# Maximizing Donor Retention: Strategies for Long-Term Engagement

By: Franky Araujo | Email: araujofranky@gmail.com



## Objective & Overview

This report aims to enhance donor engagement and increase contributions for nonprofit organizations by leveraging data-driven insights and predictive analytics. It outlines a project focused on analyzing donor lifetime value, churn, and interactions based on donor history and demographics. The report provides a comprehensive methodology, performance evaluation metrics, and actionable strategies for nonprofits to foster enduring donor relationships and maximize fundraising effectiveness.

# Data Insights

The dataset originates from donation transactions logged within a nonprofit's CRM, covering the period from 2015 to April 2024 and pertains to a small to mid-sized nonprofit organization. While invaluable for analysis, the original dataset is restricted from public access and is not available upon request, with donor identifying information conscientiously redacted to safeguard confidentiality.

## About the Data

- The dataset encompasses various data types, including datetime, float, and string, presenting a challenge in standardization and analysis. Addressing the complexities of mixed data types was crucial to reduce the likelihood of any potential misinterpretations or inconsistencies.

- Text elements associated with multiple labels pose challenges, complicating the extraction and interpretation of meaningful information and impacting the precision of analyses.

- The presence of pre-encoded values, such as segment codes and tags, necessitated careful consideration to ensure accurate interpretation. Misinterpreted encoded information could lead to flawed conclusions.

- The dataset's scope is constrained by the organization's size, reflecting varying degrees of data management practices, especially in earlier years.

- Moreover, the dataset's integrity may be affected by missing data due to early data practices (or lack thereof).

- Lastly, it's important to consider the potential impact of the Covid-19 pandemic on trends within the dataset.

- Code Snippet:

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 17356 entries, 0 to 17355
Data columns (total 22 columns):
```

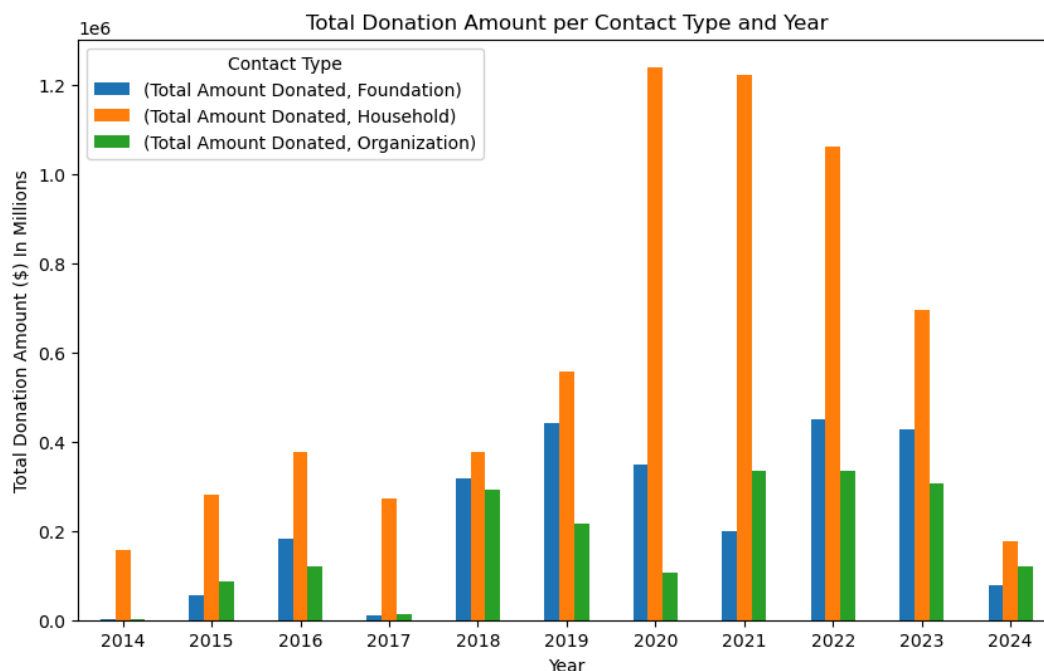| #  | Column | Non-Null Count | Dtype |
| --- | ------ | -------------- | ----- |
| 0  | Contact Id | 17356 non-null | int64 |
| 1  | Contact Primary Gender | 107 non-null | object |
| 2  | Contact Primary Age | 4 non-null | float64 |
| 3  | Contact Type | 17356 non-null | object |
| 4  | Contact Tags | 16914 non-null | object |
| 5  | Gift Date | 17356 non-null | object |
| 6  | Amount | 17356 non-null | float64 |
| 7  | Gift Type | 17356 non-null | object |
| 8  | Notes | 5980 non-null | object |
| 9  | Contact Primary Full Address | 13888 non-null | object |
| 10 | First Gift Date | 17350 non-null | object |
| 11 | First Recurring Gift Date | 10741 non-null | object |
| 12 | Recurring Gift Amount | 5328 non-null | float64 |
| 13 | Check Number | 1785 non-null | object |
| 14 | Segment Name | 10137 non-null | object |
| 15 | Campaign Name | 10137 non-null | object |
| 16 | Contact Primary Address City | 13887 non-null | object |
| 17 | Contact Primary Address State | 13778 non-null | object |
| 18 | Contact Social Score | 17354 non-null | float64 |
| 19 | Contact Primary Birth Year | 4 non-null | float64 |
| 20 | Contact Primary Birth Month | 4 non-null | float64 |
| 21 | Selected Age | 3 non-null | float64 |

