

- DATA
- PREPARATION
- VARIABLES  
SELECTION
- PARTITION
- CLASSIFICATION
- RESULTS
- SUMMARY
- CONCLUSION

- BANKRUPTCY

- BOOK CLUB

- HOME EQUITY

**DATA ANALYTICS MASTER PROGRAM**  
**PRACTICUM PROJECT**  
BY MANDY NGUYEN

Case Studies: Bankruptcy  
Charles Book Club  
Home Equity

December 08, 2021

- **BANKRUPTCY**
- **BOOK CLUB**
- **HOME EQUITY**

- **DATA**
- **PREPARATION**
- **VARIABLES  
SELECTION**
- **PARTITION**
- **CLASSIFICATION**
- **RESULTS**
- **SUMMARY**
- **CONCLUSION**



**DATA**

DATA

BANKRUPTCY

BOOK CLUB

HOME EQUITY

GOAL

Find best models to classify target variables.

OBSERVATIONS

132

4,000

5,960

VARIABLES

27

19

13

PREDICTORS

24

16

12

TARGET

“D”

“Florence”

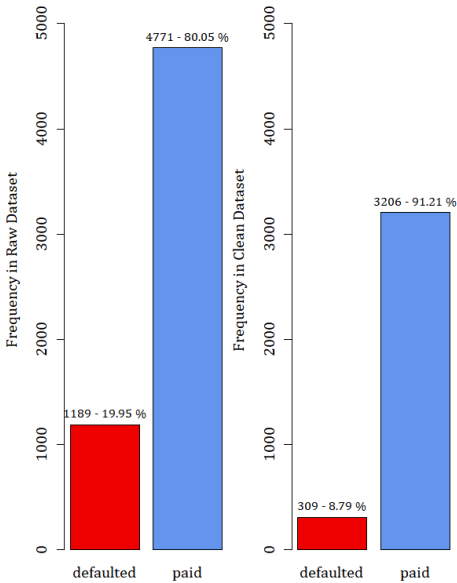
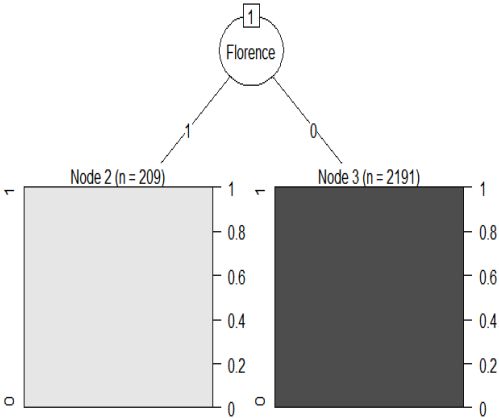
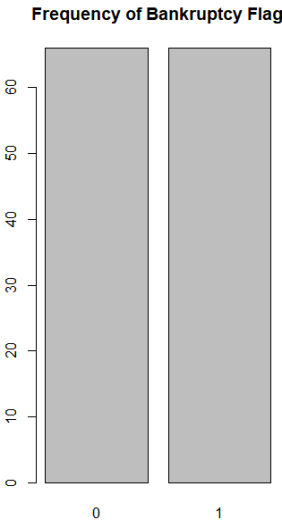
“BAD”

Class: 1 = bankrupt  
Class : 0 = healthy

Class: 1 = Yes to Florence  
Class: 0 = No to Florence

Class: 1 = Defaulted (loan)  
Class: 0 = Paid (loan)

TARGET’s CLASS  
DISTRIBUTION  
CHART



- **BANKRUPTCY**

- **BOOK CLUB**

- **HOME EQUITY**

- **DATA**

- **PREPARATION**

- **VARIABLES  
SELECTION**

- **PARTITION**

- **CLASSIFICATION**

- **RESULTS**

- **SUMMARY**

- **CONCLUSION**

**PREPARATION**

	• <b>BANKRUPTCY</b>	• <b>BOOK CLUB</b>	• <b>HOME EQUITY</b>
• <b>DATA</b>			
• <b>PREPARATION</b>			
• <b>VARIABLES SELECTION</b>	ATTRIBUTE ANALYSIS	ATTRIBUTE ANALYSIS	ATTRIBUTE ANALYSIS
	RANDOMIZE DATA SET	COMPUTE RFM SCORE	REMOVE MISSING VALUES
• <b>PARTITION</b>	SELECT VARIABLES SET	SELECT VARIABLES SET	COMBINE CLASSES
• <b>CLASSIFICATION</b>	NORMALIZE DATA	PARITION	NORMALIZE DATA
• <b>RESULTS</b>	PARTITION	BALANCE	PARTITION
• <b>SUMMARY</b>			BALANCE TRAINING
• <b>CONCLUSION</b>			

