Dan Schumacher, 5/7/24

# Enhancing Donor Prediction Accuracy for the National Veterans' Organization through Advanced Modeling Techniques

## Introduction

Non-profit organizations significantly depend on donations to fund their operations and deliver services. The National Veterans' Organization (NVO), with a vast database of over 13 million potential donors, primarily utilizes direct mail for fundraising. However, with an average donation return of $13 against a mailing cost of $13.60 for every 20 letters sent, the organization faces a net loss. This scenario underlines the critical need for a more targeted marketing strategy that reduces costs and enhances the efficiency of fundraising efforts.

## Objectives

The objective of this project is twofold:

1. Utilize various predictive modeling techniques to optimize the direct mail campaign, aiming to increase the profitability of the National Veterans' Organization.
2. Apply and showcase the modeling techniques learned in Dr. Campbell's class and beyond this semester, emphasizing practical applications in a real-world context.

## Data Sources

The data was made available on Canvas by Dr. Campbell.

## Details of the Data

The training data contains 3,000 records and is evenly balanced between Donors and Non Donors. There are 22 variables in the dataset. These include:

- Zip code category, homeowner status, number of children, income, sex, wealth category, homevalue, median family income, average family income, number of promotions sent to the individual, lifetime number of promotions received by the individual, sum of all donations given, dollar amount of last gift, months since last donation, number of months between previous 2 gifts, average gift from donor.

## Types of analysis Performed

I performed both traditional modeling and utilized an artificial neural network to generate my predictions.

Dan Schumacher, 5/7/24

## Traditional Modeling

In the traditional modeling, I started by creating dummy variables of all of the predictors. Next I removed predictor variables that were highly correlated with one another. Starting with all predictors, I would calculate the Variance Inflation Factor (Vif) of each variable, then if any variable had a Vif greater than the cut off (I used 10), I would remove the variable and recalculate the Vif without that variable. This process repeated until all Vifs were under the cut off. At the end of the process I had removed five variables
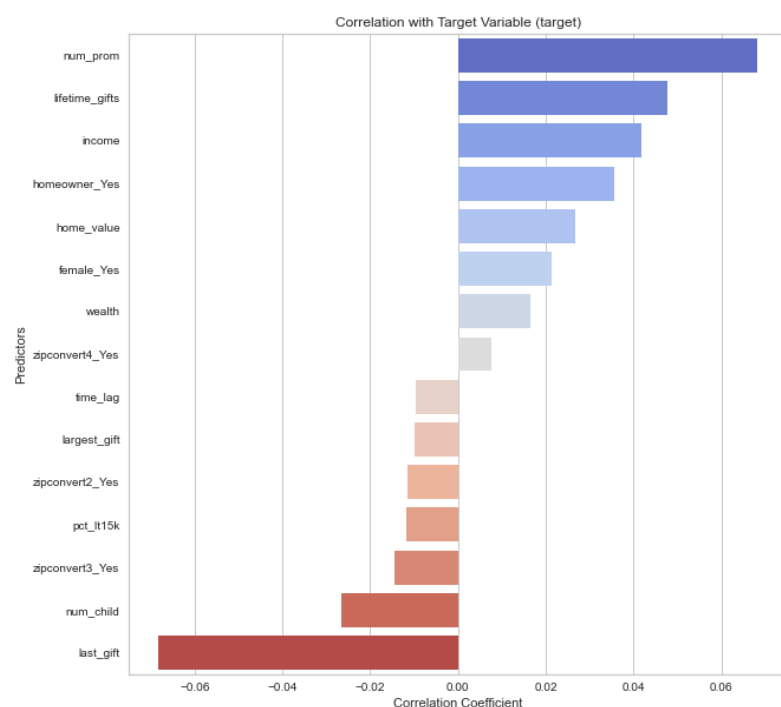
REMOVED: ['zipconvert5_Yes', 'avg_fam_inc', 'months_since_donate', 'med_fam_inc', 'avg_gift']

I fit the data to a Logistic Regression, LDA, QDA, KNN (k=1), and a Naive Bayes classifier.

The results are shown below. Our top performer on both the dev and test set was LDA. It is also noteworthy that LDA had the highest F1 score, achieving a good balance between precision and recall.

| name | tp | tn | fp | fn | acc | prec | recall | f1 | test_acc |
|---|---|---|---|---|---|---|---|---|---|
| Logistic Regression | 167 | 151 | 145 | 137 | 0.530000 | 0.535256 | 0.549342 | 0.542208 | 0.533333 |
| LDA | 163 | 166 | 149 | 122 | 0.548333 | 0.522436 | 0.571930 | 0.546064 | 0.558333 |
| QDA | 26 | 268 | 286 | 20 | 0.490000 | 0.083333 | 0.565217 | 0.145251 | 0.525000 |
| KNN | 148 | 153 | 164 | 135 | 0.501667 | 0.474359 | 0.522968 | 0.497479 | 0.475000 |
| Naïve Bayes | 38 | 263 | 274 | 25 | 0.501667 | 0.121795 | 0.603175 | 0.202667 | 0.516667 |

To get a rough idea of what variables our model is using to generate predictions, we can make a correlation chart.



Correlation with Target Variable (target)

Here we can see that the larger the bar, the more correlated the variable is with the target. Red bars show negative correlations while blue bars show positive correlations. One thought I had, to achieve superior test accuracy, was to only use the most correlated variables on each side. I tested only using variables with an absolute correlation coefficient above .05. Unfortunately, This did not improve results.

Speaking of which, it is worth mentioning that I attempted many more models, and much more complex ones. I had a total of 79 entries on the competition's website. For the more complex models I tried transforming the continuous variables by log, squared, cubed, square root, inverse, boxcox, sigmoid, sine and cosine transformations. After Vif, in the complex modeling, I had a total of 84 variables. However, they did not help. Both the dev scores and test accuracy were several percentages shorter across each category.

