

How can we accurately predict churn rates for Netflix users?

Stakeholders: Netflix , Netflix users, Investors, Competitors and Researchers

What is a churn rate?

The annual percentage rate at which customers stop subscribing to a service or employees leave a job.

$$\left(\frac{\text{Lost Customers}}{\text{Total Customers at the Start of Time Period}} \right) \times 100$$

Churn rate for this dataset is around ~80% meaning **4 out of 5** customers don't stay for a full year or more.

Why is this important?

This is important because multiple parties are involved, it helps **Netflix** know how many of their users they are retaining, it helps the **users** know that their streaming service is interested in keeping them as customers and therefore interested in enhancing user experience and lastly is important to us **researchers** to know what makes a person to stay as paying customer or not in the biggest streaming platform in the world, is it their age? Their country? Their income?

Churn rates were defined in this matter:

1. Customer stayed for $x \geq 365$ days then customer **did not churn**.
2. Customer stayed for $x < 365$ days then customer **did churn**.

Raw Dataset + Data Cleaning

	User ID	Subscription Type	Monthly Revenue	Join Date	Last Payment Date	Country	Age	Gender	Device	Plan Duration
0	1	Basic	10	15-01-22	10-06-23	United States	28	Male	Smartphone	1 Month
1	2	Premium	15	05-09-21	22-06-23	Canada	35	Female	Tablet	1 Month
2	3	Standard	12	28-02-23	27-06-23	United Kingdom	42	Male	Smart TV	1 Month
3	4	Standard	12	10-07-22	26-06-23	Australia	51	Female	Laptop	1 Month
4	5	Basic	10	01-05-23	28-06-23	Germany	33	Male	Smartphone	1 Month

Important notes about dataset:

1. The **granularity** of the data is costumer, meaning that each row represents one individual customer, therefore we need to drop user ID in the preprocessing part.
2. Since Netflix only offers subscription on a monthly basis (they don't offer multiple months or yearly subscription) we can drop Plan Duration since all entries the same. ('1 Month')
3. Instead of using Join Date and Last Payment Dates as features we transformed it to a binary feature ('Yes', 'No') called '**Churned**'

Methods

Data Preprocessing/ Data Cleaning

```
1 # CREATING A CLEANED DATAFRAME WITH ONLY USEFUL FEATURES
2 df['Last Payment Date'] = pd.to_datetime(df['Last Payment Date'], format = "%d-%m-%y")
3
4 df['Join Date'] = pd.to_datetime(df['Join Date'], format = "%d-%m-%y")
5
6 max_value = max(df['Last Payment Date'])
7 min_value = min(df['Last Payment Date'])
8
9
10 df['Churned'] = df.apply(lambda r : 'No' if (r['Last Payment Date'] - r['Join Date']).days >= 365 else 'Yes', axis=1)
11 data_training = df.drop(columns= ['Plan Duration', 'User ID', 'Join Date', 'Last Payment Date'])
12 data_training
```

Visualize						
Subscription Type	Monthly Revenue	Country	Age	Gender	Device	Churned
Basic 40%	10 - 15	United States 18%	26 - 51	Female 50.3%	Laptop 25.4%	Yes 82.2%
Standard 30.7%		Spain 18%		Male 49.7%	Tablet 25.3%	No 17.8%
Premium 29.3%		8 others 63.9%			2 others 49.2%	
0 Basic	10	United States	28	Male	Smartphone	No
1 Premium	15	Canada	35	Female	Tablet	No
2 Standard	12	United Kingdom	42	Male	Smart TV	Yes
3 Standard	12	Australia	51	Female	Laptop	Yes
4 Basic	10	Germany	33	Male	Smartphone	Yes
5 Premium	15	France	29	Female	Smart TV	No
6 Standard	12	Brazil	46	Male	Tablet	No
7 Basic	10	Mexico	39	Female	Laptop	Yes
8 Standard	12	Spain	37	Male	Smartphone	Yes
9 Premium	15	Italy	44	Female	Smart TV	Yes

1. Putting all dates into a 'readable' format using `pd.to_datetime`.
2. Using a lambda to determine churned rates.
3. Dropping 4 features: Plan Duration, User ID, Join Date, Last Payment Date.

K-means

Why K-means?

Key Points:

- **Clustering Goal:**
 - Segment Netflix users into distinct groups based on behavioral and demographic data for actionable insights.
- **Advantages of K-Means:**
 - Simplicity: Easy to implement and computationally efficient for the dataset size.
 - Interpretability: Clear cluster assignments make it straightforward to analyze user segments.
 - Scalability: Handles a medium-sized dataset (2,500 rows) effectively while maintaining speed.
- **Alignment with Project Goals:**
 - Groups users with similar attributes to understand churn patterns and identify high-risk segments.
 - Provides actionable insights for personalized marketing and retention strategies.

K-means

Data Pre-Processing

Steps Taken:

- Cleaned missing data and converted date fields to datetime.
- Applied one-hot encoding for categorical features (country, device).
- Scaled numerical features (e.g., age, revenue) using StandardScaler.