```
dtypes: float64(7), int64(1), object(14)
memory usage: 2.9+ MB
```

Addressing these challenges through comprehensive data preprocessing, exploration, and interpretation strategies was essential to maintain the integrity of the dataset and enhance the reliability of the analysis results.

# Initial Findings

Thankfully, despite any challenges, the raw data yielded valuable insights during the exploratory data analysis phase. The subsequent section will delve into the general characteristics of the data before focusing on the response variable, churn.
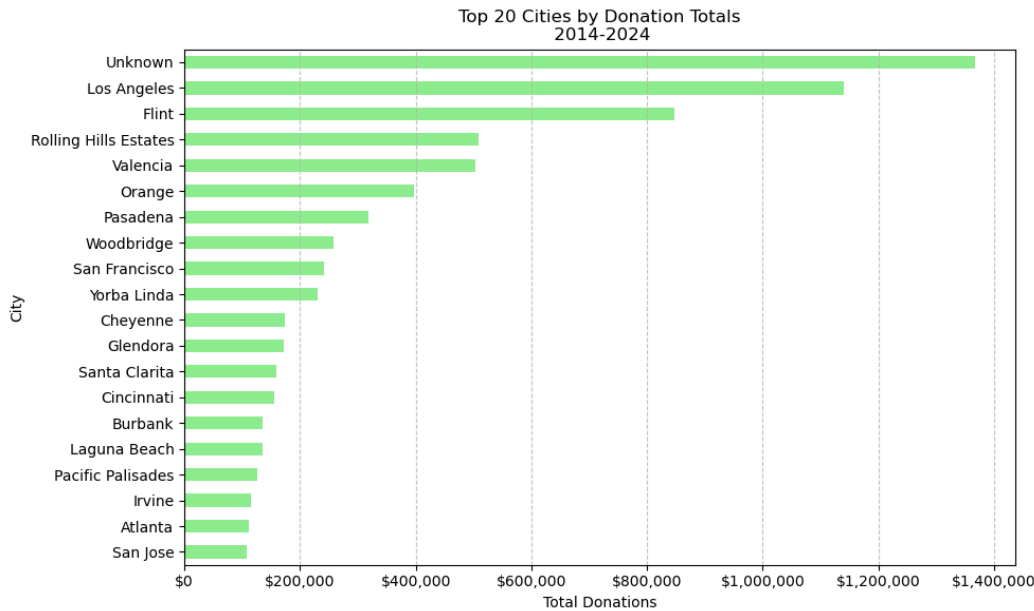
Looking at the donation amount per contract type (per year), the data shows trends across different contact types exhibit varying patterns over the years. Foundation contacts experienced a consistent upward trajectory in donations from 2014 to 2016, culminating in a peak of 182,699.24 USD in 2016. However, donations fluctuated thereafter, hitting a low point of 10,000 USD in 2017 before gradually rebounding. Household contacts demonstrated steady growth in donations, with notable surges in 2020 and 2021, marking the highest total donations of 1,222,522.84 USD in 2021. Meanwhile, Organization contacts initially displayed modest donations in 2014 but witnessed substantial increases in 2015 and 2016. Although there were fluctuations in subsequent years, donations peaked in 2018 at 292,304.59 USD. This provides guidance on the types of donors that should be focused on for fundraising.



Beyond understanding donor types, such as individual contributors, foundations, and companies, exploring the geographic locations of donors adds valuable insights to bolster marketing campaigns and event planning. Consequently, an analysis of donor cities was conducted.
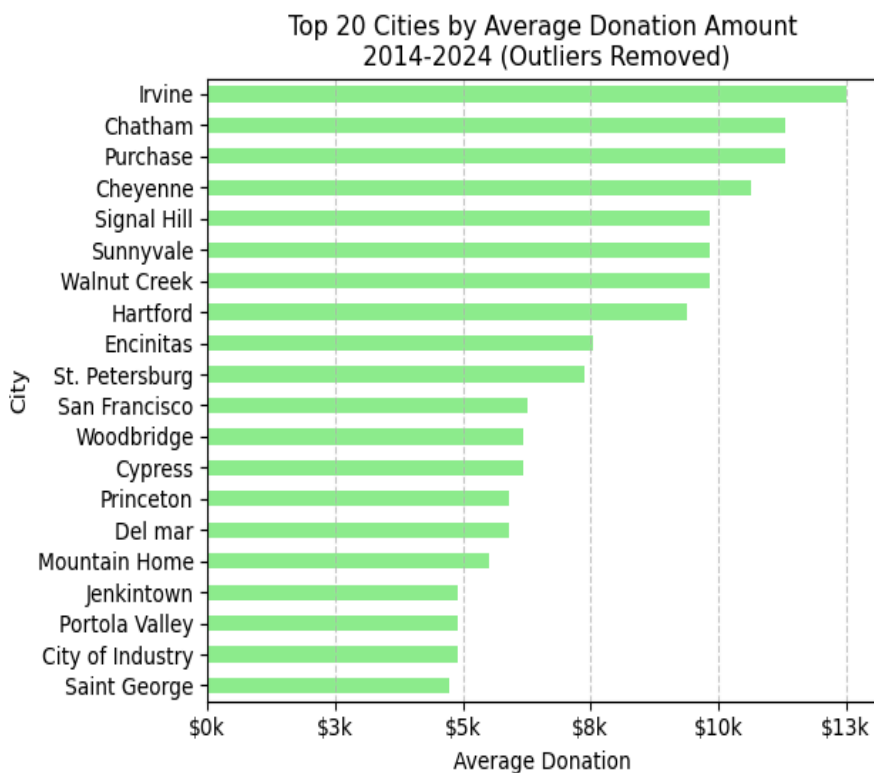
The analysis of donation data reveals intriguing insights into trends across various cities. Topping the list in total donations made is the category of "Unknown," likely attributed to missing or withheld donor city information, followed closely by Los Angeles, reflecting engagement congruent with the organization's location. Other notable cities include Flint, Rolling Hills Estates, and Valencia, indicating pockets of generosity within the dataset.
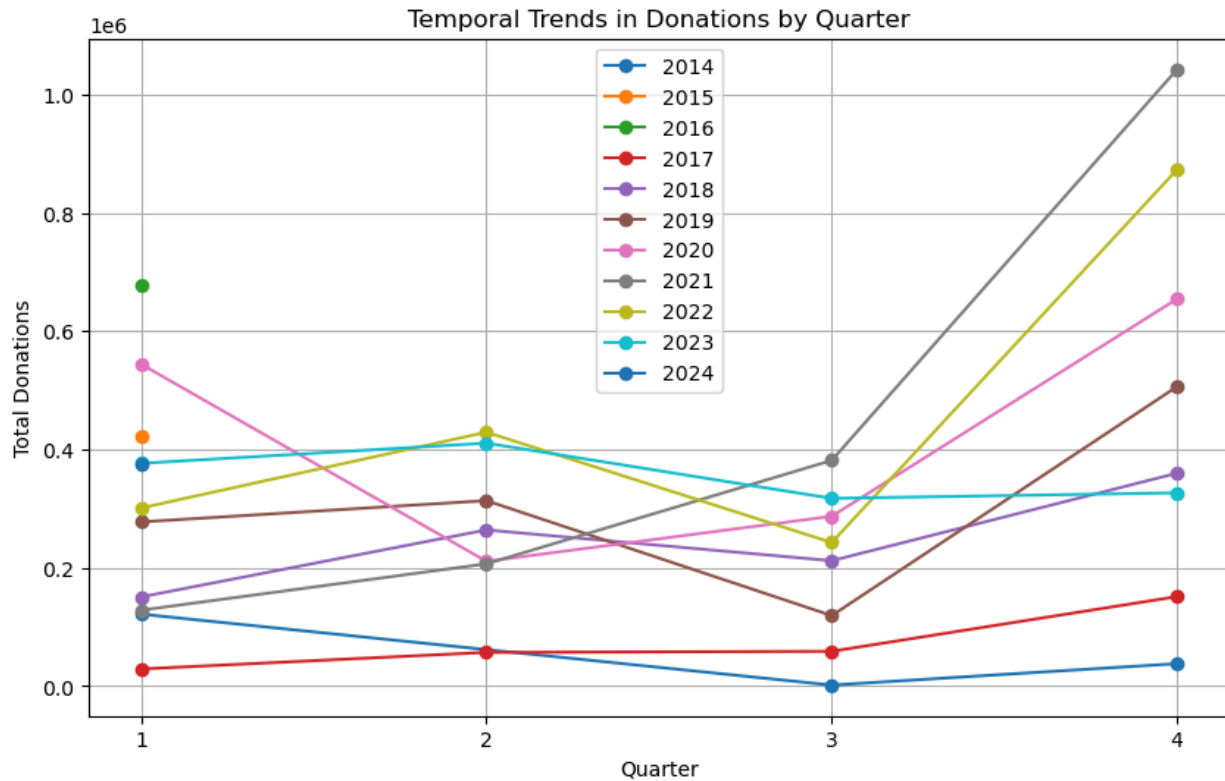
Delving into average donation amounts, Irvine emerges as a standout, boasting the highest average donation per contribution, suggesting a propensity for impactful individual giving within the community. Noteworthy contributions also come from cities like Chatham, Purchase, and Cheyenne, signaling a penchant for higher-value philanthropy in these areas. Interestingly, while Irvine leads in average donation amount, it does not rank among the top cities in total donation amounts, implying a focus on quality over quantity in charitable giving. Examining donation counts, the prevalence of "Unknown" entries underscores the need for enhanced data collection processes, while Los Angeles, Valencia, and others showcase significant community engagement through substantial numbers of recorded donations.

