



# Ruffalo Noel Levitz Project



Team 1: Jisoo Lee, Jordan Kloewer,  
Ran An, Xuan Zhang



# Today's Agenda

---

**Introduction**



**Data Gathering  
and  
Preparation**



**Modeling**



**Conclusion**

# Introduction

1. Assumptions
2. Challenges
3. Original Dataset

**RUFFALO**<sup>SM</sup>  
**NOEL LEVITZ**

# Assumptions

---



No external factors



Correct data



Use original dataset



Most important target  
variable

# Challenges

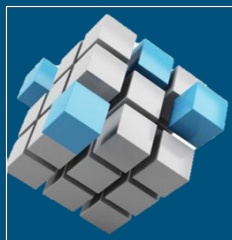


## Other Challenges

- Group members in different sections

## Data Challenges

- Ignoring columns
- Missing values
- Incomplete rows



## Modeling Challenges

- Too much time prepping data
- Unknown modeling techniques

# Original Dataset

Name: Pledge

Number/Letter: 39, AM

Description: Pledge flag

Format of Column: Single digit (1 or blank)

Column Type in R: Integer

Number of Missing/NAs in R: 244787

Questions About Column: What data is currently used that would cause someone to have a "pledge flag?"

Any person with data populated in PledgeAmt (Column BI) and PledgeTot (Column BJ) (excluding where the pledge amounts/totals = 0) should be flagged as pledges. These fields are also numeric versions of the FY17\_PLEDGE\_AMOUNT column which was character in the original data for some reason

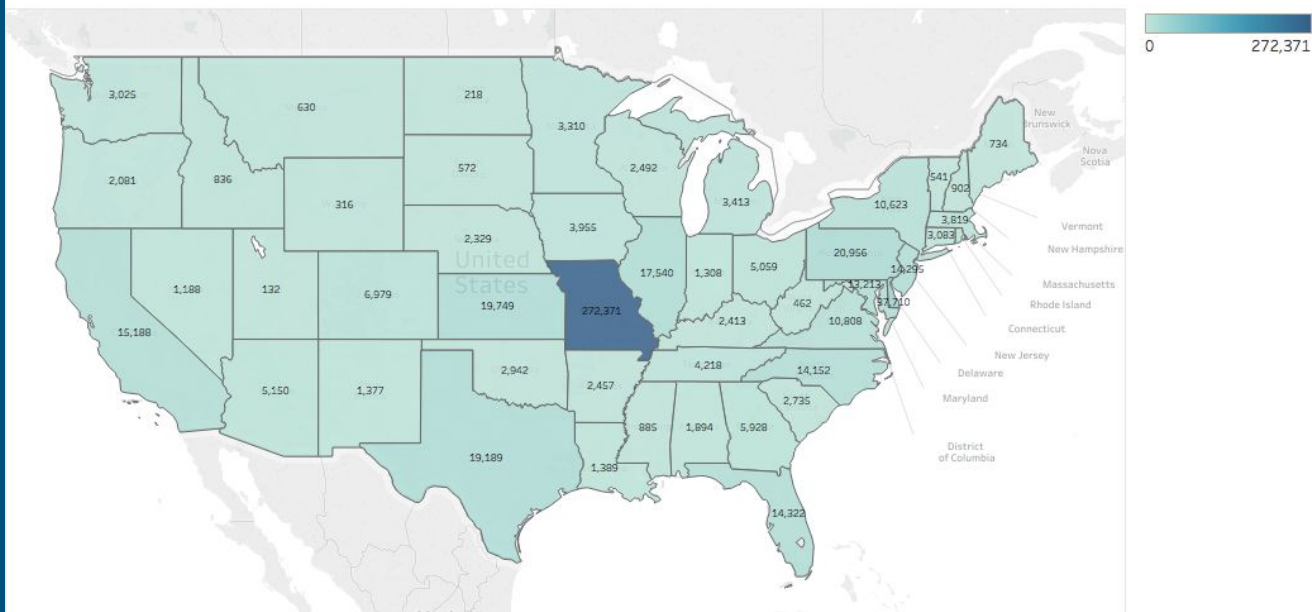
Any Similar Columns/Do Number Match Up?: Pldg

**ACTION: DELETE** - the column Pldg is the exact same but that one is formatted as a binary variable (0,1) with no blank lines

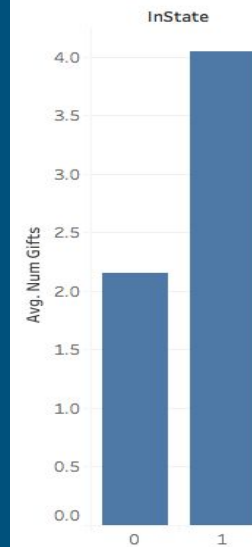
- 15,000,000+ cells
- Data Dictionary
  - Column descriptions
  - Number/text format
  - Column type in R
  - Missing values
  - Related columns
  - Questions
- Data Visualizations

# Original Dataset

Number of Gifts by State



Effect of Being In-State



Avg. Num Gifts for each InState. Color shows details about Avg. Num Gifts.

