

GROUP PROJECT

BA 706

PROFESSOR: DAVID PARENT

PREPARED BY:

SAMSON OGWUCHE: 301196569

CHARLES NTAMACK: 301209795

Introduction

The aim of the model is to provide a comprehensive analysis of the bank loan campaign and highlight the key variables that determined customers' acceptance or rejection of the loan campaign. The dataset is about the campaign carried out by a bank that was intended to get its customers to subscribe to bank loans. The dataset has 17 variables and 45,211 observations. There are several variables of interest in the dataset. Below is the dictionary to define each of the 17 variables in the dataset.

1. Age: Age in the dataset is a numeric variable that describes the ages of the customers in the loan campaign.
2. Job: This refers to the profession of the bank customers targeted for the loan campaign.
3. Marital status: This refers to the marital status of the customers contacted for the bank loan campaign.
4. Education: Education describes the level of education of the customers contacted.
5. Default: Default describes if a customer has had a credit default previously.
6. Balance: Balance relates to the average yearly euro balances in the customers' accounts.
7. Housing: This refers to whether a customer has a housing loan.
8. Loan: Loan refers if a customer has an existing personal loan.
9. Contact: This refers to the medium of communication with the customers.
10. Day: Day refers to the last contact day of the month.
11. Month: Month refers to the last contact month of the year.
12. Duration: This refers to the last contact duration in seconds.
13. Campaign: Campaign refers to the number of contacts performed during the campaign and with each customer.
14. Pdays: This refers to the number of days that passed after the customer was last contacted.
15. Previous: Previous refers to the number of campaigns performed before the current campaign and with each customer.
16. Poutcome: This refers to the outcome of the previous marketing campaigns performed.
17. Y: This is the binary variable in the dataset that refers to customers' responses which is either a yes or a no. This is our target variable in the model.

Exploring the data source:

We explored the data source to see the properties of the sample. We changed the setting of the sample method from top to random and fetch size from default to maximum. From the exploration, we observed the mean age of the customers who participated in the campaign is 40.95471 years. This implies that the loan campaign was targeted at customers who are in their youthful age.

The screenshot shows the SAS Enterprise Miner interface with a node configuration window open. The window title is "Results - Node StatExplore Diagram: Kc_house". The main pane displays a table titled "Interval Variables" with data for the "TRAIN" role. The table includes columns for Data Role, Target, Target Level, Variable, Median, Missing, Non Missing, Minimum, Maximum, Mean, Standard Deviation, Skewness, Kurtosis, and Role. The data shows various statistics for variables like duration, previous, balance, campaign, day, and age. A tooltip for the "age" variable indicates it has a value of -1. The left sidebar shows the node tree with "Kc_house" selected. The bottom status bar shows the date and time as "IT, 2023-05-10 - ... 12:13 PM" and the system as "Connected to ClassApp-049".

We have a normal range in our dataset except for age that has -1 which has been corrected by setting the age range between 16 and 100.

See the screenshot below for the normal age range of 16 to 100.

Interactive Replacement Interval Filter

Columns:		<input type="checkbox"/> Label	<input type="checkbox"/> Mining	<input type="checkbox"/> Basic	<input type="checkbox"/> Statistics
Name	Use	Limit Method	Replacement Lower Limit	Replacement Upper Limit	Replace Method
age	Default	User Specified	16	100	Missing
balance	Default	Default	.	.	Default
campaign	Default	Default	.	.	Default
day	Default	Default	.	.	Default
duration	Default	Default	.	.	Default
pdays	Default	Default	.	.	Default
previous	Default	Default	.	.	Default

Generate Summary OK Cancel

Explore - EMW56.FIMPORT_DATA

File View Actions Window

Sample Properties

Property	Value
Rows	45211
Columns	18
Library	EMW56
Modeling	FIMPORT_DATA
Type	DATA
Sample Method	Random
Fetch Size	Max
Patched Rows	45211
Random Seed	12345

Sample Statistics

Obs #	Variable ...	Label	Type	Percent Missing	Minimum	Maximum	Mean	Number o...	Mode	Perce...	Mode
1	Y	CLASS	CLASS	0	-1	46211	23906	2	88.30152NO	.	.
2	id	VAR	VAR	0	1001	46211	23906	3	40.95471	.	.
3	age	VAR	VAR	0.019907	-1	999	40.95471	1	.	.	.