## ANN

Another approach that I took was making an artificial neural network for the classification task. This network had the following parameters:
1. An input layer with both categorical and continuous features processed via embedding and batch normalization respectively
2. Multiple hidden layers with configurations of [200, 100] neurons each incorporating ReLU activation and dropout regularization set at 0.4
3. Final output layer using a softmax function for binary classification
4. Robustly scaled input variables to boost performance
5. 100 epochs of training with an early stopping patience of 5
6. Cross validation with 20 splits

This approach achieved much better dev accuracy with an OOB accuracy of 59%. However, at test time it performed practically the same as LDA, and it had the downside of being much more complex, less interpretable, and doesn't always produce the same results.

## Results

The models displayed varying levels of success:
- Accuracy Metrics: Each traditional model was evaluated based on accuracy, precision, recall, and F1-score. The LDA emerged as the most accurate model, as well as the one with the highest F1, showcasing a balance between precision and recall.
- Feature Importance: Looking at the correlation with target variable plot, we can see that the target variable is most strongly associated with number of promotions, and negatively associated with last gift. This loosely suggests that individuals who have donated recently are more likely to donate again, and people who receive a lot of promotions tend to donate more.

## Discussion

The results indicated that, with all their bells and whistles, advanced modeling techniques such as artificial neural networks do not always have better predictive power then their traditional less complicated counterparts. However, none of the variables exhibited extremely high correlations with the target, suggesting potential limitations in the dataset's explanatory power.

## Conclusion and Future Work

The study successfully demonstrated the potential predictive models to enhance donor prediction accuracy, thereby aiding the NVO in optimizing its fundraising strategies. For future work, further experimentation with ANN hyperparameters (such as learning rates, epochs, batch sizes, dropout rates) is recommended. Additionally, obtaining more comprehensive data, including demographic information, military service details, and any other variables that could potentially have a higher correlation with the target variable, would be useful. Having more variables more closely correlated to the target would significantly improve predictive power.

## Appendices

The appendices contain two different notebooks. The first is the simple modeling I performed to achieve the best results. The second includes transformed variables (test accuracy was not included, as all results were in the mid 40s). Finally there is a python script that is my artificial neural network.

## Simple Model By Hand

```
In [126]:   1  import pandas as pd
            2  pd.set_option("max_colwidth", None)
            3
            4  import pycaret
            5  import numpy as np
            6  import matplotlib.pyplot as plt
            7  from pycaret.classification import *
            8  from sklearn.model_selection import train_test_split
            9  from sklearn.metrics import accuracy_score
           10
           11  from functions.homebrew import *
           12  import numpy as np
           13  import pandas as pd
           14
           15  from sklearn.linear_model import LogisticRegression
           16  from sklearn.discriminant_analysis import LinearDiscriminantAnalysis as LDA, QuadraticDiscriminantAnalysis as QDA
           17  from sklearn.naive_bayes import GaussianNB
           18  from sklearn.neighbors import KNeighborsClassifier
           19  from sklearn.preprocessing import StandardScaler
           20  from sklearn.model_selection import train_test_split, cross_val_score
           21  from sklearn.metrics import accuracy_score
           22  from tqdm import tqdm
           23  from itertools import combinations
           24  import pickle
           25  import os
           26
           27  # If you're using statsmodels or ISLP for specific tasks, keep these imports
           28  import statsmodels.api as sm
           29  # Assuming ISLP and homebrew are custom modules specific to your project
           30  from ISLP import load_data, confusion_table
           31  from ISLP.models import ModelSpec as MS, summarize, contrast
           32  import statsmodels.api as sm
           33  from scipy import stats
```

## Helper Functions

```
In [127]:   1  def convert_confusion_matrix(df, name):
            2      """
            3      Converts a confusion matrix dataframe into a format with columns for model name, TP, TN, FP, FN.
            4
            5      Args:
            6      df (pd.DataFrame): Confusion matrix dataframe with multi-index (Truth, Predicted) and columns [0, 1].
            7
            8      Returns:
            9      pd.DataFrame: Reformatted dataframe with model evaluation metrics.
           10      """
           11      # Extracting the values from the confusion matrix
           12      tn, fp, fn, tp = df.iloc[0, 0], df.iloc[0, 1], df.iloc[1, 0], df.iloc[1, 1]
           13      acc = (tp + tn) / (tp + tn + fp + fn)
           14      prec = tp / (tp +fp)
           15      recall = tp / (tp + fn)
           16      f1 = 2 * ((prec * recall)/(prec + recall))
           17      # Creating a new dataframe with the desired format
           18      metrics_df = pd.DataFrame({
           19          "name": name,
           20          "tp": [tp],
           21          "tn": [tn],
           22          "fp": [fp],
           23          "fn": [fn],
           24          'acc': acc,
           25          'prec': prec,
           26          'recall': recall,
           27          'f1': f1
           28      })
           29
           30      return metrics_df
```

```
In [128]:   1  def format_results(df):
            2      df = np.where(df == 1, 'Donor','No Donor')
            3      return df
```

## LOAD DATA

```
In [140]:   1  df = pd.read_csv('./data/df.csv').drop('Unnamed: 0', axis=1)
```

```python
In [141]:  1  train = df[df['type'] == 'train'].drop('type',axis =1)
           2  dev = df[df['type'] == 'dev'].drop('type',axis =1)
           3  test = df[df['type'] == 'test'].drop('type',axis =1)
```