# Data Gathering and Preparation

1. Data Cleaning
2. Data Selection
3. Dataset Used For Modeling

**RUFFALO**<sup>SM</sup>  
**NOEL LEVITZ**



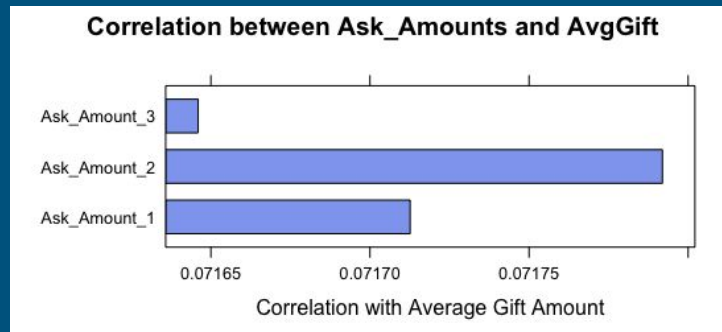
# Data Cleaning

---

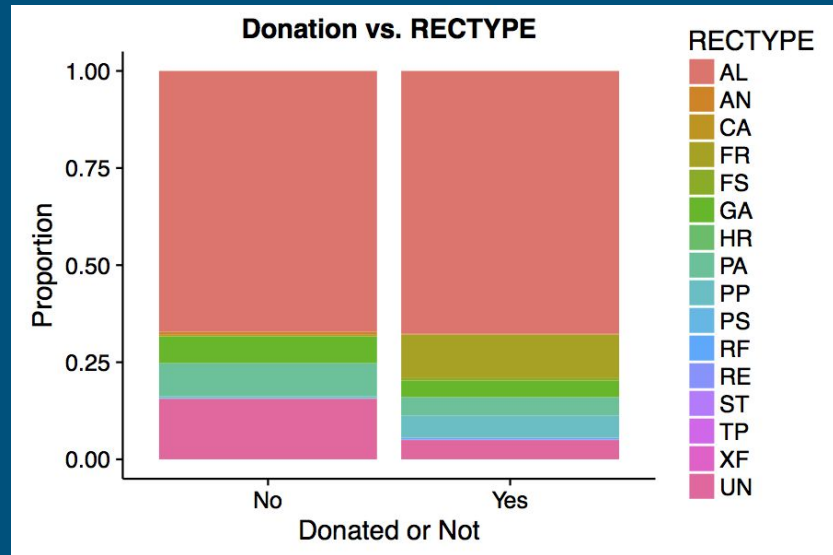
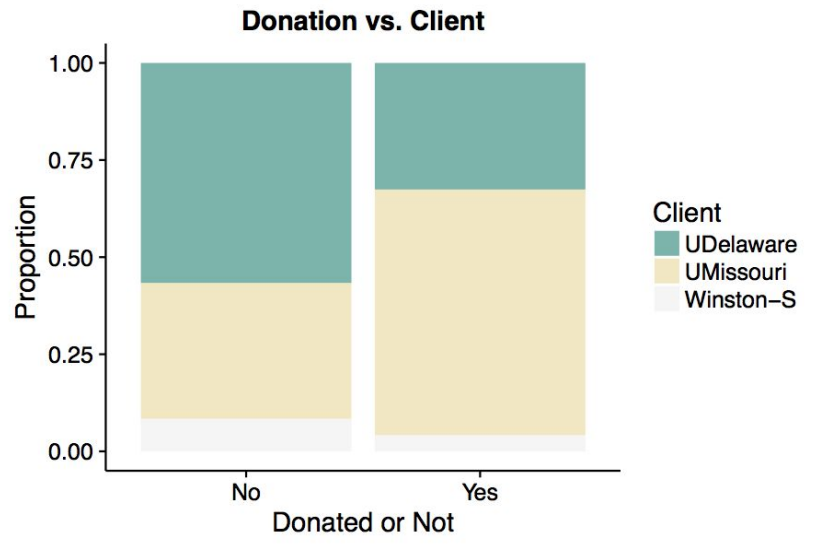
- After examination,
  - Removed errors and outliers: negative talk time and age below 3
  - Transformed some columns to improve interpretability
    - STATE: change to CRegion to reduce the number of levels
    - RecentGradYr, LastGradYr: change to YearSinceRecentGrad to have numeric values
  - Deleted 32 irrelevant and duplicate columns such as PHONE\_RESEARCH, PAYMENT\_TYPE
- Changed missing values to 0, except of AGE
  - Normally missing values of AGE are replaced to mean/median of AGE, but to avoid misleading impact of 53,161 NAs of AGE on models, separated dataset into:
    - 1) Dataset without NAs of AGE
    - 2) Dataset with replaced NAs of AGE to median

