

FINAL PRESENTATION

CUSTOMER CHURN

A DATA ANALYSIS PROJECT

BY: AHMAD SOHAIL / 64011748

	G	H	I	J	K	L	
	Age	Tenure	Balance	NumOfPro	HasCrCards	IsActiveMe	Estim
Female	42	2	0	1	1	1	101348
Female	41	1	83807.86	1	0	1	112542.8
Female	42	8	159660.8	3	1	0	113931.6
Female	39	1	0	2	0	0	93826.63
Female	43	2	125510.8	1	1	1	79084.1
Male	44	8	113755.8	2	1	0	149756.7
Male	50	7	0	2	1	1	10062.8
Female	29	4	115046.7	4	1	0	119346.9
Male	44	4	142051.1	2	0	1	74940.5
Male	27	2	134603.9	1	1	1	71725.73
Male	31	6	102016.7	2	0	0	80181.15
Male	24	3	0	2	1	0	76390
Male	34	10	0	2	1	0	2626
Male	25	5	0	2	0	0	19
	35	7	0	2	1	1	6
	45	3	143129.4	2	0	1	
	58	1	132602.9	1	1		
	24	9	0	2	1		
	45	6	0	1	0		
	24	0	0	0	1		

Table Of Content

➤ Introduction

➤ About Dataset

➤ Recap

➤ New EDA

➤ Model(s)

➤ Finetuning & Evaluation

➤ Difficulties

➤ Q&A

THYNK UNLIMITED



Introduction

- Banks collect vast amounts of data to enhance their services and meet business objectives.
- A critical aspect of this data analysis is understanding customer behavior, such as churn rates.
- For this Data Analysis Project, I was assigned to work on a model to predict Customer Churn based on a dataset from Kaggle.
- This analysis helps the bank identify at-risk customers and take proactive measures to retain them.



THYNK UNLIMITED

About Dataset

Dataset by *RADHESHYAM KOLLIPARA*
that involves **Customer data collected**
by a bank.

- 1). 10,000 records to analyze
 - 2). 10+ necessary features to predict from
 - 3). Includes non-numeric and numeric values.



	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	
1	RowNum	CustomerId	Surname	CreditScore	Geography	Gender	Age	Tenure	Balance	NumOfProducts	HasCreditCard	IsActiveMember	EstimatedSalary	Exited	Complaint	Satisfaction	Card Type	Point Earned		
2	1	1.6E+07	Hargrave	619	France	Female	42	2	0	1	1	1	101349	1	1	2	DIAMOND	464		
3	2	1.6E+07	Hill	608	Spain	Female	41	1	83807.9	1	0	1	112543	0	1	3	DIAMOND	456		
4	3	1.6E+07	Onio	502	France	Female	42	8	159661	3	1	0	113932	1	1	3	DIAMOND	377		
5	4	1.6E+07	Boni	699	France	Female	39	1	0	2	0	0	93826.6	0	0	5	GOLD	350		
6	5	1.6E+07	Mitchell	850	Spain	Female	43	2	125511	1	1	1	79084.1	0	0	5	GOLD	425		
7	6	1.6E+07	Chu	645	Spain	Male	44	8	113756	2	1	0	149757	1	1	5	DIAMOND	484		
8	7	1.6E+07	Bartlett	822	France	Male	50	7	0	2	1	1	10062.8	0	0	2	SILVER	206		
9	8	1.6E+07	Obinna	376	Germany	Female	29	4	115047	4	1	0	119347	1	1	2	DIAMOND	282		
10	9	1.6E+07	He	501	France	Male	44	4	142051	2	0	1	74940.5	0	0	3	GOLD	251		
11	10	1.6E+07	H?	684	France	Male	27	2	134604	1	1	1	71725.7	0	0	3	GOLD	342		
12	11	1.6E+07	Bearce	528	France	Male	31	6	102017	2	0	0	80181.1	0	0	3	GOLD	264		
13	12	1.6E+07	Andrews	497	Spain	Male	24	3	0	2	1	0	76390	0	0	3	GOLD	249		
14	13	1.6E+07	Kay	476	France	Female	34	10	0	2	1	0	26261	0	0	3	SILVER	119		
15	14	1.6E+07	Chin	549	France	Female	25	5	0	2	0	0	190858	0	0	3	PLATINUM	549		
16	15	1.6E+07	Scott	635	Spain	Female	35	7	0	2	1	1	65951.7	0	0	2	GOLD	318		
17	16	1.6E+07	Goforth	616	Germany	Male	45	3	143129	2	0	1	64327.3	0	0	5	GOLD	308		
18	17	1.6E+07	Romeo	653	Germany	Male	58	1	132603	1	1	0	5097.67	1	0	2	SILVER	163		
19	18	1.6E+07	Henderson	549	Spain	Female	24	9	0	2	1	1	14406.4	0	0	3	SILVER	544		
20	19	1.6E+07	Muldrow	587	Spain	Male	45	6	0	1	0	0	158685	0	0	3	PLATINUM	732		
21	20	1.6E+07	Hao	726	France	Female	24	6	0	2	1	1	54724	0	0	4	GOLD	477		
22	21	1.6E+07	McDonald	732	France	Male	41	8	0	2	1	1	170886	0	0	3	PLATINUM	568		
23	22	1.6E+07	Dellucci	636	Spain	Female	32	8	0	2	1	0	138555	0	0	2	DIAMOND	336		
24	23	1.6E+07	Gerasimov	510	Spain	Female	38	4	0	1	1	0	118914	1	1	2	DIAMOND	887		
25	24	1.6E+07	Mosman	669	France	Male	46	3	0	2	0	1	8487.75	0	0	2	SILVER	665		
26	25	1.6E+07	Yen	846	France	Female	38	5	0	1	1	1	187616	0	0	3	PLATINUM	225		
27	26	1.6E+07	Maclean	577	France	Male	25	3	0	2	0	1	124508	0	0	3	PLATINUM	659		
28	27	1.6E+07	Young	756	Germany	Male	36	2	136816	1	1	1	170042	0	0	5	DIAMOND	236		
29	28	1.6E+07	Nebechi	571	France	Male	44	9	0	2	0	0	38433.4	0	0	4	PLATINUM	448		
30	29	1.6E+07	McWilliam	574	Germany	Female	43	3	141349	1	1	1	100187	0	0	5	DIAMOND	499		
31	30	1.6E+07	Lucciano	411	France	Male	29	0	59697.2	2	1	1	53483.2	0	0	2	GOLD	343		
32	31	1.6E+07	Azikiwe	591	Spain	Female	39	3	0	3	1	0	140469	1	1	3	DIAMOND	298		
33	32	1.6E+07	Odinakach	533	France	Male	36	7	85311.7	1	0	1	156732	0	0	5	SILVER	628		
34	33	1.6E+07	Sanderson	553	Germany	Male	41	9	110113	2	0	0	81898.8	0	0	3	GOLD	611		

Recap

- Developed an initial predictive model to analyze customer churn.
- Conducted basic Exploratory Data Analysis (EDA).
- Executed data preprocessing to prepare the dataset.

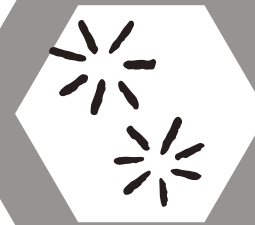


New EDA



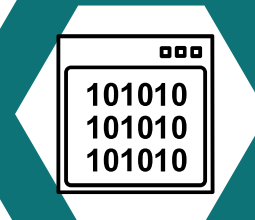
Extensive research

Extended the Textual EDA.



More visualizations

Extended the usage of plots that are helpful.



Key Observation(s)

Found Dataset's key error prone areas.