- **BANKRUPTCY**

- **BOOK CLUB**

- **HOME EQUITY**

- DATA
- PREPARATION
- VARIABLES SELECTION
- PARTITION
- CLASSIFICATION
- RESULTS
- SUMMARY
- CONCLUSION

## REMOVE MISSING VALUES

## COMBINE CLASSES

[illegible]

- **BANKRUPTCY**

- **BOOK CLUB**

- **HOME EQUITY**

- **DATA**

- **PREPARATION**

- **VARIABLES  
SELECTION**

- **PARTITION**

- **CLASSIFICATION**

- **RESULTS**

- **SUMMARY**

- **CONCLUSION**

# VARIABLES SELECTION

	<div><ul style="list-style-type: none"><li>• <b>BANKRUPTCY</b></li></ul></div>	<div><ul style="list-style-type: none"><li>• <b>BOOK CLUB</b></li></ul></div>	<div><ul style="list-style-type: none"><li>• <b>HOME EQUITY</b></li></ul></div>
<div><ul style="list-style-type: none"><li>• <b>DATA</b></li><li>• <b>PREPARATION</b></li><li>• <b>VARIABLES SELECTION</b></li><li>• <b>PARTITION</b></li><li>• <b>CLASSIFICATION</b></li><li>• <b>RESULTS</b></li><li>• <b>SUMMARY</b></li><li>• <b>CONCLUSION</b></li></ul></div>	<div>LOGISTIC REGRESSION</div> <div>CORRELATION</div>	<div>RFM</div> <div>GIVEN ATTRIBUTES</div> <div>LOGISTIC REGRESSION</div>	<div>LOGISTIC REGRESSION</div> <div>CORRELATION</div> <div>INFORMATION VALUES (SCORECARD BY WOE)</div>





- DATA
- PREPARATION
- **VARIABLES SELECTION**
- PARTITION
- CLASSIFICATION
- RESULTS
- SUMMARY
- CONCLUSION

- **BANKRUPTCY**

- **BOOK CLUB**

- **HOME EQUITY**

### **“RFM” SCORE**

Concatenate R\_code, F\_code, M\_code to creates RFM classes

#### **R\_CODE**

Recency (“Rcode” attribute):

- 0–2 months (Rcode = 1)
- 3–6 months (Rcode = 2)
- 7–12 months (Rcode = 3)
- 13 months and up (Rcode = 4)

#### **F\_CODE**

Frequency (“Fcode” attribute):

- 1 book (Fcode = 1)
- 2 books (Fcode = 2)
- 3 books and up (Fcode = 3)

#### **M\_CODE**

Monetary (“Mcode” attribute):

- \$0–\$25 (Mcode = 1)
- \$26–\$50 (Mcode = 2)
- \$51–\$100 (Mcode = 3)
- \$101–\$200 (Mcode = 4)
- \$201 and up (Mcode = 5)

Seq.	ID	M	R	F	Mcode	Rcode	Fcode	Yes_Florence	No_Florence	RFM_score
1	25	297	14	2	5	4	2	0	1	425

**RFM = “4\_5\_2”**

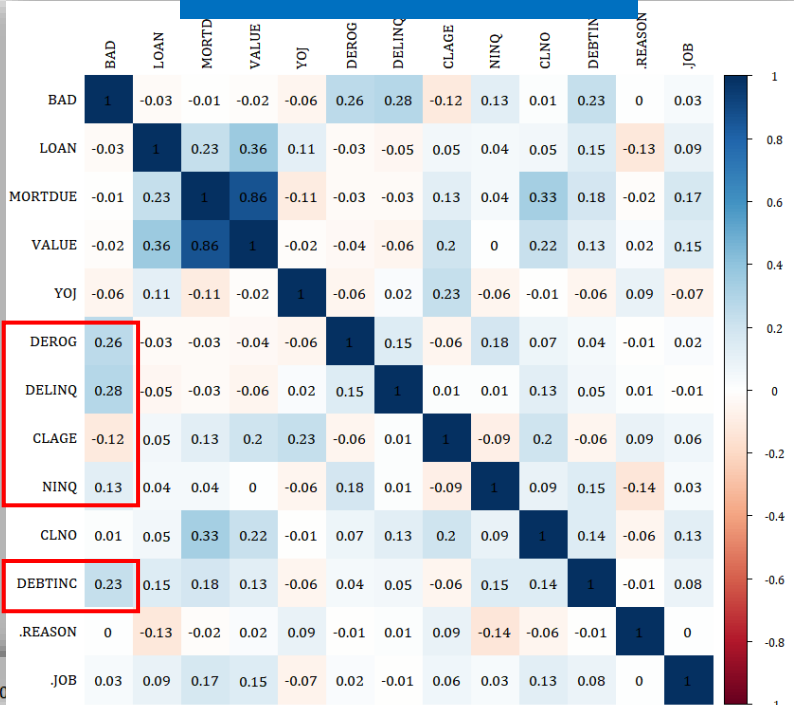
- DATA
- PREPARATION
- VARIABLES SELECTION
- PARTITION
- CLASSIFICATION
- RESULTS
- SUMMARY
- CONCLUSION