EMW56.FIMPORT_DATA

Obs #	Id	age	y
1	1001	999no	
2	1002	44no	
3	1003	33no	
4	1004	47no	
5	1005	33no	
6	1006	35no	
7	1007	28no	
8	1008	2no	
9	1009	58no	
10	1010	42no	
11	1011	41no	
12	1012	29no	
13	1013	53no	
14	1014	58no	
15	1015	57no	
16	1016	51no	
17	1017	45no	
18	1018	51no	
19	1019	60no	
20	1020	33no	
21	1021	28no	
22	1022	58no	
23	1023	32no	

age

Type here to search SAS I.T. 2022-12-06 - 3019... Balance for 301196... Enterprise Miner - ... Explore - EMW56.FL... ENG 10:27 PM US 12/12/2022

Explore - EMW56.FIMPORT_train

File View Actions Window

Sample Statistics

Obs #	Variable	Type	Percent	Minimum	Maximum	Mean	Number o.	Mode Percentage	Mode
1	contact	CLASS	0	.	.	3		64.77406	CELLULA
2	default	CLASS	0	.	.	2		98.19734	NO
3	education	CLASS	0	.	.	4		9.31629	SECONDARY
4	housing	CLASS	0	.	.	2		55.53823	YES
5	job	CLASS	0	.	.	12		21.52573	BLUE-COLLAR
6	loan	CLASS	0	.	.	2		83.97735	NO
7	marital	CLASS	0	.	.	3		66.19639	MARRIED
8	month	CLASS	0	.	.	12		30.44634	MAY
9	poutcome	CLASS	0	.	.	4		81.7478	UNKNOWN
10	r	CLASS	0	.	.	2		88.30152	NO
11	age	VAR	0	1001	46211	23006			
12	age	VAR	0.019907	-1	899	40.95471			
13	balance	VAR	0.006638	-8019	102127	1362.347			
14	campaign	VAR	0	1	63	2.763841			
15	day	VAR	0	1	31	1.9942			
16	duration	VAR	0	0	4918	258.1931			
17	pdays	VAR	0	-1	871	40.19783			
18	previous	VAR	0	0	275	0.580323			



From the pie chart above, the percentage of customers who said yes to the loan campaign is 11.6984% while that of no is 88.3015%.

Replacement:

The dataset has some missing values and to replace them, we did imputation for the missing values. The replacement node was used to perform the task of replacing the missing values by replacing all the lower limit's standard deviations with the upper limit standard deviation values.

Results - Node Replacement (3) Diagram: Bank Loan

File Edit View Window

Total Replacement Counts

Variable	Label	Role	Train
age	age	INPUT	265
balance	balance	INPUT	745
campaign	campaign	INPUT	840
day	day	INPUT	840
duration	duration	INPUT	963
previous	previous	INPUT	582

Output

```

1 *-----
2 User: 301196569
3 Date: December 11, 2022
4 Time: 23:19:37
5 *-----
6 * Training Output
7 *-----
```