## VIF

```python
In [142]:  1  dummies = pd.get_dummies(df, drop_first=True)
           2
           3  kept, removed = remove_high_vif_features(X=dummies.drop('target_No Donor', axis=1), y=dummies['target_No Donor'], vif_threshol
           4  print('REMOVED:', removed)
```
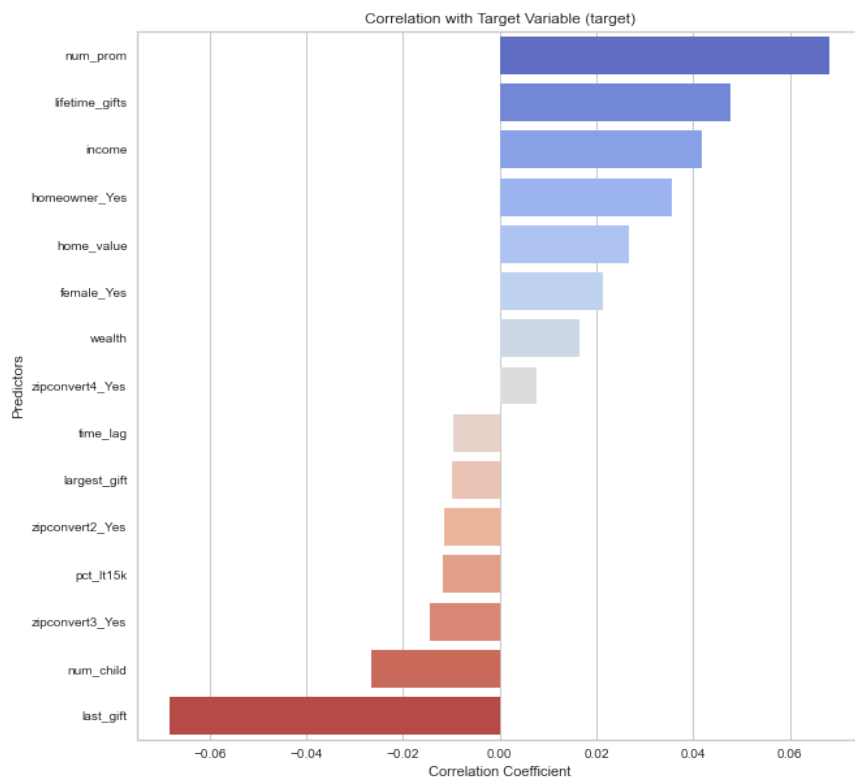
REMOVED: ['avg_fam_inc', 'months_since_donate', 'zipconvert5_Yes', 'med_fam_inc', 'avg_gift']

```python
In [144]:  1  kept['target'] = (df['target'] == 'Donor').astype(int)
```

```python
In [145]:  1  train = kept[kept['type_train'] ==1]
           2  dev = kept[(kept['type_test'] == 0) & (kept['type_train'] == 0)]
           3  test = kept[kept['type_test'] ==1]
```

```python
In [154]:  1  # for data in [train, dev, test]:
           2  #     data.drop('type_train', inplace = True, axis = 1)
           3  #     data.drop('type_test', inplace = True, axis = 1)
           4  test = test.drop('target',axis = 1)
```

```python
In [174]:  1  def cor_bars(df):
           2      corr = df.corr()
           3      # Isolating the column that represents the correlation with the target variable
           4      target_corr = corr['target'].sort_values(ascending=False)
           5
           6      # Removing the target variable from itself to avoid a perfect correlation display
           7      target_corr = target_corr.drop(labels=['target'])
           8
           9      # Plotting the correlations for visual representation
          10      plt.figure(figsize=(10, 10))
          11      sns.barplot(x=target_corr.values, y=target_corr.index, palette='coolwarm')
          12      plt.title('Correlation with Target Variable (target)')
          13      plt.xlabel('Correlation Coefficient')
          14      plt.ylabel('Predictors')
          15      plt.show()
          16  cor_bars(train)
```



## Logistic Regression

```
In [155]:   1  for col in kept.columns:
            2      print(col)
```

```
num_child
income
wealth
home_value
pct_lt15k
num_prom
lifetime_gifts
largest_gift
last_gift
time_lag
zipconvert2_Yes
zipconvert3_Yes
zipconvert4_Yes
homeowner_Yes
female_Yes
type_test
type_train
target
```

```
In [156]:   1  results_df = pd.DataFrame()
```

```
In [157]:   1  # Selecting features and target variable for training data
            2  X_train = train.drop(['target'], axis =1 )
            3  y_train = train['target']
            4  X_test = dev.drop(['target'], axis = 1)
            5  y_test = dev['target']
            6
            7  # Fitting logistic regression model
            8  glm = sm.GLM(y_train, X_train, family=sm.families.Binomial())
            9  glm = glm.fit()
           10
           11  # Summarizing results
           12  # print(results.summary())
```

```
In [158]:   1  log_preds = (glm.predict(X_test) >= 0.5).astype(int)
            2  log_acc = accuracy_score(log_preds, y_test)
            3  print(log_acc)
            4
            5  d = confusion_table(log_preds,y_test)
            6  results_df = pd.concat([results_df,convert_confusion_matrix(d, 'Logistic Regression')])
            7
            8  log_test_preds = (glm.predict(test) >= 0.5).astype(int)
            9  log_test_preds = format_results(log_test_preds)
           10
           11  save_df = pd.DataFrame(log_test_preds, columns=['values'])
           12  save_df.to_csv('./preds/log.csv', index=False)
```

```
0.53
```

