task aimed to test whether price sensitivity is correlated with customer churn and develop a metric to quantify price sensitivity.
- o **Hypothesis**:
- $H_0$: Price sensitivity is not correlated with churn.
- $H_1$: Price sensitivity is correlated with churn.
  - o **Steps Taken**:
  - o We defined **price sensitivity** as the variance in a customer's monthly charges compared to the average, reflecting potential dissatisfaction with pricing.
  - o **Results**:
- **T-test**: p-value = 0.001, confirming that churned customers have higher price sensitivity.
- **Pearson Correlation**: 0.45, indicating a moderate positive correlation between price sensitivity and churn.
- **Visual Insights**: Churned customers exhibited higher price sensitivity, and were concentrated in higher monthly charge brackets.
  - o **Key Insights**:
- **Churn Drivers**: Higher price sensitivity correlates strongly with churn.
- **High-Risk Segments**: Customers with month-to-month contracts, high charges, and volatile billing are at the highest risk.
- **Business Opportunity**: Addressing price sensitivity through strategies like discounts or stable pricing can significantly reduce churn.

# Sub task 3

## : Half-Page Summary of Key Findings and Suggestions for Data Augmentation

The goal of this analysis is to summarize key insights and recommend strategies for enhancing data to improve predictive model performance.

**Key Findings**:

- **Churn Rate**: Approximately 10.29% of customers churned, highlighting a class imbalance.
- **Key Demographic Insights**: Customers with month-to-month contracts, senior citizens, and those without dependents or a partner have higher churn rates.
- **Predictors of Churn**:
    - **Contract Type**: Short-term contracts are linked to higher churn.
    - **Price Sensitivity**: Customers with fluctuating monthly charges are more likely to churn.
    - **Internet and Support Services**: Fiber-optic customers and those lacking support services are more likely to churn.

**Model Performance**:

The Random Forest model showed strong performance with **93.6% accuracy**, **94% precision**, and **92.8% recall** for churned customers, achieving an **AUC of 0.935**.

**Price Sensitivity**:

A significant correlation between price sensitivity and churn was confirmed, underscoring the need for adjusted pricing strategies.

**Data Augmentation Suggestions**:

- **Customer Sentiment**: Incorporate feedback from surveys and complaints to detect churn signals early.
- **Usage and Billing Data**: Track service usage and support requests to identify at-risk customers.
- **Competitor and Geographical Data**: Adding competitor insights and location data will refine churn predictions.

# Sub task 4

## Feature Evaluation, Engineering, and Granularity

**Objective**:

This task aimed to assess the current features' effectiveness in predicting churn, explore new feature engineering opportunities, refine feature granularity, and eliminate unnecessary features to optimize model performance.

### 1. Evaluation of Features Against Churn

- **Contract Type**: Customers with month-to-month contracts had significantly higher churn. **Chi-square p-value** < 0.001 confirmed the strong link.
- **Price Sensitivity**: A moderate positive correlation (0.45) was found, indicating that customers with greater fluctuations in monthly charges are more likely to churn.
- **Demographics**: Senior citizens had a **1.89 times higher odds of churn**.
- **Value-Added Services**: Customers lacking tech support, online security, or streaming services were more likely to churn, with **p-values** < 0.05 for each.

### 2. Feature Engineering

- **Price Sensitivity**: Calculated by the fluctuation in monthly charges, improving recall by **4.5%**.
- **Customer Engagement Score**: Linked to retention, helping to flag potential churners.
- **Tenure Banding**: Grouping customers by tenure showed that churn probability decreases after **12 months**.
- **Service Combination Feature**: Combining services revealed that customers with "Fiber without tech support" were more likely to churn.

### 3. Granularity Improvements

- **Monthly Charges**: Grouped into buckets (low, medium, high) for better model interpretability.
- **Geographical Data**: Considered adding location clusters to address regional service quality.
- **Contract Type**: One-hot encoding was applied for better algorithm compatibility.

### 4. Removal of Unnecessary Features

- Removed features like **Customer IDs**, and highly correlated ones like **Total Charges** and **Streaming TV/Movies**, improving model efficiency.

### 5. Model Performance After Refinement

Post-feature engineering, the **Random Forest classifier** achieved:

- **Accuracy**: 93.6%
- **Recall** for churned customers: 92.8%
- **AUC**: 0.935

This represented a **5.2% improvement in recall** and a **3.6% increase in precision** for churned customers.

## 6. Insights and Recommendations

The newly engineered features, particularly **Price Sensitivity**, **Engagement Score**, and **Tenure Banding**, significantly enhanced model performance. Removing redundant features helped reduce noise, resulting in a more efficient and accurate model.