Top 20 Cities by Donation Totals
2014-2024

Top 20 Cities by Average Donation Amount
2014-2024 (Outliers Removed)

In striving to better understand fundraising trends within the nonprofit sector in Los Angeles, the analysis explored donations by quarter. This examination aligns with the sector's overall goal of focusing the majority of fundraising efforts during the holiday season (ie Q4).
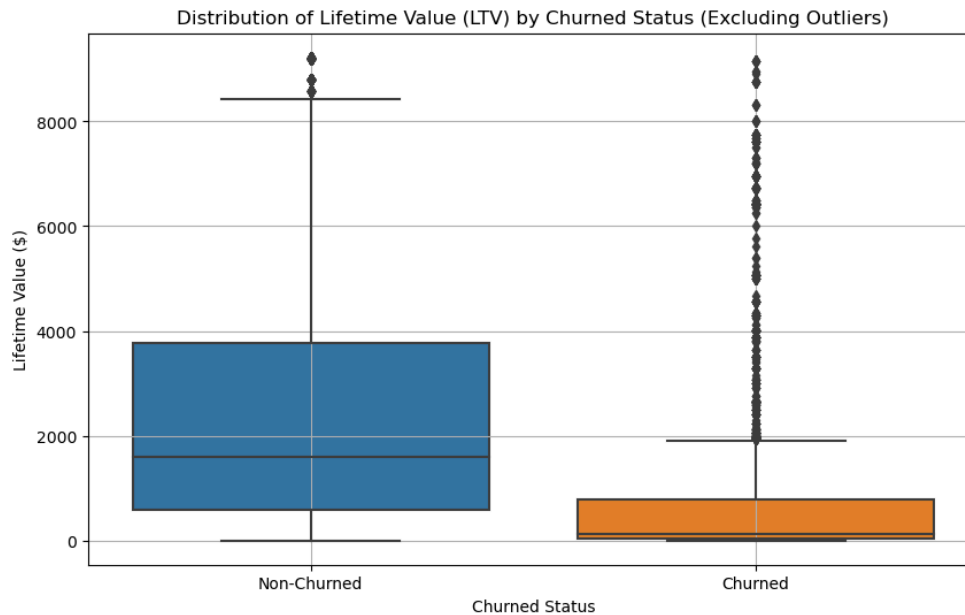


Several factors contribute to this pattern: Firstly, households often prioritize charitable giving during the holidays, motivated by altruism and tax benefits. Additionally, individuals may choose to donate towards the end of the year as they realize surplus income that can be allocated to charitable causes instead of paying it in taxes. This trend underscores the importance of strategic fundraising efforts by nonprofits to capitalize on donor generosity and maximize their impact.
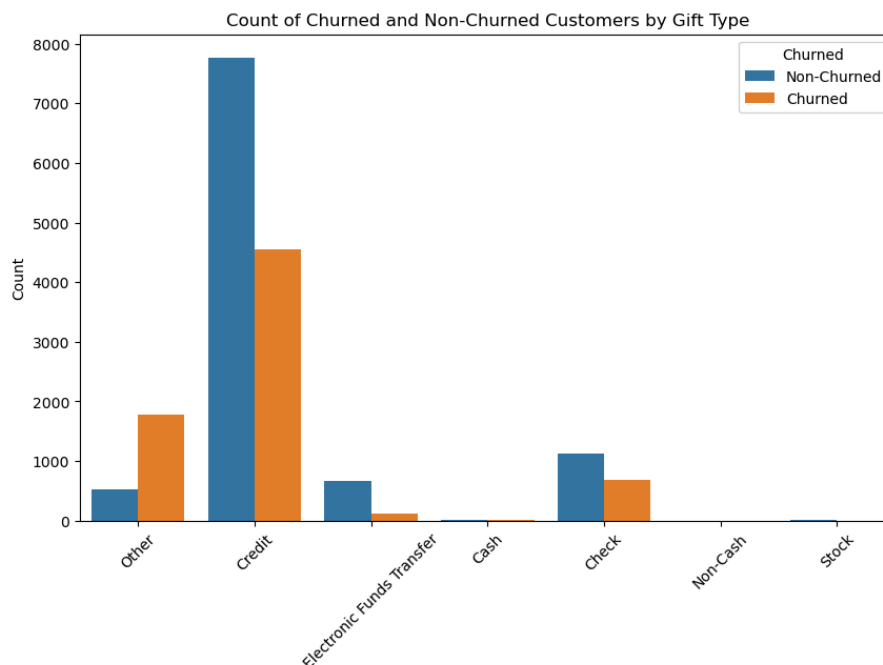
## Churn

Lastly, an exploration into churn was conducted to uncover insights regarding donors who continue to give versus those who cease giving. For this analysis, donors are classified as active if they have made a donation within the 12 months preceding March 2024, while inactive (i.e., churned) donors have not. This timeframe accounts for annual donors who might otherwise be miscategorized as churned.

The comparison between churned and non-churned donors reveals significant differences in donation behavior. Non-churned donors exhibit a higher median donation amount, indicating more consistent or higher-level contributions on average. Additionally, non-churned donors display a wider range of donation amounts, suggesting greater diversity with some individuals contributing

substantially more. In contrast, churned donors have a narrower range centered around a lower median, indicating less varied donation behavior with typically smaller contributions. However, several outliers among churned donors suggest occasional significant one-time contributions. These findings underscore the importance of donor retention strategies to maintain consistent support and encourage sustained engagement over time.



Distribution of Lifetime Value (LTV) by Churned Status (Excluding Outliers)