## LDA

```
In [159]:   1  lda = LDA(store_covariance=True)
            2  lda.fit(X_train, y_train)
            3
            4  lda_preds = lda.predict(X_test)
            5
            6  lda_acc = accuracy_score(lda_preds,y_test)
            7  print(lda_acc)
            8
            9  d = confusion_table(lda_preds,y_test)
           10  results_df = pd.concat([results_df,convert_confusion_matrix(d, 'LDA')])
           11
           12
           13  lda_test_preds = (lda.predict(test) >= 0.5).astype(int)
           14  lda_test_preds = format_results(lda_test_preds)
           15
           16  save_df = pd.DataFrame(lda_test_preds, columns=['values'])
           17  save_df.to_csv('./preds/lda.csv', index=False)
```

```
0.5483333333333333
```

# QDA.

```
In [160]:    1  qda = QDA(store_covariance=True)
             2  qda.fit(X_train, y_train)
             3
             4  qda_preds = qda.predict(X_test)
             5
             6  qda_acc = accuracy_score(qda_preds,y_test)
             7
             8  print(qda_acc)
             9
            10  d = confusion_table(qda_preds,y_test)
            11  results_df = pd.concat([results_df,convert_confusion_matrix(d, 'QDA')])
            12
            13  qda_test_preds = (qda.predict(test) >= 0.5).astype(int)
            14  qda_test_preds = format_results(qda_test_preds)
            15
            16  save_df = pd.DataFrame(qda_test_preds, columns=['values'])
            17  save_df.to_csv('./preds/qda.csv', index=False)
```

0.49

## KNN

```
In [161]:    1  df['type']
```

```
Out[161]:  0       train
           1       train
           2         dev
           3       train
           4       train
                   ...
           3115     test
           3116     test
           3117     test
           3118     test
           3119     test
           Name: type, Length: 3120, dtype: object
```

```
In [162]:    1  knn1 = KNeighborsClassifier(n_neighbors=1)
             2  knn1.fit(X_train, y_train)
             3  knn1_pred = knn1.predict(X_test)
             4  knn1_acc = accuracy_score(knn1_pred,y_test)
             5
             6  print(knn1_acc)
             7
             8  d = confusion_table(knn1_pred, y_test)
             9  results_df = pd.concat([results_df,convert_confusion_matrix(d, 'KNN')])
            10
            11  knn1_test_preds = (knn1.predict(test) >= 0.5).astype(int)
            12  knn1_test_preds = format_results(knn1_test_preds)
            13
            14  save_df = pd.DataFrame(knn1_test_preds, columns=['values'])
            15  save_df.to_csv('./preds/knn1.csv', index=False)
```

0.5016666666666667

## NB

```
In [163]:    1  nb = GaussianNB()
             2  nb.fit(X_train, y_train)
             3  nb_preds = nb.predict(X_test)
             4  nb_acc = accuracy_score(nb_preds,y_test)
             5
             6  print(nb_acc)
             7  save_df = pd.DataFrame(nb_preds, columns=['values'])
             8  save_df.to_csv('./preds/nb.csv', index=False)
             9
            10  d = confusion_table(nb_preds, y_test)
            11  results_df = pd.concat([results_df,convert_confusion_matrix(d, 'Naïve Bayes')])
            12
            13  nb_test_preds = (nb.predict(test) >= 0.5).astype(int)
            14  nb_test_preds = format_results(nb_test_preds)
            15
            16  save_df = pd.DataFrame(nb_test_preds, columns=['values'])
            17  save_df.to_csv('./preds/nb.csv', index=False)
```

0.5016666666666667

```
In [1]:      1  # results_df
```

```
In [168]:   1  test_acc = { # these are from running on the website
            2      'log': 0.5333333,
            3      'lda': 0.5583333,
            4      'qda': 0.525,
            5      'knn':  0.475,
            6      'nb': 0.5166667,
            7  }
```

```
In [169]:   1  results_df['test_acc'] = test_acc.values()
```

```
In [170]:   1  results_df
```

Out[170]:

| | name | tp | tn | fp | fn | acc | prec | recall | f1 | test_acc |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | Logistic Regression | 167 | 151 | 145 | 137 | 0.530000 | 0.535256 | 0.549342 | 0.542208 | 0.533333 |
| 0 | LDA | 163 | 166 | 149 | 122 | 0.548333 | 0.522436 | 0.571930 | 0.546064 | 0.558333 |
| 0 | QDA | 26 | 268 | 286 | 20 | 0.490000 | 0.083333 | 0.565217 | 0.145251 | 0.525000 |
| 0 | KNN | 148 | 153 | 164 | 135 | 0.501667 | 0.474359 | 0.522968 | 0.497479 | 0.475000 |
| 0 | Naïve Bayes | 38 | 263 | 274 | 25 | 0.501667 | 0.121795 | 0.603175 | 0.202667 | 0.516667 |

## By hand Complex

```
In [49]:    1  import pandas as pd
            2  pd.set_option("max_colwidth", None)
            3
            4  import pycaret
            5  import numpy as np
            6  import matplotlib.pyplot as plt
            7  from pycaret.classification import *
            8  from sklearn.model_selection import train_test_split
            9  from sklearn.metrics import accuracy_score
           10
           11  from functions.homebrew import *
           12  import numpy as np
           13  import pandas as pd
           14
           15  from sklearn.linear_model import LogisticRegression
           16  from sklearn.discriminant_analysis import LinearDiscriminantAnalysis as LDA, QuadraticDiscriminantAnalysis as QDA
           17  from sklearn.naive_bayes import GaussianNB
           18  from sklearn.neighbors import KNeighborsClassifier
           19  from sklearn.preprocessing import StandardScaler
           20  from sklearn.model_selection import train_test_split, cross_val_score
           21  from sklearn.metrics import accuracy_score
           22  from tqdm import tqdm
           23  from itertools import combinations
           24  import pickle
           25  import os
           26
           27  # If you're using statsmodels or ISLP for specific tasks, keep these imports
           28  import statsmodels.api as sm
           29  # Assuming ISLP and homebrew are custom modules specific to your project
           30  from ISLP import load_data, confusion_table
           31  from ISLP.models import ModelSpec as MS, summarize, contrast
           32  import statsmodels.api as sm
           33  from scipy import stats
```

