

Bank Marketing Project

Yufei Tang

Lei Wang

Hua Lu



Agenda



Project Introduction



Data Information



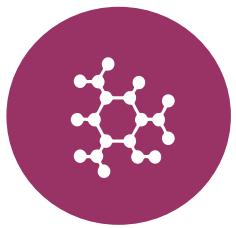
Data Preprocessing



Data Model 1 –
Logistical Regression



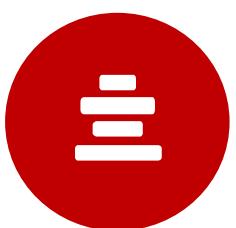
Data Model 2 –
Decision Tree



Data Model 3 –
Neural Network



Other Analyses

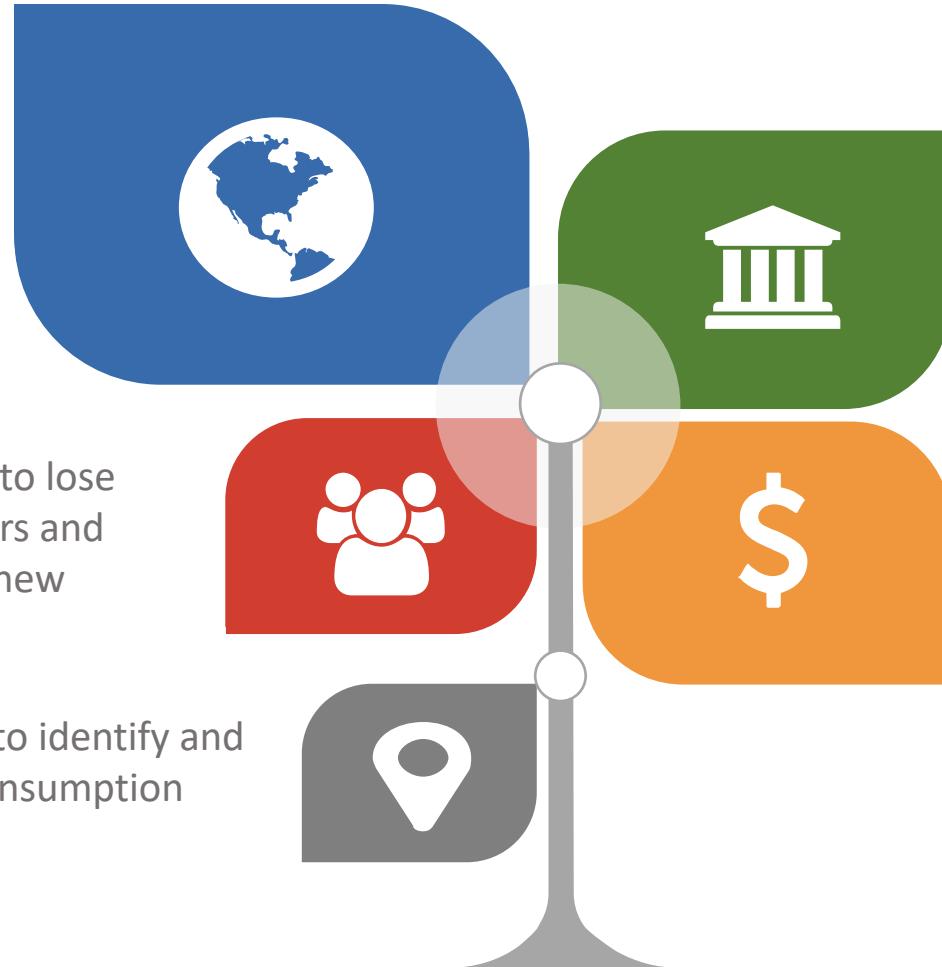


Summary

Project Introduction

Project Background

The whole national economy in the country is stagnant for a long time



Increasing competition among banks in the banking market

The Bank continues to lose its current customers and hardly to attract new customers

The bank has challenges to identify and predict customer's consumption behaviours

Customers have various investment opinions in the current market

Project Introduction

Project Goal & Benefits



Project Goal

- To distinguish potential targeted customers for our bank
- To analyse influences of each attribute on customer's final decision