• BANKRUPTCY

• BOOK CLUB

• HOME EQUITY

CORRELATION



GLM

Call:  
glm(formula = STATUS ~ LOAN + JOB + DEROG + DELINQ + CLAGE +  
NINQ + CLNO + DEBTINC, family = "binomial", data = X)

Deviance Residuals:

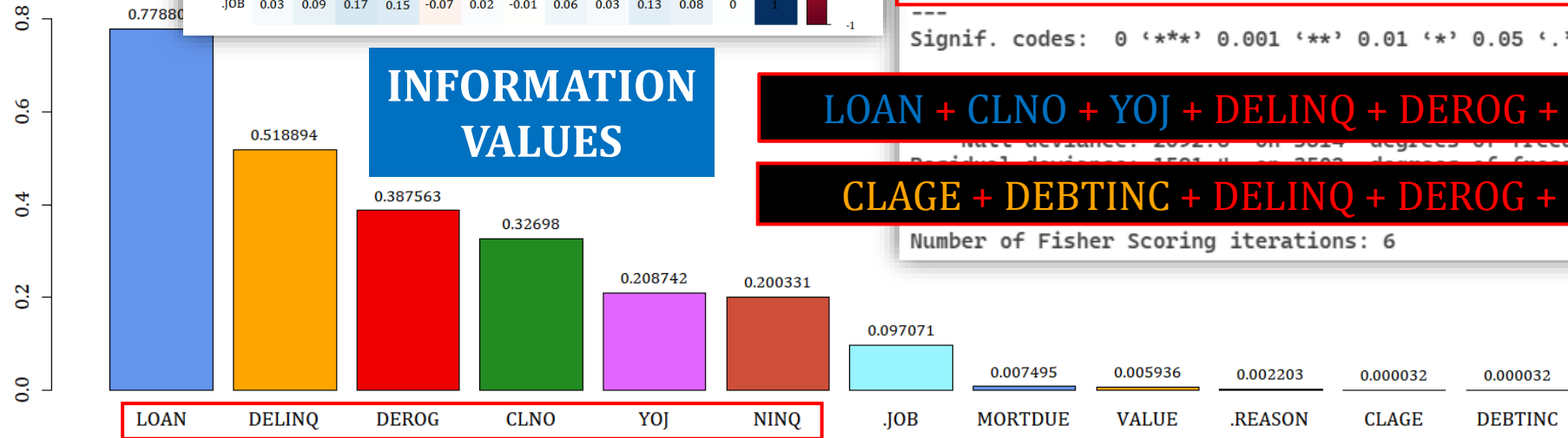
	Min	1Q	Median	3Q	Max
	-3.9354	0.1913	0.2816	0.3999	1.6910

Coefficients:

	Estimate	Std. Error	z value	Pr(> z )
(Intercept)	5.100e+00	4.661e-01	10.943	< 2e-16 ***
LOAN	1.480e-05	6.992e-06	2.117	0.034282 *
JOBOffice	6.049e-01	2.636e-01	2.295	0.021725 *
JOBOther	8.124e-02	2.065e-01	0.393	0.694035
JOBProfExe	-3.175e-02	2.346e-01	-0.135	0.892357
JOBSales	-1.321e+00	4.462e-01	-2.961	0.003063 **
JOBSelf	-7.832e-01	3.953e-01	-1.981	0.047558 *
DEROG	-7.427e-01	1.012e-01	-7.336	2.20e-13 ***
DELINQ	-7.511e-01	6.907e-02	-10.874	< 2e-16 ***
CLAGE	5.540e-03	1.036e-03	5.346	9.01e-08 ***
NINQ	-1.269e-01	3.739e-02	-3.395	0.000687 ***
CLNO	1.980e-02	7.643e-03	2.590	0.009592 **
DEBTINC	-1.044e-01	1.031e-02	-10.119	< 2e-16 ***

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

INFORMATION VALUES



LOAN + CLNO + YOJ + DELINQ + DEROG + NINQ

CLAGE + DEBTINC + DELINQ + DEROG + NINQ

Number of Fisher Scoring iterations: 6

- **BANKRUPTCY**

- **BOOK CLUB**

- **HOME EQUITY**

- **DATA**

- **PREPARATION**

- **VARIABLES  
SELECTION**

- **PARTITION**

- **CLASSIFICATION**

- **RESULTS**

- **SUMMARY**

- **CONCLUSION**

# PARTITION & CLASSIFICATION

- **DATA**
- **PREPARATION**
- **VARIABLES SELECTION**
- **PARTITION**
- **CLASSIFICATION**
- **RESULTS**
- **SUMMARY**
- **CONCLUSION**