# Data Selection

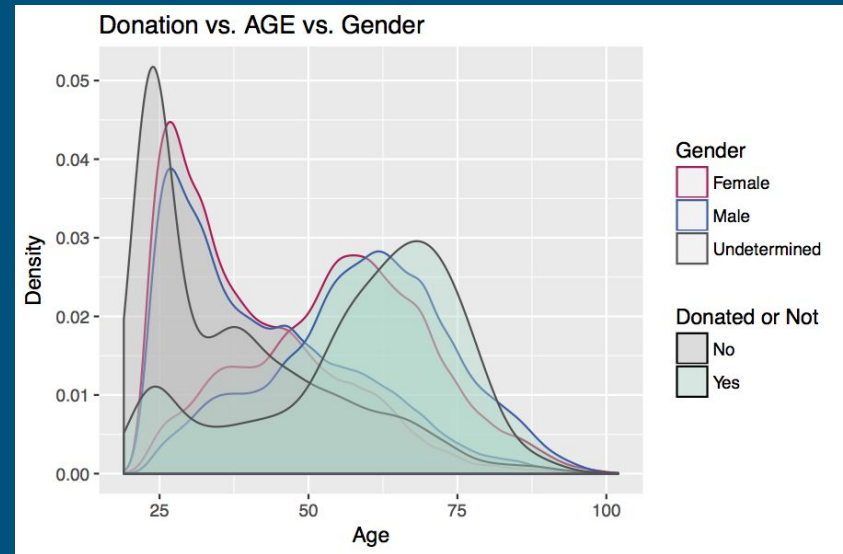
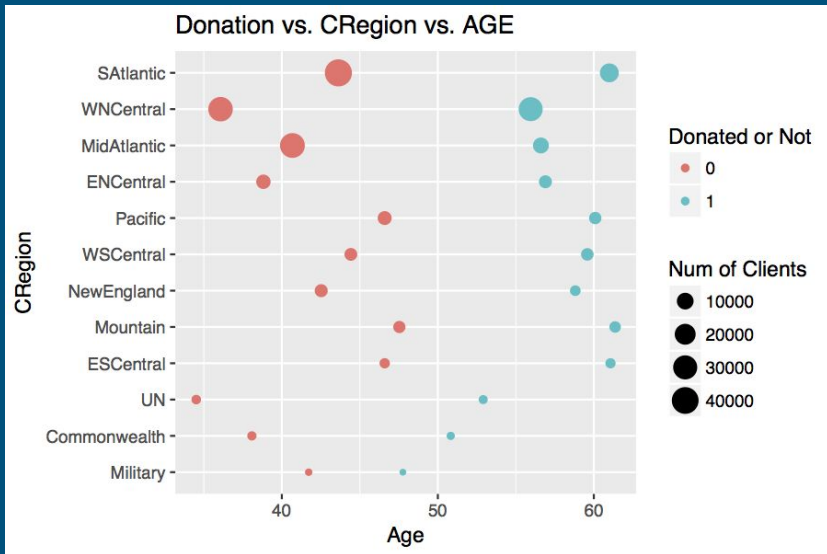
- For models to predict whether donate or not,
  - Transformed NumGifts into Donor01 to have binary values
  - Further removed 11 columns like FSTGFTA, FSTGFTD, LSTGFTD
  - Extracted 4 highly correlated columns:
    - AGE - YearSinceGrad & YearSinceRecentGrad
    - Ask\_Amount\_2 - Ask\_Amount\_1 & Ask\_Amount\_3



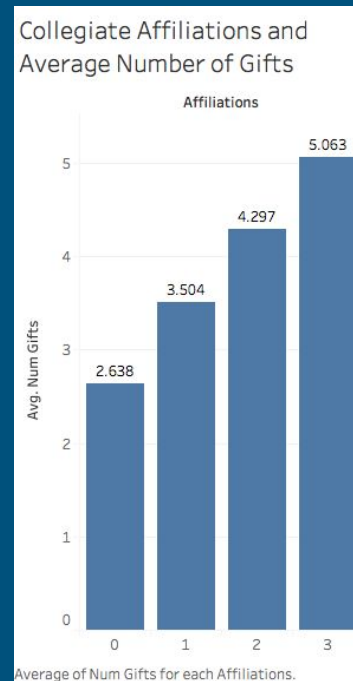
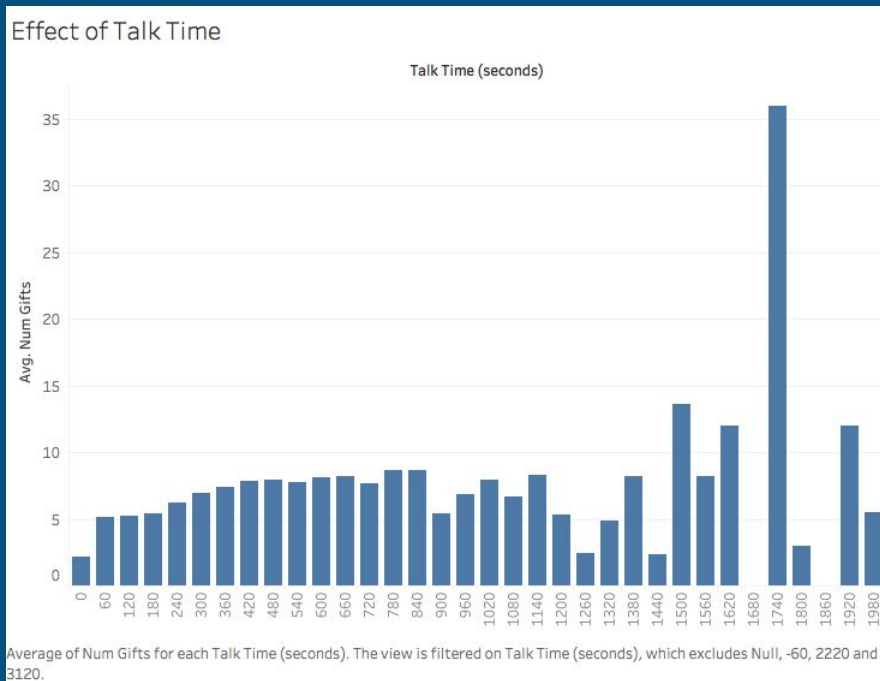
# Data Selection