In addition to this, it was also noted that non-churned donors tend to give via check, credit, or EFTs (Electronic Funds Transfers) as opposed to other methods such as cash donations, in-kind contributions, or one-time online donations. This preference for more formal and recurring payment methods suggests a stronger commitment and ongoing engagement with the organization.



Count of Churned and Non-Churned Customers by Gift Type

It's important to note that while the Jupyter notebooks offer a more extensive exploratory data analysis, the above highlights the key points leading up to preprocessing and modeling which will be discussed next.

# The Preprocessing Approach

In this section, the various preprocessing steps, model selection, and training approaches will be discussed, including:

- Data Cleaning
- Feature Engineering
- Encoding & Scaling
- Train/Test Split Approach
- Using a Baseline Model
- Model Selection & Performance Metrics

## Data Cleaning

The dataset comprises key columns such as Contact Id, Contact Type, Gift Date, Amount, Gift Type, Contact Tags, First Gift Date, and various demographic fields like Contact Primary Gender, Contact Primary Age, and Contact Primary Address. These provide detailed insights into donor behavior and demographics. However, several columns had significant missing values, including Contact Primary Gender, Contact Primary Age, Notes, Contact Primary Full Address, First Recurring Gift Date, Recurring Gift Amount, Check Number, Segment Name, Campaign Name, Contact Primary Address City, Contact Primary Address State, Contact Social Score, Contact Primary Birth Year, Contact Primary Birth Month, and Selected Age. Inconsistent data types, especially in date-related columns, required conversion to datetime formats. Low non-null values in Contact Primary Gender and Contact Primary Age led to their removal.

Specific actions were taken for each column. Columns such as Selected Age, Contact Primary Age, Contact Primary Birth Month, Contact Primary Birth Year, Contact Primary Gender, Check Number, Recurring Gift Amount, and First Recurring Gift Date were removed due to insufficient data or irrelevance. Notes were retained for potential vectorization. Missing values in Campaign Name and Segment Name were imputed with general information, while Contact Primary Address State, Contact Primary Address City, and Contact Primary Full Address were imputed using other address-related fields. Contact Tags were encoded to manage multiple values within each element.

For columns with a high percentage of missing values, strategies included imputation (filling in missing values), removal (eliminating columns or rows with excessive missing data), and encoding missingness (treating missing values as a distinct category).