## Helper Functions

```
In [50]:    1  def add_transformations(data, cont_cols):
            2      for var in cont_cols:
            3          data[f'log_{var}'] = np.log(data[var] + 1)
            4          data[f'sq_{var}'] = data[var]**2
            5          data[f'sqrt_{var}'] = np.sqrt(data[var])
            6          data[f'inv_{var}'] = 1 / (data[var] + 1)
            7          data[f'boxcox_{var}'], _ = stats.boxcox(data[var] + 1)
            8          data[f'sigmoid_{var}'] = 1 / (1 + np.exp(-data[var]))
            9          data[f'sin_{var}'] = np.sin(data[var])
           10          data[f'cos_{var}'] = np.cos(data[var])
```

```
In [51]:   1  def convert_confusion_matrix(df, name):
           2      """
           3      Converts a confusion matrix dataframe into a format with columns for model name, TP, TN, FP, FN.
           4
           5      Args:
           6      df (pd.DataFrame): Confusion matrix dataframe with multi-index (Truth, Predicted) and columns [0, 1].
           7
           8      Returns:
           9      pd.DataFrame: Reformatted dataframe with model evaluation metrics.
          10      """
          11      # Extracting the values from the confusion matrix
          12      tn, fp, fn, tp = df.iloc[0, 0], df.iloc[0, 1], df.iloc[1, 0], df.iloc[1, 1]
          13      acc = (tp + tn) / (tp + tn + fp + fn)
          14      prec = tp / (tp +fp)
          15      recall = tp / (tp + fn)
          16      f1 = 2 * ((prec * recall)/(prec + recall))
          17      # Creating a new dataframe with the desired format
          18      metrics_df = pd.DataFrame({
          19          "name": name,
          20          "tp": [tp],
          21          "tn": [tn],
          22          "fp": [fp],
          23          "fn": [fn],
          24          'acc': acc,
          25          'prec': prec,
          26          'recall': recall,
          27          'f1': f1
          28      })
          29
          30      return metrics_df
```

```
In [52]:   1  def format_results(df):
           2      df = np.where(df == 1, 'Donor','No Donor')
           3      return df
```

## LOAD DATA

```
In [53]:   1  df = pd.read_csv('./data/df.csv').drop('Unnamed: 0', axis=1)
```

```
In [54]:   1  train = df[df['type'] == 'train'].drop('type',axis =1)
           2  dev = df[df['type'] == 'dev'].drop('type',axis =1)
           3  test = df[df['type'] == 'test'].drop('type',axis =1)
```

## VIF

```
In [7]:    1  dummies = pd.get_dummies(df, drop_first=True)
           2
           3  cat_cols = [
           4      'zipconvert2_Yes', 'zipconvert3_Yes', 'zipconvert4_Yes', 'boxcox_zipconvert5_Yes',
           5      'homeowner_Yes', 'female_Yes', 'type_train', 'type_dev', 'type_test'
           6  ]
           7
           8  cont_cols = [col for col in dummies.columns if col not in cat_cols + ['target']]
           9  add_transformations(dummies, cont_cols)
          10
          11  kept, removed = remove_high_vif_features(X=dummies.drop('target_No Donor', axis=1), y=dummies['target_No Donor'], vif_threshol
          12  print('REMOVED:', removed)
```

REMOVED: ['num_child', 'income', 'cos_target_No Donor', 'sin_target_No Donor', 'sigmoid_target_No Donor', 'boxcox_target_No Dono
r', 'inv_target_No Donor', 'sqrt_target_No Donor', 'sq_target_No Donor', 'cos_zipconvert5_Yes', 'sin_zipconvert5_Yes', 'sigmoid_z
ipconvert5_Yes', 'boxcox_zipconvert5_Yes', 'inv_zipconvert5_Yes', 'sqrt_zipconvert5_Yes', 'sq_zipconvert5_Yes', 'log_zipconvert5_
Yes', 'log_income', 'sq_income', 'log_num_child', 'sq_num_child', 'sigmoid_avg_fam_inc', 'inv_num_child', 'sqrt_months_since_dona
te', 'sqrt_med_fam_inc', 'sq_wealth', 'boxcox_avg_fam_inc', 'boxcox_num_prom', 'log_last_gift', 'months_since_donate', 'log_avg_g
ift', 'inv_home_value', 'inv_avg_fam_inc', 'boxcox_time_lag', 'log_wealth', 'inv_med_fam_inc', 'log_largest_gift', 'boxcox_home_v
alue', 'sqrt_num_prom', 'log_lifetime_gifts', 'sqrt_income', 'boxcox_pct_lt15k', 'sqrt_wealth', 'sq_months_since_donate', 'sqrt_a
vg_fam_inc', 'sqrt_avg_gift', 'boxcox_med_fam_inc', 'sqrt_time_lag', 'boxcox_largest_gift', 'sqrt_pct_lt15k', 'inv_income', 'sqrt
_home_value', 'sigmoid_num_child', 'log_num_prom', 'sqrt_last_gift', 'log_months_since_donate', 'boxcox_avg_gift', 'sqrt_lifetime
_gifts', 'avg_fam_inc', 'wealth', 'largest_gift', 'med_fam_inc', 'boxcox_last_gift', 'inv_pct_lt15k', 'zipconvert5_Yes', 'num_pro
m', 'pct_lt15k', 'home_value', 'log_avg_fam_inc', 'inv_wealth', 'log_time_lag', 'boxcox_lifetime_gifts', 'sqrt_num_child', 'avg_g
ift', 'inv_time_lag', 'log_med_fam_inc', 'boxcox_income', 'inv_largest_gift', 'inv_last_gift', 'sq_avg_fam_inc', 'inv_num_prom',
'last_gift']

```
In [55]:   1  final_vars = list(kept.corr().drop('target')[np.abs(kept.corr()['target'].drop('target')) > .05].index)
```

```
In [56]:   1  regress = kept[final_vars]
```

```
In [57]:  1  regress['target'] = (df['target'] == 'Donor').astype(int)
```

```
In [72]:  1  kept = kept.drop('log_target_No Donor', axis =1)
```

```
In [59]:  1  train = kept[regress['type_train'] ==1]
          2  dev = kept[(regress['type_test'] == 0) & (regress['type_train'] == 0)]
          3  test = kept[regress['type_test'] ==1]
```

```
In [60]:  1  for data in [train, dev, test]:
          2      data.drop('type_train', inplace = True, axis = 1)
          3      data.drop('type_test', inplace = True, axis = 1)
          4  test = test.drop('target', axis =1)
```

