

✓ Predictive Modeling: the importance of data in decision making for NPOs.

Introduction

Non-Profit Organizations (NPOs) play a vital role in social, educational, and economic development, but some, especially in underdeveloped countries like Haiti, rely heavily on gut instinct and traditional methods to make decisions, which leads to poor donor retention, ineffective communication, weak donor interaction, and unsuccessful campaigns. Since NPOs rely heavily on fundraising, understanding donor behavior is a must, as one of the major challenges NPOs face is maintaining donor engagement and loyalty over time.

The primary objective of this project is to showcase the importance of using data to make informed decisions in NPOs. For that purpose, the project aims to analyze donor behavior and categorize donors for a random and fictitious organization. The insights gained throughout this analysis will foster the use of data and machine learning in decision-making, which can therefore help NPOs design more effective fundraising and communication strategies, ultimately enhancing their ability to sustain and grow their initiatives.

To achieve this, we use the [Donor Data](#) from Kaggle, which comprises six tables detailing donors, donations, campaigns, project results, and engagement activities. These data provide a comprehensive view of donor interactions and contributions, forming the basis for our analysis.

Methodology

Exploratory data Analysis

- data preparation
- business understanding
- data understanding

Major questions

These questions deepen our grasp of the business and data context, steering the analysis toward the NPOs' core concerns and enabling a thorough exploration of donor dynamics.

This part is crucial, as it serves as a compass that guides our analysis and helps address stakeholders' concerns or uncover relevant information.

Modeling/ Model Evaluation

Recommendations

Contact information

the contact information of the two data scientist who work on this project is provided.

✓ EDA – Exploratory Data Analysis

✓ Data preparation

Libraries and files importation

```
1 # libraries importation
2 import pandas as pd
3 from sklearn.model_selection import train_test_split, GridSearchCV
4 from sklearn.preprocessing import LabelEncoder, OneHotEncoder, StandardScaler
5 from sklearn.metrics import roc_auc_score, f1_score, recall_score, accuracy_score, confusion_matrix, classification_report
6 from sklearn.linear_model import LogisticRegression
7 from sklearn.preprocessing import OneHotEncoder, StandardScaler
8 from sklearn.compose import ColumnTransformer
9 from sklearn.pipeline import Pipeline
10 from sklearn.ensemble import RandomForestClassifier, RandomForestRegressor
11 from xgboost import XGBClassifier
12 import matplotlib.pyplot as plt
13 import seaborn as sns
14 import numpy as np
15
16 # files importation
17 donors= pd.read_csv('donors.csv')
18 campaigns= pd.read_csv('campaigns.csv')
```

```

19 donations= pd.read_csv('donations_linked.csv')
20 engagement_history= pd.read_csv('engagement_history.csv')
21 engagement_outcomes =pd.read_csv('engagement_outcomes.csv')
22 impact= pd.read_csv('impact.csv')

```

```

1 print(f"Number of Donors" ,donors.shape[0] ,"\nNumber donations", donations.shape[0],
2       "\nNumber of Campaigns", campaigns.shape[0], "\nEngagement_history",engagement_history.shape[0],"\nImpact",impact.shap

```

```

Number of Donors 4891
Number donations 17840
Number of Campaigns 10
Engagement_history 24238
Impact 50

```

Merging all Tables

```

1 Donors_donations= pd.merge(donors,donations , on='DonorID', how='left')
2 Donors_donations_Campaigns= pd.merge(Donors_donations, campaigns, on='CampaignID', how='left')
3 Donors_donations_Campaigns_engagement_history =pd.merge(Donors_donations_Campaigns, engagement_history, on='DonorID', how=
4 Donors_donations_Campaigns_engagement_history_impacts =pd.merge(Donors_donations_Campaigns_engagement_history,impact, on = '

```

Comment: there are 7 tables in the dataset, we merge 6 of them because donors.csv and donations.csv are the same tables.

```

1 # all tables are merged into one
2 df= Donors_donations_Campaigns_engagement_history_impacts
3

```

```

1 # the merged tables
2 df.head()
3

```

	DonorID	Age	Gender	Location	JoinDate	DonationID	DonationDate	Amount	CampaignID	CampaignName	StartDate	EndDate
0	DNR00001	56	Male	QLD	2020-10-03 22:59:27.552825	DNT002656	2015-06-16	124.99	CAMP005	Youth Initiative 005	59:27.6	59:27
1	DNR00001	56	Male	QLD	2020-10-03 22:59:27.552825	DNT002656	2015-06-16	124.99	CAMP005	Youth Initiative 005	59:27.6	59:27
2	DNR00001	56	Male	QLD	2020-10-03 22:59:27.552825	DNT002656	2015-06-16	124.99	CAMP005	Youth Initiative 005	59:27.6	59:27
3	DNR00001	56	Male	QLD	2020-10-03 22:59:27.552825	DNT002656	2015-06-16	124.99	CAMP005	Youth Initiative 005	59:27.6	59:27
4	DNR00001	56	Male	QLD	2020-10-03 22:59:27.552825	DNT002656	2015-06-16	124.99	CAMP005	Youth Initiative 005	59:27.6	59:27

```

1 print( 'Number of Records',df.shape[0])

```

```

Number of Records 440083

```

Check for null or missing values

```

1 df.isnull().sum().sort_values(ascending=False)

```

	0
CampaignID	638
Amount	638
DonationDate	638
DonationID	638
Cost	638
ActualAmount	638
Value	638
ImpactType	638
StartDate	638
CampaignName	638
EndDate	638
TargetAmount	638
Location	0
Gender	0
Age	0
DonorID	0
JoinDate	0
Channel	0
Date	0
EngagementOutcome	0

dtype: int64

```
1 # drop rows with null values
2 df[df.isnull().any(axis=1)]
3 df = df.dropna(subset=['DonationID'])
```

Dealing with Duplicates values

```
1 # drop rows with duplicates
2 duplicates = df[df.duplicated()]
3 print("Number of duplicate rows before dropping:", duplicates.shape[0])
4 df = df.drop_duplicates()
5 print("Number of duplicate rows after dropping:", df.duplicated().sum())
6 # Dropping rows duplicate donation ID
7 df = df.drop_duplicates(subset=['DonationID'])
```

Number of duplicate rows before dropping: 205
Number of duplicate rows after dropping: 0

Comment: Dropping duplicated DonationID rows is a must since we don't want to consider the same donation many time, which would biased our analysis.

```
1 # to verify if all null values were dropped run the code in this cell:
2 # df.isnull().sum().sort_values()
```

Fix the dates Format

```
1 df['JoinDate'] = pd.to_datetime(df['JoinDate'])
2 df['Date'] = pd.to_datetime(df['Date'], dayfirst=True)
3 df['DonationDate'] = pd.to_datetime(df['DonationDate'], dayfirst=True)
```

Remove unnecessary columns

```
1 df = df.drop(['StartDate', 'EndDate'], axis=1)
```

Table's shape

```
1 # Run the code in this cell to see table's shape after cleaning
2 df.shape
```

```
(17840, 18)
```

```
1 # after cleaning we save the merged dataset
2 df.to_csv('Data.csv',index=0)
```

Explore basic statistics for numerical features (mean, median, min, max)

```
1 df.describe().round(2)
```

	Age	JoinDate	DonationDate	Amount	TargetAmount	ActualAmount	Date	Value	Cost
count	17840.00	17840	17840	17840.00	17840.00	17840.00	17840	17840.00	17840.0
mean	48.44	2019-10-26 01:02:47.732196864	2020-02-09 03:17:45.470851840	110.24	62092.24	59866.38	2023-02-25 02:59:26.098654720	33.52	100.0
min	18.00	2015-04-23 22:59:27.552825	2015-01-01 00:00:00	10.00	33193.34	33179.39	2022-07-25 00:00:00	19.00	100.0
25%	33.00	2017-08-20 22:59:27.552825088	2017-07-07 00:00:00	38.90	49248.12	44689.69	2022-09-22 00:00:00	24.00	100.0
50%	48.00	2019-10-27 22:59:27.552825088	2020-02-11 00:00:00	79.15	58107.99	57113.83	2022-12-12 00:00:00	25.00	100.0
75%	64.00	2022-01-08	2022-09-11	117.50	60702.66	60470.70	2023-05-22	27.00	100.0

Business understanding

In today's fast-evolving fundraising landscape, nonprofit organisations face growing pressure to engage supporters meaningfully, demonstrate results, and build trust through transparency. Traditional, instinct-driven approaches no longer meet modern donor expectations.