## Feature Engineering

To augment the dataset for comprehensive analysis and modeling, textual information from the 'Contact Tags' and 'Notes' columns was consolidated into the new 'Tags_Notes_Combined' column. This leverages Natural Language Processing (NLP) techniques to capture nuanced donor characteristics and engagement insights embedded within the text data. By merging these columns, we ensure a better representation of donor attributes, enabling deeper analysis and more accurate predictive modeling.

Additionally, numerical attributes were extracted from datetime columns like 'Gift Date,' 'First Gift Date,' and 'Last Gift Date.' This process involves breaking down datetime information into discrete components such as year, month, day, hour, minute, and second. By capturing temporal nuances in numerical form, we provide a more granular perspective on donor behavior over time, facilitating deeper insights and predictive modeling. These preprocessing steps collectively enhance the dataset's analytical potential, paving the way for robust analysis and actionable insights.

## Encoding & Scaling

One-hot encoding was applied to several variables characterized by a low number of unique values, including 'Gift Type,' 'Campaign Name,' and 'Contact Type.' Additionally, a pragmatic approach was adopted for 'Contact Primary Address State,' encoding whether the donor resided in California or not, given the dataset's California-centric nature. Through one-hot encoding, categorical variables were transformed into binary vectors, enhancing their utility in subsequent analysis and modeling tasks.

For text-based variables such as 'Tags_Notes_Combined,' 'Contact Primary Full Address,' 'Contact Primary Address City,' and 'Segment Name,' an NLP (Natural Language Processing) approach was employed. This involved tokenizing and preprocessing the text data before training Word2Vec models on each column's text content. The Word2Vec models were then used to vectorize the text data, generating numerical representations of the textual information. These vectorized values were stored in a new DataFrame, replacing the original text columns, facilitating the integration of textual information into the dataset for subsequent analysis and modeling tasks.

Standard Scaling was exclusively applied to the continuous variables to prevent features with larger magnitudes from exerting undue influence during model training. Since class labels lack numerical magnitudes and text vectors from Word2Vec are inherently normalized during vectorization, they were unaltered by the scaling process. The continuous variables identified for scaling, such as 'Contact Social Score,' 'Donor Tenure Years,' 'Amount,' 'Average Amount,' and 'LTV,' were standardized using the StandardScaler fitted on the training set. This ensured consistency and compatibility across the dataset while preserving the integrity of categorical variables and text vectors.

