# DTSA 5510 Final Project

# GitHub repo URL

https://github.com/k-khuu/DTSA-5510-Final-Project/tree/master

# Introduction and Problem Description

Customer attrition is a major concern in the financial industry, where retaining credit card users is essential for long-term profitability. Identifying customers who are likely to disengage allows businesses to take proactive steps to improve retention and build stronger relationships.

This project explores a dataset available on Kaggle containing information about credit card holders, including demographic details, account activity, and whether they have stopped using their cards. The objective is to uncover patterns that may indicate early signs of churn.

To tackle this, I applied two unsupervised learning techniques: **K-Means Clustering** and **Agglomerative Clustering**. These methods help group customers based on shared behaviors and characteristics, revealing potential segments with varying levels of attrition risk. I also used XGBoost as a supervised model to benchmark predictive performance and assess feature relevance.

Throughout the process, I ran multiple iterations of each model, experimenting with different feature sets and tuning parameters to see what combinations produced the most meaningful insights. It was a hands-on exploration of how clustering and classification can complement each other in understanding customer behavior

# Import required modules

```
import warnings
warnings.filterwarnings("ignore", category=RuntimeWarning) # Suppress
minor runtime warnings

!pip install gower # Install Gower package

# Data manipulation
import pandas as pd
import numpy as np
from itertools import permutations

# Visualization
import altair as alt
import seaborn as sns
import matplotlib.pyplot as plt

# Preprocessing
from sklearn.preprocessing import LabelEncoder, StandardScaler
```

```
from sklearn.model selection import train test split
# Clustering
from sklearn.cluster import KMeans, AgglomerativeClustering
from scipy.cluster.hierarchy import dendrogram
# Tree-based models
from sklearn.tree import DecisionTreeClassifier
from sklearn import tree
# Evaluation metrics
from sklearn.metrics import (
    accuracy_score,
    balanced accuracy score,
    confusion matrix,
    precision recall fscore support
)
# Gradient boosting
import xgboost as xgb
# Altair configuration
alt.data transformers.disable max rows()
# Distance metric for mixed data types
import gower
Requirement already satisfied: gower in
/usr/local/lib/python3.11/dist-packages (0.1.2)
Requirement already satisfied: numpy in
/usr/local/lib/python3.11/dist-packages (from gower) (1.26.4)
Requirement already satisfied: scipy in
/usr/local/lib/python3.11/dist-packages (from gower) (1.15.3)
Requirement already satisfied: mkl fft in
/usr/local/lib/python3.11/dist-packages (from numpy->gower) (1.3.8)
Requirement already satisfied: mkl random in
/usr/local/lib/python3.11/dist-packages (from numpy->gower) (1.2.4)
Requirement already satisfied: mkl umath in
/usr/local/lib/python3.11/dist-packages (from numpy->gower) (0.1.1)
Requirement already satisfied: mkl in /usr/local/lib/python3.11/dist-
packages (from numpy->gower) (2025.2.0)
Requirement already satisfied: tbb4py in
/usr/local/lib/python3.11/dist-packages (from numpy->gower) (2022.2.0)
Requirement already satisfied: mkl-service in
/usr/local/lib/python3.11/dist-packages (from numpy->gower) (2.4.1)
Requirement already satisfied: intel-openmp<2026,>=2024 in
/usr/local/lib/python3.11/dist-packages (from mkl->numpy->gower)
(2024.2.0)
Requirement already satisfied: tbb==2022.* in
/usr/local/lib/python3.11/dist-packages (from mkl->numpy->gower)
```

```
(2022.2.0)
Requirement already satisfied: tcmlib==1.* in
/usr/local/lib/python3.11/dist-packages (from tbb==2022.*->mkl->numpy-
>gower) (1.4.0)
Requirement already satisfied: intel-cmplr-lib-rt in
/usr/local/lib/python3.11/dist-packages (from mkl_umath->numpy->gower)
(2024.2.0)
Requirement already satisfied: intel-cmplr-lib-ur==2024.2.0 in
/usr/local/lib/python3.11/dist-packages (from intel-openmp<2026,>=2024->mkl->numpy->gower) (2024.2.0)
```

# **Exploratory Data Analysis**

Kaggle link: https://www.kaggle.com/datasets/sakshigoyal7/credit-card-customers

# **Data Description**

This is a dataset of credit card customers. It consists of over 10,000 customers mentioning their age, salary, marital\_status, credit card limit, credit card category, etc. There are 23 customer features/columns.