✓
0s

df.info()

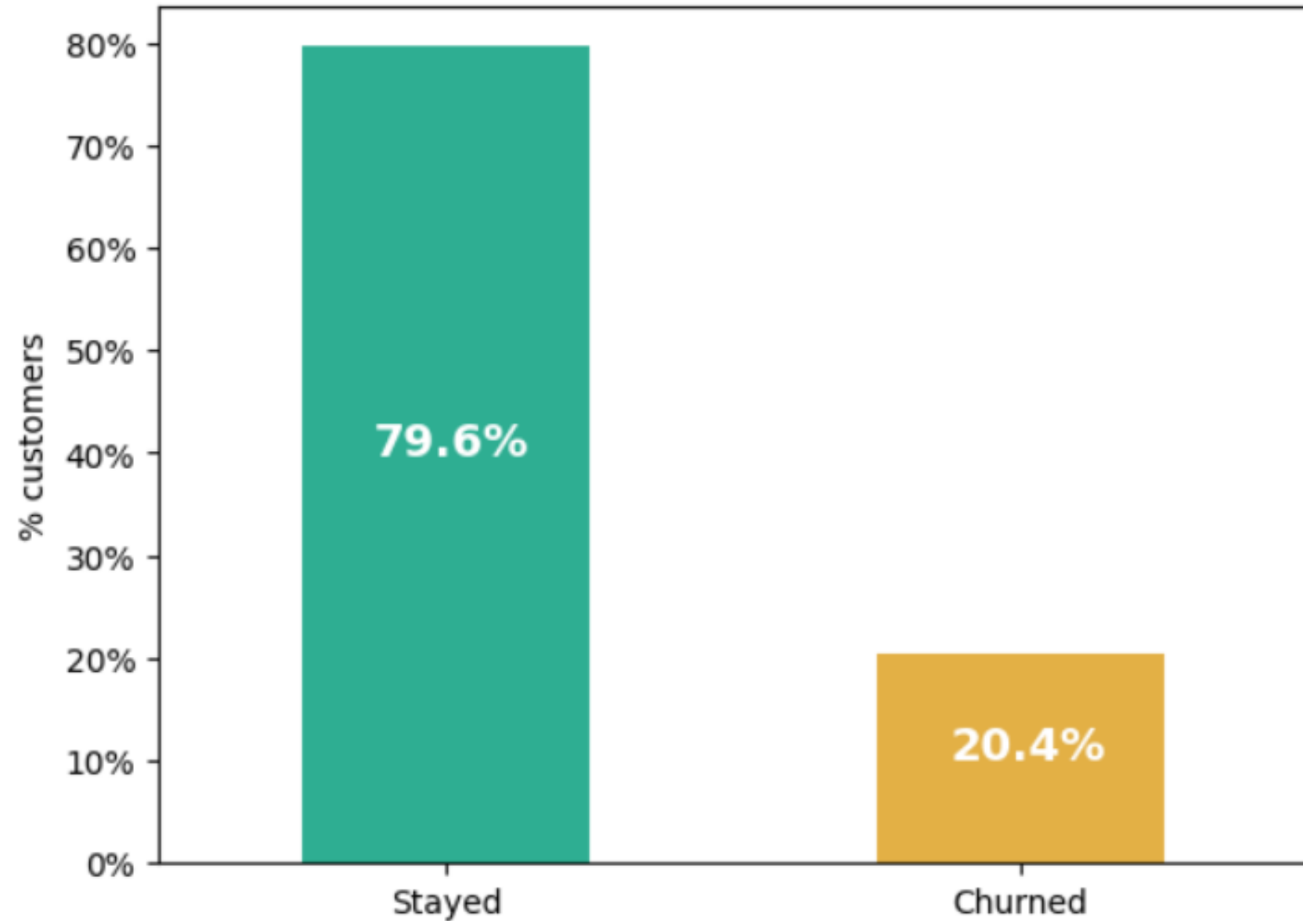
>>> `<class 'pandas.core.frame.DataFrame'>`
RangeIndex: 10000 entries, 0 to 9999
Data columns (total 18 columns):
Column Non-Null Count Dtype
--- -
0 RowNumber 10000 non-null int64
1 CustomerId 10000 non-null int64
2 Surname 10000 non-null object
3 CreditScore 10000 non-null int64
4 Geography 10000 non-null object
5 Gender 10000 non-null object
6 Age 10000 non-null int64
7 Tenure 10000 non-null int64
8 Balance 10000 non-null float64
9 NumOfProducts 10000 non-null int64
10 HasCrCard 10000 non-null int64
11 IsActiveMember 10000 non-null int64
12 EstimatedSalary 10000 non-null float64
13 Exited 10000 non-null int64
14 Complain 10000 non-null int64
15 Satisfaction Score 10000 non-null int64
16 Card Type 10000 non-null object
17 Point Earned 10000 non-null int64
dtypes: float64(2), int64(12), object(4)
memory usage: 1.4+ MB

✓
s

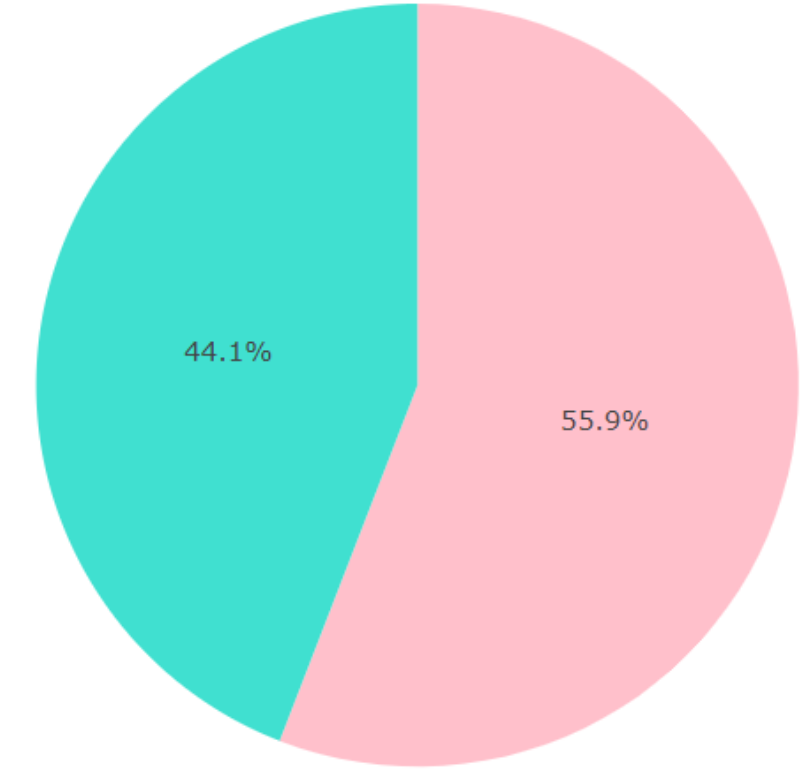
df.isnull().sum()

>>> RowNumber 0
CustomerId 0
Surname 0
CreditScore 0
Geography 0
Gender 0
Age 0
Tenure 0
Balance 0
NumOfProducts 0
HasCrCard 0
IsActiveMember 0
EstimatedSalary 0
Exited 0
Complain 0
Satisfaction Score 0
Card Type 0
Point Earned 0
dtype: int64