## Train/Test Split Approach

The data was then divided into training and testing sets using a test size of 40% and a random state of 42 for reproducibility. A larger test set was chosen due to the relatively smaller size of the dataset, ensuring a more robust evaluation of model performance.

## Baseline Model

To establish a baseline for comparison, predictions were obtained from the Dummy Classifier `(DummyClassifier(strategy = 'stratified'))` . The accuracy, precision, recall, F1-score, and ROC AUC were calculated to assess the performance of the dummy classifier.

- **Accuracy**: The dummy classifier achieved an accuracy of approximately 52.25%. This indicated that around 52.25% of the model's predictions were correct.
- **Precision**: With a precision of approximately 41.80%, the dummy classifier correctly identified around 41.80% of the positive predictions among all positive instances.
- **Recall**: The recall, or sensitivity, of the classifier was approximately 43.08%. This suggested that around 43.08% of the actual positive instances were correctly identified by the model.
- **F1-score**: The F1-score, which balanced precision and recall, was approximately 42.43%. A higher F1-score indicated better overall performance.
- **ROC AUC**: The ROC AUC (Receiver Operating Characteristic Area Under the Curve) of the classifier was approximately 50.83%. This metric evaluated the classifier's ability to distinguish between positive and negative instances, with a value close to 0.5 indicating performance similar to random guessing.

These baseline performance metrics provided a reference point for evaluating the performance of more sophisticated classifiers in subsequent analyses

## Model Selection & Performance Metrics

The dataset underwent evaluation with multiple classification models to determine their performance in predicting churn. Six models were assessed: RandomForest, GradientBoosting, DecisionTree, SVM, KNN, and ANN.

**Accuracy**: GradientBoosting and DecisionTree exhibited almost perfect accuracy scores, while SVM showed the lowest accuracy, indicating a struggle with the dataset.
**Precision**: GradientBoosting demonstrated the highest precision, indicating reliability in minimizing false positives, whereas KNN displayed the lowest precision.
**Recall**: GradientBoosting nearly captured all positive instances, while ANN struggled with identifying true positives, showing the lowest recall.