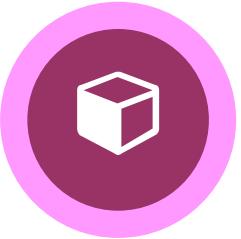
Project Benefits



To predict whether a customer will subscribe a term deposit in our bank or not



To contact with most potential consumers effectively



To control and reduce the marketing cost

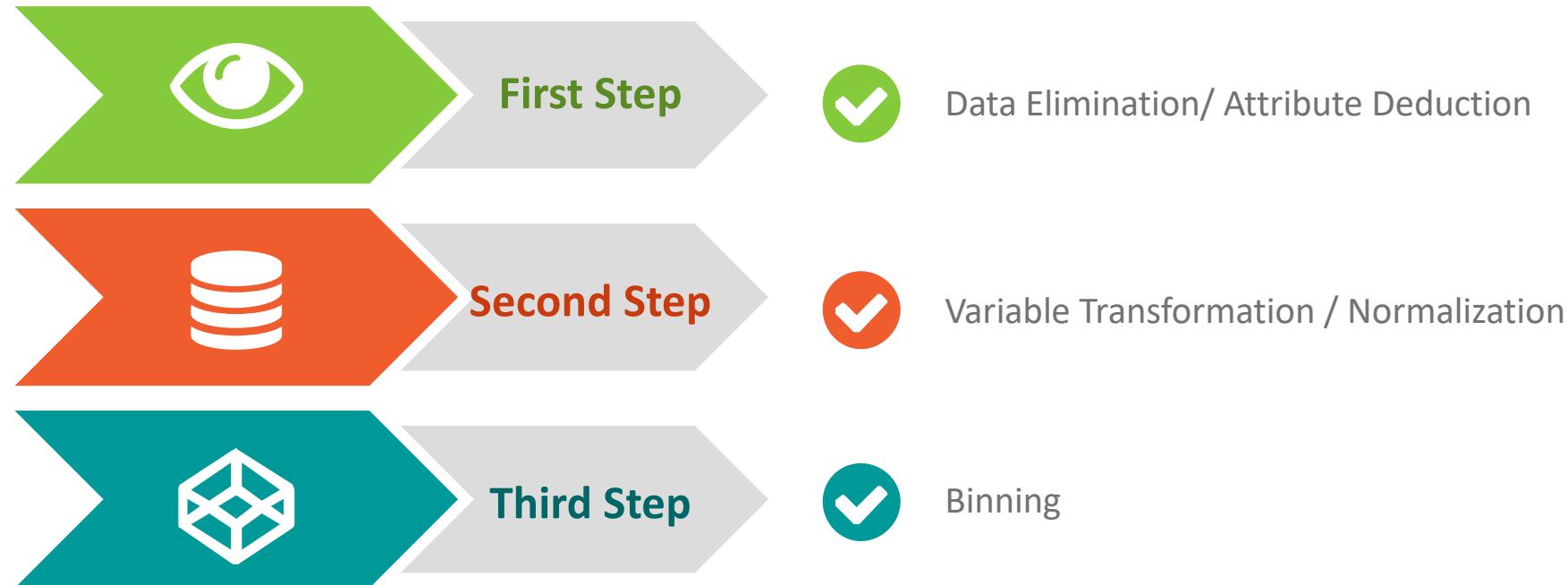


Data Information



Data Preprocessing

Preprocessing Steps



Data Preprocessing

Data Elimination/ Attribute Deduction

Missing Data Elimination:

Why

- 1) Large scale(over 50%) data missing in “contact” and “poutcom” columns
- 2) Occasional data(less than 5%) missing in “education” column

How

Removed “contact” and “poutcom” columns
Deleted rows with missing date in “education” column

Attribute Reduction:

Why

- 1) Information gained after result is obtained--- “duration” column

How

Removed this column



Data Preprocessing

Variable Transformation / Normalization



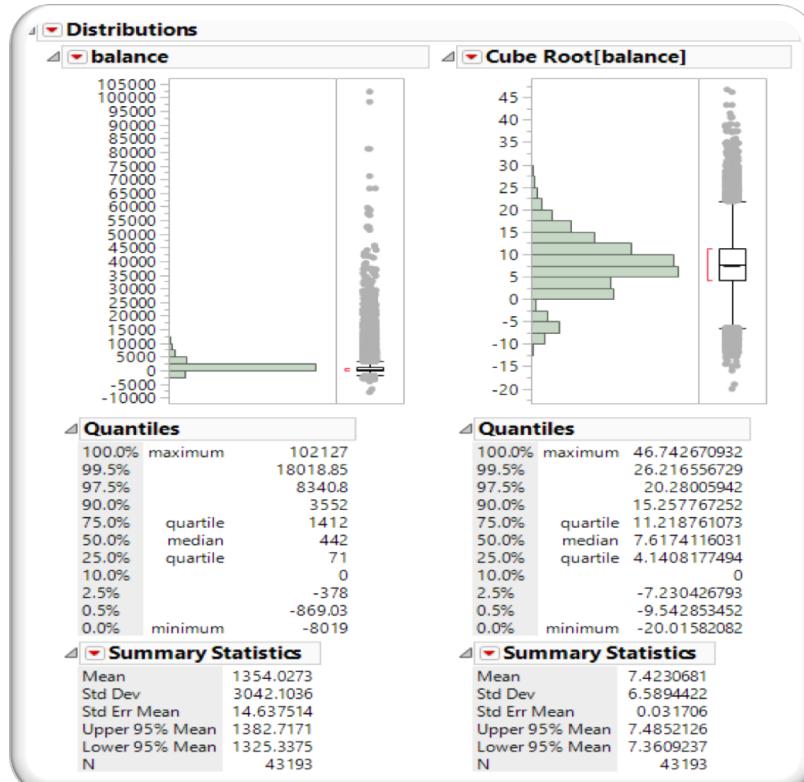
The distribution of “balance” is left skewed and the value range is too large

Why



Using cube root function to transform this variable into a normal distribution

How



Data Preprocessing

Binning



- 1) Making the dataset easier to handle and interpret
- 2) Avoid over emphasize on specific continuous variables

Why



Binning categories under “job” and “month”, based on the Row percentage result of each category as a function of Yes/No

How

job	Count Old Values (11)	New Values (5)
9216 management		ATMS
7355 technician		
5000 admin.		
1540 self-employed		
9278 blue-collar		HEBSer
4004 services		
1411 entrepreneur		
1195 housemaid		
2145 retired		retired
775 student		student
1274 unemployed		unemployed

month	Count Old Values (12)	New Values (5)
2820 apr		feb,apr
2533 feb		
6037 aug		jan,jun,aug,nov
4980 jun		
3842 nov		
1318 jan		
448 mar		mar
13192 may		may,jul
6601 jul		
690 oct		sep,oct,dec
532 sep		
200 dec		