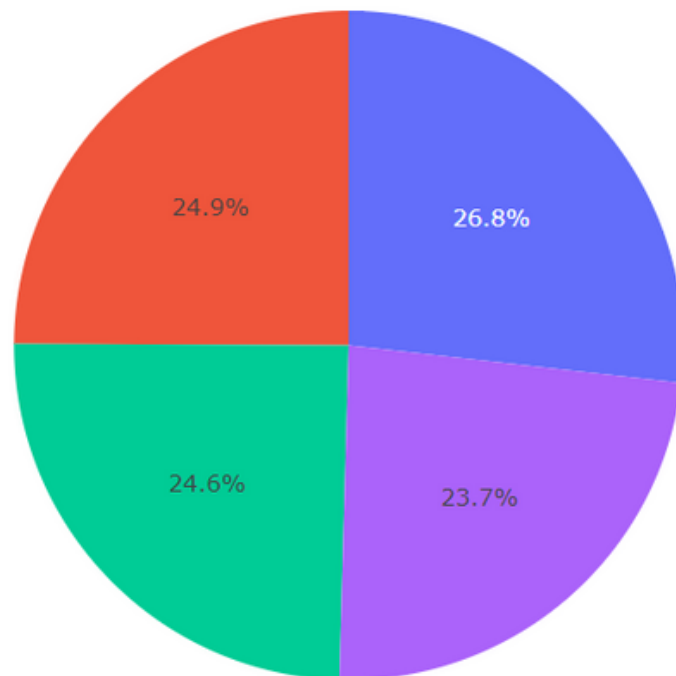
Customers by Churn



Churned percentage based on Gender

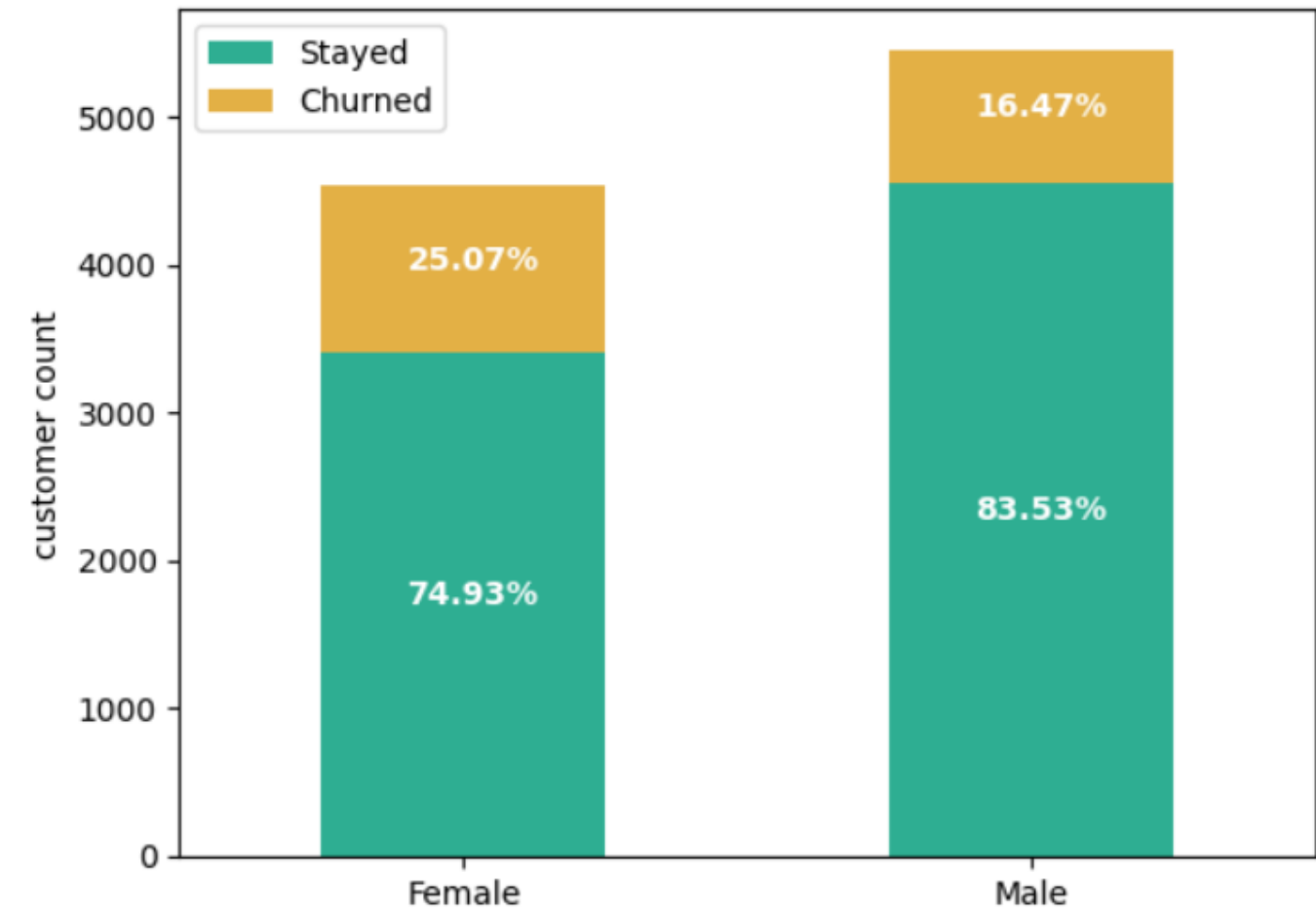


Churned percentage based on Card Type

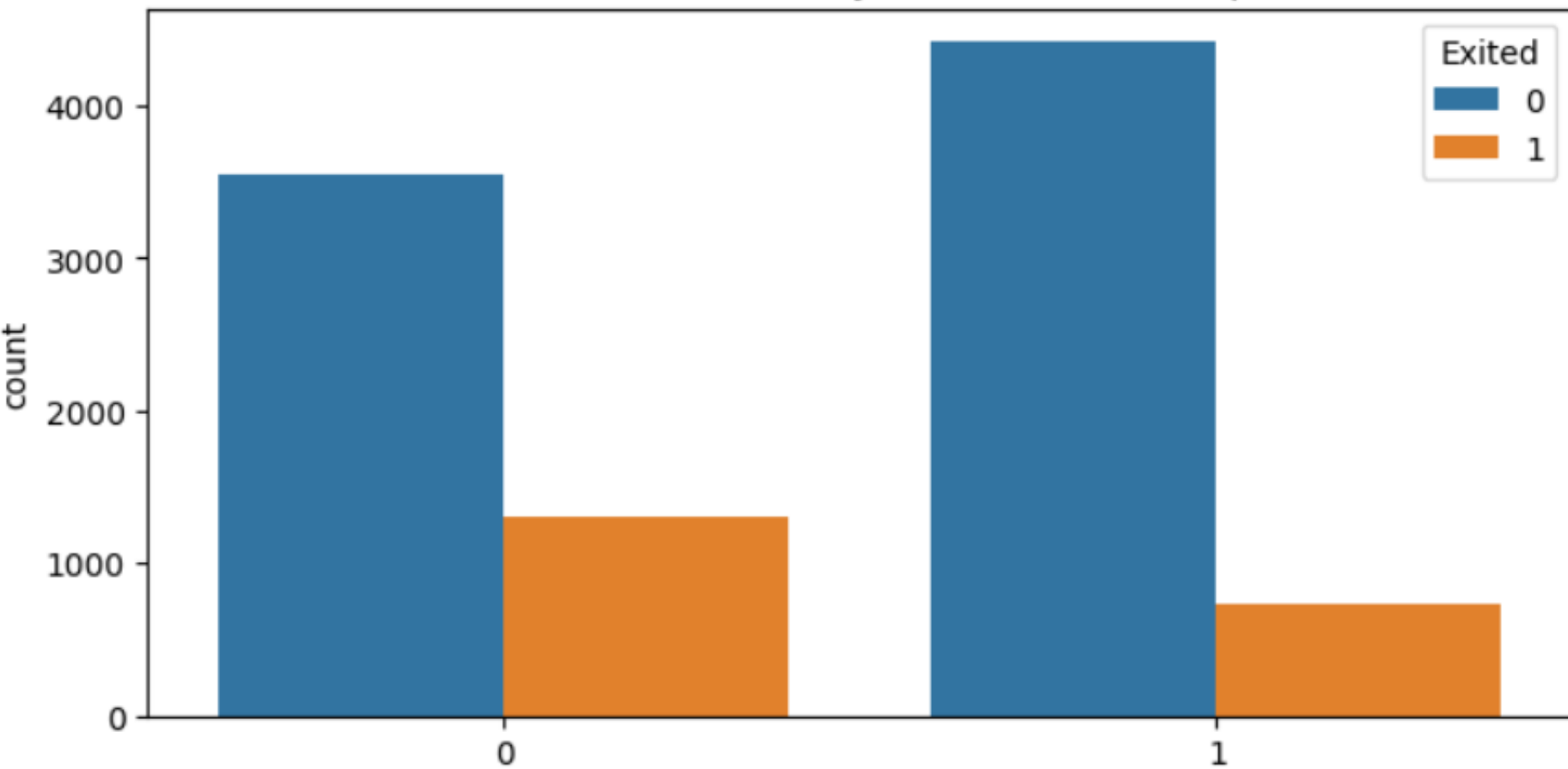