Variable Summary

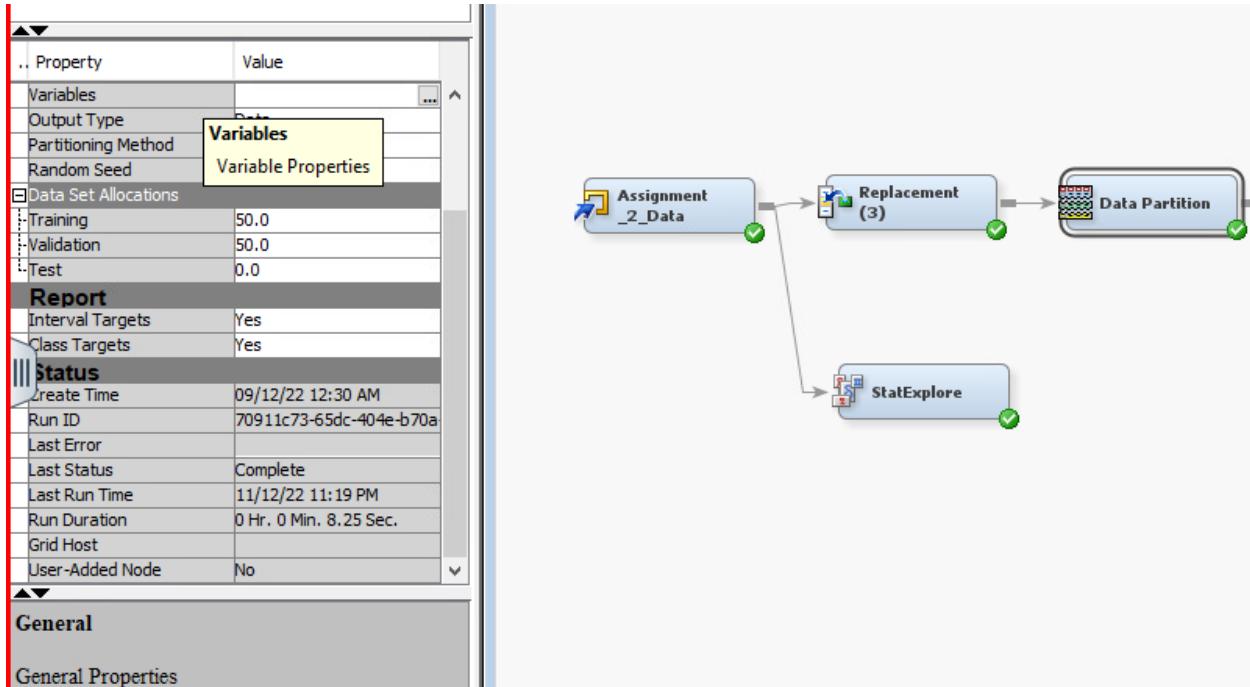
Role	Measurement Level	Frequency Count
INPUT	INTERVAL	6
INPUT	BINARY	9
REJECTED	INTERVAL	1
TARGET	BINARY	1

Interval Variables

Variable	Replace Variable	Lower limit	Upper Limit	Label	Limits Method	Replacement Method	Lower Replacement Value	Upper Replacement Value
age	REP_age	6.337283	75.57215	age	STDDEV	COMPUTED	6.337283	75.57215
balance	REP_balance	-7772.21	10486.9	balance	STDDEV	COMPUTED	-7772.21	10486.9
campaign	REP_campaign	-4.1625	20.2579	campaign	STDDEV	COMPUTED	-4.1625	20.2579
day	REP_day	-3.16101	40.7735	day	STDDEV	COMPUTED	-3.16101	40.7735
duration	REP_duration	-514.42	1030.747	duration	STDDEV	COMPUTED	-514.42	1030.747
previous	REP_previous	-6.33	7.490647	previous	STDDEV	COMPUTED	-6.33	7.490647

Data Participation:

We partitioned the data into two sets – training and validation. The training set is 50% of the dataset while validation is also 50% of the dataset.



StatExplore.

We used the StatExplore update to reject two of the variables that are not relevant to the model and also used it to explore the skewness of the variables. The variables we rejected are day and pdays. The variable ‘day’ in the dataset represents the last contact day of the month while the variable ‘plays’ represents the last day that passed after the customer was contacted from the previous campaign. The two variables were dropped because the dataset was skewed with the two variables, day and pday. The result was not good because of the two variables in the dataset; hence we dropped them.