Data Model 1 – Logistic Regression

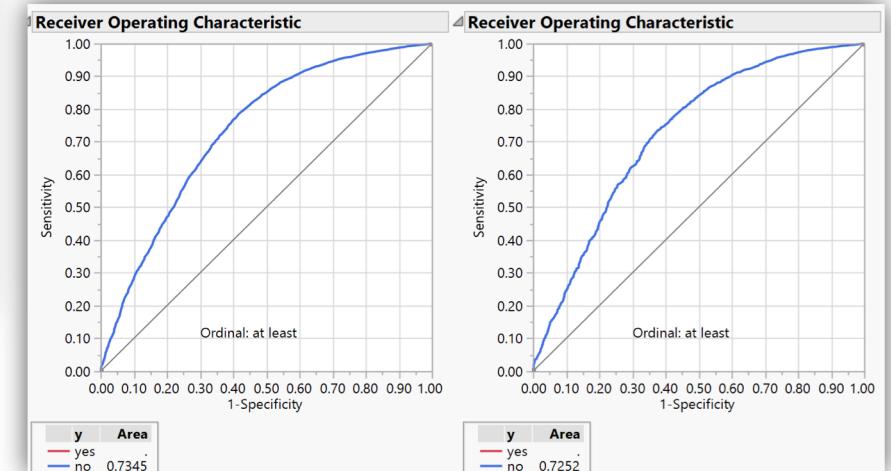
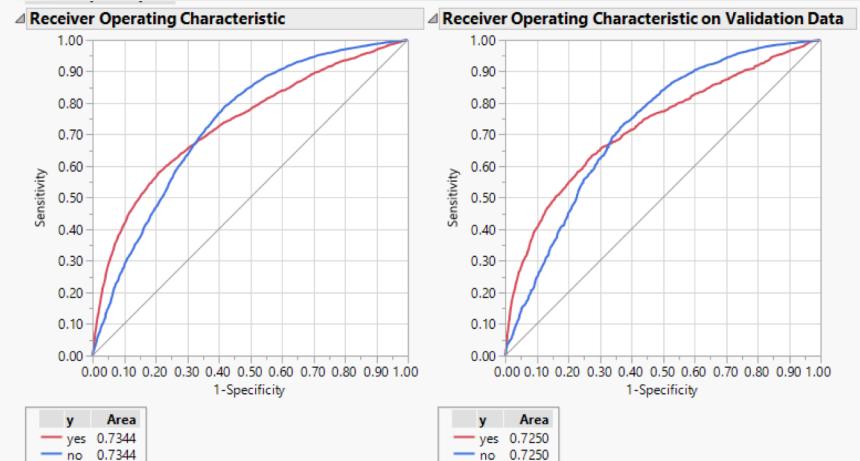


Data Model 1 – Logistic Regression

Model Evaluation

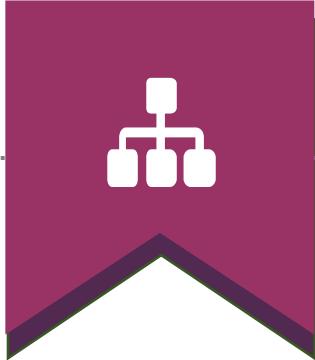


According to the ROC curves of Nominal Logistic and Stepwise Regression, the model has stable and predictive ability.



Data Model 1 – Logistic Regression

Misclassification Rate



Misclassification
Rate

Nominal Logistic Regression: 0.1158 in validation dataset

Stepwise Regression: 0.1158 in validation dataset

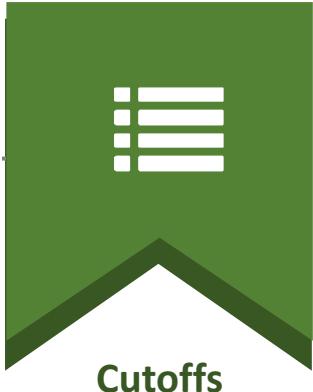
No significant performance difference in making predictions

Fit Details			
Measure	Training	Validation	Definition
Entropy RSquare	0.1138	0.1106	$1 - \text{Loglike}(\text{model})/\text{Loglike}(0)$
Generalized RSquare	0.1531	0.1493	$(1 - (L(0)/L(\text{model}))^{(2/n)})/(1 - L(0)^{(2/n)})$
Mean -Log p	0.3179	0.3214	$\sum -\text{Log}(p[j])/n$
RMSE	0.3028	0.3041	$\sqrt{\sum (y[j] - p[j])^2/n}$
Mean Abs Dev	0.1834	0.1841	$\sum y[j] - p[j] /n$
Misclassification Rate	0.1154	0.1158	$\sum (p[j] \neq p_{\text{Max}})/n$
N	32377	10816	n

Fit Details			
Measure	Training	Validation	Definition
Entropy RSquare	0.1139	0.1107	$1 - \text{Loglike}(\text{model})/\text{Loglike}(0)$
Generalized RSquare	0.1533	0.1494	$(1 - (L(0)/L(\text{model}))^{(2/n)})/(1 - L(0)^{(2/n)})$
Mean -Log p	0.3178	0.3214	$\sum -\text{Log}(p[j])/n$
RMSE	0.3028	0.3041	$\sqrt{\sum (y[j] - p[j])^2/n}$
Mean Abs Dev	0.1833	0.1841	$\sum y[j] - p[j] /n$
Misclassification Rate	0.1154	0.1158	$\sum (p[j] \neq p_{\text{Max}})/n$
N	32377	10816	n

Data Model 1 – Logistic Regression

Cutoffs



Opportunity costs: High opportunity cost which is the profit we could earn, if we predict a False Negative.

Predicted Yes VS. Yes				
Cutoff	0.5	0.4	0.3	0.2
PredictY/Y	8.99%	9.04%	22.71%	36.99%