# Data Selection



# Data Selection



# Data Selection

---

- Left 14 columns with 252,725 rows
  - Predictors: SCHOOL, AGE, RECTYPE, GENDER, NumDegrees, SuccCont, AffilCount, client, Ask\_Amount\_2, InState, TotAttempts, TALK\_TIME, and CRegion
  - Target: Donor01

# Dataset Used for Modeling

---

- Had multiple datasets that we could use:
  - Depending on AGE:
    - trainset\_rmAGE - Dataset without NAs of AGE
    - trainset\_medAGE - Dataset with replaced NAs of AGE to median of AGE (45)
  - Depending on client, (no Winston-S due to too many incomplete rows):
    - trainset\_UD - Dataset for UDelaware
    - trainset\_UM - Dataset for UMissouri
  - Depending on balancing methods that we used\*:
    - trainset\_under - Dataset with undersampling target variable
    - trainset\_both - Dataset with over and undersampling target variable

\*undersampling: reduce the number of majority class in target variable

\*oversampling: increase the number of minority class in target variable

# Modeling

1. Hypotheses
2. Modeling Techniques
3. Results
4. Evaluation Criteria

**RUFFALO**<sup>SM</sup>  
**NOEL LEVITZ**



# Hyphotheses

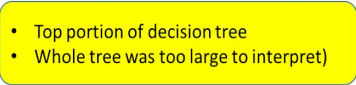
---

- Based on the findings:
  - Chance to donate
    - People living in certain areas or in-state tend to donate more than those living in other areas or out-of-state
    - Older individuals are more likely to donate than younger individuals
    - People with more call attempts tend to donate more than those with less call attempts
    - People with longer talk time have a higher chance of donating
    - Males have a higher chance to donate

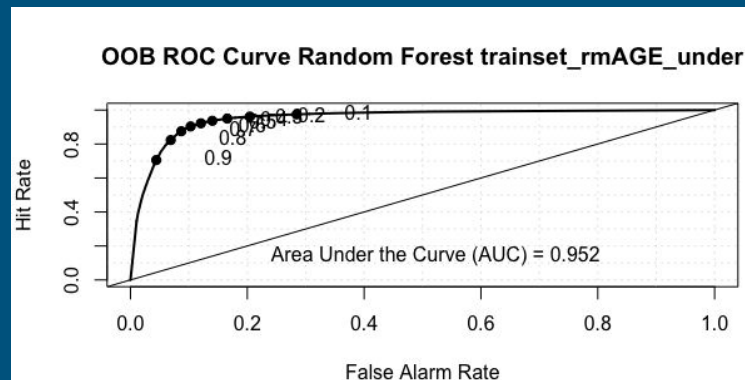
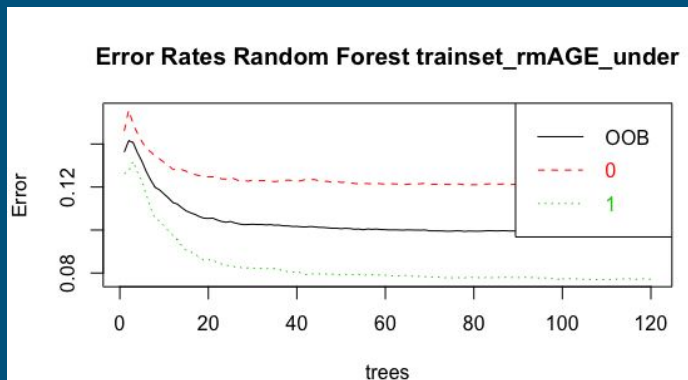
# Modeling Techniques Explored

---

- Decision Tree:
  - Used R-packages to help pick optimum parameters for min split, min bucket, max depth and complexity values
- Random Forest: advanced decision tree
  - Number of trees: the larger the better, but the longer to compute, results stops getting significantly better beyond a critical number
  - Number of variables: the lower the greater reduction of variance, but also the greater increase in bias
- Support Vector Machine (SVM):
  - Kernel: Radial Basis (rbfdot)
  - Options: The parameter has been set from 0-1.
- Used Rattle and R-packages to build three models