Adopting a data-informed approach is not just a technical upgrade—it's a cultural shift that aligns mission, message, and measurement. With the right mindset and existing tools, even small organisations can use data to boost engagement and impact.

Delivering measurable change increasingly depends on how well nonprofits engage and retain supporters. As donor expectations rise and attention spans shrink, sustaining long-term relationships becomes a major challenge. [McKinsey Global Institute](#) says that data-driven organisations are 23 times more likely to acquire customers, 6 times as likely to retain customers and 19 times as likely to be profitable, showing the importance of including data in decision making.

Our objectives focus on helping NPOs reduce donor churn, strengthen their relationships with donors, and improve alignment between fundraising and impact, with data driven analysis. We aim to help them engage donors and inspire new strategies for attraction and retention.

Data Understanding

A synthetic dataset was created to reflect the operational characteristics of a nonprofit organisation, focused on youth and community welfare. It includes relational tables for donors, donations, campaigns, impact outcomes, and engagement activities. These tables are structured to support a complete engagement analysis from start to finish.

Although the data is fictitious, it reflects realistic patterns based on [sector](#) / [research](#). These patterns include donor attrition, variations in donation behavior, and inconsistent campaign performance. The goal is to encourage the use of data in nonprofit environments and to provide the NPO community with a clear example of what a data-informed approach can accomplish. This simulation allows for open analysis while maintaining privacy and ethical standards.

The structure of the dataset makes it possible to demonstrate how nonprofit engagement can be improved through analytics without relying on sensitive real-world data.

Major questions

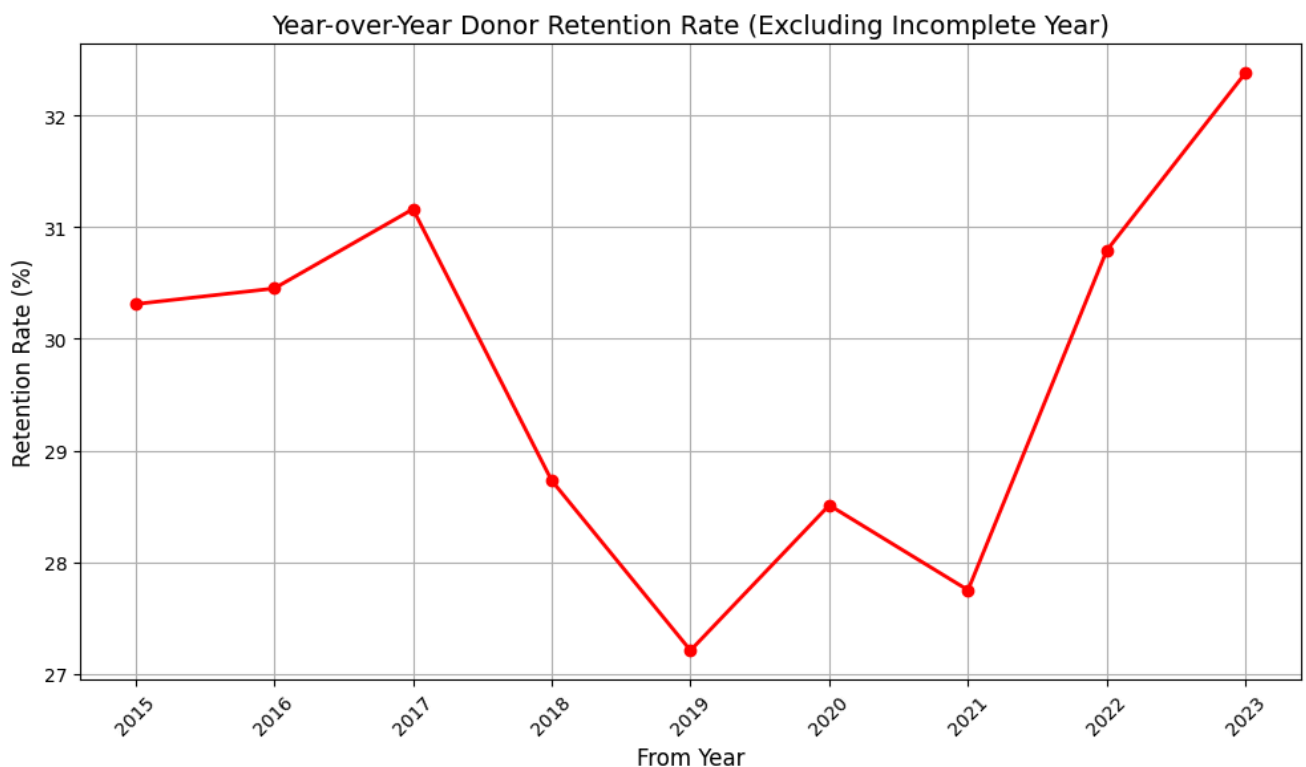
this section is a follow up of the two previous sections. it builds our business and data understanding and allow us to leverage the data to get meaningful and pertinent information. We will also tell the purpose of each Newfound information, and by the same occasion foster a data informed approached in NPOs.

How has donor retention evolved over time ?

```

1 # conversion
2 df['DonationDate'] = pd.to_datetime(df['DonationDate'], dayfirst=True)
3
4 # Extract donation year
5 df['DonationYear'] = df['DonationDate'].dt.year
6 years = sorted(df['DonationYear'].unique())
7 retention_stats = []
8
9 for i in range(len(years) - 1):
10     year_current = years[i]
11     year_next = years[i + 1]
12
13     donors_current = set(df[df['DonationYear'] == year_current]['DonorID'])
14     donors_next = set(df[df['DonationYear'] == year_next]['DonorID'])
15
16     retained_donors = donors_current & donors_next
17     retention_rate = len(retained_donors) / len(donors_current) * 100 if donors_current else 0
18
19     retention_stats.append({
20         "From Year": year_current,
21         "To Year": year_next,
22         "Donors in Year": len(donors_current),
23         "Retained Donors": len(retained_donors),
24         "Retention Rate (%)": round(retention_rate, 2)
25     })
26
27 retention_df = pd.DataFrame(retention_stats)
28
29 # Exclude the last year from the plot as it has incomplete data
30 retention_df_plot = retention_df[:-1]
31
32 plt.figure(figsize=(10, 6))
33 plt.plot(retention_df_plot["From Year"].astype(str), retention_df_plot["Retention Rate (%)"], marker='o', color = "red", 1
34 plt.title("Year-over-Year Donor Retention Rate (Excluding Incomplete Year)", fontsize=14)
35 plt.xlabel("From Year", fontsize=12)
36 plt.ylabel("Retention Rate (%)", fontsize=12)
37 plt.xticks(rotation=45)
38 plt.grid(True)
39 plt.tight_layout()
40 # to save the image
41 plt.savefig('1st_image.png')
42 # to see the image
43 plt.show()
44

```



Comment: This line graph illustrates the percentage of donors retained from one year to the next over a ten-year period. While retention fluctuates around 30%, occasional dips highlight the importance of targeted re-engagement strategies and sustained communication with donor segments.