Purpose: Ensure the dataset was prepared for clustering and machine learning.

```
1 netflix_data['Join Date'] = pd.to_datetime(netflix_data['Join Date'], format='%d-%m-%y')
2 netflix_data['Last Payment Date'] = pd.to_datetime(netflix_data['Last Payment Date'], format='%d-%m-%y')
3
4 duplicates_count = netflix_data.duplicated().sum()
5
6 updated_data_types = netflix_data.dtypes
7 updated_data_types, duplicates_count
```



This output has been hidden. [Show it.](#)

```
1 encoded_data = pd.get_dummies(netflix_data, columns=[
2 |     'Subscription Type', 'Country', 'Gender', 'Device', 'Plan Duration'], drop_first=True)
3
4 encoded_data.head()
```



This output has been hidden. [Show it.](#)

```
1 from sklearn.preprocessing import StandardScaler
2
3 scaler = StandardScaler()
4 numerical_columns = ['Age', 'Monthly Revenue']
5 encoded_data[numerical_columns] = scaler.fit_transform(encoded_data[numerical_columns])
6 final_data = encoded_data.drop(columns=['User ID', 'Join Date', 'Last Payment Date'])
7
8 final_data.head()
```



This output has been hidden. [Show it.](#)

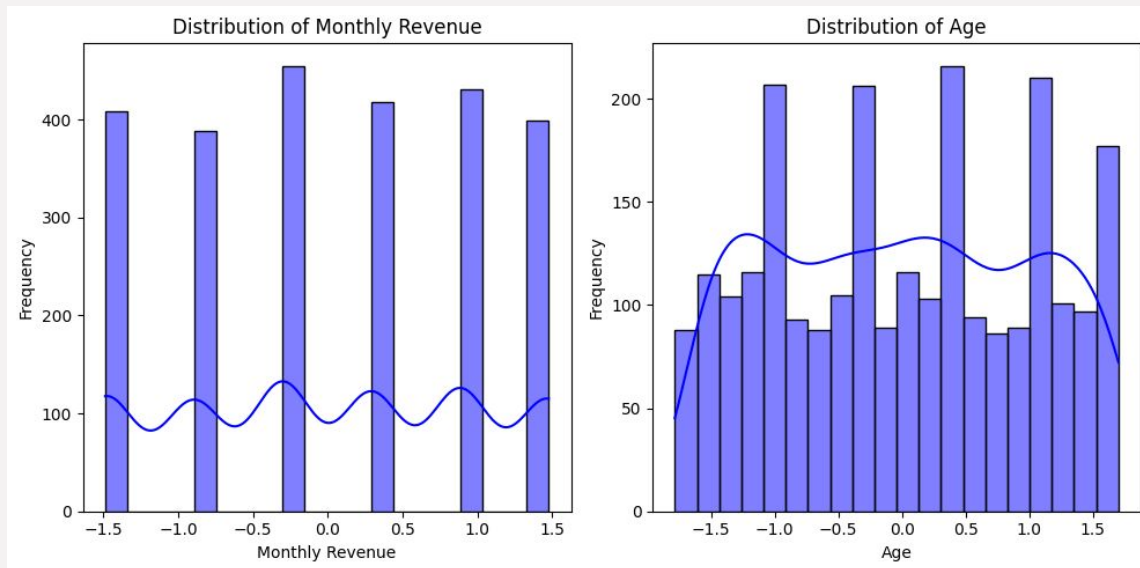
K-means

Data Visualizations

- **Visual Insights:**

- Histograms: Highlighted distributions of age and monthly revenue.
- Heatmap: Showed relationships between features.
- PCA and t-SNE: Revealed distinct user clusters and dimensionality reduction.

- **Implication:** Identified separable user groups, supporting clustering feasibility.



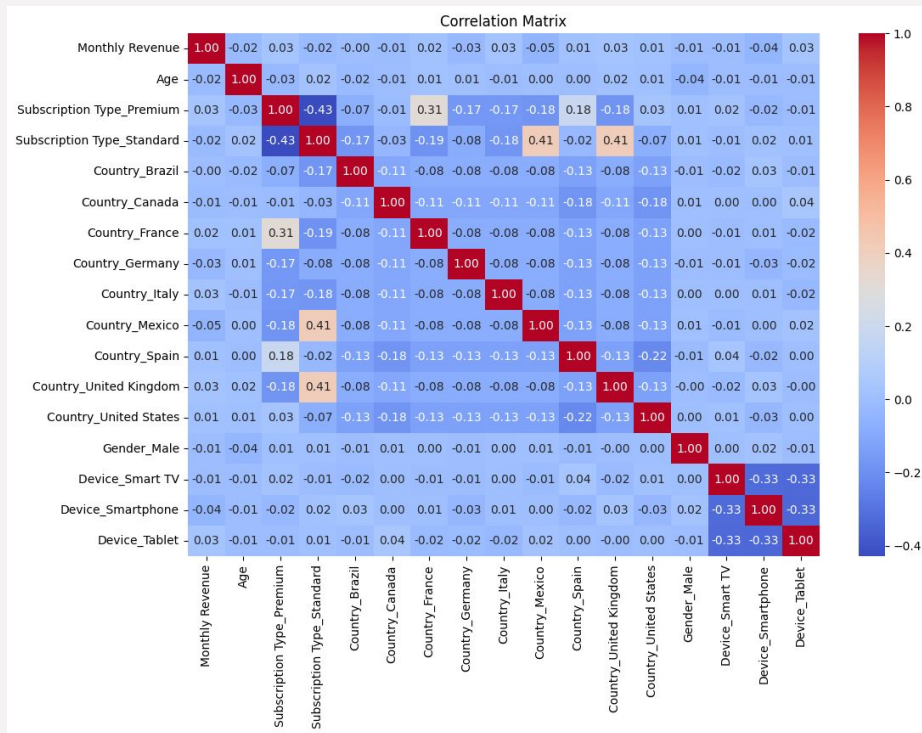
K-means

Data Visualizations

- **Visual Insights:**

- Histograms: Highlighted distributions of age and monthly revenue.
- Heatmap: Showed relationships between features.
- PCA and t-SNE: Revealed distinct user clusters and dimensionality reduction.