# Results -Random Forest

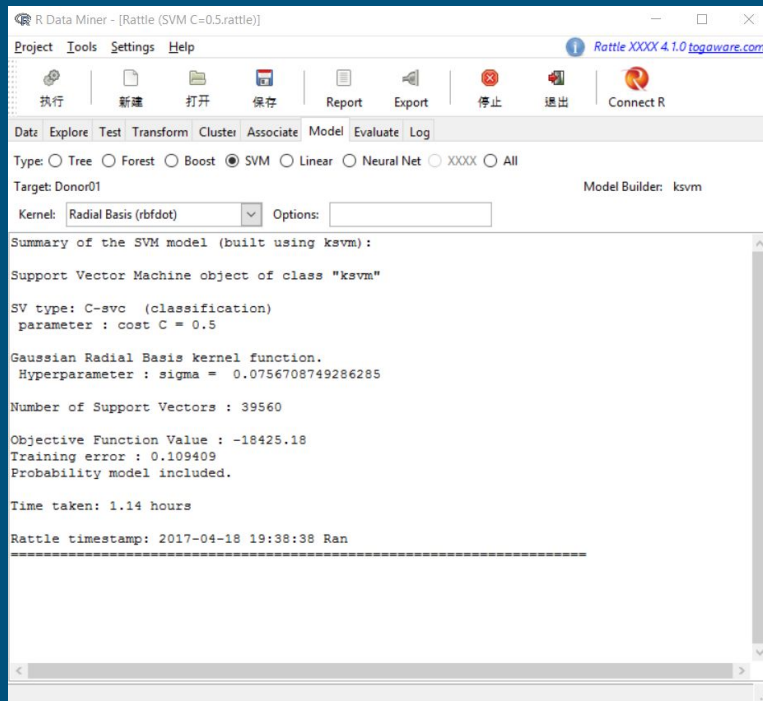
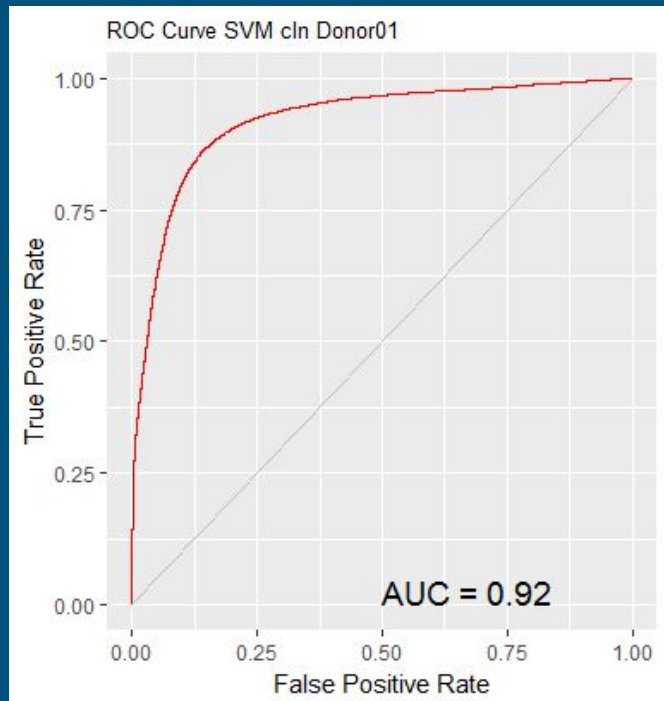


-----  
Tree 1 Rule 34 Node 2796 Decision 1

```
1: TALK TIME <= 41.5
2: CRegion IN ("Military", "Pacific", "ESCentral", "WSCentral", "Mountain", "WNCentral", "ENCentral")
3: SCHOOL IN ("Agriculture and Natural Resources", "Applied Science", "Continuing Education", "Educa
4: client IN ("UDelaware")
5: AffilCount <= 0.5
6: NumDegrees > 0.5
7: SuccCont IN ("0")
8: TotAttempts > 2.5
9: Ask_Amount_2 > 67.5
10: TotAttempts <= 7
11: TotAttempts <= 4.5
12: GENDER IN ("M", "U")
13: CRegion IN ("Military", "Pacific", "ESCentral", "Mountain", "NewEngland", "Commonwealth", "SAtla
```

-----

# Results -SVM



# Evaluation Criteria - Overall Dataset

Dataset	rmAGE	medAGE	rmAGE_under	rmAGE_both
Tree	15, 5, 15, 0.0001	21, 7, 10, 0.0001	21, 7, 10, 0.0001	21, 7, 10, 0.0001
FPR	0.0816	0.0707	0.1077	0.1071
FNR	0.1311	0.1639	0.1089	0.1061
F	0.8695	0.8457	0.8464	0.8462
Forest	120, 3	120, 3	30, 3	30, 3
FPR	0.0827	0.0693	0.1190	0.1118
FNR	0.1239	0.1402	0.0812	0.0882
F	0.8623	0.8605	0.8575	0.8594
SVM	-	0.5	-	-
FPR	-	0.0984	FPR: false positive rate (predicted 0 was really 1) FNR: false negative rate (predicted 1 was really 0) F score: overall accuracy for binary target	
FNR	-	0.2034		
F	-	0.8444		