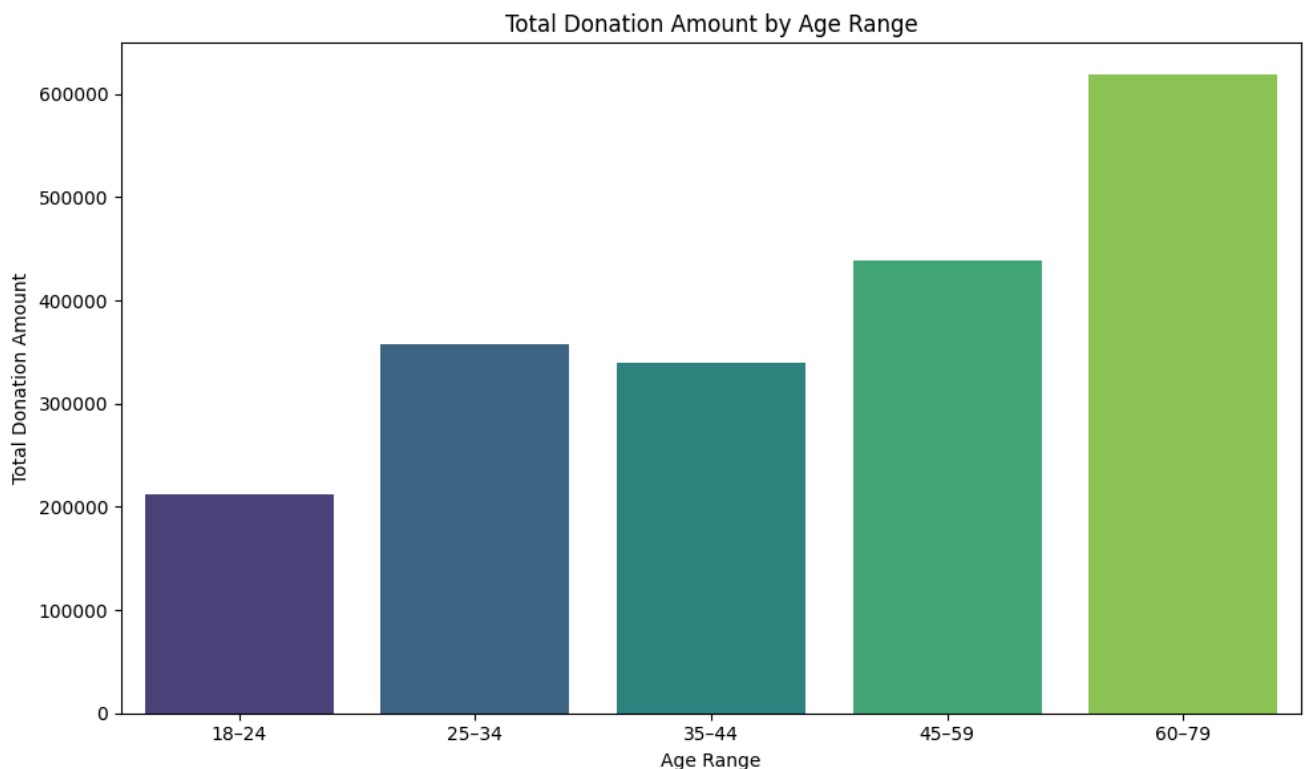
Any NPO will want this information to position themselves regarding donor retention. By knowing how they retain donors, they can determine whether or not to foster interaction with donors, build more campaigns or any other inspired-data decision that could improve donor retention.

What is the donation frequency per age range ?

```
1 # Find out Maximum and Minimum age
2 print("the oldest person are ", df.Age.min())
3 print("the youngest person are ", df.Age.max())
```

```
the oldest person are 18
the youngest person are 79
```

```
1 # Define age ranges
2
3 # Define bins and matching labels
4 bins = [18, 24, 35, 45, 59, 79]
5 labels = ['18-24', '25-34', '35-44', '45-59', '60-79'] # 5 labels for 5 intervals
6
7 # Create a new column with age ranges
8 df['AgeRange'] = pd.cut(df['Age'], bins=bins, labels=labels, right=True, include_lowest=True)
9
10 # Calculate total donation amount per age range
11 donation_by_age_range = df.groupby('AgeRange', observed=True)['Amount'].sum().reset_index()
12
13 # Plot the bar chart
14 plt.figure(figsize=(10, 6))
15 sns.barplot(x='AgeRange', y='Amount', data=donation_by_age_range, palette='viridis')
16
17 plt.title('Total Donation Amount by Age Range')
18 plt.xlabel('Age Range')
19 plt.ylabel('Total Donation Amount')
20 plt.tight_layout()
21 # to save the image
22 plt.savefig('2nd_image.png')
23 # to see the image
24 plt.show()
25
```

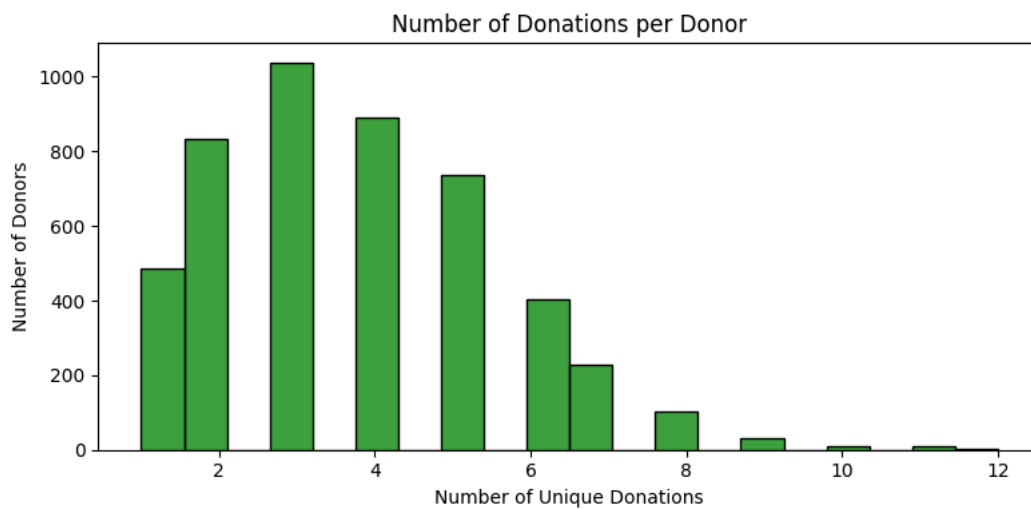


Comment: Donors who are older are likely to donate more than young individuals in the range of 25 to 34 years old. The results match [real world data](#).

Any NPO would want this information to know which age range to target the most and which one to sensitize to donate more, whether it's through campaign messages or any other means.