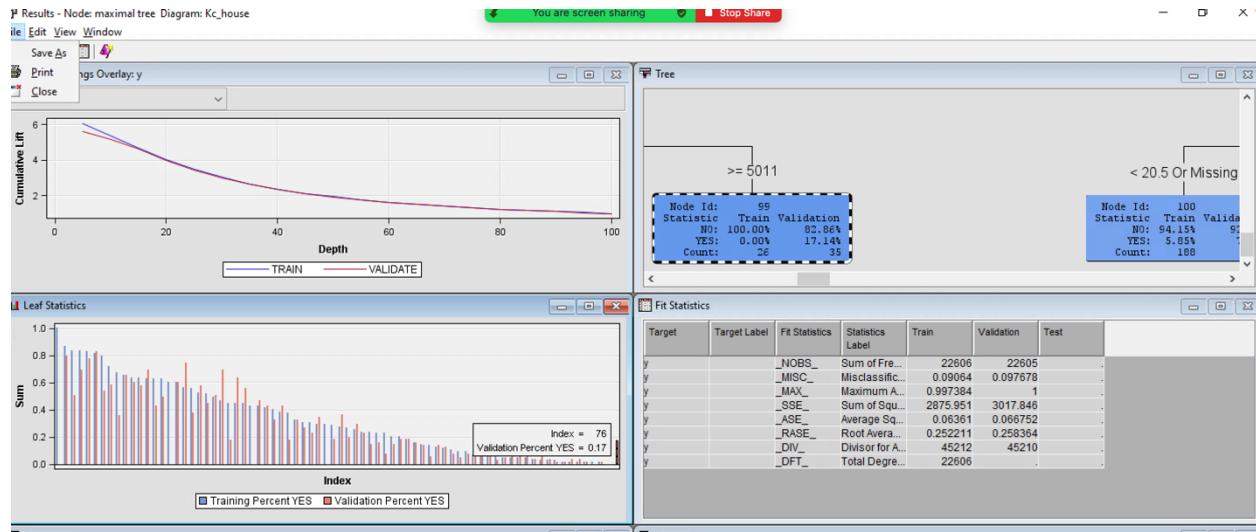
The screenshot shows the 'Variables - Stat' dialog box. At the top, there is a search bar with '(none)' and a dropdown menu. Below it are several filter buttons: 'not', 'Equal to', and a date range selector. To the right are 'Apply' and 'Reset' buttons. Underneath these are three checkboxes: 'Label', 'Mining', and 'Basic'. A 'Statistics' checkbox is located further down on the right. The main area is a table with columns: Name, Use, Report, Role, and Level. The table lists 20 variables:

Name	Use	Report	Role	Level
Id	Default	No	ID	Nominal
age	Default	No	Input	Interval
balance	Default	No	Input	Interval
campaign	Default	No	Input	Interval
contact	Default	No	Input	Nominal
day	Default	No	Rejected	Interval
default	Default	No	Input	Nominal
duration	Default	No	Input	Interval
education	Default	No	Input	Nominal
housing	Default	No	Input	Nominal
job	Default	No	Input	Nominal
loan	Default	No	Input	Nominal
marital	Default	No	Input	Nominal
month	Default	No	Input	Nominal
pdays	Default	No	Rejected	Interval
poutcome	Default	No	Input	Nominal
previous	Default	No	Input	Interval
y	Default	No	Target	Binary

At the bottom of the dialog box are buttons for 'Explore...', 'Update Path', 'OK', and 'Cancel'.