## Logistic Regression

```
In [61]:  1  results_df = pd.DataFrame()
```

```
In [62]:  1  # Selecting features and target variable for training data
          2  X_train = train.drop(['target'], axis =1 )
          3  y_train = train['target']
          4  X_test = dev.drop(['target'], axis = 1)
          5  y_test = dev['target']
          6
          7  # Fitting logistic regression model
          8  glm = sm.GLM(y_train, X_train, family=sm.families.Binomial())
          9  glm = glm.fit()
         10
         11  # Summarizing results
         12  # print(results.summary())
```

```
In [63]:  1  log_preds = (glm.predict(X_test) >= 0.5).astype(int)
          2  log_acc = accuracy_score(log_preds, y_test)
          3  print(log_acc)
          4
          5  d = confusion_table(log_preds,y_test)
          6  results_df = pd.concat([results_df,convert_confusion_matrix(d, 'Logistic Regression')])
          7
          8  log_test_preds = (glm.predict(test) >= 0.5).astype(int)
          9  log_test_preds = format_results(log_test_preds)
         10
         11  save_df = pd.DataFrame(log_test_preds, columns=['values'])
         12  save_df.to_csv('./preds/log.csv', index=False)
```

```
0.5116666666666667
```

### LDA

```
In [64]:  1  lda = LDA(store_covariance=True)
          2  lda.fit(X_train, y_train)
          3
          4  lda_preds = lda.predict(X_test)
          5
          6  lda_acc = accuracy_score(lda_preds,y_test)
          7  print(lda_acc)
          8
          9  d = confusion_table(lda_preds,y_test)
         10  results_df = pd.concat([results_df,convert_confusion_matrix(d, 'LDA')])
         11
         12
         13  lda_test_preds = (lda.predict(test) >= 0.5).astype(int)
         14  lda_test_preds = format_results(lda_test_preds)
         15
         16  save_df = pd.DataFrame(lda_test_preds, columns=['values'])
         17  save_df.to_csv('./preds/lda.csv', index=False)
```

```
0.5066666666666667
```

## QDA.

```
In [65]:   1  qda = QDA(store_covariance=True)
           2  qda.fit(X_train, y_train)
           3
           4  qda_preds = qda.predict(X_test)
           5
           6  qda_acc = accuracy_score(qda_preds,y_test)
           7
           8  print(qda_acc)
           9
          10  d = confusion_table(qda_preds,y_test)
          11  results_df = pd.concat([results_df,convert_confusion_matrix(d, 'QDA')])
          12
          13  qda_test_preds = (qda.predict(test) >= 0.5).astype(int)
          14  qda_test_preds = format_results(qda_test_preds)
          15
          16  save_df = pd.DataFrame(qda_test_preds, columns=['values'])
          17  save_df.to_csv('./preds/qda.csv', index=False)
```

0.49166666666666664

## KNN

```
In [66]:   1  knn1 = KNeighborsClassifier(n_neighbors=1)
           2  knn1.fit(X_train, y_train)
           3  knn1_pred = knn1.predict(X_test)
           4  knn1_acc = accuracy_score(knn1_pred,y_test)
           5
           6  print(knn1_acc)
           7
           8  d = confusion_table(knn1_pred, y_test)
           9  results_df = pd.concat([results_df,convert_confusion_matrix(d, 'KNN')])
          10
          11  knn1_test_preds = (knn1.predict(test) >= 0.5).astype(int)
          12  knn1_test_preds = format_results(knn1_test_preds)
          13
          14  save_df = pd.DataFrame(knn1_test_preds, columns=['values'])
          15  save_df.to_csv('./preds/knn1.csv', index=False)
```

0.4816666666666667

## NB

```
In [67]:   1  nb = GaussianNB()
           2  nb.fit(X_train, y_train)
           3  nb_preds = nb.predict(X_test)
           4  nb_acc = accuracy_score(nb_preds,y_test)
           5
           6  print(nb_acc)
           7  save_df = pd.DataFrame(nb_preds, columns=['values'])
           8  save_df.to_csv('./preds/nb.csv', index=False)
           9
          10  d = confusion_table(nb_preds, y_test)
          11  results_df = pd.concat([results_df,convert_confusion_matrix(d, 'Naïve Bayes')])
          12
          13  nb_test_preds = (nb.predict(test) >= 0.5).astype(int)
          14  nb_test_preds = format_results(nb_test_preds)
          15
          16  save_df = pd.DataFrame(nb_test_preds, columns=['values'])
          17  save_df.to_csv('./preds/nb.csv', index=False)
```