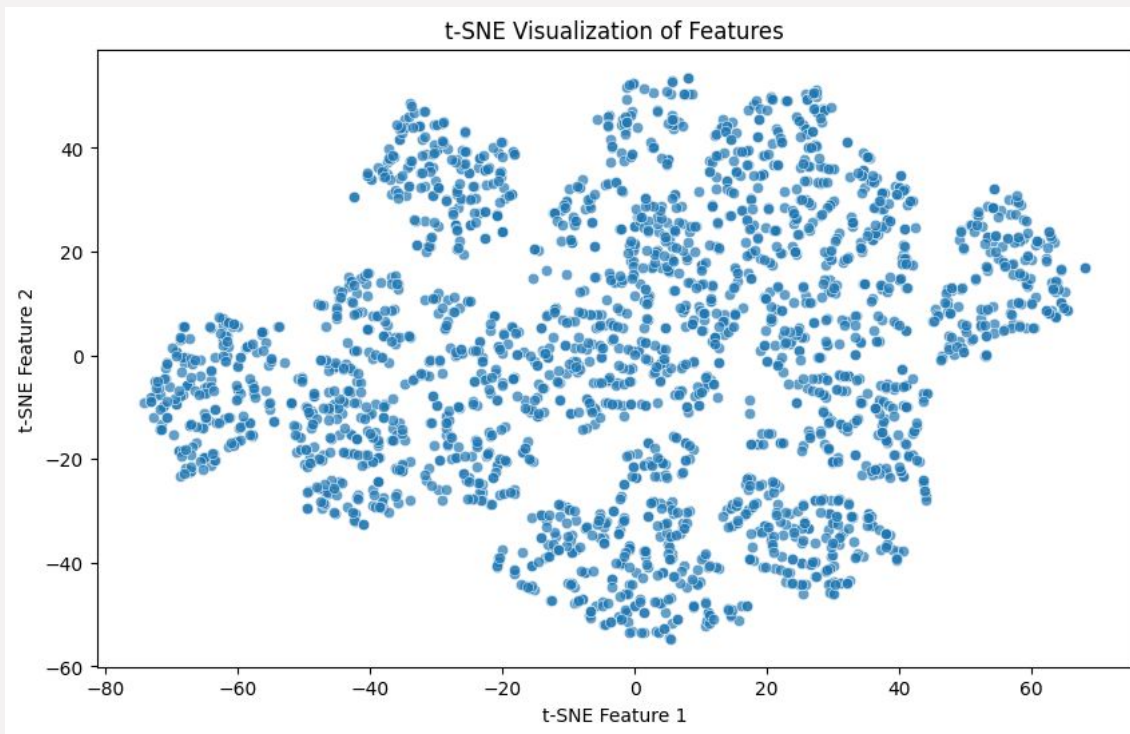
- **Implication:** Identified separable user groups, supporting clustering feasibility.



K-means

Data Visualizations

- **Visual Insights:**
 - Histograms: Highlighted distributions of age and monthly revenue.
 - Heatmap: Showed relationships between features.
 - PCA and t-SNE: Revealed distinct user clusters and dimensionality reduction.
- **Implication:** Identified separable user groups, supporting clustering feasibility.



K-means

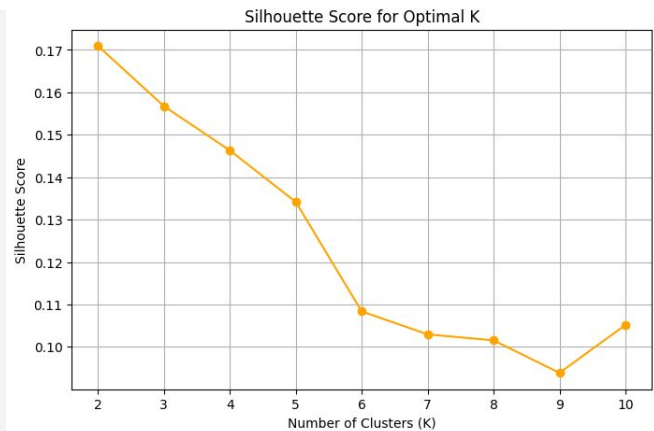
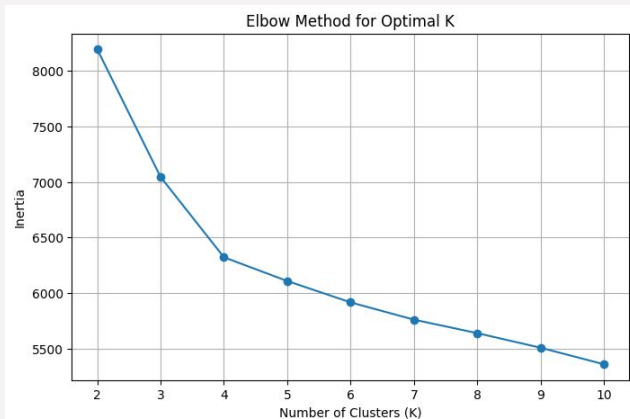
Approach:

- Used Elbow Method and Silhouette Scores to determine optimal clusters.
- Chose 4 clusters for balance between interpretability and separability.

Results:

- Clear segmentations identified through PCA visualizations.
- Clusters characterized by attributes like revenue and engagement.