- DIAMOND
- PLATINUM
- SILVER
- GOLD

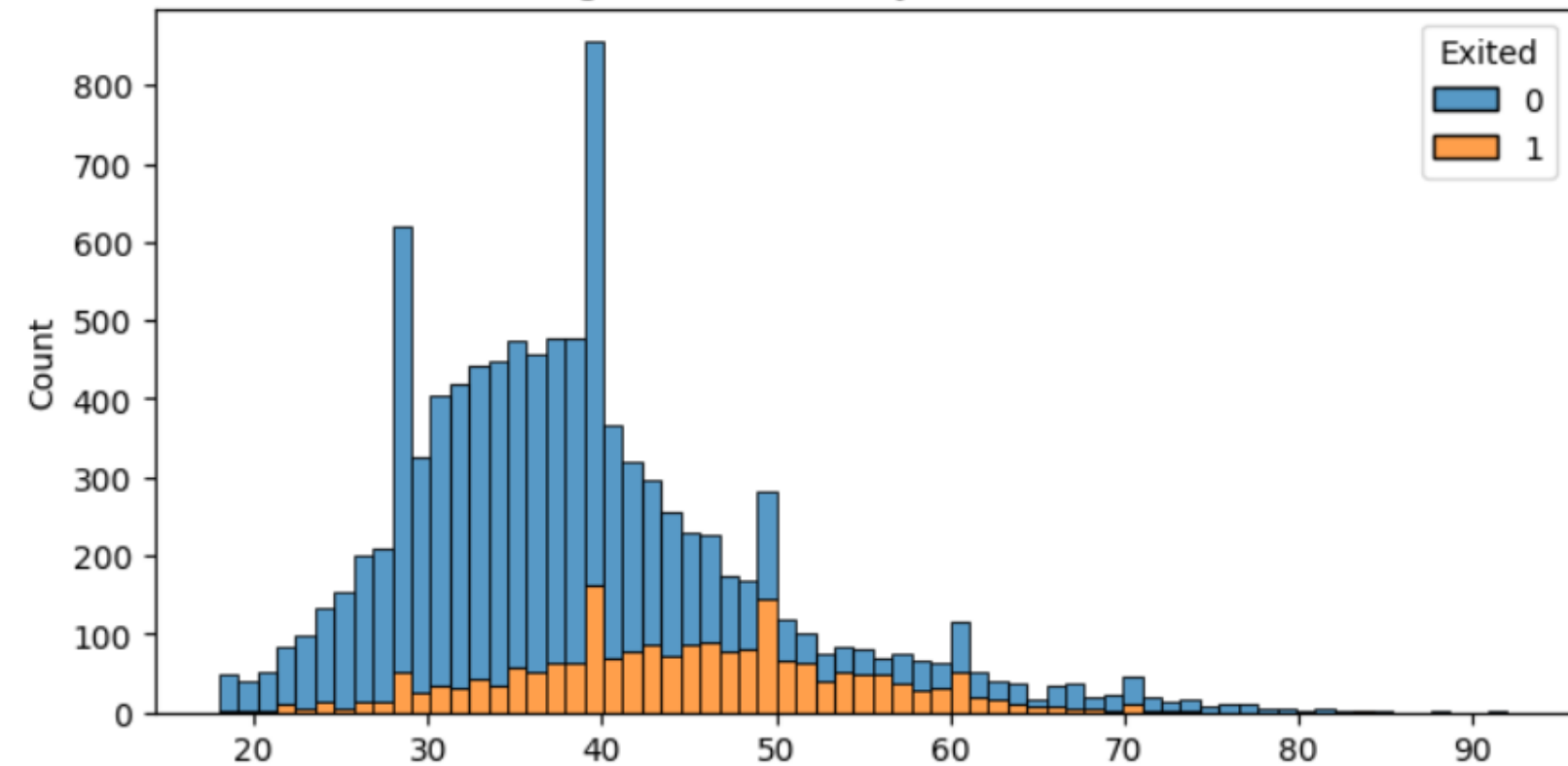
Gender Distribution and Churn rate



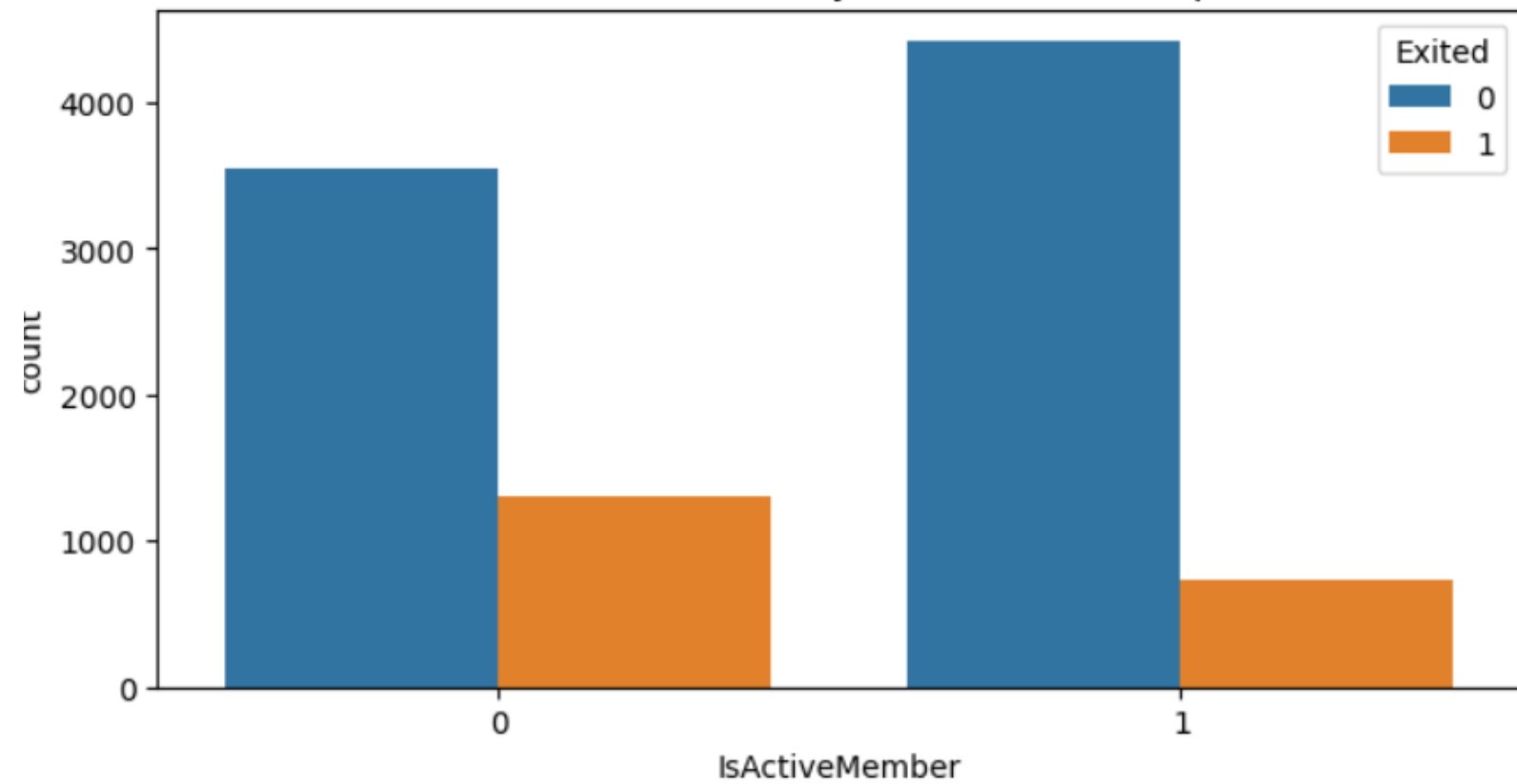
Churn Distribution by Active Membership



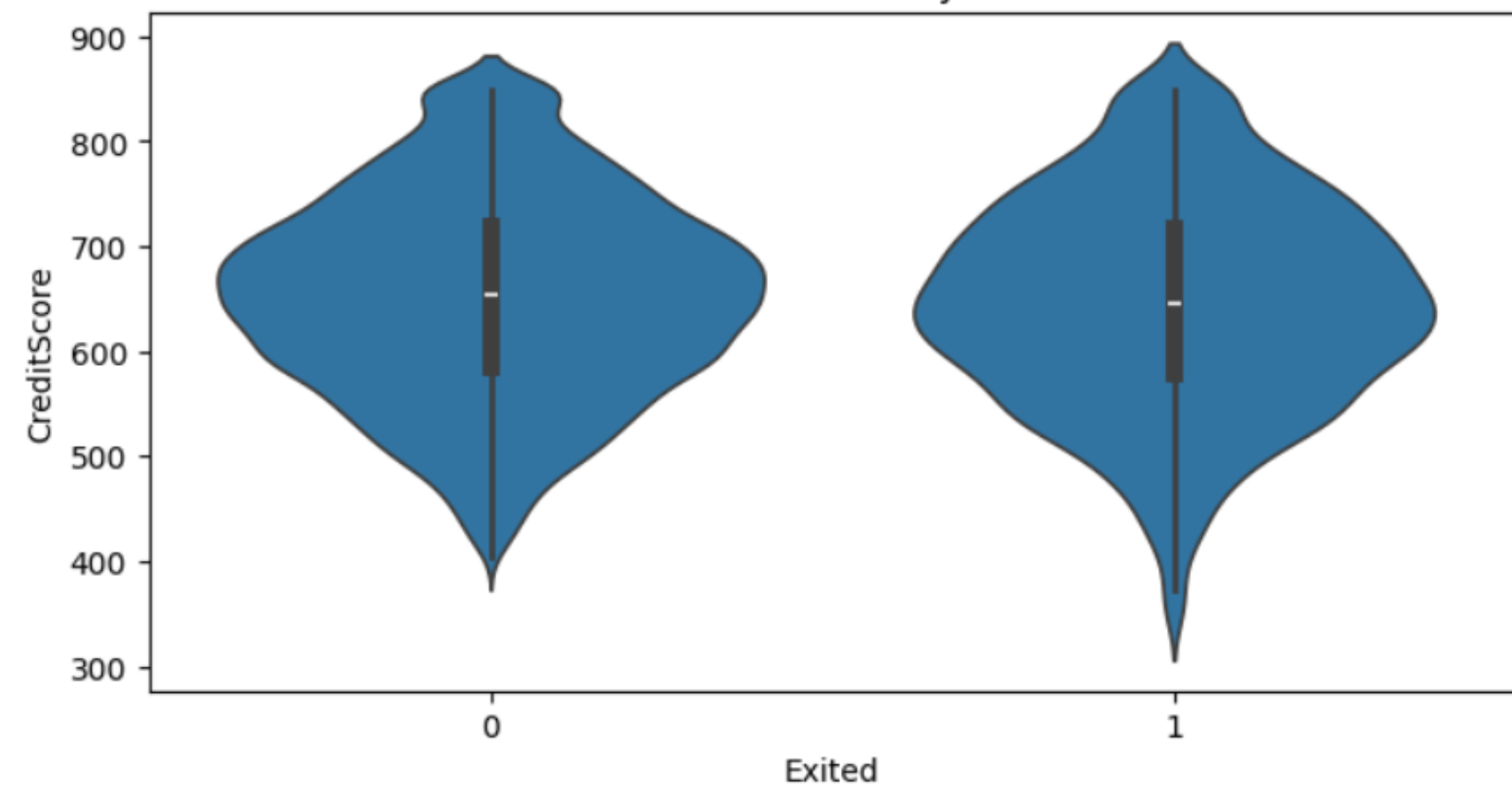