The dataset captures a variety of customer attributes relevant to credit card usage and retention analysis. These features are organized into three key groups:

Demographic information includes:

- Age
- Gender
- Number of Dependents
- Education Level

#### Account-related data covers:

- Card Category
- Tenure with the Bank (Months on Book)
- Credit Limit
- Total Revolving Balance

#### Target label:

 A binary indicator specifying whether the customer is active or has attrited (i.e., stopped using their credit card) The dataset contains both categorical variables (such as Gender, Education Level, and Card Category) and numerical variables (such as Age, Tenure, and Credit Limit), making it suitable for mixed-type analysis in clustering, classification, and churn prediction

# Load Data and Analyzing the Dataset

In this step, I load the dataset from a CSV file and perform initial preprocessing. This includes assigning appropriate categorical data types and explicitly ordering any ordinal features to ensure correct handling during analysis.

Additionally, I remove three columns: CLIENTNUM, along with two others that show little relevance to the modeling objectives and do not contribute meaningful information.

```
# Load and prepare dataset
def load data(path):
    return pd.read csv(path)
def assign categories(df, categorical map):
    for col, order in categorical map.items():
        df[col] = pd.Categorical(df[col], categories=order,
ordered=bool(order))
    return df
def drop columns(df, columns):
    return df.drop(columns=columns)
# File path
csv path = "/kaggle/input/credit-card-customers/BankChurners.csv"
# Categorical configuration
categorical config = {
    "Education Level": [
        "Unknown", "Uneducated", "High School", "College",
        "Graduate", "Post-Graduate", "Doctorate"
    "Income Category": [
        "Unknown", "Less than $40K", "$40K - $60K",
        "$60K - $80K", "$80K - $120K", "$120K +"
    ],
    "Gender": None,
    "Marital Status": None,
    "Card Category": None,
    "Attrition Flag": None
}
# Columns to remove
irrelevant columns = [
    "Attrition Flag",
    "CLIENTNUM",
"Naive Bayes Classifier Attrition Flag Card Category Contacts Count 12
mon Dependent count Education Level Months Inactive 12 mon 1",
"Naive Bayes Classifier Attrition Flag Card Category Contacts Count 12
```

```
_mon_Dependent_count_Education_Level_Months Inactive 12 mon 2"
# Pipeline
data raw = load data(csv path)
data prepped = assign categories(data raw.copy(), categorical config)
labels = data prepped[["Attrition Flag"]]
label classes = labels["Attrition Flag"].unique()
data final = drop columns(data prepped, irrelevant columns)
# Display the first few rows of the cleaned dataset
data final.head()
                         Dependent count Education Level Marital Status
   Customer Age Gender
0
             45
                      М
                                        3
                                              High School
                                                                  Married
             49
                                        5
                                                 Graduate
1
                                                                   Single
2
             51
                                                 Graduate
                                                                  Married
3
             40
                                                                  Unknown
                                              High School
             40
                      М
                                        3
                                               Uneducated
                                                                  Married
  Income Category Card Category
                                  Months on book
Total Relationship_Count
                                               39
      $60K - $80K
                            Blue
5
1
   Less than $40K
                            Blue
                                               44
6
2
                            Blue
                                               36
     $80K - $120K
4
3
   Less than $40K
                            Blue
                                               34
3
4
      $60K - $80K
                            Blue
                                               21
5
                            Contacts Count 12 mon
                                                   Credit Limit \
   Months Inactive 12 mon
0
                         1
                                                 3
                                                          12691.0
                         1
                                                 2
1
                                                           8256.0
2
                         1
                                                 0
                                                           3418.0
3
                         4
                                                 1
                                                           3313.0
4
                         1
                                                 0
                                                           4716.0
   Total Revolving Bal
                                          Total Amt Chng Q4 Q1 \
                         Avg_Open_To_Buy
0
                                 11914.0
                    777
                                                           1.335
```