How many time donors are likely to donate ?

```
1 # import necessary library
2 import matplotlib.pyplot as plt
3 import seaborn as sns
4
5 # Group by DonorID to get donation count
6 donation_counts = donations.groupby('DonorID')['DonationID'].count()
7
8 # Plot
9 plt.figure(figsize=(8, 4))
10 sns.histplot(donation_counts, bins=20, kde=False, color='green')
11 plt.title('Number of Donations per Donor')
12 plt.xlabel('Number of Unique Donations')
13 plt.ylabel('Number of Donors')
14 plt.tight_layout()
15 # to save the image
16 plt.savefig('3rd_image.png')
17 # to see the image
18 plt.show()
```



Comment: Most donors donate around 2 to 4 times. We can observe that the Number of Donors decreases for higher number of donations, showing moderate donor retention.

This type of information is very useful because some NPOs might want to set their new objectives to higher number of donations per donor.

What is the donation amount repartition per gender?

```
1 df.Age.value_counts()
2
```

count

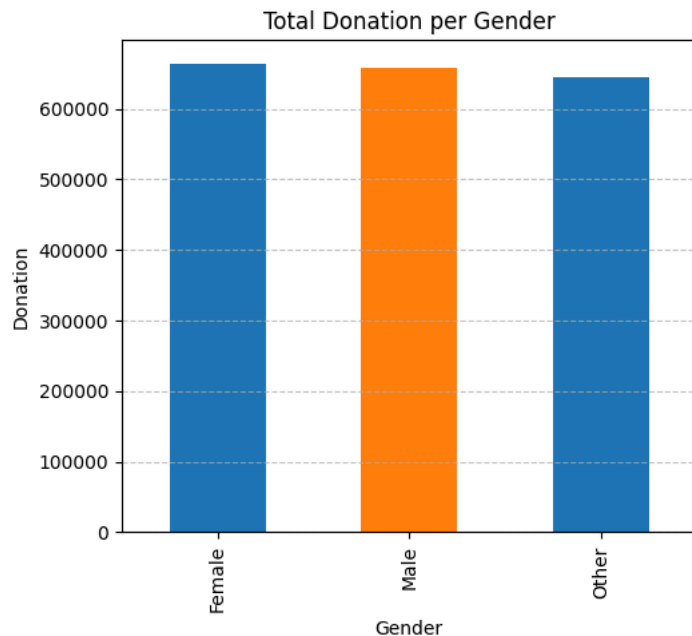
Age

45	370
40	351
53	349
33	349
43	342
...	...
70	236
51	231
28	225
67	212
60	210

62 rows × 1 columns

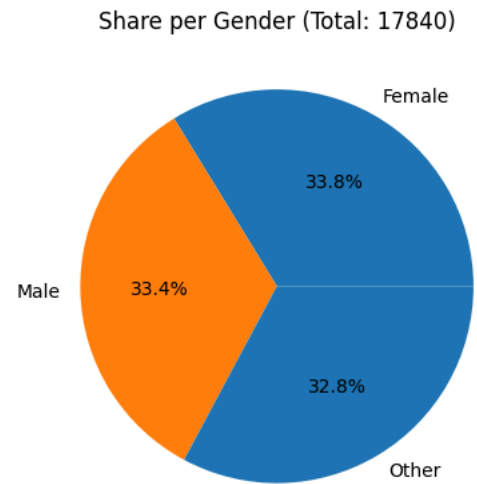
dtype: int64

```
1 donation_by_gender = df.groupby('Gender')['Amount'].sum()
2
3 # Kalkile total donation tout ansanm
4 total_donation = df['Amount'].count()
5
6 # subplots (2 graf)
7 fig, axes = plt.subplots(1, 2, figsize=(10, 5))
8
9 # 1 Bar chart donation by gender
10 donation_by_gender.plot(kind='bar', ax=axes[0], color=['#1f77b4', '#ff7f0e'])
11 axes[0].set_title('Total Donation per Gender')
12 axes[0].set_xlabel('Gender')
13 axes[0].set_ylabel('Donation')
14 axes[0].grid(axis='y', linestyle='--', alpha=0.7)
15
16 # 2 Pie chart
17 donation_by_gender.plot(kind='pie', autopct='%1.1f%%', ax=axes[1], ylabel='', colors=['#1f77b4', '#ff7f0e'])
18 axes[1].set_title(f'Share per Gender (Total: {total_donation})')
19
20
21 plt.tight_layout()
22 # to save the image
23 plt.savefig('4rth_image.png')
24 # to see the image
25 plt.show()
26 donation_by_gender
```

Amount	
Gender	
Female	663950.26
Male	657135.77
Other	645596.74

dtype: float64



Comment: For this fictitious dataset, gender is almost equally distributed, but in some regions, NPOs might observe some imbalance and, therefore, make it their objective to understand why and do something about it.

How can the different type of donors be classified ?

```

1 import numpy as np
2
3 # Convert donation date to datetime
4 df['DonationDate'] = pd.to_datetime(df['DonationDate'])
5
6 # Create donor-level summary
7 today = pd.to_datetime('today')
8 donor_summary = df.groupby('DonorID').agg(
9     FirstDonationDate=('DonationDate', 'min'),
10    LastDonationDate=('DonationDate', 'max'),
11    Frequency=('DonationID', 'count'),
12    Monetary=('Amount', 'sum')
13 ).reset_index()
14 donor_summary['Recency'] = (today - donor_summary['LastDonationDate']).dt.days
15
16 # Merge with full donor list
17 donor_full_summary = pd.merge(donors[['DonorID']], donor_summary, on='DonorID', how='left')
18 donor_full_summary['Recency'] = donor_full_summary['Recency'].fillna(np.inf)
19 donor_full_summary['Frequency'] = donor_full_summary['Frequency'].fillna(0)
20 donor_full_summary['Monetary'] = donor_full_summary['Monetary'].fillna(0)
21 # Assign segments
22 def assign_segment(row):
23     r = row['Recency']
24     f = row['Frequency']
25     m = row['Monetary']
26
27     if np.isinf(r) or pd.isna(r):
28         return 'Never Donated'
29     elif r <= 365 and f == 1:
30         return 'New Donors'
31     elif r <= 365 and f >= 4 and m >= 750:
32         return 'Champions'
33     elif r <= 1095 and f >= 4:
34         return 'Loyal Donors'
35     elif r <= 1095 and 2 <= f <= 4 and m >= 750:
36         return 'High Value Potentials'
37     elif r > 1095 and f >= 2 and m >= 250:

```

```

38         return 'Lapsed but Valuable'
39     elif 365 < r <= 1095 and f <= 2:
40         return 'At Risk'
41     elif r > 1095 and f == 1:
42         return 'Lost or Inactive'
43     elif f <= 2 and m < 250:
44         return 'Low Frequency'
45     else:
46         return 'Misc Donors'
47
48 donor_full_summary['Segment'] = donor_full_summary.apply(assign_segment, axis=1)
49
50 # Return segment summary for confirmation
51 segment_summary = donor_full_summary['Segment'].value_counts().reset_index()
52 segment_summary.columns = ['Segment', 'Count']
53 segment_summary

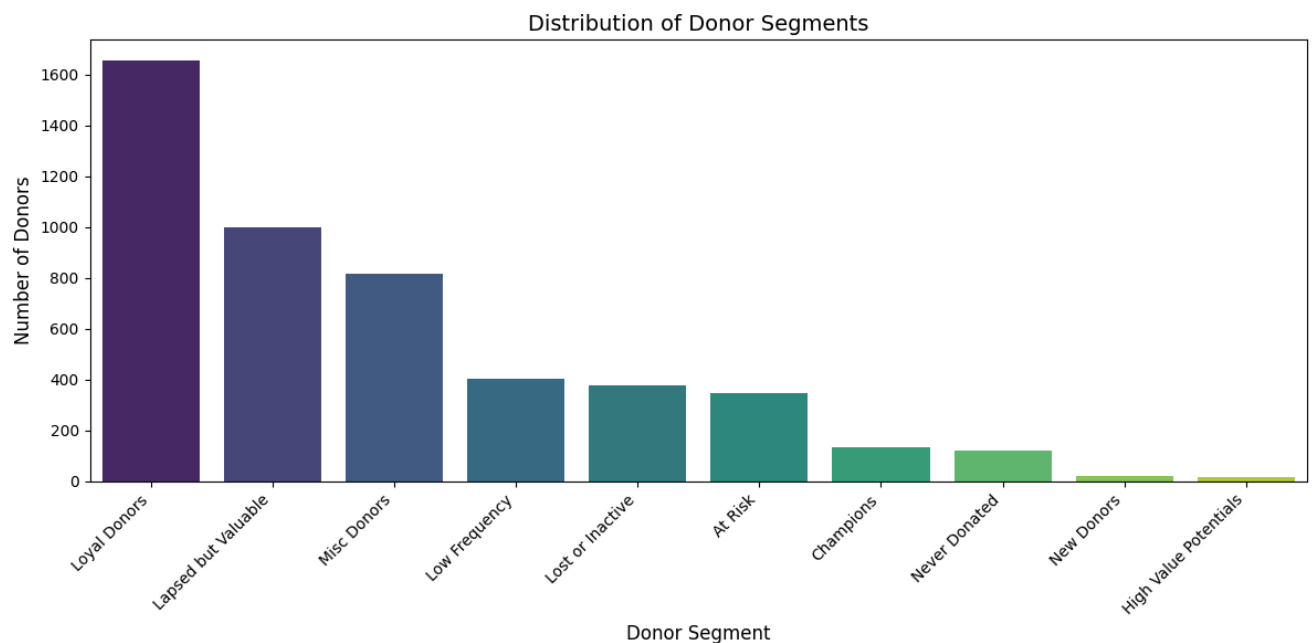
```

	Segment	Count
0	Loyal Donors	1653
1	Lapsed but Valuable	998
2	Misc Donors	818
3	Low Frequency	405
4	Lost or Inactive	376
5	At Risk	347
6	Champions	133
7	Never Donated	122
8	New Donors	20
9	High Value Potentials	19

```

1 plt.figure(figsize=(12, 6))
2 sns.barplot(x='Segment', y='Count', data=segment_summary, palette='viridis')
3 plt.title('Distribution of Donor Segments', fontsize=14)
4 plt.xlabel('Donor Segment', fontsize=12)
5 plt.ylabel('Number of Donors', fontsize=12)
6 plt.xticks(rotation=45, ha='right')
7 plt.tight_layout()
8 # to save the image
9 plt.savefig('7th_image.png')
10 # to see the image

```



Comments: Not every donor contributes the same. Therefore, donors have to be interacted with differently. Some donors may receive messages of recognition and heartfelt remerciements, while others may need to be re-engaged toward donating again. Knowing the different donor segments is crucial information for NPOs.

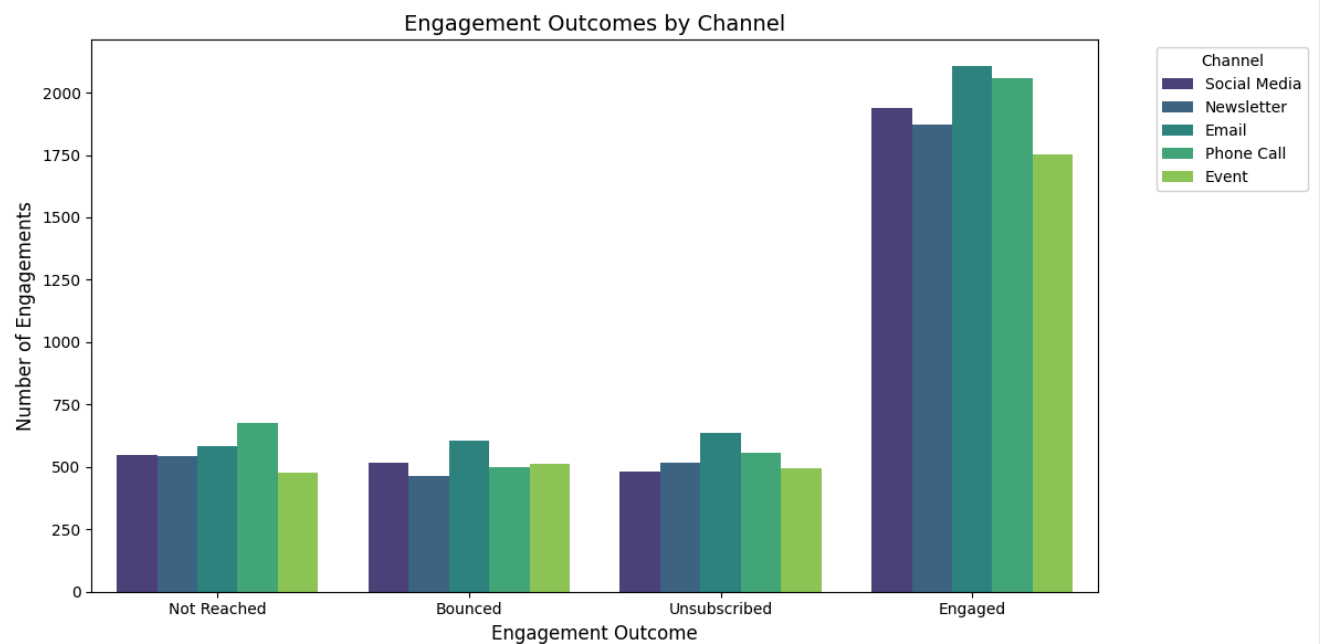
what channel results to the best engagement outcome?

```
1 df.EngagementOutcome.value_counts()
```

	count
EngagementOutcome	
Engaged	9723
Not Reached	2831
Unsubscribed	2688
Bounced	2598

dtype: int64

```
1 # Create a countplot showing outcomes across channels
2 plt.figure(figsize=(12, 6))
3 sns.countplot(data=df, x='EngagementOutcome', hue='Channel', palette='viridis')
4 plt.title('Engagement Outcomes by Channel', fontsize=14)
5 plt.xlabel('Engagement Outcome', fontsize=12)
6 plt.ylabel('Number of Engagements', fontsize=12)
7 plt.xticks(rotation=0)
8 plt.legend(title='Channel', bbox_to_anchor=(1.05, 1), loc='upper left')
9 plt.tight_layout()
10 # to save the image
11 plt.savefig('6xth_image.png')
12 # to see the image
13 plt.show()
```



```
1 df.Channel.value_counts()
```

	count
Channel	
Email	3931
Phone Call	3790
Social Media	3486
Newsletter	3398
Event	3235

dtype: int64

Comment: Channels seem to have approximately the same engagement outcome in the context of this analysis. However, some NPOs might have a certain channel that leads to better results. For example, in the context of Haiti, some NPOs may want to use social media instead of email, since it is more widely used.