K-Means Clustering



K-means

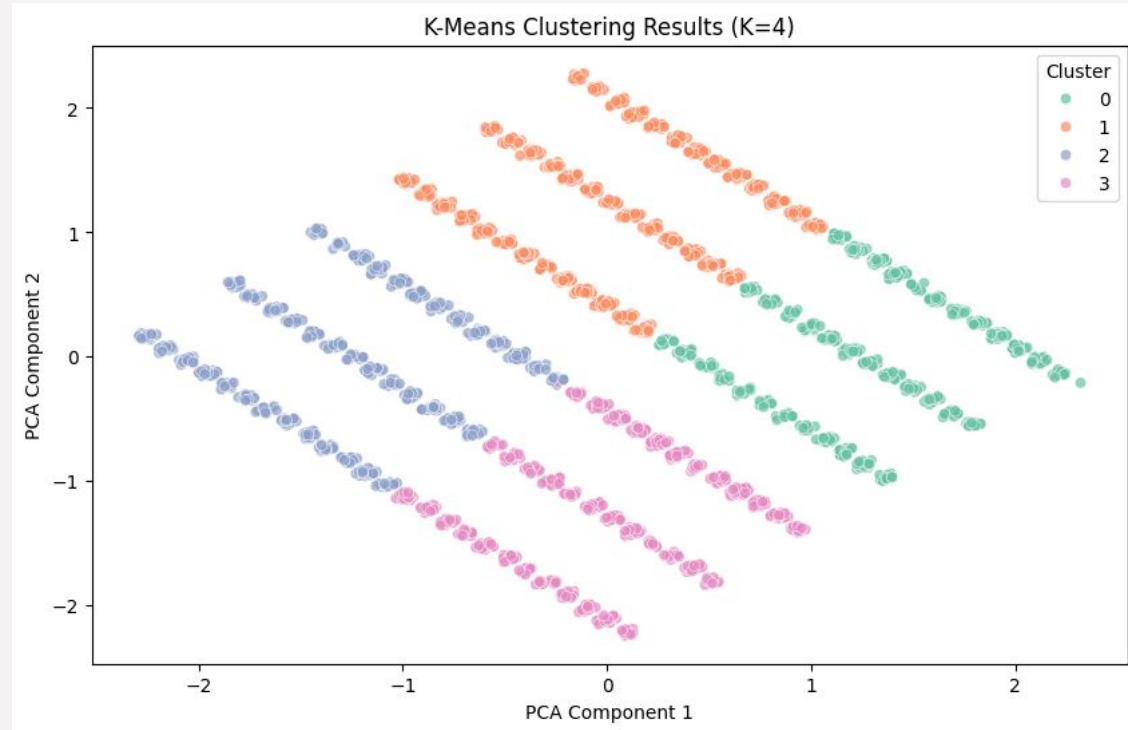
Approach:

- Used Elbow Method and Silhouette Scores to determine optimal clusters.
- Chose 4 clusters for balance between interpretability and separability.

Results:

- Clear segmentations identified through PCA visualizations.
- Clusters characterized by attributes like revenue and engagement.

K-Means Clustering



K-means

Methodology:

- Integrated Isolation Forest for churn detection.
- Analyzed cluster-level churn prediction rates.

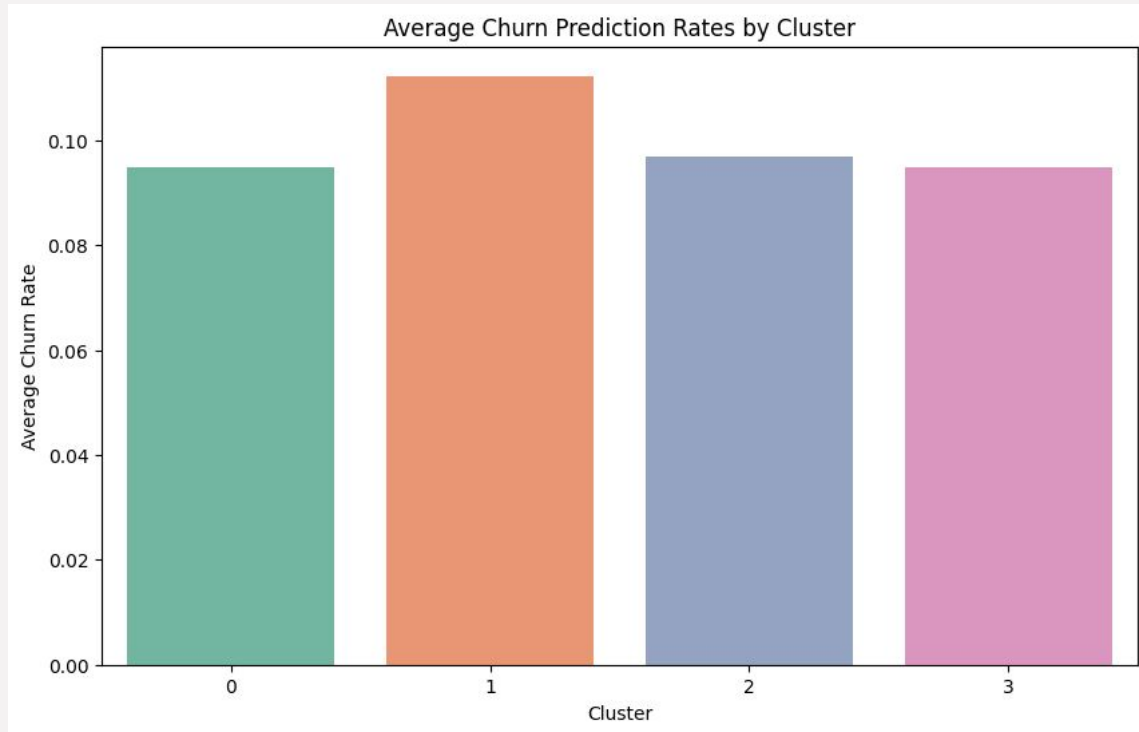
Findings:

- Cluster 1 exhibited the highest churn rate (~10.9%).
- Other clusters showed lower churn rates (~9.5%).

Business Insight:

- Cluster 1 indicates at-risk customers needing targeted retention strategies.

Churn Prediction



K-means

Was K-means Successful?

Key Takeaways:

- K-Means effectively segmented users into meaningful groups.
- Combined clustering and isolation models identified churn risk areas.

Strengths:

- K-Means successfully segmented the user base into meaningful clusters with clear demographic and behavioral differences.
- These clusters enable targeted analysis of churn risk, which can inform retention strategies.