**6**

**SUBSETS**

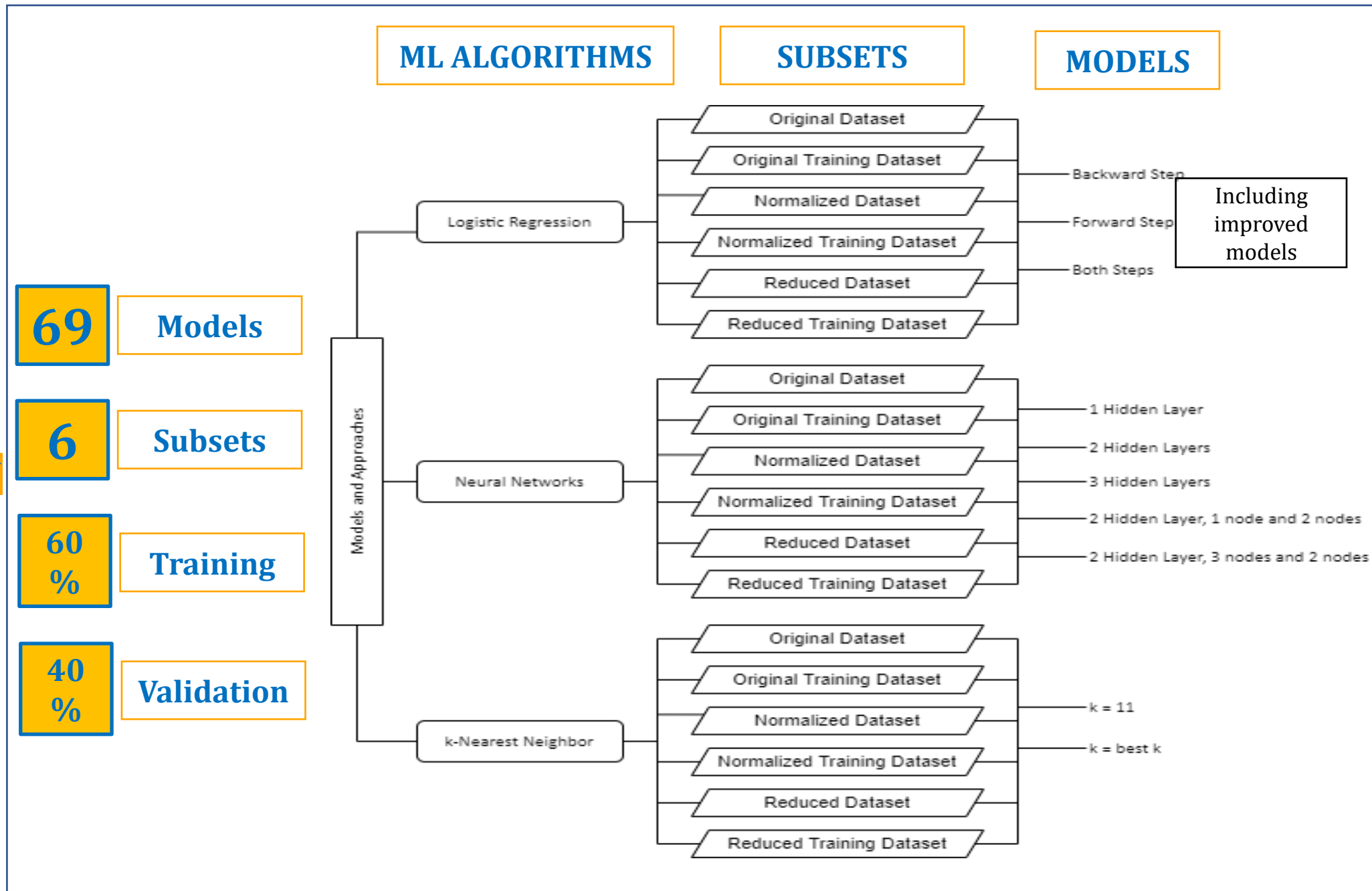
- **Original Dataset:** Original Randomized dataset with all predictors
- **Original Training Dataset:** 60% Original Dataset
- **Normalized Dataset:** Original Randomized dataset with all predictors normalized by Min-Max method
- **Normalized Training Dataset:** 60% Normalized Dataset
- **Reduced Dataset:** Original Randomized dataset with 4 predictors R9, R10, R17, R20 (selected from correlation plot)
- **Reduced Training Dataset:** 60% Reduced Dataset