Modeling

We will develop a concrete and practical plan, along with an ML pipeline, that aligns with the objective of reducing churn and unsubscriptions, strengthening donor relationships, and improving alignment between collection and impact — with clear, actionable ML deliverables.

```
1 from sklearn.model_selection import train_test_split
2 from sklearn.compose import ColumnTransformer
3 from sklearn.preprocessing import OneHotEncoder, StandardScaler
4 from sklearn.pipeline import Pipeline
5 from sklearn.metrics import accuracy_score, roc_auc_score, confusion_matrix, classification_report
6 from xgboost import XGBClassifier
7 #import shap
8 import matplotlib.pyplot as plt
9 import warnings
10 warnings.filterwarnings('ignore')
```

Model 1

below are the many steps taken to build the model.

Feature Engineering

In this part we will generate new features that will help us create our model. these features will be the criterion that will help us predict donor loyalty.

Splitting the Data The model will be based on historical data to predict the loyalty of a donor.

we will split the historical data into two parts: one for the train set and one for the test set. the train set will contain all historical data that are three years before 2025.

```
1 cutoff_date = df['DonationDate'].max() - pd.Timedelta(days=365)
2 # Historical data (before cutoff)
3 history = df[df['DonationDate'] <= cutoff_date]
4 # Future data (after cutoff)
5 future = df[df['DonationDate'] > cutoff_date]
```

Basic RFM Calculation RFM stands for : Recency, Frequency, Monetary

we will also create the target variable Loyalty indicating if the donor is still active.

```
1 # RFM calculation
2 rfm = history.groupby('DonorID').agg(
3     Recency=('DonationDate', lambda x: (cutoff_date - x.max()).days),
4     Frequency=('DonationID', 'count'),
5     Monetary=('Amount', 'sum'),
6     AvgAmount=('Amount', 'mean'),
7
8 ).reset_index()
9
10
```

```

11 # Loyalty = 1 if the donor donated after cutoff_date
12 rfm['Loyalty'] = (rfm['DonorID'].isin(future['DonorID'])).astype(int)

```

Additional Behavioral Variables

```

1 # Variance of donation amounts
2 var_df = history.groupby('DonorID')['Amount'].std().fillna(0).reset_index()
3 var_df.rename(columns={'Amount': 'DonationVariance'}, inplace=True)
4
5 # Average interval between donations (days)
6 interval_df = history.groupby('DonorID')['DonationDate'].apply(
7     lambda x: x.sort_values().diff().dt.days.mean()
8 ).fillna(0).reset_index(name='AvgInterval')
9
10 # Donor tenure
11 tenure_df = history.groupby('DonorID')['DonationDate'].min().reset_index()
12 tenure_df['Tenure'] = (cutoff_date - tenure_df['DonationDate']).dt.days
13 tenure_df.drop(columns='DonationDate', inplace=True)
14
15 # Merge
16 rfm = rfm.merge(var_df, on='DonorID', how='left')
17 rfm = rfm.merge(interval_df, on='DonorID', how='left')
18 rfm = rfm.merge(tenure_df, on='DonorID', how='left')
19

```

Finds the channel that the donor uses most often.

that will help us Capture the donor's preferred communication/interaction channel, which often correlates with loyalty.

```

1 # Most frequent channel
2 fav_channel_df = history.groupby('DonorID')['Channel'].agg(
3     lambda x: x.mode()[0] if len(x.mode()) > 0 else 'unknown'
4 ).reset_index(name='FavChannel')
5
6 # Number of different channels used
7 diversity_df = history.groupby('DonorID')['Channel'].nunique().reset_index(name='ChannelDiversity')
8
9 # Counts per channel type (integers, e.g., Online, Event, Phone)
10 channel_counts = history.pivot_table(
11     index='DonorID',
12     columns='Channel',
13     values='DonationID',
14     aggfunc='count',
15     fill_value=0
16 ).reset_index()
17
18 # Merge all channel features into the main RFM dataframe
19 rfm = rfm.merge(fav_channel_df, on='DonorID', how='left')
20 rfm = rfm.merge(diversity_df, on='DonorID', how='left')
21 rfm = rfm.merge(channel_counts, on='DonorID', how='left')
22
23 # Group by DonorID to get one AgeRange and Gender per donor (most frequent value)
24 age_df = df.groupby('DonorID')['AgeRange'].agg(lambda x: x.mode()[0]).reset_index()
25 #gender = df.groupby('DonorID')['Gender'].agg(lambda x: x.mode()[0]).reset_index()
26
27 # Merge into the main RFM dataset
28 rfm = rfm.merge(age_df, on='DonorID', how='left')
29 #rfm = rfm.merge(gender, on='DonorID', how='left')
30
31 # Display updated RFM
32 rfm.head(1)
33

```

	DonorID	Recency	Frequency	Monetary	AvgAmount	Loyalty	DonationVariance	AvgInterval	Tenure	FavChannel	ChannelDiversity
0	DNR00001	880	3	377.67	125.89	0	2.035166	1174.5	3229	Social Media	

Prepare Features We select the numerical and categorical features for the model and Separates features (X) from the target (y).

```

1 channel_cols = [c for c in channel_counts.columns if c != 'DonorID']
2
3 features = [
4     'Recency', 'Frequency', 'Monetary', 'AvgAmount',
5     'DonationVariance', 'AvgInterval', 'Tenure',
6     'ChannelDiversity', 'FavChannel', 'AgeRange'

```

```

7 ] + channel_cols
8
9 X = rfm[features]
10 y = rfm['Loyalty']
11
12
13 numeric_features = [
14     'Recency', 'Frequency', 'Monetary', 'AvgAmount',
15     'DonationVariance', 'AvgInterval', 'Tenure',
16     'ChannelDiversity'
17 ] + channel_cols
18
19 categorical_features = ['FavChannel', 'AgeRange']

```

Split train / test

```

1 ##Split train / test
2
3 X_train, X_test, y_train, y_test = train_test_split(
4     X, y, test_size=0.2, random_state=42, stratify=y
5 )
6
7

```

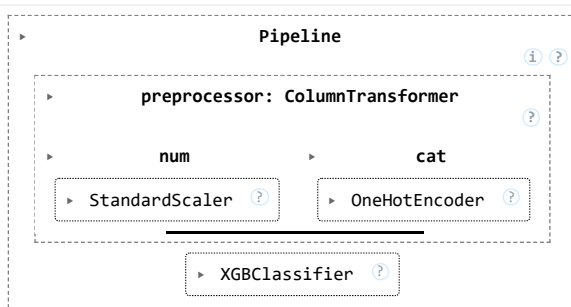
Chosen model:

The model used is XGBoost Classifier, selected for its ability to handle non-linear relationships and imbalanced datasets.

```

1 preprocessor = ColumnTransformer([
2     ('num', StandardScaler(), numeric_features),
3     ('cat', OneHotEncoder(handle_unknown='ignore'), categorical_features)
4 ])
5
6 ratio = y_train.value_counts()[0] / y_train.value_counts()[1]
7 clf_pipeline = Pipeline([
8     ('preprocessor', preprocessor),
9     ('classifier', XGBClassifier(
10         n_estimators=300,
11         learning_rate=0.05,
12         max_depth=5,
13         subsample=0.8,
14         colsample_bytree=0.8,
15         scale_pos_weight=ratio,
16         eval_metric='logloss',
17         random_state=42
18     ))
19 ])
20
21 clf_pipeline.fit(X_train, y_train)

```



Model Evaluation

```

1 from sklearn.metrics import accuracy_score, roc_auc_score, confusion_matrix
2 # Prédiction sur le test set
3 y_pred = clf_pipeline.predict(X_test)
4 y_proba = clf_pipeline.predict_proba(X_test)[:, 1]
5
6 # Accuracy
7 accuracy = accuracy_score(y_test, y_pred)
8 print("Accuracy:", accuracy)
9
10 # ROC-AUC
11 roc_auc = roc_auc_score(y_test, y_proba)
12 print("ROC-AUC:", roc_auc)
13