Data Model 1 – Logistic Regression

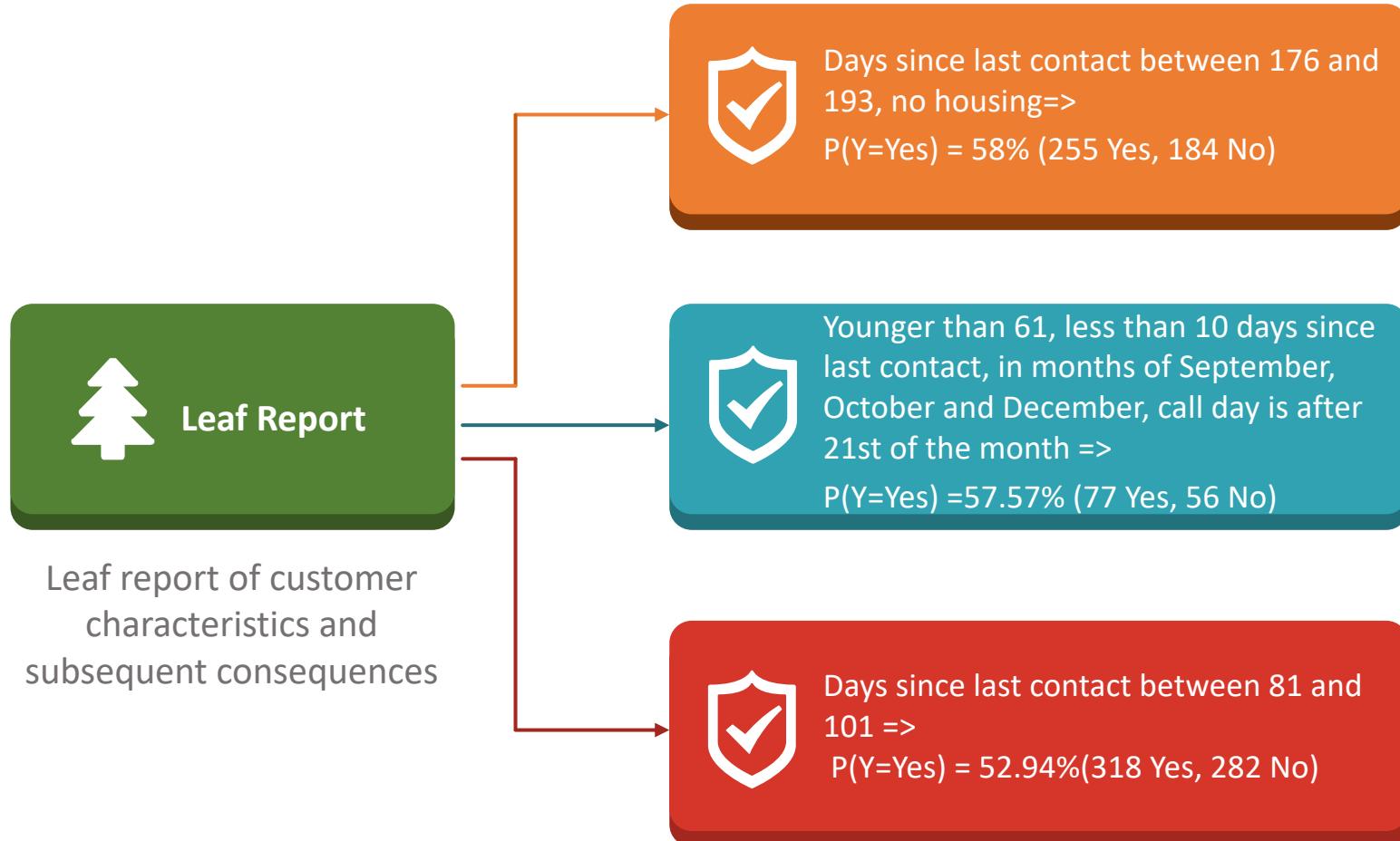
Conclusion

- 1 It is difficult to make an accurate prediction using logistic regression
- 2 We can compare the effect of variables in Nominal Logistic and Stepwise Regression to identify and decide which variables to retain.
- 3 The performances of Nominal Logistic and Stepwise Regression are alike in terms of making predictions



Data Model 2 – Decision Tree

Leaf Report



Fit Details

Measure	Training	Validation	Definition
Entropy RSquare	0.1618	0.1438	$1 - \text{Loglike}(\text{model})/\text{Loglike}(0)$
Generalized RSquare	0.2141	0.1918	$(1 - L(0)/L(\text{model}))^{(2/n)} / (1 - L(0)^{(2/n)})$
Mean -Log p	0.3007	0.3094	$\sum -\text{Log}(p_{ij})/n$
RMSE	0.2929	0.2984	$\sqrt{\sum (y_{ij} - \hat{p}_{ij})^2/n}$
Mean Abs Dev	0.1716	0.1754	$\sum y_{ij} - \hat{p}_{ij} /n$
Misclassification Rate	0.1116	0.1163	$\sum (p_{ij} \neq p_{\text{Max}})/n$
N	32377	10816	n

Confusion Matrix

		Training		Validation	
		Actual	Predicted Count	Actual	Predicted Count
Actual	Predicted	no	yes	no	yes
		no	27927 697	no	9281 267
yes	no	2917	836	yes	991 277

Data Model 2 – Decision Tress

Conclusion



Easier to interpret the results, target groups could be prioritized for decision makers



Depicted specific characteristics of certain groups



Dispersed groups and large number of leaves make it difficult for the bank to initiate marketing campaigns



Data Model 3 – Neural Networks

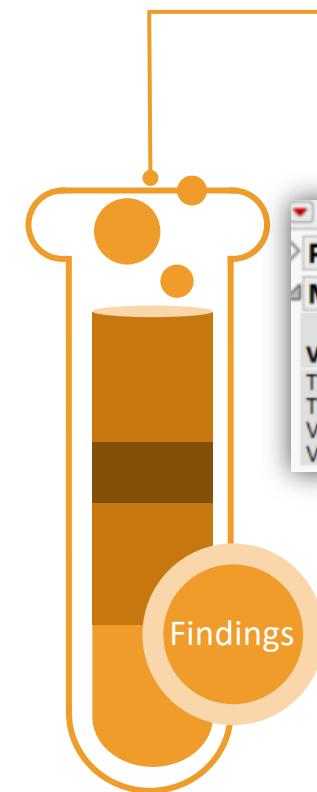
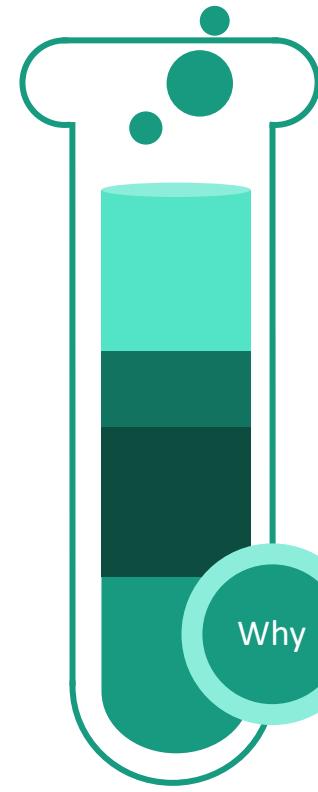
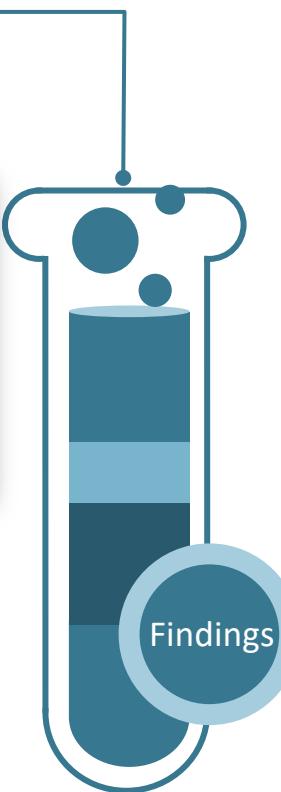
Neural Networks