Limitations:

- K-Means alone does not inherently predict churn; it groups users based on feature similarity.
- The addition of an unsupervised model like Isolation Forest helped assign churn probabilities to these clusters, extending the functionality of K-Means.

Decision Trees

Decision Trees are a foundational supervised learning algorithm, used here for churn prediction.

Conducted hyperparameter tuning with GridSearchCV:

- Tested different values for `max_depth`, `min_samples_split`, `min_samples_leaf`, and `criterion`.
- Best parameters identified: `{'criterion': 'gini', 'max_depth': 5, 'min_samples_leaf': 10, 'min_samples_split': 2}`.

Achieved an accuracy of 66%, indicating room for improvement in predictive performance.

Decision Trees are interpretable, allowing us to pinpoint specific decision paths influencing churn.

```
1 from sklearn.model_selection import GridSearchCV
2
3 params = {
4     'max_depth': [5, 10, 15, None],
5     'min_samples_split': [2, 10, 20],
6     'min_samples_leaf': [1, 5, 10],
7     'criterion': ['gini', 'entropy']
8 }
9
10 grid_search = GridSearchCV(DecisionTreeClassifier(random_state=42), param_grid=params, cv=5, scoring='accuracy')
11 grid_search.fit(X_train, y_train)
12
13 print("Best parameters for decision trees:", grid_search.best_params_)
14 best_model = grid_search.best_estimator_
```


Best parameters for decision trees: {'criterion': 'gini', 'max_depth': 5, 'min_samples_leaf': 10, 'min_samples_split': 2}

```
1 from sklearn.tree import DecisionTreeClassifier
2 from sklearn.metrics import accuracy_score, classification_report
3 from sklearn.model_selection import train_test_split
4
5 X = pd.get_dummies(data_training.drop(columns=['Churned']))
6 y = data_training['Churned']
7
8 y = y.map({'Yes':1, 'No':0})
9 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42, stratify=y)
10
11 dt_model = DecisionTreeClassifier(max_depth=5, random_state=42, class_weight='balanced')
12 dt_model.fit(X_train, y_train)
13
14 y_pred = dt_model.predict(X_test)
15 accuracy = accuracy_score(y_test, y_pred)
16
17 print(f"Decision tree accuracy: {accuracy:.2f}")
```


Decision tree accuracy: 0.66

Random Forest

Random Forest builds on Decision Trees by combining multiple trees for better accuracy and robustness.

Hyperparameter tuning involved a larger search space for:

- Best parameters: {'max_depth': 10, 'max_features': 'sqrt', 'min_samples_leaf': 1, 'min_samples_split': 2, 'n_estimators': 100}.

Improved accuracy to **82.2%**, showing the value of ensemble methods for predicting churn.

Random Forest reduces overfitting while providing feature importance scores for deeper insights.

```
1 param_grid = {
2     'n_estimators': [100, 200, 300],
3     'max_depth': [5, 10, None],
4     'min_samples_split': [2, 10],
5     'min_samples_leaf': [1, 5, 10],
6     'max_features': ['sqrt', 'log2', None]
7 }
8
9 grid_search_rf = GridSearchCV(
10     RandomForestClassifier(random_state=42), param_grid, cv=5, scoring='accuracy'
11 )
12 grid_search_rf.fit(X_train, y_train)
13
14 print("Best parameters for random forest:", grid_search_rf.best_params_)
15 best_rf_model = grid_search_rf.best_estimator_
```

Best parameters for random forest: {'max_depth': 10, 'max_features': 'sqrt', 'min_samples_leaf': 1, 'min_samples_split': 2, 'n_estimators': 100}

```
1 from sklearn.ensemble import RandomForestClassifier
2
3 rf_model = RandomForestClassifier(n_estimators=100, max_depth=10, random_state=42, min_samples_leaf=1, min_samples_split=2)
4 rf_model.fit(X_train, y_train)
5 y_pred_rf = rf_model.predict(X_test)
6
7 print("Random forest accuracy:", accuracy_score(y_test, y_pred_rf))
```