# Sub task 5

## Model Development, Evaluation, and Business Implications

*Objective:*

The goal is to develop and evaluate a predictive model to understand how well the engineered features can predict customer churn. This task also involves documenting the advantages and disadvantages of using Random Forest, evaluating the appropriateness of the chosen metrics, and connecting the model's performance to business outcomes like profits or savings.

### 1. Classification vs. Regression: Why Classification?

The task at hand is a binary classification problem where the goal is to predict whether a customer will churn (1) or stay (0).

Since churn prediction does not involve continuous outcomes, regression methods are unsuitable for this task, making classification the correct approach.

### 2. Model Performance Metrics

**Metrics Used and Their Rationale**:

- **Accuracy**:
    - Measures the overall correctness of the model's predictions.
    - **Result**: The Random Forest model achieved **93.6% accuracy**, indicating strong performance across all classes.
- **Recall (Sensitivity)**:
    - Focuses on the model's ability to correctly identify churners, which is critical in this scenario as we want to catch as many churners as possible.
    - **Result**: Recall for the churn class reached **92.8%**, ensuring fewer at-risk customers are overlooked.
    - **Reason**: Prioritizing recall ensures that most churners are flagged for intervention, preventing potential revenue loss.
- **Precision**:
    - Measures how many of the predicted churners actually end up churning.
    - **Result**: Precision for churners was **94.1%**, which minimizes false positives and prevents unnecessary interventions.
- **AUC (Area Under the ROC Curve)**:

- Evaluates the trade-off between true positive rate (recall) and false positive rate (1 - specificity).
- **Result**: The model achieved an **AUC of 0.935**, indicating excellent discriminatory ability and a robust performance overall.

## 3. Advantages and Disadvantages of Random Forest

**Advantages**:

- **High Accuracy**:

Random Forest outperformed Logistic Regression, showing notable improvements in recall and precision, making it a good choice for this type of classification task.

- **Feature Importance**:

Random Forest provides interpretable **feature importance scores**, helping to identify which factors most contribute to churn. Key features include **Contract Type**, **Price Sensitivity**, **Tenure Bands**, and **Engagement Score**.

- **Robustness**:

It handles imbalanced data effectively, especially with class-weight adjustments, making it suitable for churn prediction where churned customers are a smaller proportion of the dataset.

- **Non-Parametric**:

The model does not make assumptions about the relationship between features and target, which makes it flexible for capturing complex patterns.

**Disadvantages**:

- **Computational Cost**:

Training the Random Forest model with hyperparameter tuning (like grid search) was computationally intensive, which may slow down the process, especially with larger datasets.

- **Lack of Explainability**: While feature importance can be derived, the decision-making process within Random Forest is less transparent compared to simpler models like Logistic Regression, which might pose challenges in some business settings.

## 4. Evaluation of Model Performance

**Where the Model Underperformed**:

- **False Positives**:

The model predicted some customers as churners who were not at risk. While the financial cost of these false positives is manageable, it increases the workload for interventions, which could lead to unnecessary retention efforts.

- **Edge Cases**:

Some customers with **mid-range tenure** and **moderate price sensitivity** were misclassified. These edge cases could be fine-tuned for better performance.

**Performance Justification**: Overall, the model's performance is highly satisfactory, especially with critical business metrics like recall and precision. After incorporating engineered features and hyperparameter tuning, the model achieved **93.6% accuracy** and **92.8% recall**, which ensures it is reliable for real-world applications.

## 5. Financial Implications

**Assumptions for Financial Analysis**:

- **Average Revenue per Customer (ARPU)**: ₹1,000 per month.
- **Churn Rate**: 20% based on exploratory analysis.
- **Retention Rate with Model**: 15% (a 5% reduction in churn).
- **Cost per Intervention**: ₹100 (discounts, offers, etc.).

**Savings Calculation**:

- **Total Customers at Risk of Churning**:

Given a **20% churn rate**, approximately **20,000 churners** are estimated in the business context.

- **Revenue Loss Without Model**:

Without intervention, the churners represent a **₹240 million** loss annually (₹1,000 × 20,000 churners × 12 months).

- **Savings from Reduced Churn**:

By retaining **1,000 customers** (5% reduction in churn), the business saves **₹12 million** annually (₹1,000 × 1,000 retained customers × 12 months).

- **Cost of Interventions**:

Intervention costs for retaining 1,000 customers amount to ₹100,000 (₹100 × 1,000).

- **Net Savings**: After accounting for intervention costs, the net savings is **₹11.9 million** annually (₹12 million - ₹0.1 million).