$\begin{array}{cccccccccccccccccccccccccccccccccccc$
2517 796.0 1.405
0 4716.0 2.175
otal Trans Amt Total Trans Ct Total Ct Chng Q4 Q1
Jtilization Ratio
$1\overline{1}44$ 42 1.625
1
1291 33 3.714
)
1887 20 2.333
)
)
816 28 2.500

#### Initial Observations

Upon inspecting the dataset, we observe a mix of ordinal, nominal, and numerical features. The dataset contains 10,127 records, and all columns are complete—no missing values are present.

A review of the numerical variables reveals no apparent outliers or inconsistencies, suggesting that the data is clean and reliable at this stage. Therefore, no further cleaning or imputation is required before proceeding.

It's worth noting that during model development, additional preprocessing steps are applied to convert categorical features into numerical representations, enabling their use in machine learning algorithms.

```
# Generate a summary of column types, non-null counts, and unique
values
summary = pd.DataFrame({
    "Data Type": data final.dtypes,
    "Non-Null Count": data_final.count(),
    "Unique Values": data_final.nunique()
})
# Display the summary sorted by column name
summary.sort index()
                         Data Type
                                     Non-Null Count
                                                     Unique Values
Avg Open To Buy
                           float64
                                              10127
                                                              6813
Avg Utilization Ratio
                           float64
                                              10127
                                                               964
Card Category
                                              10127
                                                                  4
                          category
Contacts Count 12 mon
                                                                  7
                             int64
                                              10127
                                              10127
                                                              6205
Credit Limit
                           float64
Customer Age
                             int64
                                              10127
                                                                 45
Dependent count
                             int64
                                              10127
                                                                  6
```

```
Education Level
                                                10127
                                                                    7
                           category
Gender
                                                                    2
                                                10127
                           category
Income Category
                           category
                                                10127
                                                                    6
Marital Status
                                               10127
                                                                    4
                           category
Months Inactive 12 mon
                              int64
                                               10127
                                                                    7
Months on book
                                               10127
                              int64
                                                                   44
Total Amt Chng Q4 Q1
                            float64
                                                                 1158
                                               10127
Total Ct Chng \overline{Q}4 \overline{Q}1
                            float64
                                               10127
                                                                  830
Total Relationship Count
                              int64
                                               10127
                                                                    6
Total Revolving Bal
                              int64
                                               10127
                                                                 1974
Total Trans Amt
                              int64
                                                10127
                                                                 5033
Total Trans Ct
                              int64
                                                10127
                                                                  126
# Generate descriptive statistics for numerical columns only
numeric summary = data final.select dtypes(include="number").agg(
    ["count", "mean", "std", "min", "median", "max"]
).transpose()
# Rename columns for clarity
numeric summary.columns = ["Count", "Mean", "Std Dev", "Min",
"Median", "Max"]
# Display the summary
numeric summary
```

### Churn versus Existing Customers

The target variable for this classification task is the Attrition Flag, which indicates whether a customer has attrited. A quick inspection reveals that the class distribution is notably imbalanced. To address this, we apply undersampling techniques during model development to help ensure more balanced training and improve predictive performance.

```
# Bar chart showing distribution of churn labels
alt.Chart(labels).mark_bar(color="#5276A7").encode(
    x=alt.X("Attrition_Flag:N", title="Customer Status"),
    y=alt.Y("count()", title="Number of Customers")
).properties(
    width=350,
    height=300,
    title="Distribution of Churn vs Existing Customers"
)
alt.Chart(...)
```