Random forest accuracy: 0.822

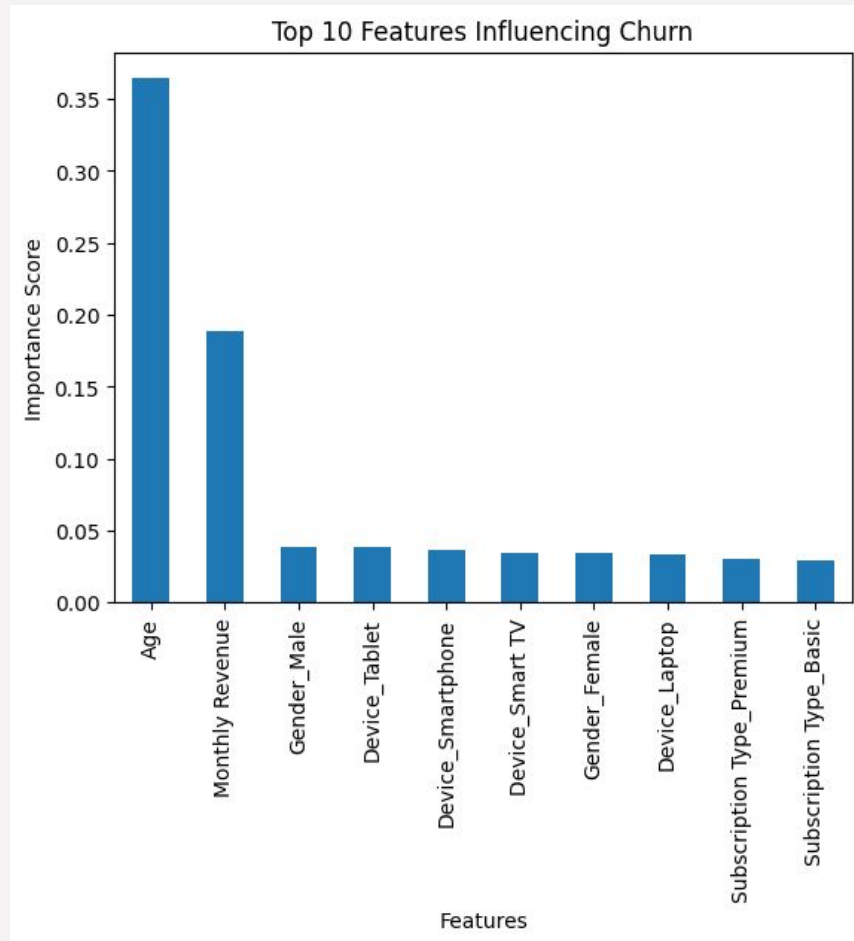
Top 10 Features Influencing Churn Rate

Feature importance from the Random Forest model highlights the key drivers of churn:

- **Age** is the most significant factor, followed by **Monthly Revenue**.
- Device type and subscription type also have notable contributions, indicating differences in user behavior.

Insights suggest that user demographics (e.g., Age) and revenue patterns are critical churn predictors.

Understanding feature importance informs targeted interventions and prioritizes areas for improvement.



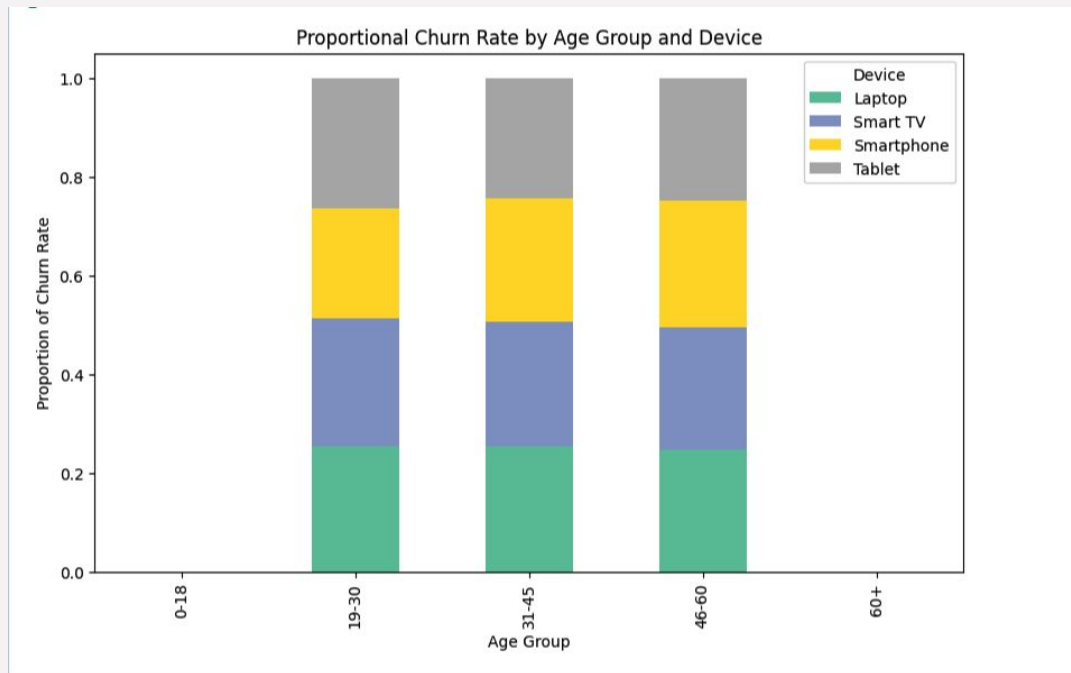
Churn Rate by Age Group and Device

This visualization explores churn distribution across **age groups** and **devices**:

- Mid-age users (31-45) show higher churn proportions on **Smartphones**.
- Older users (46-60) exhibit balanced churn across **Tablets** and **Laptops**.

Device-specific churn highlights potential usability or engagement issues.

Insight: Retention strategies should be customized based on age and device usage patterns.