```

```

14 # Confusion matrix
15 cm = confusion_matrix(y_test, y_pred)
16 print("Confusion Matrix:\n", cm)
17
18 # Visualize confusion matrix
19 plt.figure(figsize=(8, 6))
20 sns.heatmap(cm, annot=True, fmt='d', cmap='Blues', xticklabels=['No Future Donation', 'Future Donation'], yticklabels=['No
21 plt.xlabel('Predicted')
22 plt.ylabel('Actual')
23 plt.title('Confusion Matrix')
24 # to save plot
25 plt.savefig('8th_image.png')
26 # to see plot
27 plt.show()

```

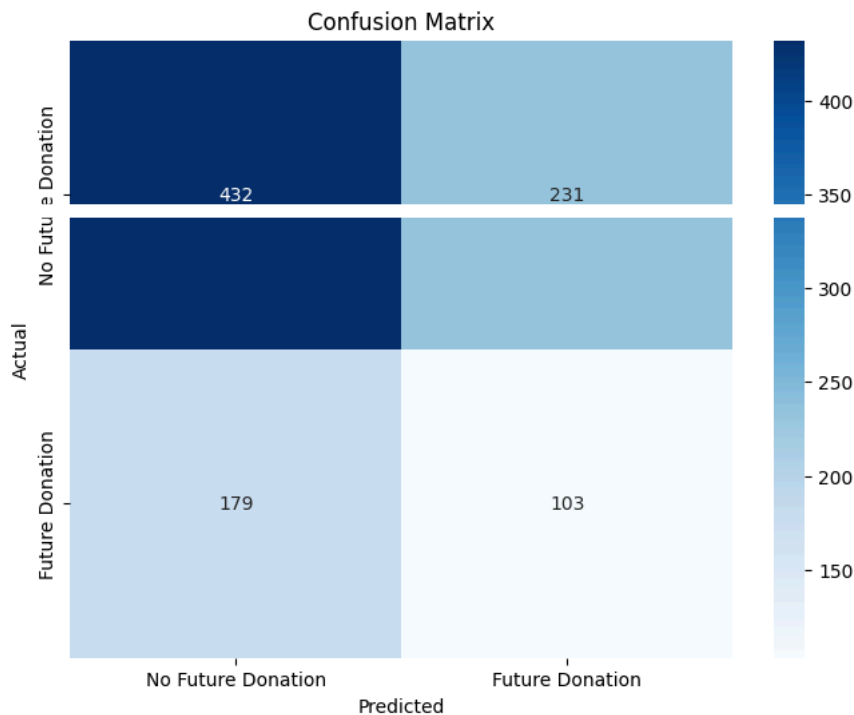
Accuracy: 0.5661375661375662

ROC-AUC: 0.5091915107559662

Confusion Matrix:

```
[[432 231]
```

```
[179 103]]
```



Interpretation

The model correctly predicts donor loyalty about 57% of the time, slightly better than random (50%), but not yet reliable for decision-making. It has limited ability to distinguish between loyal and non-loyal donors (0.5 = random). It needs improvement. The model correctly identifies many non-loyal donors (432), but struggles with loyal donors (103 correctly, 179 missed).

```
1 print("\nClassification Report:\n", classification_report(y_test, y_pred))
```

```

Classification Report:
              precision    recall  f1-score   support

     0       0.71         0.65         0.68         663
     1       0.31         0.37         0.33         282

 accuracy          0.57
 macro avg         0.51         0.51         0.51         945
 weighted avg         0.59         0.57         0.58         945

```

Overview of Model's results

```

1 rfm['Fidelite_Prob'] = clf_pipeline.predict_proba(rfm[features])[0, 1]
2
3 rfm['Fidelite_Segment'] = pd.cut(
4     rfm['Fidelite_Prob'],
5     bins=[0, 0.33, 0.66, 1],
6     labels=['Weak Loyalty ', 'Average Loyalty', 'Strong Loyalty']
7

```

```

8 )
9
10 print("\n ♦ Quick overview on the result upon all the Datateset :\n ")
11 print(rfm[['DonorID', 'Fidelite_Prob', 'Fidelite_Segment']].head())
12

```

♦ Quick overview on the result upon all the Datateset :

	DonorID	Fidelite_Prob	Fidelite_Segment
0	DNR00001	0.341485	Average Loyalty
1	DNR00002	0.504820	Average Loyalty
2	DNR00003	0.149082	Weak Loyalty
3	DNR00004	0.508528	Average Loyalty
4	DNR00005	0.583386	Average Loyalty

▼ Model 2 : Random Forest

```

1 rf_pipeline = Pipeline([
2     ('preprocessor', preprocessor),
3     ('classifier', RandomForestClassifier(
4         n_estimators=300,
5         max_depth=None,
6         min_samples_split=5,
7         min_samples_leaf=2,
8         class_weight='balanced',
9         random_state=42,
10        n_jobs=-1
11    ))
12 ])
13
14 # Convert target variables to integer type
15 y_train = y_train.astype(int)
16 y_test = y_test.astype(int)
17
18 # Train the model
19 rf_pipeline.fit(X_train, y_train)
20
21 # Predictions
22 y_pred = rf_pipeline.predict(X_test)
23 y_proba = rf_pipeline.predict_proba(X_test)[:, 1]
24
25 # Evaluation
26 print("Accuracy:", accuracy_score(y_test, y_pred))
27 print("ROC-AUC:", roc_auc_score(y_test, y_proba))
28 cm = confusion_matrix(y_test, y_pred)
29 print("Confusion Matrix:\n", cm)
30 print("\nClassification Report:\n", classification_report(y_test, y_pred))
31 plt.figure(figsize=(8, 6))
32 sns.heatmap(cm, annot=True, fmt='d', cmap='Blues', xticklabels=['No Future Donation', 'Future Donation'], yticklabels=['No
33 plt.xlabel('Predicted')
34 plt.ylabel('Actual')
35 plt.title('Confusion Matrix')
36 # to save plot
37 plt.savefig('9th_image.png')
38 # to see plot
39 plt.show()

```

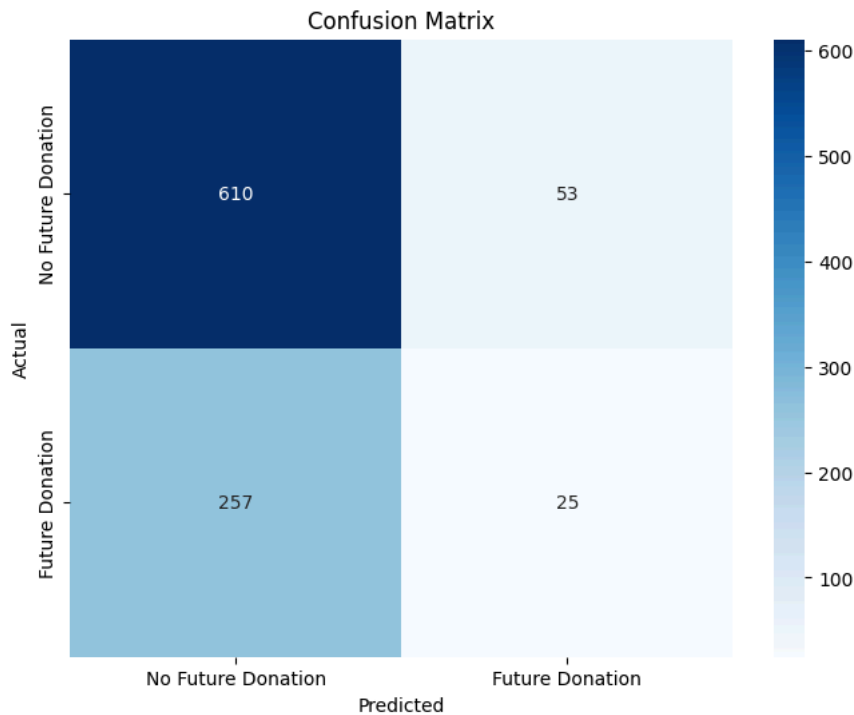

Accuracy: 0.671957671957672
ROC-AUC: 0.5066482675994566

Confusion Matrix:

```
[[610  53]
 [257  25]]
```

Classification Report:

	precision	recall	f1-score	support
0	0.70	0.92	0.80	663
1	0.32	0.09	0.14	282
accuracy			0.67	945
macro avg	0.51	0.50	0.47	945
weighted avg	0.59	0.67	0.60	945