# Customer Breakdown by Gender

The gender distribution appears reasonably balanced, which is consistent with expectations for a broad customer base. There's no indication of skew or irregularity

```
# Bar chart showing gender distribution
alt.Chart(data_final).mark_bar(color="#D65F5F").encode(
    x=alt.X("Gender:N", title="Gender"),
    y=alt.Y("count()", title="Customer Count")
).properties(
    width=350,
    height=300,
    title="Gender Breakdown of Credit Card Customers"
)
alt.Chart(...)
```

#### Breakdown of Customer Education Levels

The distribution of education levels among customers shows a clear concentration in a few categories. Most customers have either a Graduate or High School education, followed by a smaller group with unknown education status. College, Post-Graduate, and Doctorate levels are less represented. Overall, the spread is moderately skewed, but not unexpected for a general consumer credit dataset.

```
# Display sorted counts of education levels
education counts =
data_final["Education_Level"].value counts(sort=True)
education counts
Education Level
Graduate
                 3128
High School
                 2013
                 1519
Unknown
Uneducated
                 1487
College
                 1013
Post-Graduate
                516
Doctorate
                  451
Name: count, dtype: int64
# Bar chart showing distribution of education levels
alt.Chart(data_final).mark_bar(color="#F2C94C").encode(
    x=alt.X("Education Level:N", title="Education Level",
sort=education counts.index.tolist()),
    y=alt.Y("count()", title="Number of Customers")
).properties(
    width=350,
    height=300,
    title="Customer Distribution by Education Level"
)
alt.Chart(...)
```

# Breakdown of Customer Income Categories

The customer base is concentrated in the lower income brackets, with the largest group earning less than 40K. Moderate representation follows in the 40K to 60K and 60K to 80K ranges. Higher income categories, including 80K to 120K and 120K and above, are less common. A notable portion of customers have their income listed as "Unknown." Overall, the distribution leans toward lower to middle income levels.

```
# Display sorted counts of income categories
income counts =
data final["Income Category"].value counts(ascending=False)
income counts
Income Category
Less than $40K
                  3561
$40K - $60K
                  1790
$80K - $120K
                  1535
$60K - $80K
                  1402
Unknown
                  1112
$120K +
                  727
Name: count, dtype: int64
# Bar chart showing distribution of income categories
alt.Chart(data final).mark bar(color="#76B7B2").encode(
    x=alt.X("Income_Category:N", title="Income Category", sort="-y"),
    v=alt.Y("count()", title="Number of Customers")
).properties(
    width=350,
    height=300,
    title="Customer Distribution by Income Category"
)
alt.Chart(...)
```

# Customer Segmentation by Card Category

The distribution of Card Category is heavily skewed, which suggests it may offer limited value for modeling or analysis without further transformation.

```
# Bar chart showing distribution of card categories
alt.Chart(data_final).mark_bar(color="#9B51E0").encode(
    x=alt.X("Card_Category:N", title="Card Category", sort="-y"),
    y=alt.Y("count()", title="Number of Customers")
).properties(
    width=350,
    height=300,
    title="Customer Distribution by Card Category"
)
alt.Chart(...)
```

# Age Distribution of Customers

The age distribution of customers follows a roughly bell-shaped curve, with the highest concentration in the mid-40s to early 50s range. This segment represents the peak of the customer base, while both younger and older age groups appear less represented. The overall spread suggests a mature customer population, with most individuals falling between their late 30s and mid-60s.

```
# Histogram of customer age distribution
alt.Chart(data_final).mark_bar(color="#6FCF97").encode(
    x=alt.X("Customer_Age:Q", bin=alt.Bin(maxbins=20), title="Age
Range"),
    y=alt.Y("count()", title="Number of Customers")
).properties(
    width=350,
    height=300,
    title="Distribution of Customer Age"
)
alt.Chart(...)
```