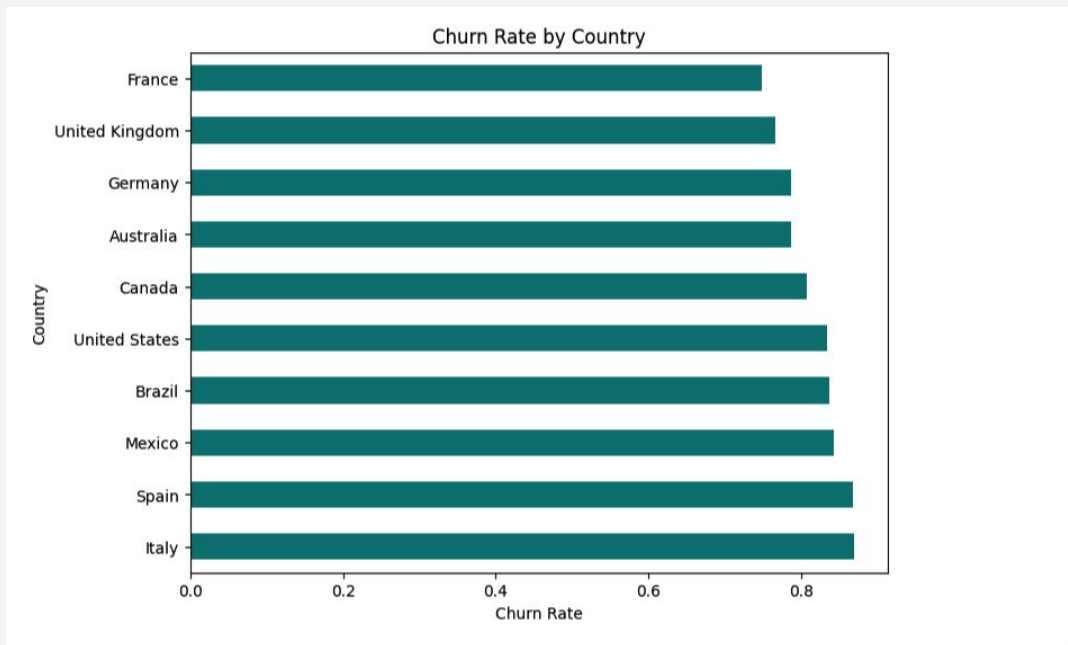
Churn Rate by Country

Horizontal bar chart compares churn rates across countries:

- Countries like **Italy**, **Spain**, and **Mexico** have the highest churn rates.
- Countries like **France** and the **United Kingdom** show relatively lower churn.

Regional differences suggest cultural, content, or pricing factors influencing churn.

Insights: support localized strategies, such as tailored content or competitive pricing, in high-churn regions.



Neural Networks

```
1 # NEURAL NETWORK 1
2 X = data_training[['Age']].values
3 y = data_training['Churned']
4
5 X_train, X_test, Y_train, Y_test = train_test_split(X, y, test_size=0.2)
6
7
8 le = LabelEncoder()
9
10 Y_train = le.fit_transform(Y_train)
11
12 Y_test = le.fit_transform(Y_test)
13
14
15
16 #neural network
17 clf = MLPClassifier(hidden_layer_sizes = 100, solver='lbfgs', activation = 'relu')
18
19 clf.fit(X_train, Y_train)
20 yp = clf.predict(X_test)
21
22 accuracy = accuracy_score(Y_test, yp)
23 accuracy
```



0.886

```
1 #NEURAL NETWORK USING CONTRY AS A FEATURE
2
3 #put the data in binary
4
5 X = pd.get_dummies(data_training['Country'])
6 y = data_training['Churned']
7
8 X_train, X_test, Y_train, Y_test = train_test_split(X, y, test_size=0.2)
9
10
11 le = LabelEncoder()
12
13 Y_train = le.fit_transform(Y_train)
14
15 Y_test = le.fit_transform(Y_test)
16
17
18 #sample_weights = compute_sample_weight(class_weight= 'balanced', y=Y_train)
19 #neural network
20 clf = MLPClassifier(hidden_layer_sizes = (10,), solver='lbfgs', activation = 'relu')
21
22 clf.fit(X_train, Y_train)
23 yp = clf.predict(X_test)
24
25 accuracy = accuracy_score(Y_test, yp)
26 accuracy
```



0.884

Starting out with one-feature Neural Networks using only Age in NN1 and Country in NN2 (Country needed to be one-hot-encoded since it's a nominal feature)

Some annotations:

-Hidden layer size didn't seemed to increase/decrease accuracy of the model.

-After some trails it was determined optimal train/test split was 80/20.

Neural Networks

```
1 #NEURAL NETWORK USING Subscription Type AND Monthly Revenue
2
3 from sklearn.neural_network import MLPClassifier
4 from sklearn.preprocessing import LabelEncoder
5
6 from sklearn.metrics import accuracy_score
7 from sklearn.metrics import classification_report
8
9
10
11 #put the data in binary
12
13 X1= pd.get_dummies(data_training['Subscription Type'])
14 X2 = pd.get_dummies(data_training['Monthly Revenue'])
15 #X3 = pd.get_dummies(data_training['Age'])
16 #X4 = pd.get_dummies(data_training['Subscription Type'])
17 X = np.concatenate([X1, X2], axis=1)
18
19 y = data_training['Churned']
20
21 X_train, X_test, Y_train, Y_test = train_test_split(X, y, test_size=0.2)
22
23
24 le = LabelEncoder()
25
26 Y_train = le.fit_transform(Y_train)
27
28 Y_test = le.fit_transform(Y_test)
29
30
31
32 #neural network
33 clf = MLPClassifier(hidden_layer_sizes = (10,), solver='lbfgs', activation = 'relu')
34
35 clf.fit(X_train, Y_train)
36 yp = clf.predict(X_test)
37
38 accuracy = accuracy_score(Y_test, yp)
39 accuracy
```

0.822

```
1
2 #NEURAL NETWORK USING COUNTRY AND GENDER
3
4 from sklearn.neural_network import MLPClassifier
5 from sklearn.preprocessing import LabelEncoder
6
7 from sklearn.metrics import accuracy_score
8 from sklearn.metrics import classification_report
9
10
11
12 #put the data in binary
13
14 X1= pd.get_dummies(data_training['Country'])
15 X2 = pd.get_dummies(data_training['Gender'])
16 #X3 = pd.get_dummies(data_training['Age'])
17 #X4 = pd.get_dummies(data_training['Subscription Type'])
18 X = np.concatenate([X1, X2], axis=1)
19
20 y = data_training['Churned']
21
22 X_train, X_test, Y_train, Y_test = train_test_split(X, y, test_size=0.2)
23
24
25 le = LabelEncoder()
26
27 Y_train = le.fit_transform(Y_train)
28
29 Y_test = le.fit_transform(Y_test)
30
31
32
33 #neural network
34 clf = MLPClassifier(hidden_layer_sizes = (10,), solver='lbfgs', activation = 'relu')
35
36 clf.fit(X_train, Y_train)
37 yp = clf.predict(X_test)
38
39 accuracy = accuracy_score(Y_test, yp)
40 accuracy
```

0.818

Started to concatenate features (after one them), after some trials found that 3 or more features overfitted the model leading to accuracy rates of around 70-75%, therefore the course of action was to leave it to two.