# BANKRUPTCY

# BOOK CLUB

# HOME EQUITY

- DATA
- PREPARATION
- VARIABLES SELECTION
- PARTITION
- CLASSIFICATION
- RESULTS
- SUMMARY
- CONCLUSION



- **BANKRUPTCY**
- **BOOK CLUB**
- **HOME EQUITY**

- **CLASSIFICATION**

Models	Accuracy	95% CI	No Info. Rate	P-Value	Kappa	Sensitivity	Specificity
<b>LOGISTIC REGRESSION MODELS</b>							
<b>Improved Original Dataset (step = backward = both)</b> Retained Ratios: R3, R5, R6, R9, R10, R16, R17, R18, R22, R23, R24	0.8679	(0.7466, 0.9452)	0.5283	1.709e-07	0.7333	0.9286	0.8000
<b>Normalized Dataset</b>	0.9434	(0.8434, 0.9882)	0.6604	1.006e-06	0.8721	0.9714	0.8889
<b>Reduced Dataset</b>							
<b>Reduced Dataset</b>	0.9245	(0.8179, 0.9791)	0.6604	6.766e-06	0.8271	0.9714	0.8333
<b>Improved Reduced Dataset (step = backward = both)</b> Retained Ratios: R9, R10, R17	0.9245	(0.8179, 0.9791)	0.6604	6.766e-06	0.8271	0.9714	0.8333
<b>NEURAL NETWORKS</b>							
<b>Improved Original Dataset (hidden: layer = 1; node = 3)</b>	0.9623	(0.8702, 0.9954)	0.5283	2.356e-12	0.9243	0.9643	0.9600
<b>Normalized Dataset (hidden: layer = 1; node = 1)</b>	0.9623	(0.8702, 0.9954)	0.6604	1.104e-07	0.9159	0.9714	0.9444
<b>Improved Normalized Dataset (hidden: layer = 1; node = 3)</b>	0.9811	(0.8993, 0.9995)	0.6604	7.945e-09	0.9585	0.9714	1.0000
<b>Reduced Dataset (hidden: layer = 1; node = 1)</b>	0.9623	(0.8702, 0.9954)	0.6604	1.104e-07	0.9159	0.9714	0.9444
<b>K-NEAREST NEIGHBOR</b>							
<b>Improved Original Dataset (k = 1)</b>	1	(0.9328, 1)	0.5283	2.055e-15	1	1.0000	1.0000
<b>Improved Original Dataset (k = 3)</b>	0.8113	(0.6803, 0.9056)	0.5283	1.709e-05	0.6198	0.8571	0.7600
<b>Improved Normalized Dataset (k = 1)</b>	1	(0.9328, 1)	0.6604	2.812e-10	1	1.0000	1.0000
<b>Improved Normalized Dataset (k = 2)</b>	0.8868	(0.7697, 0.9573)	0.6604	0.0001552	0.7665	0.8286	1.0000
<b>Improved Reduced Dataset (k = 1)</b>	1	(0.9328, 1)	0.6604	2.812e-10	1	1.0000	1.0000
<b>Improved Reduced Dataset (k = 3)</b>	0.9057	(0.7934, 0.9687)	0.6604	3.58e-05	0.7979	0.8857	0.9444

	<div><div>BANKRUPTCY</div></div>	<div><div>BOOK CLUB</div></div>	<div><div>HOME EQUITY</div></div>		
DATA					
PREPARATION					
VARIABLES SELECTION					
PARTITION					
CLASSIFICATION	<div><div>SIGNIFICANT CLASS: "1"</div><div>CUTOFF (C): 0.5</div></div>	<div><div>RFM Score</div><div>C &gt; mean</div><div>ALL</div></div>	<div><div>K-Nearest Neighbors</div><div>K14</div><div>SET 1</div></div>	<div><div>LOGISTIC REGRESSION</div><div>C0.5</div><div>C0.5</div><div>C0.5</div><div>ALL</div><div>SET 2</div><div>SET 3</div></div>	
RESULTS	ACCURACY	0.81	0.92	0.91	0.91
	SENSITIVITY	0.00	0.00	0.04	0.02
SUMMARY	SPECIFICITY	0.87	1.00	0.99	1.00
CONCLUSION	<div><div>SET 1: PREDICTORS: "FirstPurch", "Related.Purchase"; TARGETS: "Yes_Florence", "No_Florence"</div><div>SET 2: "ArtBks", "ItalArt", "ItalAtlas", "ItalCook"</div><div>SET 3: "Rcode", "Fcode", "Mcode"</div></div>				



- **HOME EQUITY**

- 0.45**

**SET 1:** PREDICTORS: "FirstPurch", "Related.Purchase"; TARGETS: "Yes\_Florence", "No\_Florence"  
**SET 2:** "ArtBks", "ItalArt", "ItalAtlas", "ItalCook"  
**SET 3:** "Rcode", "Fcode", "Mcode"