# Evaluation Criteria - Client Specific

Dataset	UDelaware (UD)	UMissouri (UM)	UM_under	UM_both
Tree	15, 5, 5, 0.0001	15, 5, 15, 0.0001	90, 30, 10, 0.0001	90, 30, 5, 0.0001
FPR	0.0110	0.5014	0.4866	0.7421
FNR	0.1556	0.1420	0.8543	0.6128
F	0.8943	0.5674	0.1279	0.2487
Forest	60, 3	90, 3	30, 3	30, 3
FPR	0.0115	0.2879	0.8137	0.8239
FNR	0.0969	0.1133	0.8401	0.8127
F	0.9268	0.7045	0.1251	0.1251
SVM	0.5	0.5	-	-
FPR	0.0186	0.2199	FPR: false positive rate (predicted 0 was really 1) FNR: false negative rate (predicted 1 was really 0) F score: overall accuracy for binary target	
FNR	0.2232	0.1447		
F	0.9628	0.5745		

# Conclusion

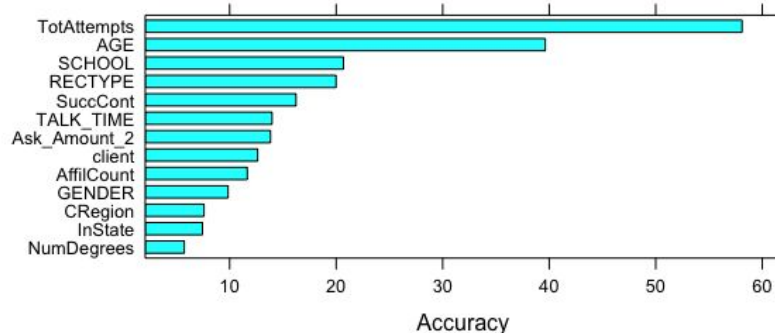
1. Final Best Model
2. Insights
3. Recommendations
4. Final Takeaways

**RUFFALO**<sup>SM</sup>  
**NOEL LEVITZ**

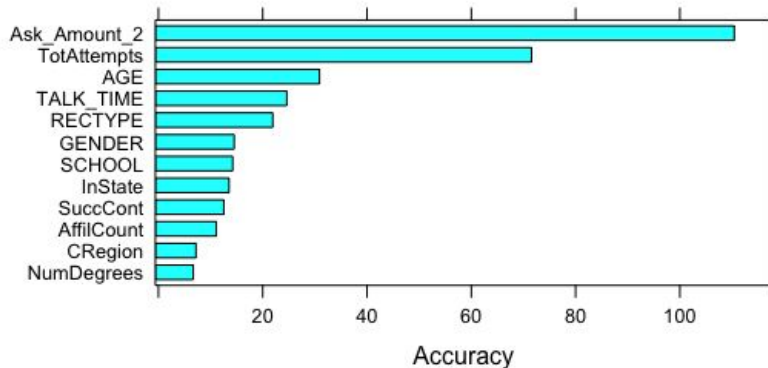


# Final Best Model

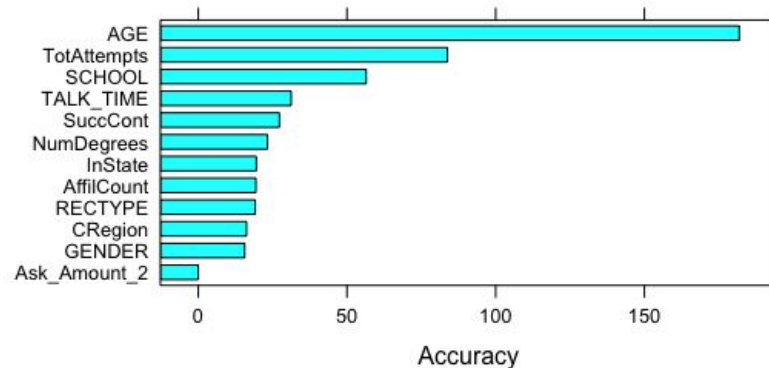
Variable Importance Random Forest trainset\_rmAGE\_under