0.495

```
In [68]:   1  results_df
```

Out[68]:

|   | name | tp | tn | fp | fn | acc | prec | recall | f1 |
|---|------|-----|-----|-----|-----|----------|----------|----------|----------|
| 0 | Logistic Regression | 159 | 148 | 153 | 140 | 0.511667 | 0.509615 | 0.531773 | 0.520458 |
| 0 | LDA | 156 | 148 | 156 | 140 | 0.506667 | 0.500000 | 0.527027 | 0.513158 |
| 0 | QDA | 26 | 269 | 286 | 19 | 0.491667 | 0.083333 | 0.577778 | 0.145658 |
| 0 | KNN | 140 | 149 | 172 | 139 | 0.481667 | 0.448718 | 0.501792 | 0.473773 |
| 0 | Naïve Bayes | 18 | 279 | 294 | 9 | 0.495000 | 0.057692 | 0.666667 | 0.106195 |

```
In [69]:  1  # test_acc = {
          2  #     'log': 0.5333333,
          3  #     'lda': 0.5583333,
          4  #     'qda': 0.525,
          5  #     'knn':  0.475,
          6  #     'nb': 0.5166667,
          7  # }
```

```
In [70]:  1  # results_df['test_acc'] = test_acc.values()
```

```
In [71]:  1  results_df
```

Out[71]:

| | name | tp | tn | fp | fn | acc | prec | recall | f1 | test_acc |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | Logistic Regression | 159 | 148 | 153 | 140 | 0.511667 | 0.509615 | 0.531773 | 0.520458 | 0.533333 |
| 0 | LDA | 156 | 148 | 156 | 140 | 0.506667 | 0.500000 | 0.527027 | 0.513158 | 0.558333 |
| 0 | QDA | 26 | 269 | 286 | 19 | 0.491667 | 0.083333 | 0.577778 | 0.145658 | 0.525000 |
| 0 | KNN | 140 | 149 | 172 | 139 | 0.481667 | 0.448718 | 0.501792 | 0.473773 | 0.475000 |
| 0 | Naïve Bayes | 18 | 279 | 294 | 9 | 0.495000 | 0.057692 | 0.666667 | 0.106195 | 0.516667 |

```python
 1  # ============================================================================
 2  # SET-UP AND IMPORTS
 3  # ============================================================================
 4  import os
 5  import torch
 6  import torch.nn as nn
 7  import numpy as np
 8  import pandas as pd
 9  import matplotlib.pyplot as plt
10  from torch.utils.data import DataLoader, TensorDataset, ConcatDataset
11  from sklearn.model_selection import KFold
12  from sklearn.preprocessing import StandardScaler
13
14  # Set directory and device setup
15  os.chdir('/home/dan/FUNDRAISING')
16  device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
17  print(f"Using device: {device}")
18
19  # ============================================================================
20  # PREPARE DATA
21  # ============================================================================
22  df = pd.read_csv('./data/df.csv', index_col='Unnamed: 0')
23  df['zipconvert'] = df[['zipconvert2', 'zipconvert3', 'zipconvert4', 'zipconvert5']].apply(
24      lambda row: 'zc' + str(row.idxmax()[-1]) if pd.notna(row.idxmax()) else 'unknown', axis=1
25  )
26  test = df[df['type']=='test']
27
28  df = df[df['type']!='test']
29
30  df.drop(['zipconvert2', 'zipconvert3', 'zipconvert4', 'zipconvert5', 'type'], axis=1, inplace=True)
31  test.drop(['zipconvert2', 'zipconvert3', 'zipconvert4', 'zipconvert5', 'type'], axis=1, inplace=True)
32
33  cat_cols = ['homeowner', 'female', 'zipconvert', 'wealth', 'income', 'num_child']
34  cont_cols = [col for col in df.columns if col not in cat_cols + ['target']]
35  emb_sizes = [(df[col].astype('category').cat.codes.max() + 1, min(50, (df[col].nunique() + 1) // 2)) for col in cat_cols]
36
37  # Function to process datasets
38  def process_data(data):
39      cats = np.stack([data[col].astype('category').cat.codes.values for col in cat_cols], axis=1)
40      conts = np.stack([data[col].values for col in cont_cols], axis=1)
41      scaler = StandardScaler()
42
43      conts = scaler.fit_transform(conts)
44      y = data['target'].map({'Donor': 1, 'No Donor': 0}).values
45      return torch.tensor(cats, dtype=torch.int64), torch.tensor(conts, dtype=torch.float), torch.tensor(y, dtype=torch.long)
46
47
48  cats, conts, targets = process_data(df)
49  test_cats, test_conts, test_targets = process_data(test)
50
51  dataset = TensorDataset(cats, conts, targets)
52
53  test_dataset = TensorDataset(test_cats, test_conts, test_targets)
54  test_loader = DataLoader(test_dataset, batch_size=2048, shuffle=False)
55
56  # ============================================================================
57  # DEFINE A TABULAR MODEL
58  # ============================================================================
59  class TabularModel(nn.Module):
60      def __init__(self, emb_sizes, n_cont, out_sz, layers, p=0.5):
61          super().__init__()
62          self.embeds = nn.ModuleList([nn.Embedding(ni, nf) for ni, nf in emb_sizes])
63          self.emb_drop = nn.Dropout(p)
64          self.bn_cont = nn.BatchNorm1d(n_cont)
65
66          # Calculate the total embedding output size
67          total_emb_size = sum(nf for _, nf in emb_sizes)
68          n_in = total_emb_size + n_cont  # Total input size for the first linear layer
69
70          layerlist = nn.ModuleList()
71          for output_size in layers:
72              layerlist.append(nn.Linear(n_in, output_size))
73              layerlist.append(nn.ReLU())
74              layerlist.append(nn.BatchNorm1d(output_size))
75              layerlist.append(nn.Dropout(p))
76              n_in = output_size  # Update n_in to the output size of the current layer
77
78          layerlist.append(nn.Linear(layers[-1], out_sz))  # Final output layer
79          self.layers = nn.Sequential(*layerlist)
80
81      def forward(self, x_cat, x_cont):
82          embeddings = [e(x_cat[:, i]) for i, e in enumerate(self.embeds)]
83          x = torch.cat(embeddings, 1)
84          x = self.emb_drop(x)
```

```python
 85            x_cont = self.bn_cont(x_cont)
 86            x = torch.cat([x, x_cont], 1)
 87            return self.layers(x)
 88
 89  # Initialize model
 90  model = TabularModel(emb_sizes, len(cont_cols), 2, [200, 100], p=0.4)
 91  model = model.to(device)
 92
 93
 94  #endregion
 95  #region # EARLY STOPPING
 96  # ============================================================================
 97  # EARLY STOPPING
 98  # ============================================================================
 99  class EarlyStopping:
100      def __init__(self, patience=5, verbose=False, delta=0):
101          self.patience = patience
102          self.verbose = verbose
103          self.delta = delta
104          self.best_score = None
105          self.early_stop = False
106          self.counter = 0
107
108      def __call__(self, val_loss, model):
109          score = -val_loss
110
111          if self.best_score is None:
112              self.best_score = score
113          elif score < self.best_score + self.delta:
114              self.counter += 1
115              if self.verbose:
116                  print(f'EarlyStopping counter: {self.counter} out of {self.patience}')
117              if self.counter >= self.patience:
118                  self.early_stop = True
119          else:
120              self.best_score = score
121              self.counter = 0
122
123
124  # ============================================================================
125  # CROSS-VALIDATION SETUP
126  # ============================================================================
127  def calculate_accuracy(y_pred, y_true):
128      y_pred_classes = torch.argmax(y_pred, dim=1)
129      correct = (y_pred_classes == y_true).float()  # convert into float for division
130      acc = correct.sum() / len(correct)
131      return acc
132
133  kf = KFold(n_splits=20, shuffle=True, random_state=42)
134  early_stopping = EarlyStopping(patience=5, verbose=True)
135  results = []
136
137  for fold, (train_idx, val_idx) in enumerate(kf.split(dataset)):
138      train_subsampler = torch.utils.data.SubsetRandomSampler(train_idx)
139      val_subsampler = torch.utils.data.SubsetRandomSampler(val_idx)
140      train_loader = DataLoader(dataset, batch_size=32, sampler=train_subsampler)
141      val_loader = DataLoader(dataset, batch_size=32, sampler=val_subsampler)
142
143      model = TabularModel(emb_sizes, len(cont_cols), 2, [100,200], p=0.4).to(device)
144      criterion = nn.CrossEntropyLoss()
145      optimizer = torch.optim.Adam(model.parameters(), lr=1e-5)
146
147      for epoch in range(50):  # Adjust as needed
148          model.train()
149          total_loss, total_acc = 0, 0
150          for cats, conts, y in train_loader:
151              cats, conts, y = cats.to(device), conts.to(device), y.to(device)
152              optimizer.zero_grad()
153              outputs = model(cats, conts)
154              loss = criterion(outputs, y)
155              acc = calculate_accuracy(outputs, y)
156              total_loss += loss.item()
157              total_acc += acc.item()
158              loss.backward()
159              optimizer.step()
160          avg_train_loss = total_loss / len(train_loader)
161          avg_train_acc = total_acc / len(train_loader)
162
163          model.eval()
164          val_loss, val_acc = 0, 0
165          with torch.no_grad():
166              for cats, conts, y in val_loader:
167                  cats, conts, y = cats.to(device), conts.to(device), y.to(device)
168                  outputs = model(cats, conts)
169                  loss = criterion(outputs, y)
170                  acc = calculate_accuracy(outputs, y)
171                  val_loss += loss.item()
```

```python
                    val_acc += acc.item()
            avg_val_loss = val_loss / len(val_loader)
            avg_val_acc = val_acc / len(val_loader)

            print(f'Fold {fold+1}, Epoch {epoch+1}, Train Loss: {avg_train_loss:.4f}, Train Acc: {avg_train_acc:.4f}, Val Loss: {avg_val_loss:.4f}, Val Acc: {avg_val_acc:.4f}')

            # Call early stopping
            early_stopping(avg_val_loss, model)
            if early_stopping.early_stop:
                print("Early stopping")
                break

    results.append((avg_val_loss, avg_val_acc))

average_val_loss = sum(x[0] for x in results) / len(results)
average_val_acc = sum(x[1] for x in results) / len(results)
print(f'Average Validation Loss: {average_val_loss:.4f}, Average Validation Accuracy: {average_val_acc:.4f}')

preds = []
model.eval()

with torch.no_grad():
    for cats, conts, y in test_loader:
        cats, conts, y = cats.to(device), conts.to(device), y.to(device)
        output = model(cats, conts)
        predicted = output.argmax(dim=1)  # Ensure you use dim=1
        # print(predicted)
        preds.append(predicted.cpu())  # Move predictions to CPU
# print(type(preds[0]))
# Concatenate all batch predictions into a single tensor
# preds = torch.cat(preds)
preds = preds[0].tolist()
# print(preds)
final_preds = []
donor=1
no_donor=1
for i in preds:
    if i == 1:
        final_preds.append('Donor')
        donor +=1
    else:
        final_preds.append('No Donor')
        no_donor+=1
print(f'% Donor = {donor / (donor+no_donor)}')


save_df = pd.DataFrame(final_preds, columns=['values'])
save_df.to_csv('./preds/preds.csv', index=False)
```