### Number of Dependents per Account Holder

The distribution of dependents per account holder peaks at three dependents, indicating this is the most common household size among customers. Zero and two dependents follow closely behind, while one, four, and five dependents show progressively lower counts. Overall, the data suggests a mix of small to mid-sized households, with three dependents being the most typical.

```
# Bar chart showing distribution of dependent counts
alt.Chart(data_final).mark_bar(color="#F2994A").encode(
    x=alt.X("Dependent_count:0", title="Number of Dependents"),
    y=alt.Y("count()", title="Customer Count")
).properties(
    width=350,
    height=300,
    title="Distribution of Dependents per Customer"
)
alt.Chart(...)
```

#### Feature Selection Based on Attrition Correlation

To get a clearer picture of what influences customer churn, I looked at how each feature correlates with Attrition\_Flag. A handful of variables stood out with strong signals:

- Total\_Relationship\_Count
- Months\_Inactive\_12\_mon

- Contacts\_Count\_12\_mon
- Total\_Revolving\_Bal
- Total\_Amt\_Chng\_Q4\_Q1
- Total\_Trans\_Ct

These are the ones most closely tied to attrition behavior and are likely to add real value to predictive models. The rest did not show much correlation and may not contribute meaningfully without further engineering.

To test this, I will run two sets of models:

- One using just the correlated features
- Another using all available features

This way, we can see whether a leaner input set improves performance or if the extra data adds nuance worth keeping.

# Analysis (Model Building and Training)

# Models - Unsupervised

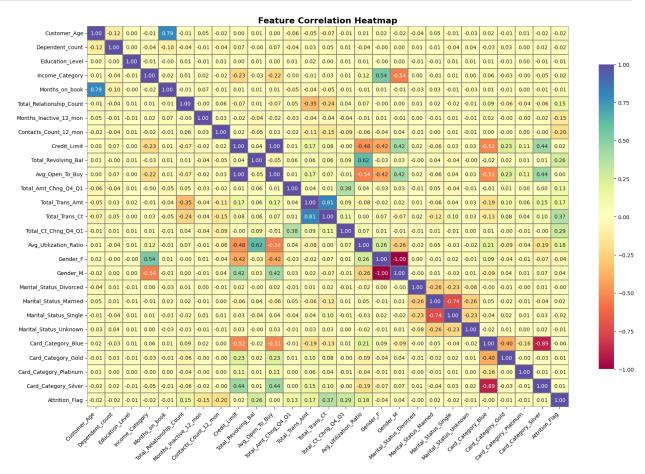
To explore customer segmentation without relying on labeled data, I tested two clustering approaches KMeans and Agglomerative Clustering across seven configurations. Each variation tweaks the input features or preprocessing strategy to uncover different structural patterns in the data:

- 1. KMeans using all available columns
- 2. KMeans using only the features most correlated with attrition
- 3. Agglomerative Clustering with all columns, converting categorical features to numeric codes
- 4. Agglomerative Clustering with all columns and ward linkage for tighter cluster formation
- 5. Agglomerative Clustering on an undersampled dataset to balance class representation
- 6. Agglomerative Clustering using only the correlated features, with categorical data numerically encoded
- 7. Agglomerative Clustering using Gower distance to better handle mixed data types