Accuracy was increased to 81% for NN3 (using Country and Gender as features) and 82% (using Subscription Type and Monthly Revenue as features)

Neural Networks

```
33 history=model.fit(X, y, epochs=100, validation_split=0.2, batch_size=32)
34 accuracy = model.evaluate(X_test, Y_test, verbose=0)[1]
35 accuracy
36 0.86100
37 Epoch 86/100
38 63/63 [=====] - 0s 3ms/step - loss: 0.5217 - accuracy: 0.7790 - val_loss: 0.2591 - val_accuracy: 0.992
39 Epoch 87/100
40 63/63 [=====] - 0s 3ms/step - loss: 0.5217 - accuracy: 0.7790 - val_loss: 0.2582 - val_accuracy: 0.992
41 Epoch 88/100
42 63/63 [=====] - 0s 4ms/step - loss: 0.5217 - accuracy: 0.7790 - val_loss: 0.2559 - val_accuracy: 0.992
43 Epoch 89/100
44 63/63 [=====] - 0s 3ms/step - loss: 0.5219 - accuracy: 0.7790 - val_loss: 0.2545 - val_accuracy: 0.992
45 Epoch 90/100
46 63/63 [=====] - 0s 3ms/step - loss: 0.5218 - accuracy: 0.7790 - val_loss: 0.2563 - val_accuracy: 0.992
47 Epoch 91/100
48 63/63 [=====] - 0s 3ms/step - loss: 0.5217 - accuracy: 0.7790 - val_loss: 0.2618 - val_accuracy: 0.992
49 Epoch 92/100
50 63/63 [=====] - 0s 3ms/step - loss: 0.5217 - accuracy: 0.7790 - val_loss: 0.2600 - val_accuracy: 0.992
51 Epoch 93/100
52 63/63 [=====] - 0s 3ms/step - loss: 0.5219 - accuracy: 0.7790 - val_loss: 0.2588 - val_accuracy: 0.992
53 Epoch 94/100
54 63/63 [=====] - 0s 3ms/step - loss: 0.5217 - accuracy: 0.7790 - val_loss: 0.2610 - val_accuracy: 0.992
55 Epoch 95/100
56 63/63 [=====] - 0s 3ms/step - loss: 0.5217 - accuracy: 0.7790 - val_loss: 0.2534 - val_accuracy: 0.992
57 Epoch 96/100
58 63/63 [=====] - 0s 3ms/step - loss: 0.5218 - accuracy: 0.7790 - val_loss: 0.2568 - val_accuracy: 0.992
59 Epoch 97/100
60 63/63 [=====] - 0s 3ms/step - loss: 0.5216 - accuracy: 0.7790 - val_loss: 0.2593 - val_accuracy: 0.992
61 Epoch 98/100
62 63/63 [=====] - 0s 3ms/step - loss: 0.5217 - accuracy: 0.7790 - val_loss: 0.2582 - val_accuracy: 0.992
63 Epoch 99/100
64 63/63 [=====] - 0s 3ms/step - loss: 0.5217 - accuracy: 0.7790 - val_loss: 0.2569 - val_accuracy: 0.992
65 Epoch 100/100
66 63/63 [=====] - 0s 3ms/step - loss: 0.5218 - accuracy: 0.7790 - val_loss: 0.2543 - val_accuracy: 0.992
```

0.8360000252723694

```
1 #NEURAL NETWORK USING COUNTRY AND AGE AS A FEATURE, THIS TIME USING A SEQUENTIAL MODEL AND USING EPOCHS
2
3 from sklearn.neural_network import MLPClassifier
4 from sklearn.preprocessing import LabelEncoder
5
6 from sklearn.metrics import accuracy_score
7 from sklearn.metrics import classification_report
8 import tensorflow as tf
9
10 from tensorflow.keras.layers import Dense
11
12
13 X1= pd.get_dummies(data_training['Country'])
14 X2 = pd.get_dummies(data_training['Gender'])
15 X3 = pd.get_dummies(data_training['Age'])
16 #X4 = pd.get_dummies(data_training['Subscription Type'])
17 X = np.concatenate([X1, X3], axis=1)
18
19 y = pd.get_dummies(data_training['Churned'])
20
21 X_train, X_test, Y_train, Y_test = train_test_split(X, y, test_size=0.2)
22
23
24 model = Sequential()
25 #model.add(Dropout(0.2))
26
27 model.add(Dense(10, activation='relu', input_dim = X.shape[1]))
28 model.add(Dense(2, activation='softmax'))
29
30
31
32 model.compile(loss='categorical_crossentropy', optimizer='adam', metrics=['accuracy'])
33
34 history=model.fit(X, y, epochs=100, validation_split=0.25, batch_size=32)
35
36 accuracy = model.evaluate(X_test, Y_test, verbose=0)[1]
37 accuracy
```

This time I did a more advanced Neural Network using two features but importing Keras to initialize a Sequential model adding a dense layer with 2 neurons and another dense layer with 10 neurons, with relu and softmax activation features, and finally did a validation split of 25% and got the accuracy to almost 85%

Results and Conclusions

1. Although one of our initial assumptions was that all three methods used (NN, Random Forests and K-means) would lead us to successful results only Random Forests and NN performed well while K-means didn't.
2. The K-Means clustering effectively segmented users into distinct groups, with Cluster 1 identified as the highest churn risk (10.9%) due to unique demographic and behavioral attributes, enabling targeted retention strategies, while additional models like Isolation Forest enhanced churn prediction accuracy.
3. Actionable insights: address regional disparities by optimizing content offerings, pricing models, and support services in high-churn countries like Italy and Spain, while reinforcing successful practices in low-churn regions like France and the UK.