Combining different amount of nodes in layers and different activation Functions to generalize or recognize unrealized patterns

Finding 1

Comparison of Activation functions(TanH, Linear, Gaussian)

Model Comparison		
Predictors		
Measures of Fit for y		
Validation	Creator	.2 .4 .6 .8
Training	Neural	
Training	Neural	
Training	Neural	
Validation	Neural	
Validation	Neural	
Validation	Neural	
RMSE	Mean	Misclassification Rate
0.2956	0.1747	0.1122
0.3027	0.1834	0.1153
0.2977	0.1774	0.1126
0.2996	0.1768	0.1153
0.3042	0.1842	0.1155
0.3010	0.1792	0.1156



Finding 2

Comparison of Number of Nodes/ Layers (3/3, 3/0)

Model Comparison		
Predictors		
Measures of Fit for y		
Validation	Creator	.2 .4 .6 .8
Training	Neural	
Training	Neural	
Validation	Neural	
Validation	Neural	
RMSE	Mean	Misclassification Rate
0.2956	0.1747	0.1122
0.2972	0.1774	0.1152
0.2996	0.1768	0.1153
0.3015	0.1800	0.1168

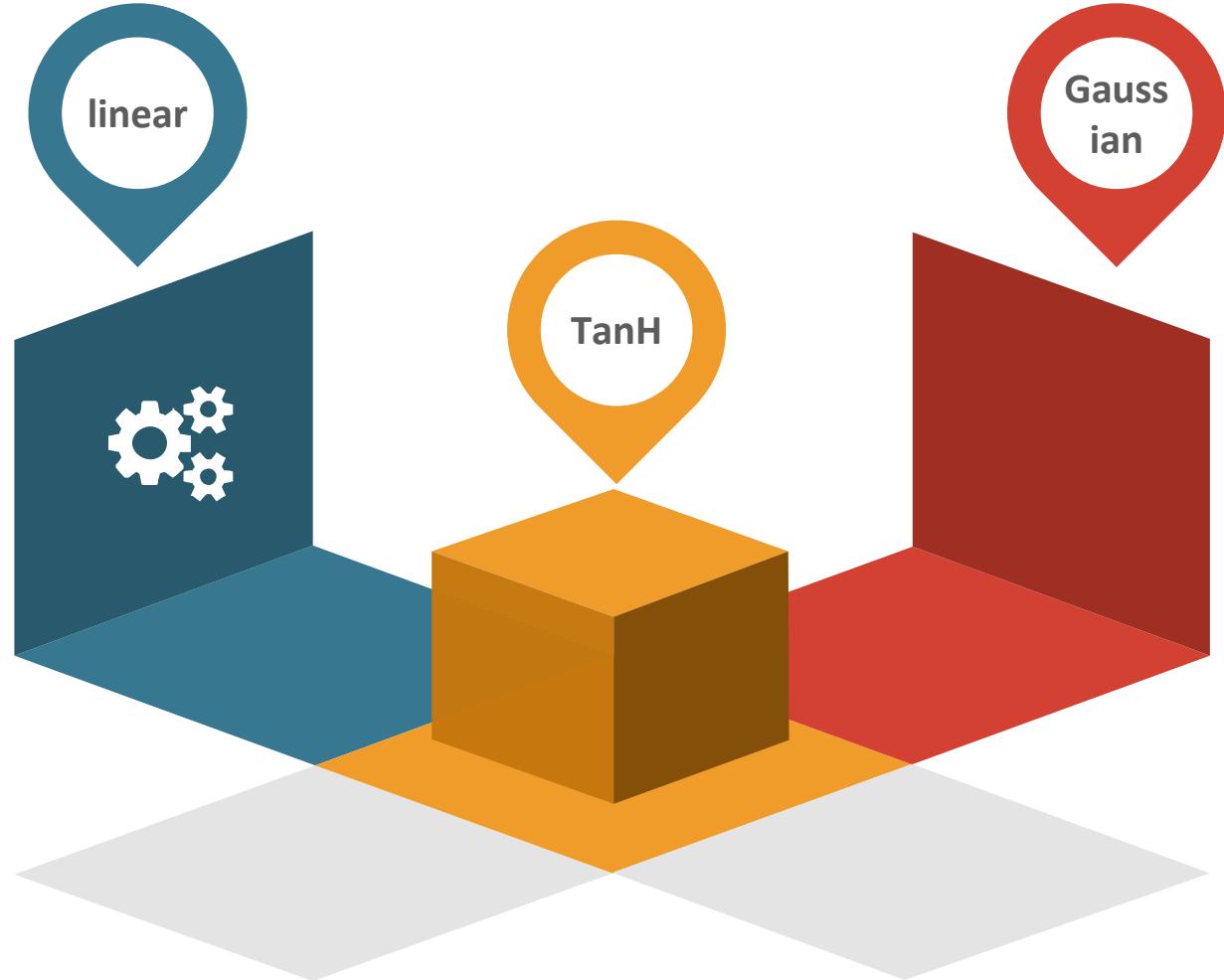
TanH has the lowest Misclassification Rate