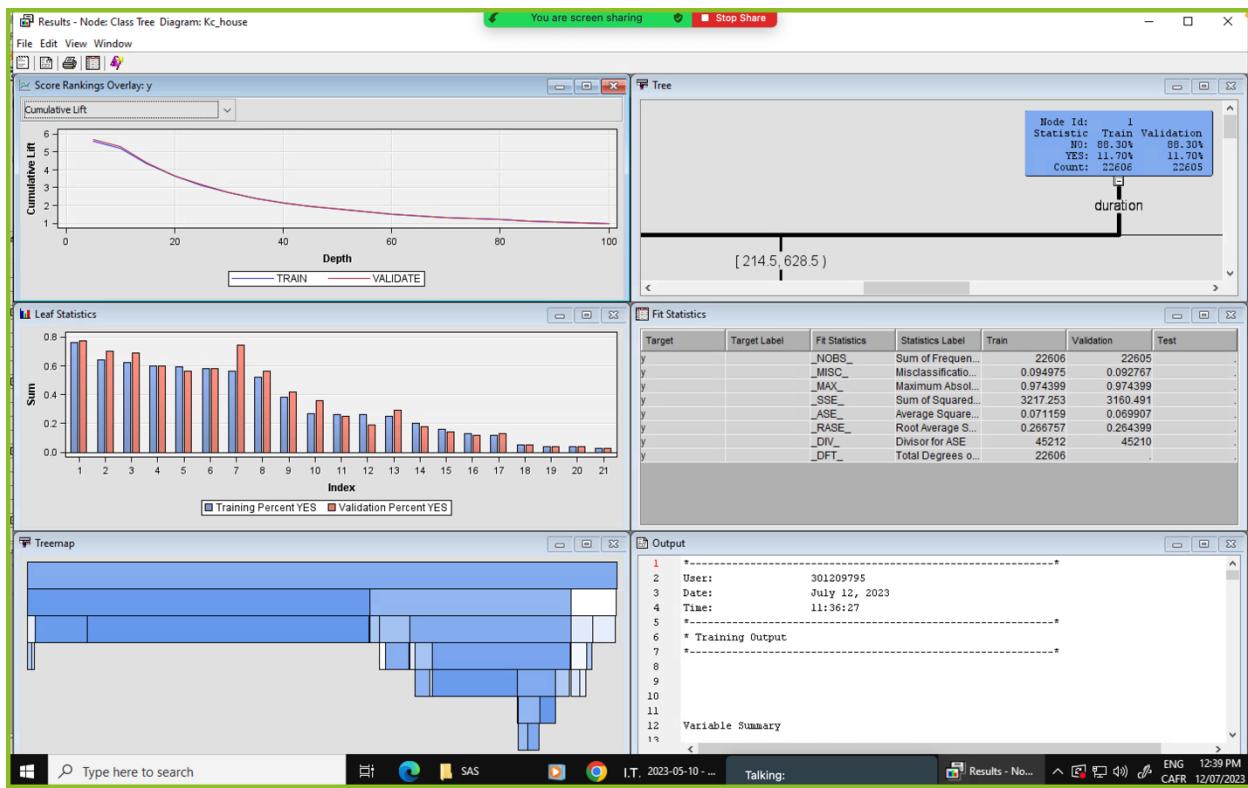
Decision Tree statistics

We built three decision trees using the average squared error (ASE) and misclassification as the assessment measures. Our first decision tree is the Maximum decision tree and we used ‘largest’ as the subtree method with maximum branch of 3. The Maximum decision tree has 76 leaf statistics and an ASE of **0.066752**. However, our target was to get an optimal ASE and so we have to run the optimal decision tree.



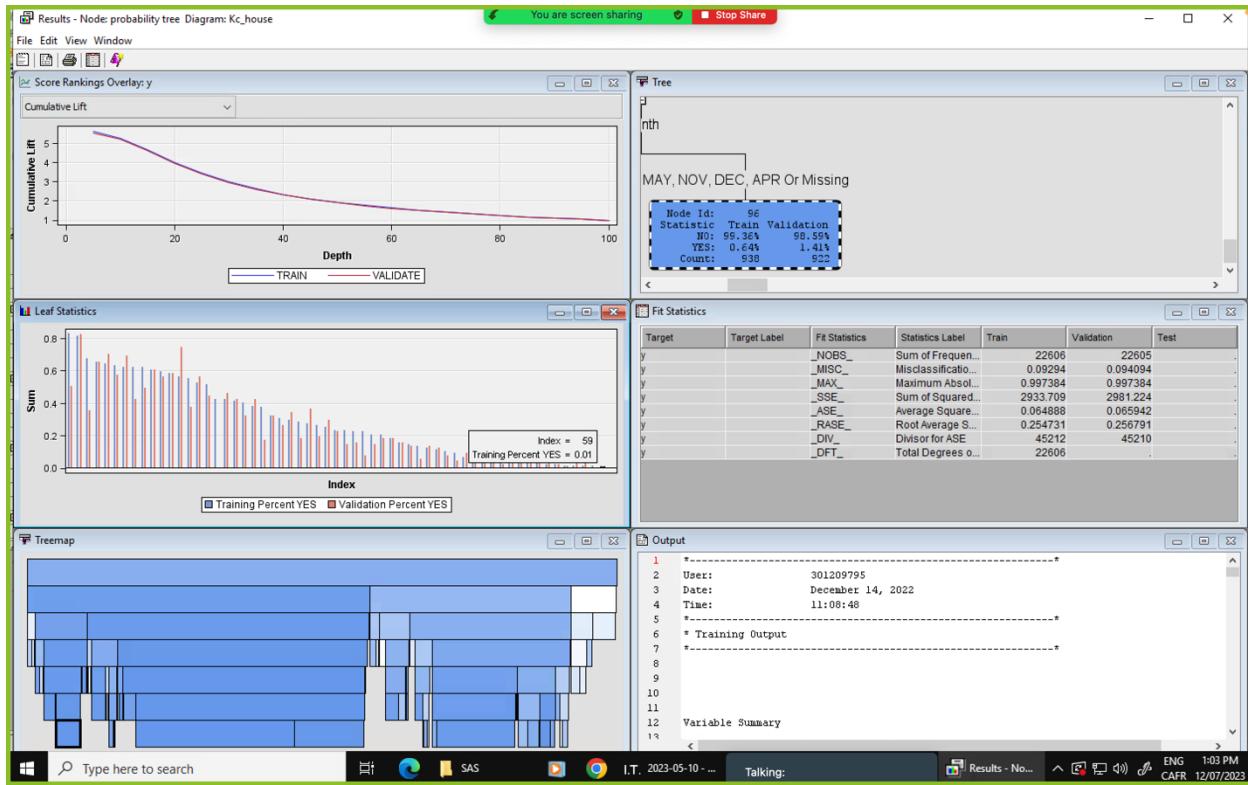
Class Decision Tree.