Interpretation

The model correctly predicts about 67% of loyal donors. That's higher than the XGBoost model (≈ 0.57). We can therefore say that this one does better overall at predicting loyal donors

However, accuracy alone is misleading when classes are imbalanced (as is the case here).

Model 3 : regression task

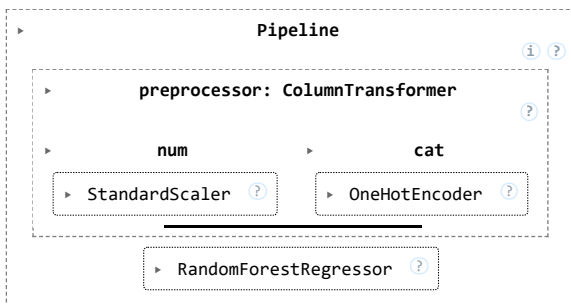
this model is build to try to predict future donation amount.

```
1 from sklearn.preprocessing import OneHotEncoder
2 from sklearn.compose import ColumnTransformer
3 from sklearn.pipeline import Pipeline
4 from sklearn.preprocessing import StandardScaler
5
6 # Aggregate donor-level data
7 dataset = df.groupby('DonorID').agg({
8     'Age': 'first',
9     'Gender': 'first',
10    'JoinDate': 'first',
11    'Location': 'first',
12    'Amount': ['sum', 'mean', 'count'],
13    'Channel': lambda x: x.mode()[0],
14    'EngagementOutcome': lambda x: x.mode()[0],
15    'Value': 'first'
16 }).reset_index()
17
18 dataset.columns = ['DonorID', 'Age', 'Gender', 'JoinDate', 'Location', 'TotalAmount', 'AvgAmount', 'DonationCount', 'Channel', 'Enga
19
20 # 2. Create the Target
21 # target 2: Future average donation amount (regression)
22 dataset['FutureDonation'] = np.random.randint(0,2,size=len(dataset)) # Randomly assign whether donor will donate (0/1)
```

```

23 dataset['FutureAmount'] = np.where(dataset['FutureDonation']==1,          # Future donation amount
24                                     np.random.uniform(10,500,size=len(dataset)),
25                                     0)
26
27 # Target variable: 'Value'
28 y = dataset['FutureAmount']
29
30 # Select relevant features
31 features = ['Age', 'Gender', 'Location', 'JoinDate', 'TotalAmount', 'AvgAmount', 'DonationCount', 'Channel', 'EngagementOutcome', 'V
32 X = dataset[features]
33
34 # Define numeric and categorical columns
35 numeric_features = ['Age', 'TotalAmount', 'DonationCount']
36 categorical_features = ['Gender', 'Channel', 'EngagementOutcome']
37
38 # Create preprocessor
39 preprocessor = ColumnTransformer(
40     transformers=[
41         ('num', StandardScaler(), numeric_features),          # Standardize numeric features
42         ('cat', OneHotEncoder(handle_unknown='ignore'), categorical_features) # One-hot encode categorical features
43     ])
44
45 from sklearn.model_selection import train_test_split
46
47 # Split the dataset
48 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
49
50 from sklearn.ensemble import RandomForestRegressor
51
52 # Create the model
53 model = Pipeline(steps=[
54     ('preprocessor', preprocessor),          # Apply preprocessing
55     ('regressor', RandomForestRegressor(random_state=42)) # Train Random Forest Regressor
56 ])
57
58 # Train the model
59 model.fit(X_train, y_train)
60

```



This pipeline predicts the expected future donation amount for each donor, enabling nonprofits to plan fundraising strategies and target high-value donors more effectively.

```

1 from sklearn.metrics import mean_squared_error, r2_score
2
3 # Predict on the test set
4 y_pred = model.predict(X_test)
5
6 # Calculate the Mean Squared Error (MSE)
7 mse = mean_squared_error(y_test, y_pred)
8 rmse = mse ** 0.5 # Root Mean Squared Error (RMSE)
9
10 # Calculate the coefficient of determination R²
11 r2 = r2_score(y_test, y_pred)
12
13 print(f"RMSE: {rmse}")
14 print(f"R²: {r2}")
15

```

RMSE: 171.94733378480254
R²: -0.13073355442002632

These metrics allow you to quantify the performance of the regression model in predicting donors' future donation amounts. The MSE measures the average squared difference between predicted and actual values. The RMSE is the square root of the MSE, providing an error expressed in the same units as the target variable (FutureAmount). A lower RMSE indicates better predictive performance.

Model Performance

RMSE \approx 165.5 On average, the model's predictions for future donation amounts deviate by about 165 units from the actual values. Since donations vary widely, this represents a relatively large error, indicating that the predictions are not very precise.

$R^2 \approx -0.10$ An R^2 value below 0 indicates that the model performs worse than a simple baseline that predicts the mean donation for all donors. This suggests that the Random Forest model fails to capture meaningful patterns from the current features for predicting future donation amounts.

Testing model With a new data.

```
1 new_data = {
2     'Age': [45],
3     'Gender': ['Male'],
4     'Location': ['Ouest, Haiti'],
5     'TotalAmount': [1700.00],
6     'Channel': ['Email'],
7     'ImpactType': ['Education'],
8     'DonationCount': [4],
9     'EngagementOutcome': ['none']
10 }
11 new_df = pd.DataFrame(new_data)
```

```
1
2 # future value
3 predicted_value = model.predict(new_df)
4 print(f" Future value expected : {predicted_value[0]}")
```

Future value expected : 225.47299810821997

Recommendations

The use of data plays a crucial role in decision-making. Data can be used for inferential anysis, predictive anlysis, etc... data allows organisation to make better and more informed decision overall.

✚ Contact information 1

Bellow: The Contact Information Of the Two Data scientist in charge of the project.

- First Name: Haender Michael
- Last Name : Jean Louis
- Email : michaelhaenderjeanlouis@gmail.com
- Phone : +509 31742772
- whatsapp : +509 33905060
- LinkedIn : https://www.linkedin.com/in/michael-haender-jean-louis-4b7320316?utm_source=share&utm_campaign=share_via&utm_content=profile&utm_medium=ios_app

✚ Contact information 2

- First Name : John Widno
 - Last Name : DORCY
 - Email : dorcyjohnwidno97@gmail.com
 - Phone : +509 38 05 8388
 - LinkedIn : <https://www.linkedin.com/in/john-widno-dorcy-19399a216/>
-