The misclassification rate of TanH with 3/0 of nodes is slightly lower

Data Model 3 – Neural Networks

Conclusion

By using TanH function with 3/0 nodes which has the highest rate of predicting “Yes”, the Bank could make predictions based on Neural Network method



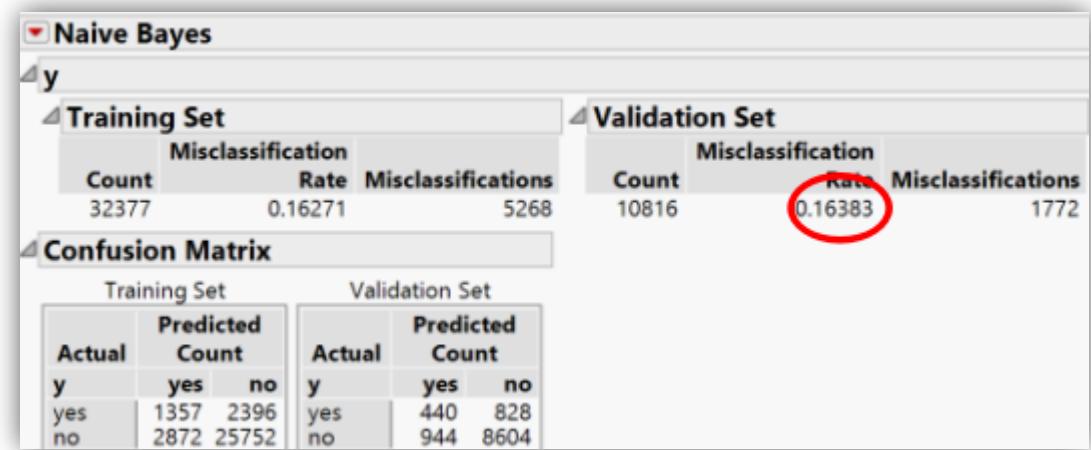
Other Analyses– Naïve Bayes



It is time effective to obtain a predictive model



Comparison basis to find best model to solve the problem



Conclusion

01 Number One

With best performance in predicting probability of people who would accept the offer

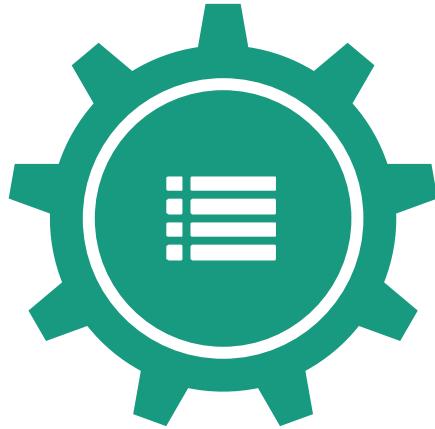
02 Number Two

With noticeably higher Misclassification in exchange

03 Number Three

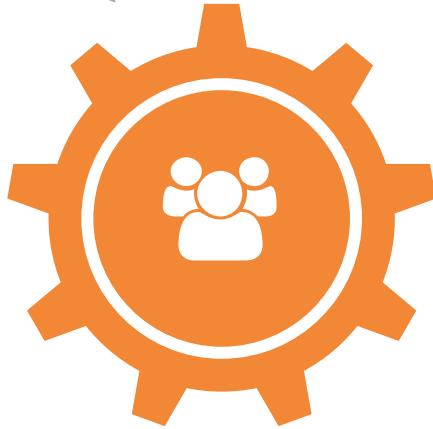
The bank should use this model with caution

Other Analyses— Hierarchical Clustering



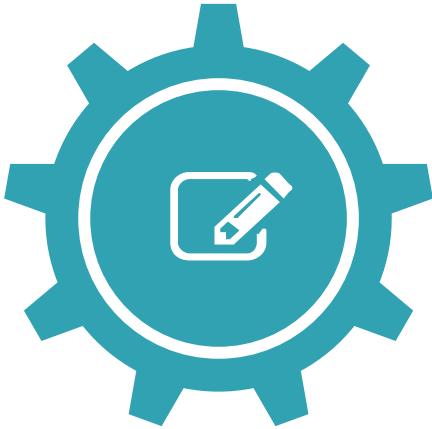
Reason 1

To segment customers and find factors that may affect their decisions



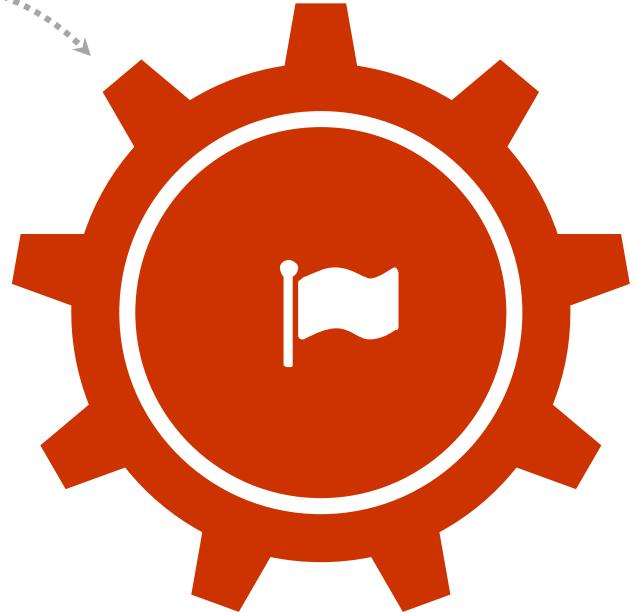
Reason 2

To find target groups of customers that have higher probability to accept the offer