Age Distribution by Churn Status



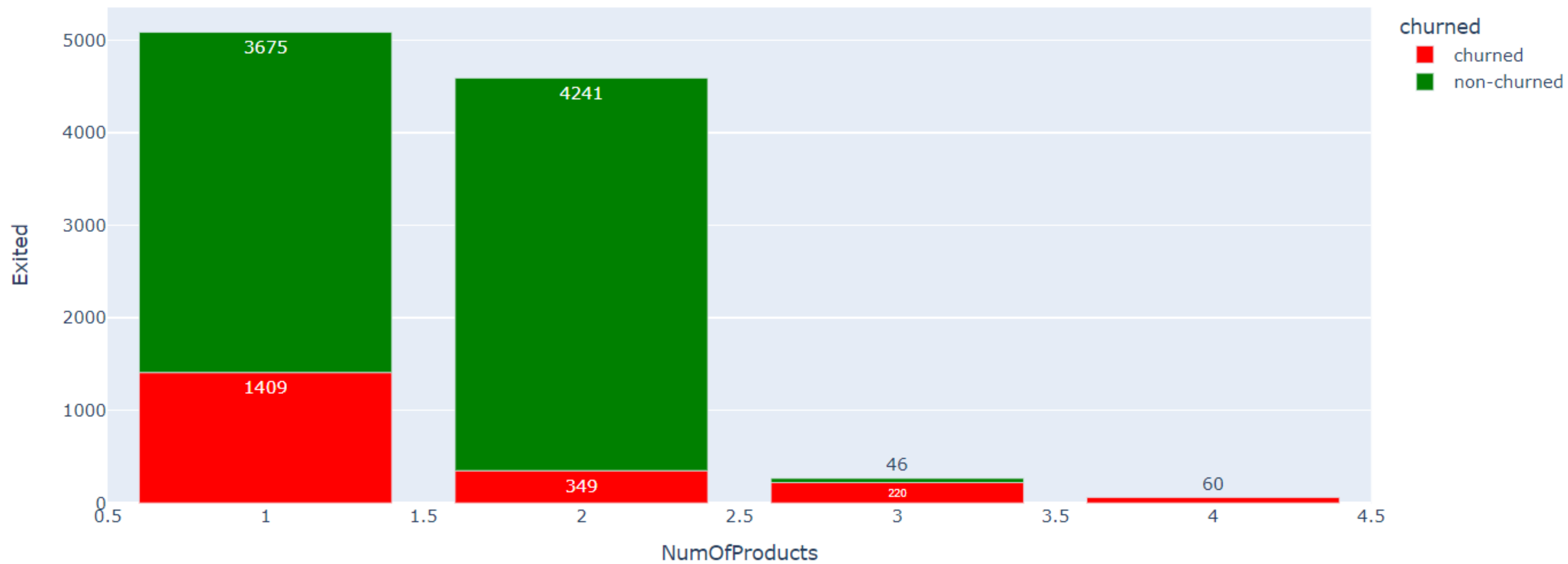
Churn Distribution by Active Membership



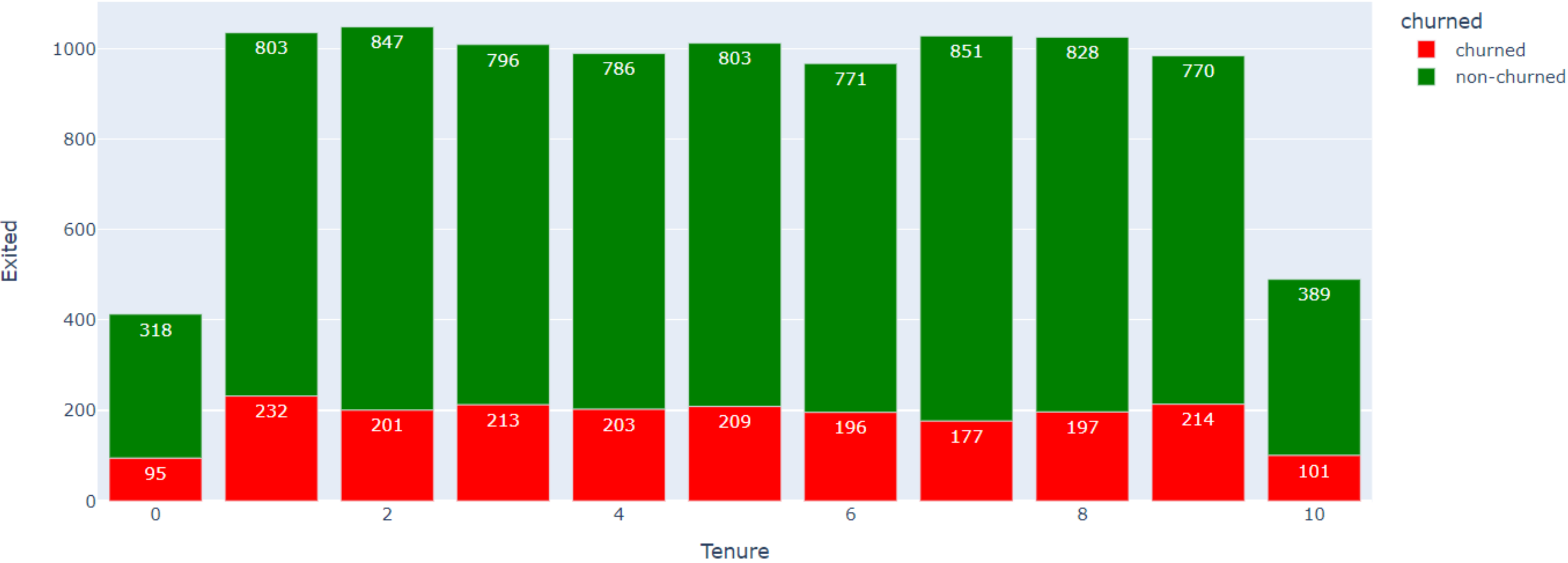
Credit Score Distribution by Churn Status