For each setup, I evaluated clustering quality using accuracy, precision, recall, and F1 score; giving a clearer sense of which configurations best separate customer behaviors.

```
# Try different label orders to find the best match with actual labels
def find best label mapping(true labels, predicted labels,
class names, metric="accuracy"):
    def evaluate mapping(mapping):
        remapped = true labels.map(mapping)
        scores = precision_recall_fscore_support(remapped,
predicted labels, average="weighted")
        return {
            "accuracy": accuracy_score(remapped, predicted_labels),
            "precision": scores[0],
            "recall": scores[1],
            "f1": scores[2]
        }
    top score = float("-inf")
    top metrics = {}
    top mapping = {}
    for combo in permutations(range(len(class names))):
        current map = {name: idx for name, idx in zip(class names,
combo)}
        current metrics = evaluate mapping(current map)
        if current metrics[metric] > top score:
            top score = current metrics[metric]
            top metrics = current_metrics
            top mapping = {idx: name for idx, name in zip(combo,
class names)}
    return combo, top score, top metrics, top mapping
# Extract column groups from config
ordinal cols = [col for col, order in categorical config.items() if
order is not Nonel
nominal cols = [col for col, order in categorical config.items() if
order is None and col != "Attrition Flag"]
# Encode ordinal features
for col in ordinal cols:
    data final[col] = LabelEncoder().fit transform(data final[col])
# One-hot encode nominal features
data_final = pd.get_dummies(data_final, columns=nominal cols)
# Compute correlation matrix
correlations = corr data.corr().round(2)
# Set up the plot
plt.figure(figsize=(18, 12))
```

```
sns.heatmap(
    correlations,
    annot=True,
    fmt=".2f",
    cmap="Spectral",
                              # Vibrant diverging color palette
    center=0.
                               # Centered at zero for better contrast
    linewidths=0.3,
    linecolor="gray"
    cbar kws={"shrink": 0.8}
                               # Smaller color bar
)
plt.title("Feature Correlation Heatmap", fontsize=16,
fontweight="bold")
plt.xticks(rotation=45, ha="right")
plt.yticks(rotation=0)
plt.tight_layout()
plt.show()
```



```
key_features_for_attrition = [
    "Total_Relationship_Count",
    "Months_Inactive_12_mon",
    "Contacts_Count_12_mon",
```

```
"Total Revolving Bal"
    "Total Amt Chng Q4 Q1",
    "Total Trans Ct"
1
# Combine features and labels for sampling
df = pd.concat([data final, labels], axis=1)
majority = df[df["Attrition Flag"] == "Existing Customer"]
minority = df[df["Attrition Flag"] == "Attrited Customer"]
# Downsample majority class to match minority count
majority sampled = majority.sample(n=len(minority), random state=42)
df balanced = pd.concat([majority sampled, minority]).sample(frac=1,
random state=42)
X resampled = df balanced.drop("Attrition Flag", axis=1)
y resampled = df balanced["Attrition Flag"]
# Scale all features
scaler all = StandardScaler()
scaled all = scaler all.fit transform(data final)
scaled all resampled = scaler all.fit transform(X resampled)
# Scale selected correlative features
scaler subset = StandardScaler()
scaled subset =
scaler subset.fit transform(data final[key features for attrition])
scaled subset resampled =
scaler subset.fit transform(X resampled[key features for attrition])
```

#### KMeans with All Available Columns

KMeans clustering across the full feature set yielded modest results, with an F1 score of 0.755.

```
'recall': 0.792436062012442,
'f1': 0.7559451980077926}
```

#### KMeans with Correlation-Selected Features

KMeans delivered comparable results when limited to the correlation-selected features, reaching an F1 score of 0.571.

```
# Fit KMeans model on selected correlative features
kmeans_subset = KMeans(n_clusters=2, n_init="auto", random_state=42)
subset_cluster_labels = kmeans_subset.fit_predict(scaled_subset)

# Evaluate clustering performance against true labels
_, _, subset_metrics, label_map_subset = label_permute_compare(
        labels["Attrition_Flag"], subset_cluster_labels, label_classes,
"accuracy"
)

subset_metrics
{'accuracy': 0.5083440308087291,
    'precision': 0.666806152582151,
    'recall': 0.5083440308087291,
    'f1': 0.5713161095441398}
```