Reason 3

To help the bank to design marketing strategies toward target customers



How

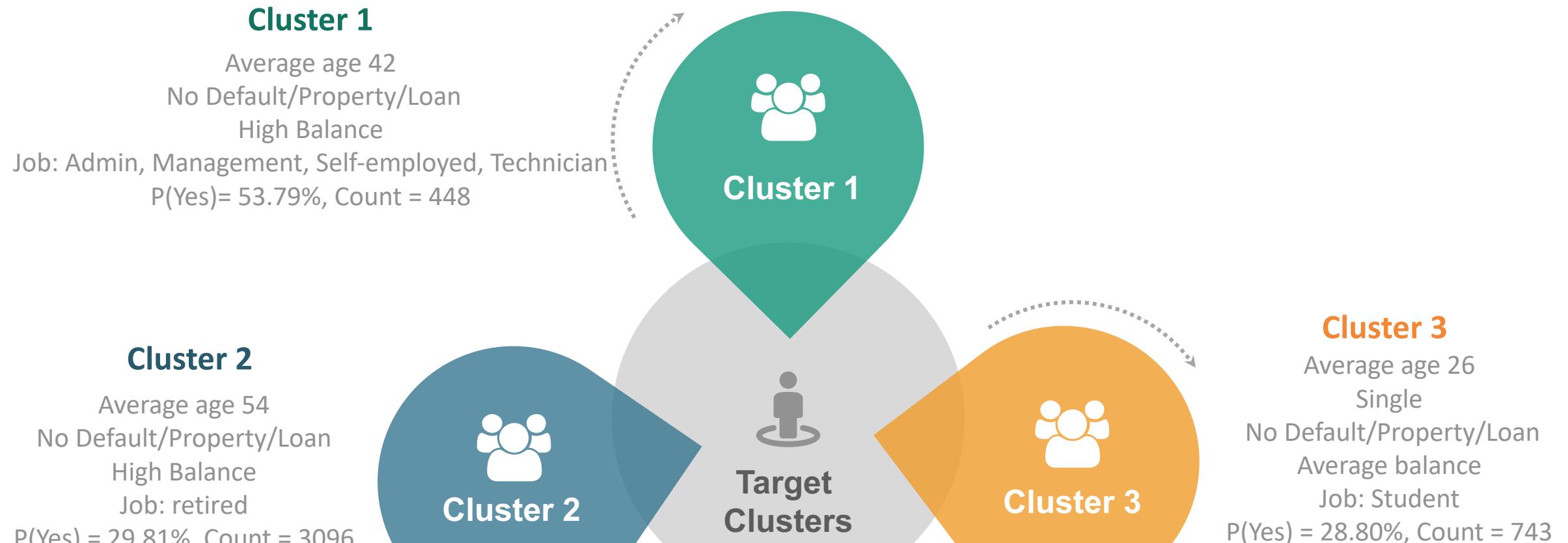
Start from a larger number of clusters, and gradually approach the best number where the objective of the clustering is met

Other Analyses– Hierarchical Clustering

Conclusion

-  01 Out Eight customer clusters we proposed to study, we identified 3 of them as target groups
-  02 The three target groups have symbolic characteristics to label and to distinguish from the other
-  03 The counts of the three groups are not large, efforts could be made to enlarge target groups by formulating a larger customer basis
-  04 Marketing campaigns can be devised according to the target clusters to expand cluster size with high value or to improve converting rates or potential customers