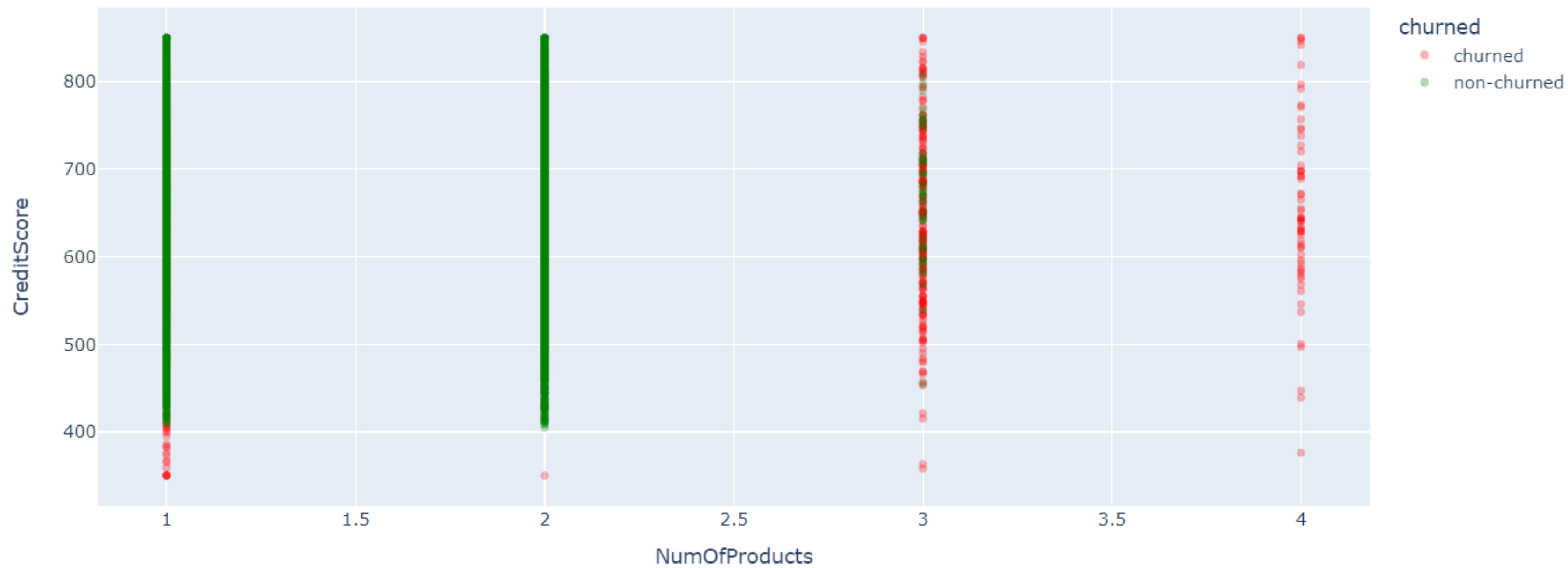
churned & non-churned based on number of purchased products



churned & non-churned based on Tenure

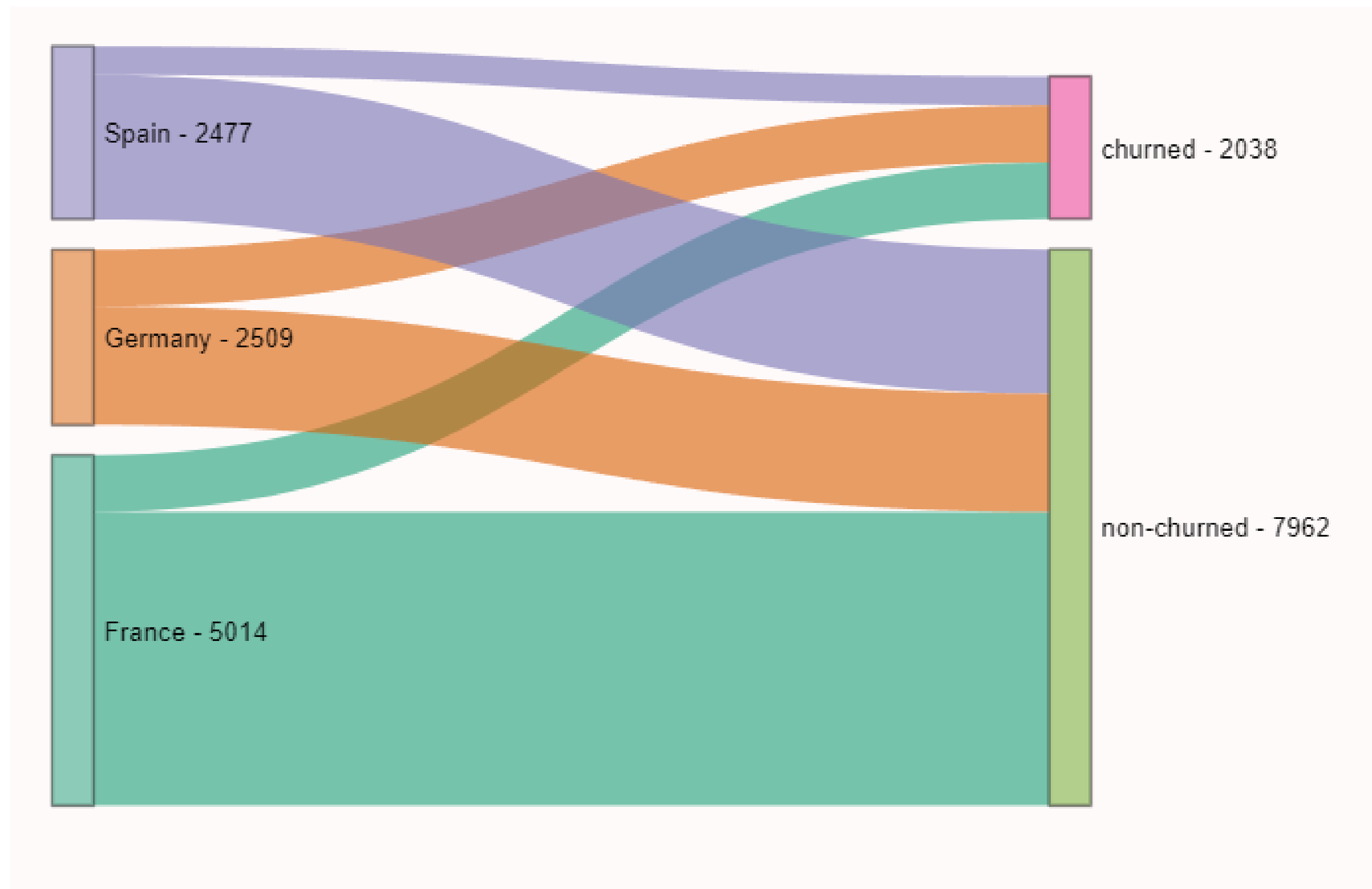


churned & non-churned based on the number of purchased products

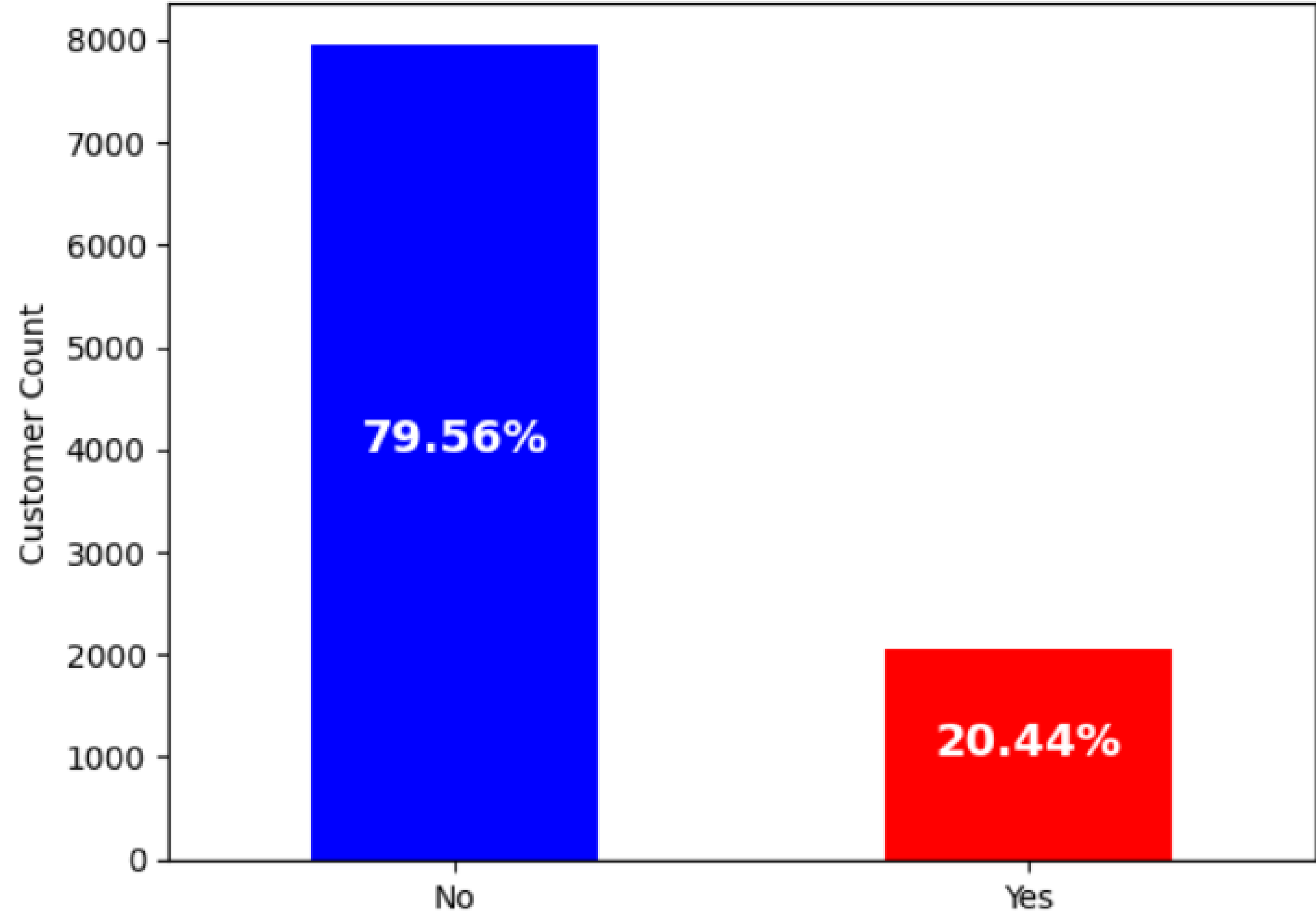


`aget` is deprecated and will be removed in version 1.4, use `transform_reference` instead.