Variable Importance Random Forest trainset\_UD



Variable Importance Random Forest trainset\_UM



# Insights

---

## Important Factors

- Age
- Total Attempts
- Talk Time

## Insignificant Factors

- Number of degrees
- Client's location
- Gender

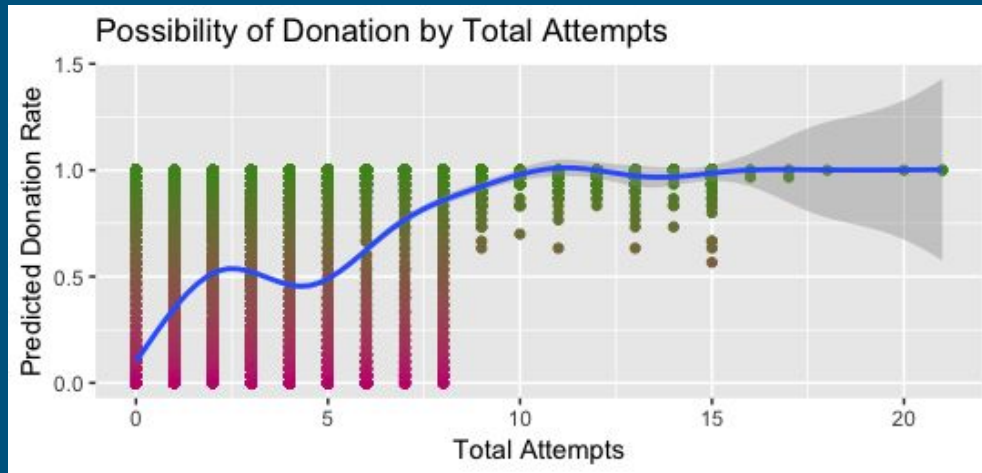
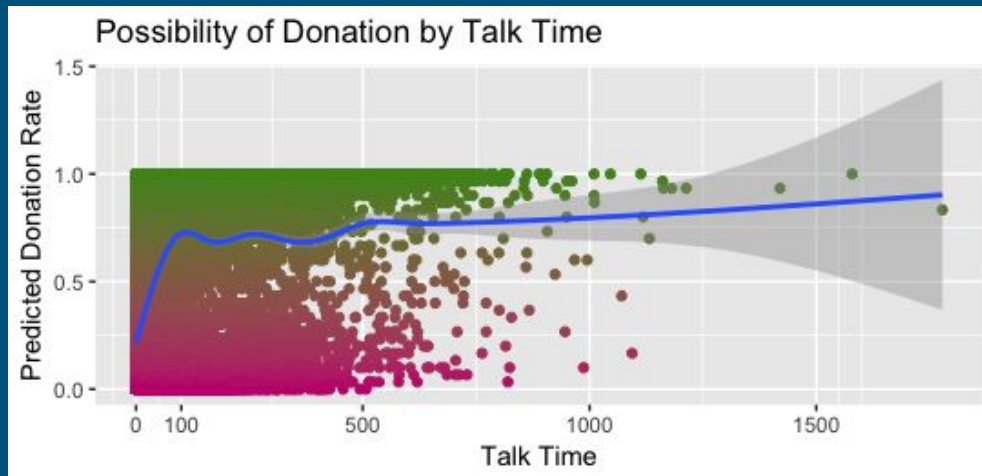
## Strategic Recommendations