- BANKRUPTCY

- BOOK CLUB

- HOME EQUITY

6

Subsets

Variables Sets

Variables

Partitions

TRAINING 70%

VALIDATION 30%

IMBALANCED

BALANCED

ALL

SET 1

SET 2

LOAN  
MORTDUE  
VALUE  
REASON  
JOB  
YOJ  
DEROG  
DELINQ  
CLAGE  
NINQ  
CLNO  
DEBTINC

LOAN  
  
YOJ  
DEROG  
DELINQ  
  
NINQ  
CLNO

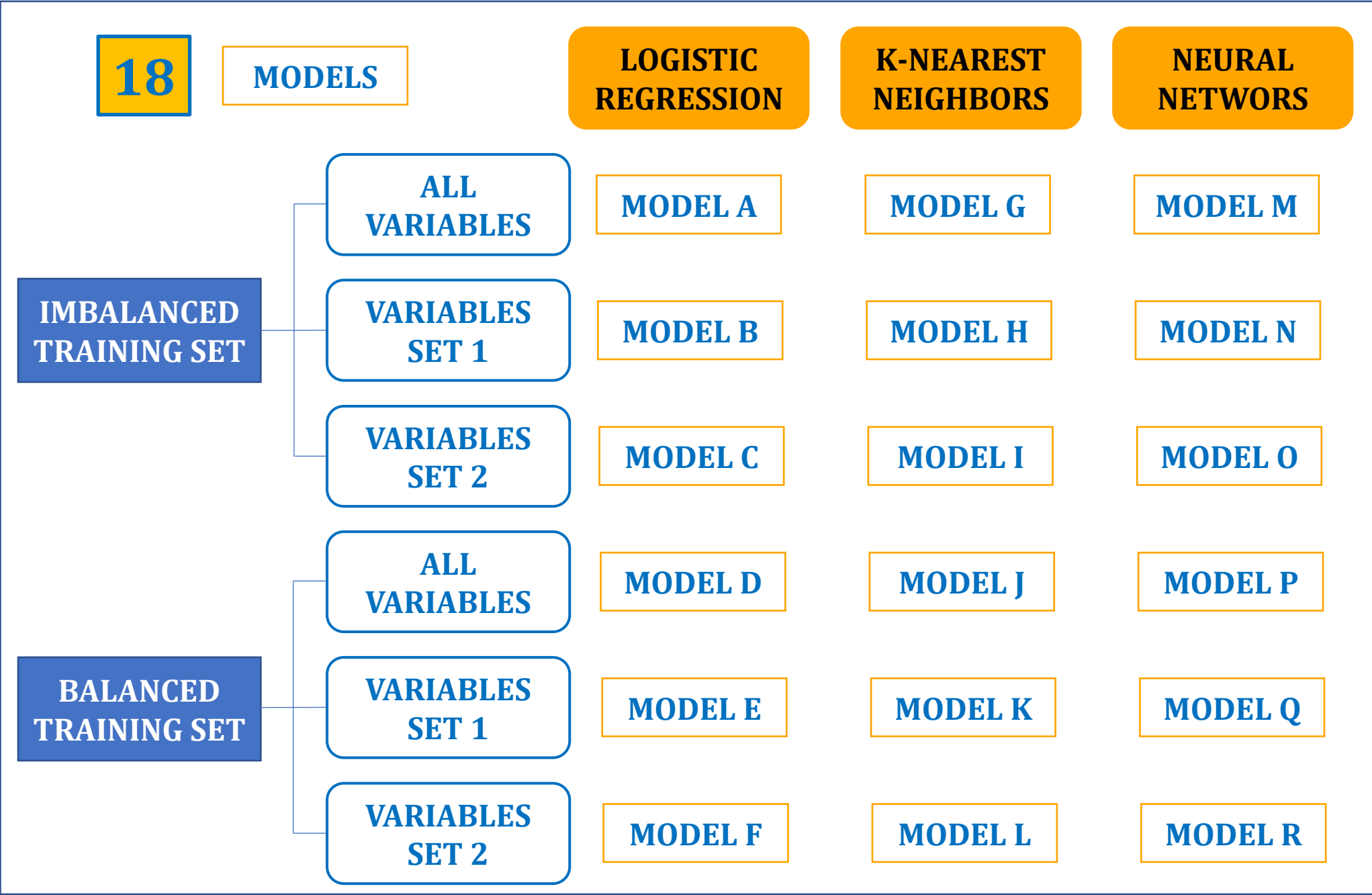
DEROG  
DELINQ  
CLAGE  
NINQ  
  
DEBTINC

- DATA
- PREPARATION
- VARIABLES SELECTION
- PARTITION
- CLASSIFICATION
- RESULTS
- SUMMARY
- CONCLUSION

- BANKRUPTCY

- BOOK CLUB

- HOME EQUITY



- **HOME EQUITY**

- BEST MODEL: B: IMBALANCED TRAINING SET**  
**SET 1: LOAN, JOY, CLAGE, DEBTINC, NINQ, CLNO**

- DATA
- PREPARATION
- VARIABLES SELECTION
- PARTITION
- **CLASSIFICATION**
- RESULTS
- SUMMARY
- CONCLUSION

• BANKRUPTCY			• BOOK CLUB			• <b>HOME EQUITY</b>		