**F1 Score**: GradientBoosting achieved the highest F1 score, indicating a good balance between precision and recall, while ANN displayed the lowest, indicating significant imbalance.

**ROC AUC**: GradientBoosting excelled in distinguishing between classes, whereas SVM struggled, reflecting in its lower AUC.

**Cross-validation score**: DecisionTree showed perfect cross-validation scores, indicating highly reliable performance across different subsets of data. SVM, on the other hand, displayed inconsistent scores, reflecting instability.

**Cross-validation score average**: DecisionTree stood out with a perfect average performance, while SVM had the lowest, indicating overall weaker performance.

Based on the results, the GradientBoosting model was selected for further application and testing, with potential fine-tuning to enhance its performance.

Looking ahead, several pivotal next steps can be considered based on the analysis results. These may include conducting hyperparameter tuning, performing additional feature importance analysis to identify influential features, exploring ensemble methods, and delving into model interpretability techniques for decision-making.

By integrating these next steps into the workflow, the performance and efficacy of the GradientBoosting model can be further refined, ultimately facilitating more accurate predictions and actionable insights.

# From Analysis to Action: Effective Strategies

Based on this project, the following strategies have been crafted to empower nonprofits in Los Angeles to enhance donor engagement, optimize fundraising efforts, and foster long-term sustainability. These strategies leverage data-driven insights and tailored approaches to cultivate meaningful connections with donors, maximize fundraising impact, and mitigate churn risk.

1. Segmented Donor Engagement

   Strategy: Utilize segmentation to tailor engagement strategies for different donor groups based on their giving patterns, preferences, and demographics.

   Implementation: Develop communication plans, fundraising appeals, and stewardship efforts tailored to the unique characteristics and preferences of each donor segment. For example, donors who prefer recurring donations may receive targeted messaging highlighting the impact of sustained giving, while major donors may be invited to exclusive events or provided with customized acknowledgment and recognition.

   Impact: By engaging donors in a more personalized and meaningful way, nonprofits can foster stronger connections, increase donor loyalty, and ultimately drive higher levels of giving and support over time.

2. Strategic Fundraising Timing