Flow of churned and non-churned clients based on country



Amount of Complaining Customers



Exited		
Complain		
0	7952	4
1	10	2034

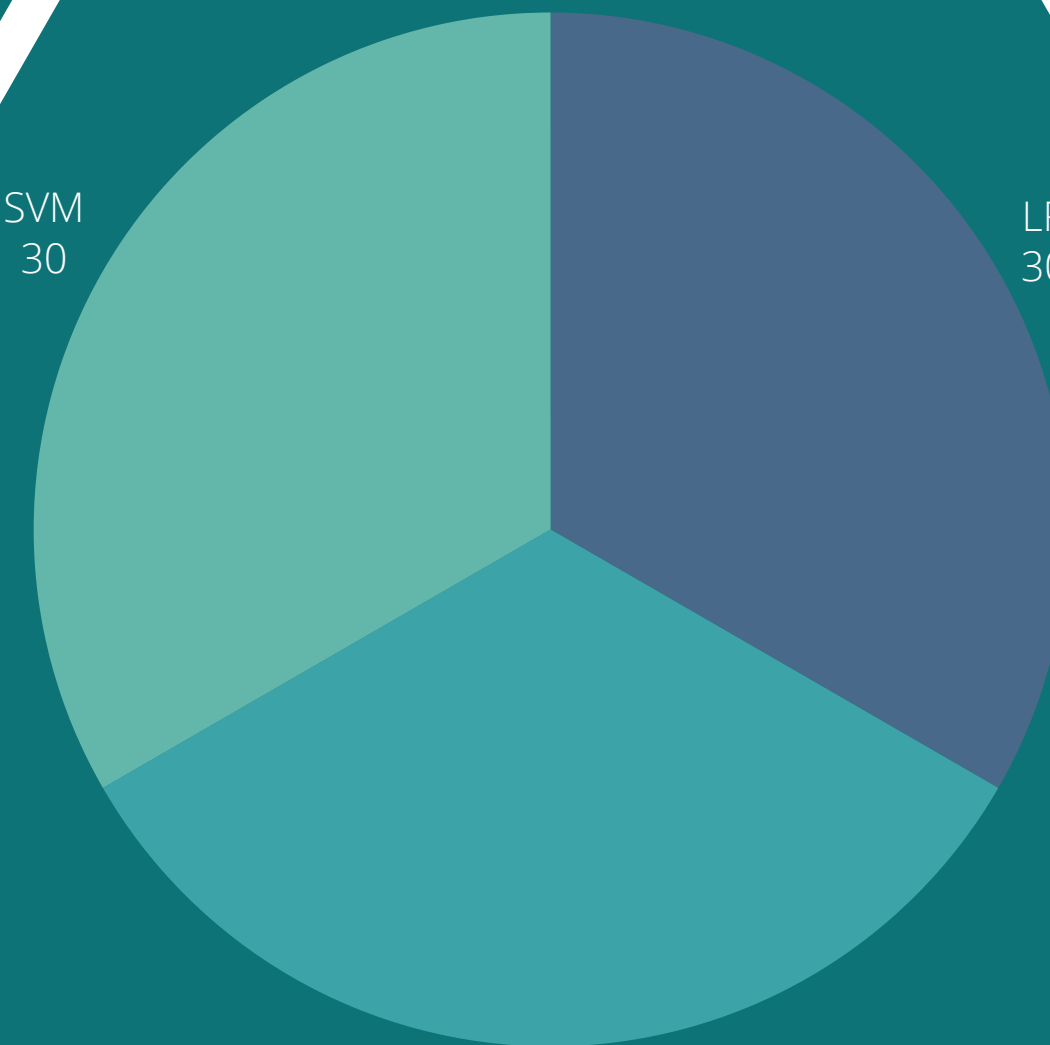
Model (s)

I explored three algorithms for this dataset: logistic regression, Support Vector Machine (SVM), and Random Forest Classifier. Each model underwent meticulous hyperparameter tuning to enhance accuracy and ensure and other evaluation metrics. Addressed potential overfitting and bias to uphold the integrity of the predictions.

SVM
30

LR
30

Random Forest Tree Classifier
30



First Model: Random Forest Tree Classifier

[+ Code](#)[+ Text](#)

✓
11s



```
from sklearn.ensemble import RandomForestClassifier
from sklearn import metrics
from sklearn.model_selection import train_test_split

#one hot encoding
bank_dummies = pd.get_dummies(df)

y = bank_dummies['Exited'].values
X = bank_dummies.drop(columns = ['Exited'])

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

model_rf = RandomForestClassifier(n_estimators=1000, oob_score= True, random_state=77, max_leaf_nodes=30)

model_rf.fit(X_train, y_train)

preds = model_rf.predict(X_test)

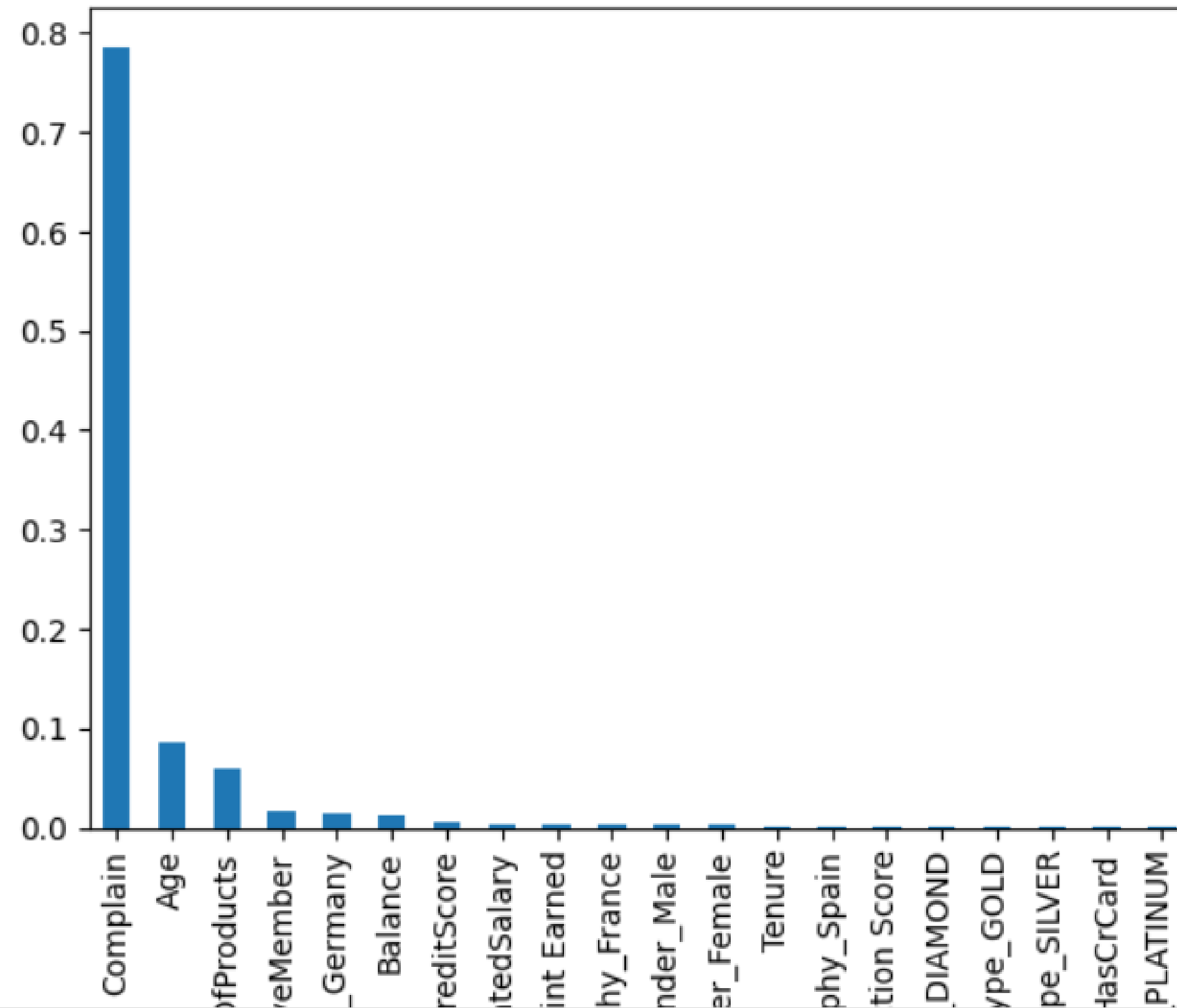
metrics.accuracy_score(y_test, preds)
```