<div>K-Nearest Neighbors</div> <div>SIGNIFICANT CLASS: "PAID"</div> <div>K: various</div>			IMBALANCED TRAINING SET			BALANCED TRAINING SET		
			K5	K5	K5	K15	K7	K15
			G	H	I	J	K	L
			ALL	SET 1	SET 2	ALL	SET 1	SET 2
ACCURACY			0.92	0.91	0.92	0.73	0.83	0.84
SENSITIVITY			0.99	0.99	0.99	0.78	0.87	0.91
SPECIFICITY			0.08	0.21	0.10	0.35	0.41	0.40
BEST MODEL: I:			IMBALANCED TRAINING SET					
			K5					
			SET 2: DEROG, DELINQ, CLAGE, NINQ, DEBTINC					

- **HOME EQUITY**

- ## CONCLUSION

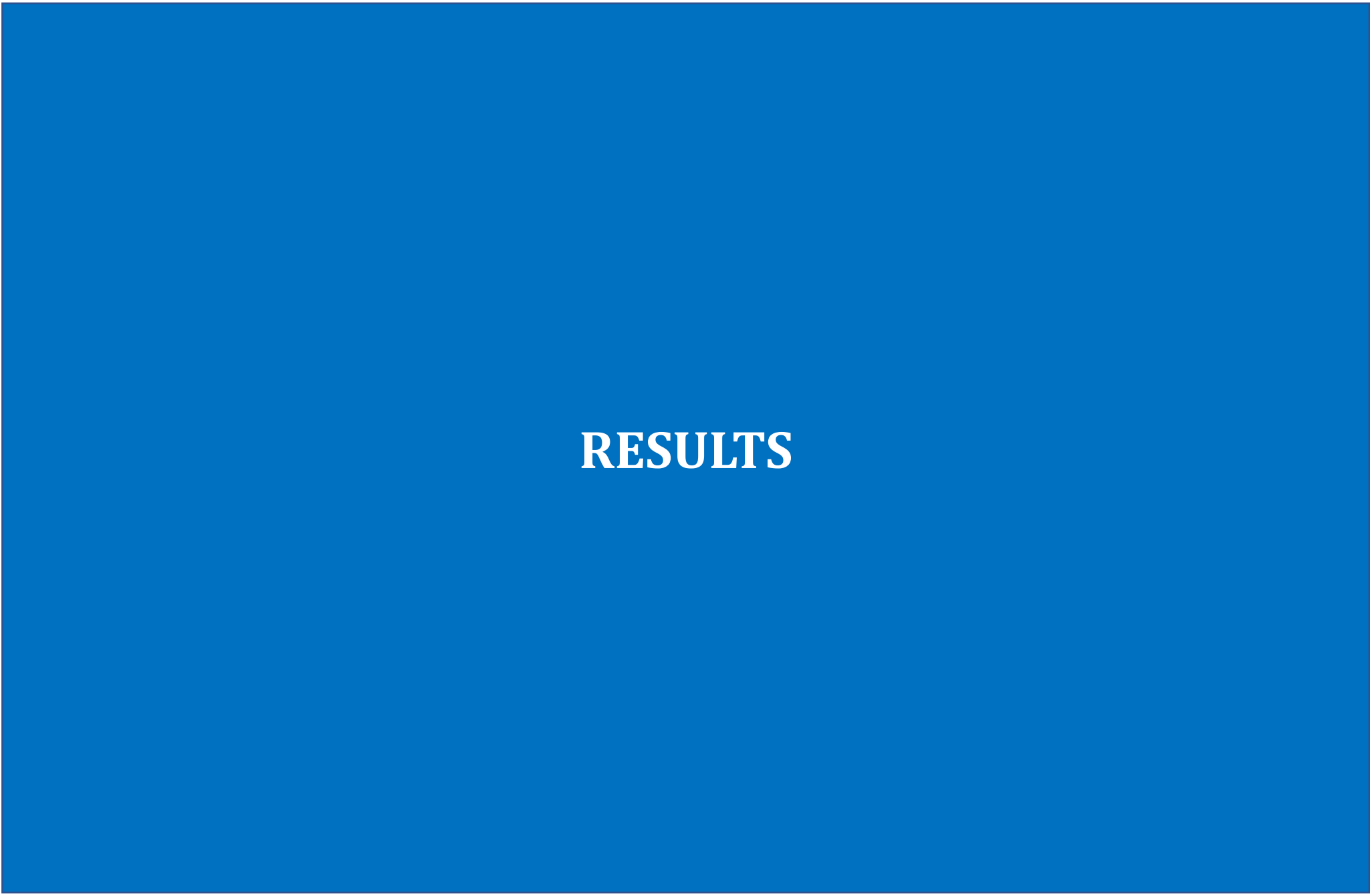
NEURAL NETWORKS	IMBALANCED TRAINING SET			BALANCED TRAINING SET		
	L(1)	L(1,1)	K5	K15	K7	K15
	M	N	O	P	Q	R
	ALL	SET 1	SET 2	ALL	SET 1	SET 2
ACCURACY	0.92	0.92	0.62	0.08	0.08	0.15
SENSITIVITY	0.99	1.00	0.66	0.00	0.00	0.11
SPECIFICITY	0.00	0.00	0.2	1.00	1.00	0.61
BEST MODEL: N: IMBALANCED TRAINING SET (1,1): 2 LAYERS, 1 NODE EACH LAYER SET 1: LOAN, JOY, CLAGE, DEBTINC, NINQ, CLNO						