   Strategy: Analyze seasonal donation trends and geographic preferences to strategically time fundraising campaigns and maximize fundraising impact.

   Implementation: Use insights from the analysis to identify peak donation periods, such as year-end giving or seasonal holidays, and concentrate fundraising efforts during these times. Additionally, target regions or cities with higher average donation amounts or engagement levels for more focused outreach and solicitation efforts.

   Impact: By aligning fundraising efforts with periods of heightened donor activity and focusing resources on areas with the greatest potential for donor engagement, nonprofits can optimize fundraising outcomes and achieve greater efficiency in their fundraising campaigns.

3. Proactive Donor Retention

   Strategy: Implement proactive donor retention strategies informed by churn analysis to prevent donor attrition and cultivate long-term donor relationships.

   Implementation: Develop targeted retention initiatives, such as personalized stewardship programs, donor appreciation events, and milestone acknowledgments, to foster ongoing engagement and loyalty among donors. Additionally, identify key factors contributing to donor churn, such as communication gaps or lack of engagement, and address these issues proactively to mitigate churn risk.

   Impact: By prioritizing donor retention and implementing proactive strategies to address churn risk factors, nonprofits can maintain a stable donor base, increase donor lifetime value, and minimize the need for costly donor acquisition efforts, ultimately driving greater sustainability and impact for their organization.

# Closing Remarks

In closing, this report underscores the role of data-driven strategies in empowering nonprofit organizations to maximize their fundraising effectiveness and donor engagement. By harnessing the insights derived from comprehensive data analysis, nonprofits can not only optimize their fundraising outcomes but also cultivate enduring relationships with their supporters. As nonprofits continue to navigate an evolving landscape, integrating these data-driven strategies into their organizational frameworks will be essential for driving sustainable growth, fostering community impact, and advancing their missions in a rapidly changing world.

# Project Links

**GitHub Repository**

- [Link to Your GitHub Repository](Link to Your GitHub Repository)

**Jupyter Notebooks**

1. [Data Wrangling Initial Exploration & Cleaning](Data Wrangling Initial Exploration & Cleaning)
2. [Exploratory Data Analysis](Exploratory Data Analysis)
3. [Pre-processing, Training Data Development, and Modeling](Pre-processing, Training Data Development, and Modeling)