Other Analyses– Hierarchical Clustering



Summary

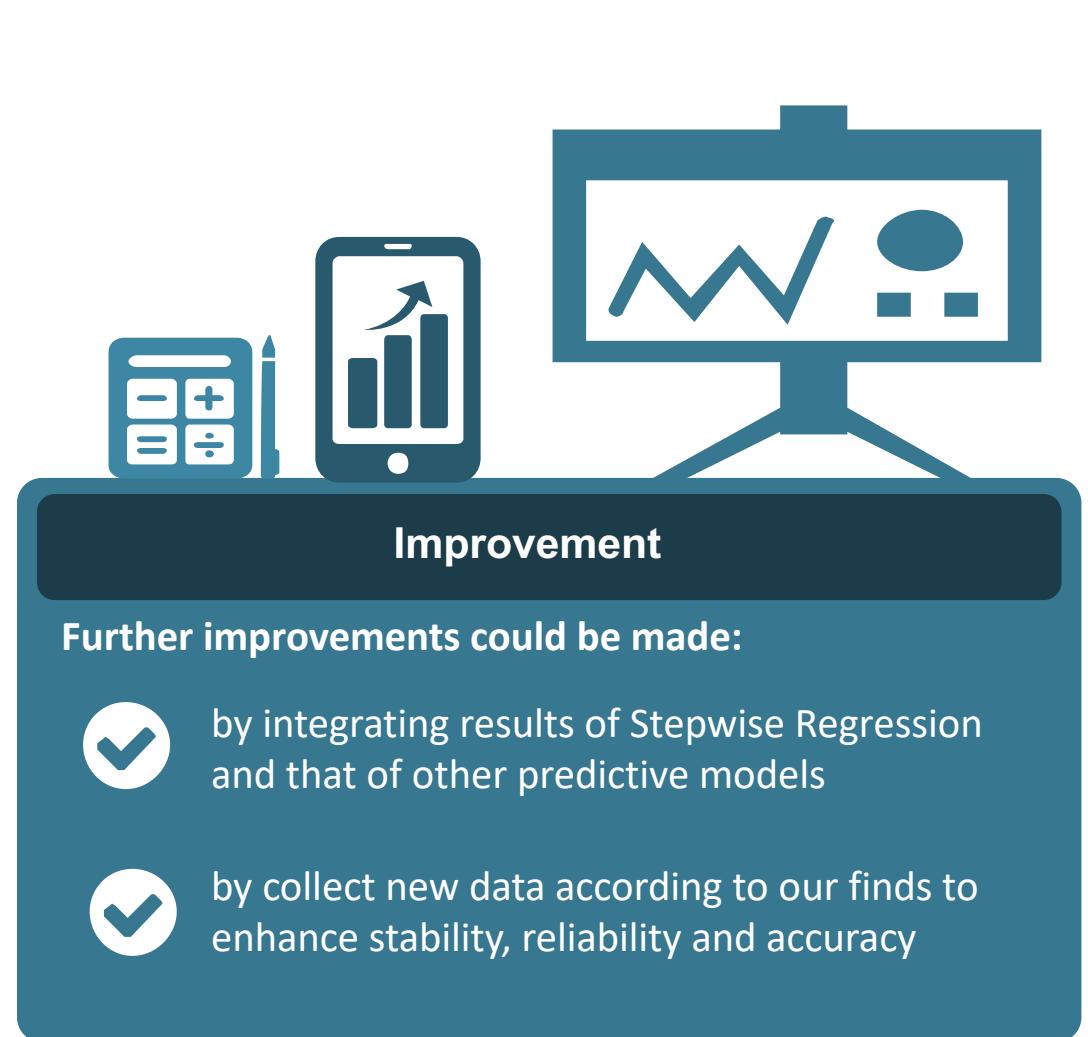
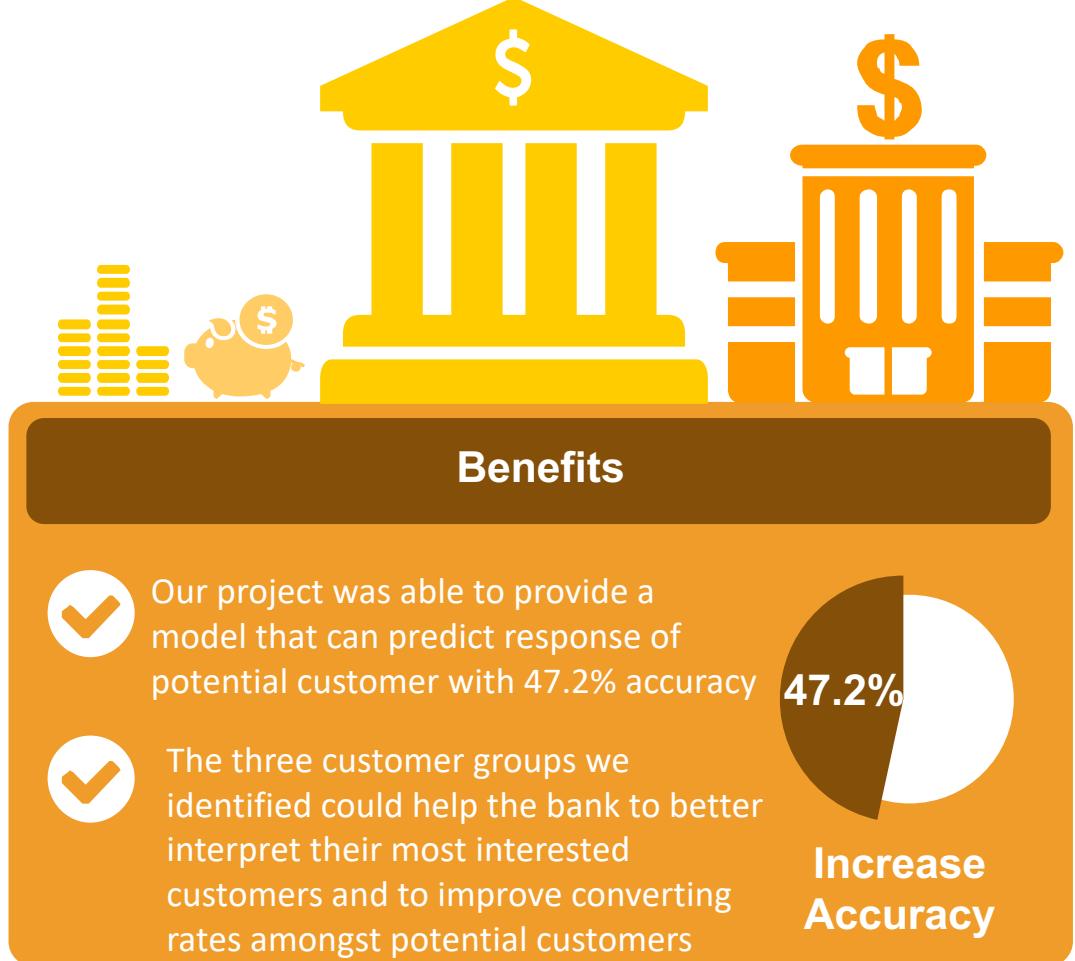
Models Comparison

Model	Misclassification	Predict "Yes"
Logistic Regression	0.1158	8.99%
Decision Tree	0.1163	21.85%
Neural Networks	0.1153	14.9%
Naïve Bayes	0.1638	34.7%

Predicted Yes VS. Yes (Naïve Bayes)				
Cut-off	0.5	0.4	0.3	0.2
P(Yes) VS. Yes	34.7%	38.26%	40.9%	44.4%
Predicted Yes VS. Yes (Logistic Regression)				
Cut-off	0.5	0.4	0.3	0.2
P(Yes) VS. Yes	8.99%	9.04%	22.71%	36.99%
Predicted Yes VS. Yes (Decision Tree)				
Cut-off	0.5	0.4	0.3	0.2
P(Yes) VS. Yes	21.85%	32.66%	43.23%	47.2%
Predicted Yes VS. Yes (Neural Network)				
Cut-off	0.5	0.4	0.3	0.2
P(Yes) VS. Yes	14.9%	29.47%	31.71%	41.76%

- Neural Networks has the lowest misclassification
- Naive Bayes has a significantly better accuracy in predicting "Yes"
- Decision Tree Model is recommended, taking into consideration of opportunity costs

Summary



Thank you !