- DATA
- PREPARATION
- VARIABLES SELECTION
- PARTITION
- CLASSIFICATION
- RESULTS
- SUMMARY
- CONCLUSION

		<div>LOGISTIC REGRESSION</div> <div>IMBALANCED</div>	<div>K-Nearest Neighbors</div> <div>IMBALANCED</div>	<div>NEURAL NETWORKS</div> <div>IMBALANCED</div>
	SIGNIFICANT CLASS: "PAID"			
	CUTOFF: 0.5	<div>B</div> <div>SET 1</div>	<div>K5</div> <div>I</div> <div>SET 2</div>	<div>L(1,1)</div> <div>N</div> <div>SET 1</div>
	ACCURACY	0.92	0.92	0.92
	SENSITIVITY	0.99	0.99	1.00
	SPECIFICITY	0.09	0.10	0.00
	<div>MODEL B: IMBALANCED TRAING SET - SET 1: LOAN, JOY, CLAGE, DEBTINC, NINQ, CLNO</div> <div>MODEL I: IMBALANCED TRAING SET - K5 - SET 2: DEROG, DELINQ, CLAGE, NINQ, DEBTINC</div> <div>MODEL N: IMBALANCED TRAING SET - (1,1): 2 LAYERS, 1 NODE EACH LAYER - SET 1: LOAN, JOY, CLAGE, DEBTINC, NINQ, CLNO</div>			

- **BANKRUPTCY**
- **BOOK CLUB**
- **HOME EQUITY**

- **DATA**
- **PREPARATION**
- **VARIABLES  
SELECTION**
- **PARTITION**
- **CLASSIFICATION**
- **RESULTS**
- **SUMMARY**
- **CONCLUSION**



**RESULTS**



	• RESULTS	• BANKRUPTCY	• BOOK CLUB	• HOME EQUITY
• ALGORITHM		NEURAL NETWORKS	LOGISTICS REGRESSION	K-Nearest Neighbors
• TRAINING SET		FULL ORIGINAL RANDOMIZED	BALANCED TRAINING	IMBALANCED TRAINING
• TARGET		“D” Class: 1 = bankrupt Class : 0 = healthy	“Florence” Class: 1 = Yes to Florence Class: 0 = No to Florence	“BAD” Class: 1 = Defaulted (loan) Class: 0 = Paid (loan)
• PREDICTORS		ALL	SET 3 “Rcode”, “Fcode”, “Mcode”	SET 2: DEROG, DELINQ, CLAGE, NINQ, DEBTINC
• CUTOFF		0.5	0.5	0.5
• OTHERS		LAYER 1; NODE 3		K5

- **BANKRUPTCY**

- **BOOK CLUB**

- **HOME EQUITY**

- **DATA**

- **PREPARATION**

- **VARIABLES  
SELECTION**

- **PARTITION**

- **CLASSIFICATION**

- **RESULTS**

- **SUMMARY**

- **CONCLUSION**

# CONCLUSION & KEY LEARNING

DATA

GOAL

Find best models to classify target variables.

OBSERVATIONS

132

4,000

5,960

VARIABLES

27

19

13

PREDICTORS

24

16

12

TARGET

“D”

“Florence”

“BAD”

Class: 1 = bankrupt  
Class : 0 = healthy

Class: 1 = Yes to Florence  
Class: 0 = No to Florence

Class: 1 = Defaulted (loan)  
Class: 0 = Paid (loan)

TARGET’s CLASS  
DISTRIBUTION  
CHART