0.9985

```
#Finding important features and plotting them based on their affect on the model
forest_importances = pd.Series(model_rf.feature_importances_, index=X_train.columns).sort_values(ascending=False)
forest_importances.plot(kind='bar')

#Complain causing bias for our dataset
```

<Axes: >



Finetuning

- Varied hyperparameters that will be showed in demo for higher evaluation score.
- Used different techniques for overfitting.



```
model_rf.fit(X_train, y_train)

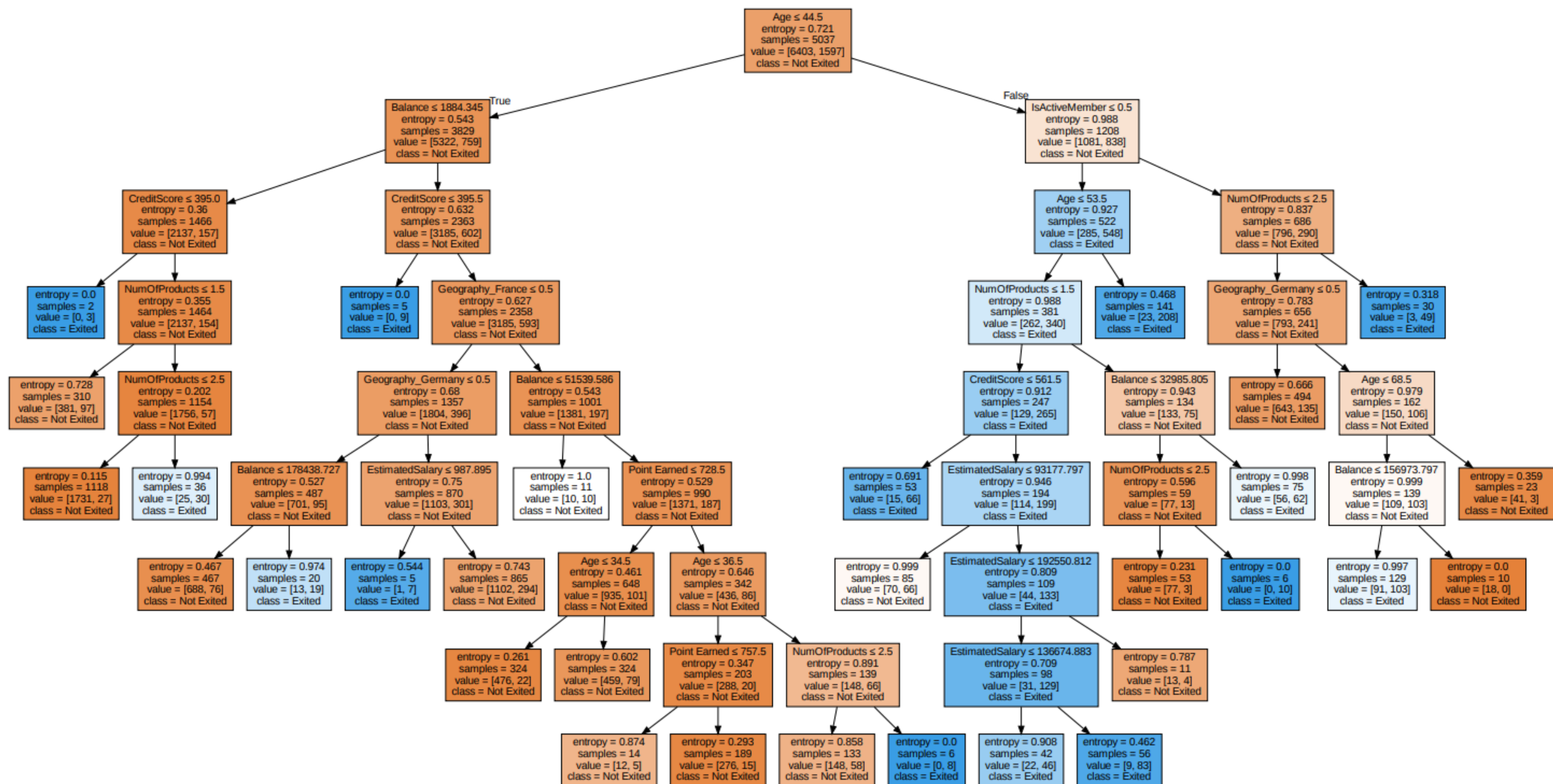
preds = model_rf.predict(X_test)

metrics.accuracy_score(y_test, preds)
```

0.84575

```
# Calculating the accuracy
accuracy = accuracy_score(y_test, preds)
print(f"Accuracy with entropy: {accuracy}")
```

Accuracy with entropy: 0.847



Second Model: Logistics Regression

```
▶ from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import train_test_split
from sklearn.metrics import classification_report, confusion_matrix, accuracy_score
from sklearn.preprocessing import StandardScaler
import pandas as pd

# Drop 'Complain' and set 'Exited' as the target variable
X = bank_dummies.drop(columns=['Exited', 'Complain'])
y = bank_dummies['Exited'].values

# Split the dataset into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=0)

# Feature scaling for better performance of the logistic regression
scaler = StandardScaler()
X_train = scaler.fit_transform(X_train)
X_test = scaler.transform(X_test)

logistic_classifier = LogisticRegression(random_state=0)
logistic_classifier.fit(X_train, y_train)
```

```

[[1528   67]
 [ 309   96]]
precision    recall  f1-score   support

     0       0.83     0.96     0.89     1595
     1       0.59     0.24     0.34      405

 accuracy          0.81     2000
 macro avg       0.71     0.60     0.61     2000
weighted avg       0.78     0.81     0.78     2000

```

Accuracy of the model: 0.812

```

# Apply SMOTE to oversample the minority class
smote = SMOTE(random_state=42)
X_smote, y_smote = smote.fit_resample(X, y)

```




```
[[3023  231]
 [ 410 2706]]
```

	precision	recall	f1-score	support
0	0.88	0.93	0.90	3254
1	0.92	0.87	0.89	3116
accuracy			0.90	6370
macro avg	0.90	0.90	0.90	6370
weighted avg	0.90	0.90	0.90	6370

Accuracy: 0.8993720565149137

Third Model: SVM

✓
0s

```
▶ from sklearn.svm import SVC
  from sklearn.model_selection import train_test_split
  from sklearn.metrics import accuracy_score
  from imblearn.over_sampling import SMOTE
  from imblearn.pipeline import Pipeline
  from sklearn.preprocessing import StandardScaler
```

```
  # Calculate accuracy
  accuracy = accuracy_score(y_test, predictions)
  print(classification_report(y_test, predictions))
  print(f"SVM Accuracy without SMOTE: {accuracy}")
```



	precision	recall	f1-score	support
0	1.00	1.00	1.00	1607
1	1.00	1.00	1.00	393
accuracy			1.00	2000
macro avg	1.00	1.00	1.00	2000
weighted avg	1.00	1.00	1.00	2000

SVM Accuracy without SMOTE: 0.999

```
# Split the dataset into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Train the pipeline
pipeline.fit(X_train, y_train)

# Make predictions
predictions = pipeline.predict(X_test)

# Calculate accuracy
accuracy = accuracy_score(y_test, predictions)
print(f"SVM Accuracy with SMOTE: {accuracy}")
```



SVM Accuracy with SMOTE: 0.8035

Difficulties

- Finding out clear correlation between Complain and Exit. (Dataset flaw)
- Encoding process.
- Model finetuning.



THYNK UNLIMITED

THANKS.
Q&A