- Target people in the 55-75 age group

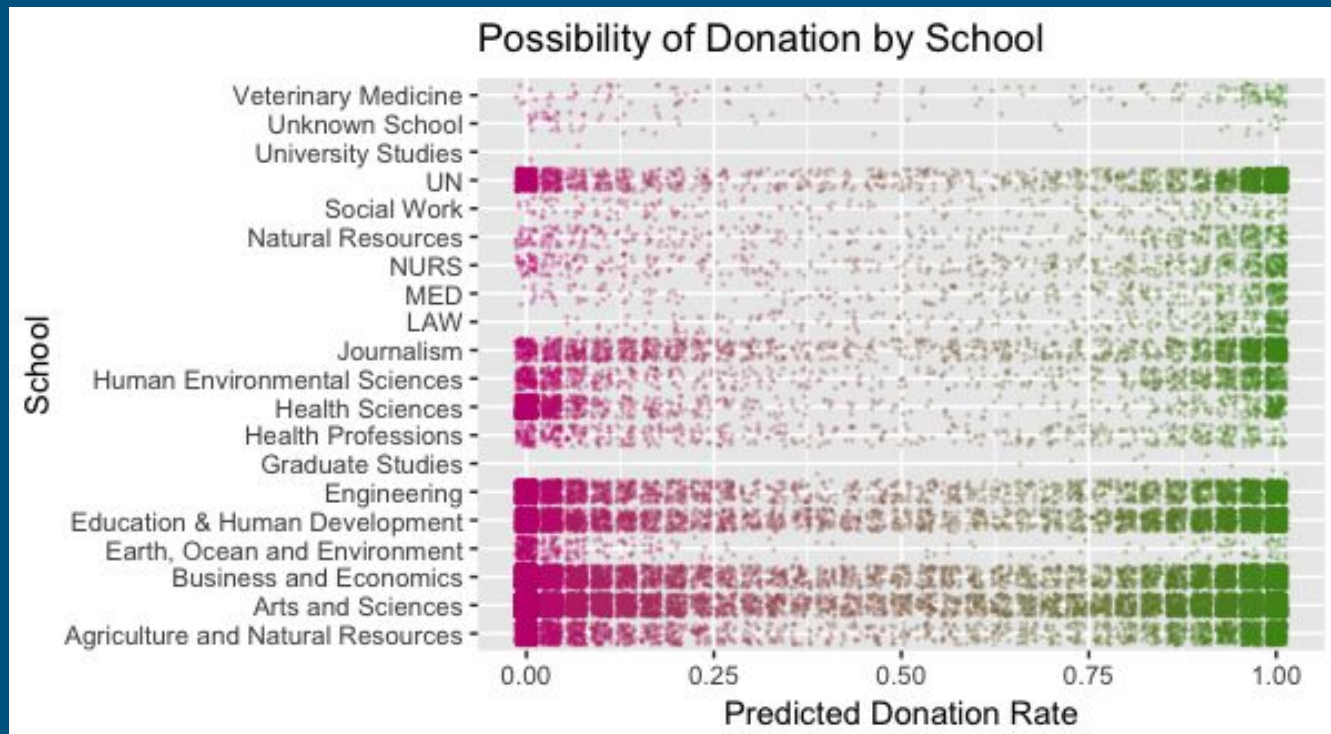
## Tactical Recommendations

- Delaware clients: set optimal ask amount
- Missouri clients: focus on older people

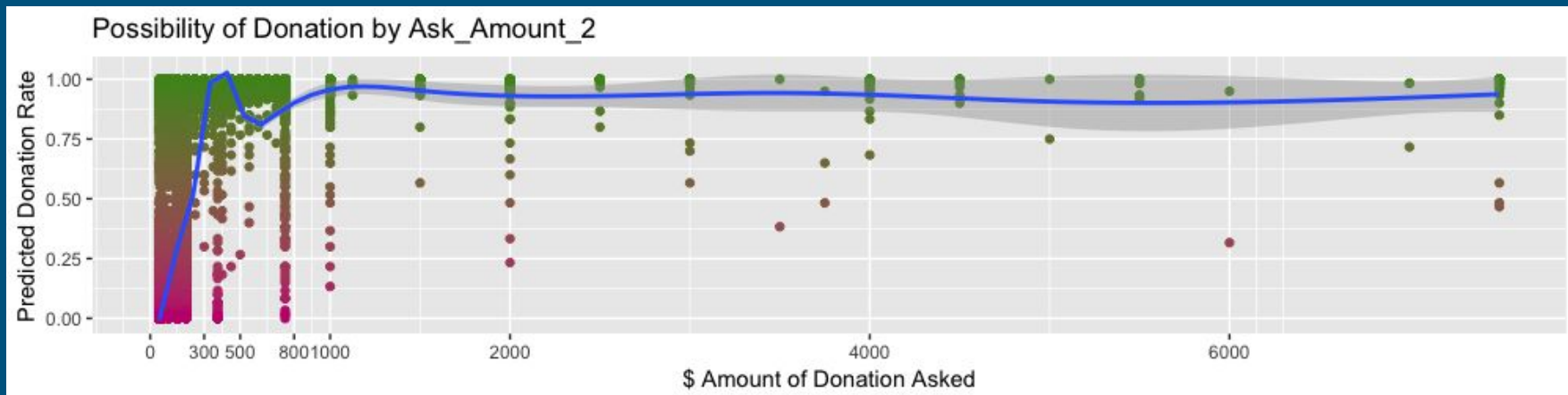
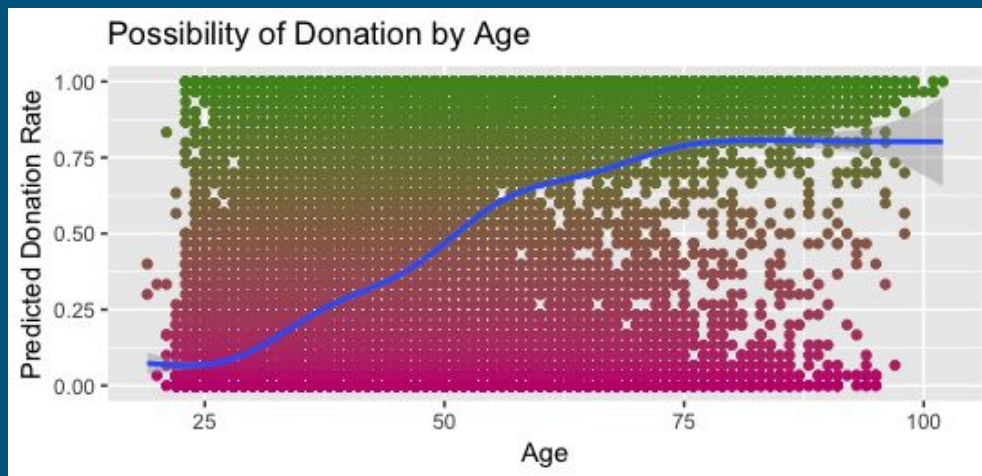
# Insights



# Insights



# Insights



# Final Takeaways

THANK YOU!