# Agglomerative Clustering with All Available Columns

Agglomerative Clustering using all columns, with categorical features numerically encoded, showed stronger performance. This approach achieved an F1 score of 0.766.

```
# Fit Agglomerative Clustering on full feature set
agg_model_all = AgglomerativeClustering(
    n_clusters=2,
    metric="euclidean",
    linkage="complete",
    compute_distances=True
)
agg_cluster_labels = agg_model_all.fit_predict(scaled_all)
# Evaluate clustering performance against true labels
_, _, agg_metrics_all, label_map_agg = label_permute_compare(
    labels["Attrition_Flag"], agg_cluster_labels, label_classes,
"accuracy"
)
agg_metrics_all
```

```
{'accuracy': 0.8383529179421348,
 'precision': 0.7448055607507813,
 'recall': 0.8383529179421348,
 'f1': 0.7664726575589648}
```

### Agglomerative Clustering with All Available Columns + Ward Linkage

Agglomerative Clustering with Ward linkage, applied to the full feature set including numerically encoded categorical variables, delivered comparable performance, yielding an F1 score of 0.756.

```
# Fit Agglomerative Clustering using Ward linkage on full feature set
agg model ward = AgglomerativeClustering(
    n clusters=2,
    linkage="ward",
    compute distances=True
)
ward cluster labels = agg model ward.fit predict(scaled all)
# Evaluate clustering performance against true labels
_, _, ward_metrics, label_map_ward = label_permute_compare(
   labels["Attrition Flag"], ward cluster labels, label classes,
"accuracy"
)
ward metrics
{'accuracy': 0.7925348079391725,
 'precision': 0.7293850521403183,
 'recall': 0.7925348079391724,
 'f1': 0.7560039370982169}
```

# Agglomerative Clustering with All Available Columns + Undersampled Dataset

Agglomerative Clustering on the undersampled, balanced dataset resulted in notably weaker performance, with an F1 score of just 0.384.

### Agglomerative Clustering with Correlation-Selected Features

Agglomerative Clustering using only the correlation-selected features showed slightly lower performance compared to the full-feature model, with an F1 score of 0.748.

```
# Fit Agglomerative Clustering with complete linkage on correlative
features
agg model subset = AgglomerativeClustering(
    n clusters=2,
    linkage="complete",
    compute distances=True
)
subset cluster labels = agg model subset.fit predict(scaled subset)
# Evaluate clustering performance against true labels
_, _, agg_metrics_subset, label_map_subset = label_permute compare(
    labels["Attrition Flag"], subset cluster labels, label classes,
"accuracy"
agg metrics subset
{'accuracy': 0.8032981139527995,
 precision': 0.702351576184927,
 'recall': 0.8032981139527995,
 'f1': 0.748687833815765}
```

# Agglomerative Clustering with Gower Distance

This approach performed similarly to the Euclidean-based method with numerically encoded categorical data, yielding a comparable F1 score of 0.762.

```
# Convert categorical columns to string for Gower compatibility
for col in ordinal_cols + nominal_cols:
    data_original[col] = data_original[col].astype(str)
```

```
# Compute Gower distance matrix
gower dist matrix = gower.gower matrix(data original)
# Fit Agglomerative Clustering using precomputed Gower distances
agg model gower = AgglomerativeClustering(
    n clusters=2,
    metric="precomputed",
    linkage="average"
)
gower cluster labels = agg model gower.fit predict(gower dist matrix)
# Evaluate clustering performance
_, _, gower_metrics, label_map_gower = label_permute_compare(
   labels["Attrition Flag"], gower cluster labels, label classes,
"accuracy"
gower metrics
{'accuracy': 0.8124814851387381,
 'precision': 0.7308330405617858,
 'recall': 0.8124814851387381,
 'f1': 0.7623152727247017}
```

### Models - Supervised

To evaluate how unsupervised clustering compares to predictive modeling, I used XGBoost, a tree-based supervised learning algorithm known for its performance and interpretability. The dataset was split with 20 percent reserved for testing. Each model was trained on the remaining data and assessed using consistent metrics including accuracy, precision, recall, and F1 score. This provided a clear benchmark for comparing supervised learning outcomes with the patterns identified through clustering.

#### XGBoost with All Available Columns

XGBoost outperformed all unsupervised approaches, achieving a notably high F1 score of 0.968.

```
# Encode target labels for classification
encoded_labels = label_encoder.fit_transform(labels["Attrition_Flag"])
# Split data into training and testing sets with stratified sampling
X_train, X_test, y_train, y_test = train_test_split(
    scaled_all,
    encoded_labels,
    test_size=0.2,
    stratify=encoded_labels,
    random_state=42
)
```

```
# Convert datasets into XGBoost's DMatrix format
dtrain = xgb.DMatrix(X train, label=y train)
dtest = xgb.DMatrix(X test, label=y test)
# Define model parameters
params = {
    "objective": "multi:softmax", # Multiclass classification
                        # Number of target c
# Maximum tree depth
    "num class": 5,
                                   # Number of target classes
    "max depth": 10,
    "eta": 1,
                                   # Learning rate
    "eval metric": "merror"  # Evaluation metric:
classification error
}
rounds = 20
# Train the XGBoost model
model = xgb.train(params, dtrain, rounds)
# Generate predictions on the test set
predictions = model.predict(dtest)
# Compute evaluation metrics
precision, recall, f1, _ = precision_recall_fscore_support(y_test,
predictions, average="weighted")
metrics = {
    "accuracy": accuracy_score(y_test, predictions),
    "precision": precision,
    "recall": recall.
    "f1": f1
}
metrics
{'accuracy': 0.9684106614017769,
 precision': 0.9684106614017769,
 'recall': 0.9684106614017769,
 'f1': 0.9684106614017769}
```

#### XGBoost with Correlation-Selected Features

Restricting XGBoost to the correlation-selected features resulted in a modest drop in performance, yielding an F1 score of 0.913.

```
encoded labels,
    test size=0.2,
    stratify=encoded labels,
    random state=42
)
# Prepare data for XGBoost
dtrain = xgb.DMatrix(X train, label=y train)
dtest = xgb.DMatrix(X_test, label=y_test)
# Model parameters
params = {
    "objective": "multi:softmax",
    "num_class": 5,
    "max depth": 10,
    "eta": 1,
    "eval metric": "merror"
rounds = 20
# Train and predict
model = xgb.train(params, dtrain, rounds)
predictions = model.predict(dtest)
# Evaluate
precision, recall, f1, _ = precision_recall_fscore support(y test,
predictions, average="weighted")
metrics = {
    "accuracy": accuracy score(y test, predictions),
    "precision": precision,
    "recall": recall,
    "f1": f1
}
metrics
{'accuracy': 0.9155972359328727,
 'precision': 0.9128463891113534,
 'recall': 0.9155972359328727,
 'f1': 0.9138648365985037}
```

# Discussion / Conclusion

This project explored different approaches to predicting customer attrition, starting with unsupervised clustering methods like KMeans and Agglomerative Clustering. I tested several configurations across both algorithms and compared their results to a supervised model using XGBoost.

XGBoost clearly outperformed the clustering models. While clustering helped uncover some structure in the data, it was not nearly as effective for classification. The supervised model delivered much stronger results across all metrics.

One surprising outcome was that using all available features consistently led to better performance than limiting the models to only the most correlated columns. I initially expected that focusing on the strongest predictors would improve accuracy, but the full feature set proved more valuable. This suggests that even weaker individual features may contribute meaningful signal when combined.

Agglomerative Clustering performed especially poorly on the undersampled dataset. This likely reflects how sensitive clustering is to data volume and distribution. With more records, even a balanced subset might have produced better results.

Going forward, I would like to explore clustering in contexts where it is better suited, such as customer segmentation or behavioral grouping. For classification tasks like churn prediction, supervised models are clearly the more reliable choice.