Our second decision tree is the Class decision tree using misclassification as an assessment measure. The subtree method used for Class decision tree is assessment. The maximum number of branches is three (3). Our second decision tree has 21 leaf statistics and an ASE of **0.069907**



Probability Decision Tree.

Our third decision tree is the Probability decision tree using average squared error as an assessment measure. The subtree method used for Probability decision tree is assessment. The maximum number of branches is three (3). Our third decision tree has 59 leaf statistics and an ASE of **0.065942**



Our key variables under the Probability decision tree are:

```

61
62      Variable Importance
63
64
65
66      Variable          Number of          Ratio of
67      Name       Label     Splitting    Validation
68
69      duration           9            1.0000        1.0000
70      poutcome          2            0.6905        0.7440
71      month             10           0.6041        0.6049
72      day               9            0.3903        0.3034
73      contact            3            0.2100        0.1907
74      REP_age   Replacement: age  3            0.1961        0.0827
75      housing            4            0.1891        0.1640
76
77
78
79      Tree Leaf Report
80
81
82      Node          Training          Validation
83      Id           Depth   Observations   Percent   Percent
84
85      106          6           6116        0.00      6108      0.00
86      118          6           3010        0.03      3061      0.03
87      107          6           2679        0.01      2756      0.01
88      <=            <           <=>        < .01>    < .01>    < .01>

```

Decision Tree Model Conclusion:

Using average squared error (ASE) as our model assessment measures, our best decision tree is the Probability decision tree with an ASE of **0.065942**. This is the best because it has the lowest ASE.

Imputation:

Under imputation, we changed the indicator variable type to unique and the role to input. See the screenshot below. This step was performed to impute the missing value(s) in preparation for regression and neural networks. The only missing value for age was replaced with the mean of the non-missing variables.

Enterprise Miner - BA706 CLASS WORK1

File Edit View Actions Options Window Help

BA706 CLASS WORK1

Data Sources

Diagrams

Bank Loan

Organic2

Organics

PREDICTIVE ANALYTICS

Semester Project

Supermarket_sales

Model Packages

Property Value

Tuning Parameters

Tree Imputation

Score

Hide Original Variables Yes

Indicator Variables

Type Unique

Source Imputed Variables

Role Input

Report

Validation and Test Data No

Distribution of Missing No

Status

Create Time 09/12/22 2:10 AM

Run ID 7dde0cd3-cfcc-478b-b53d-

Last Error

Last Status Complete

Last Run Time 11/12/22 11:22 PM

Run Duration 0 Hr, 0 Min, 7.50 Sec.

Grid Host

User-Added Node No

General

General Properties

.. Property Value

Tuning Parameters

Tree Imputation

Score

Hide Original Variables Yes

Indicator Variables

Type Unique

Source Imputed Variables

Role Input

Report

Validation and Test Data No

Distribution of Missing No

Status

Create Time 09/12/22 2:10 AM

Run ID 7dde0cd3-cfcc-478b-b53d-

Last Error

Last Status Complete

Last Run Time 11/12/22 11:22 PM

Run Duration 0 Hr, 0 Min, 7.50 Sec.

Grid Host

User-Added Node No

General

General Properties

We then imputed the missing value for the age to **40.9**. See the screenshot below.

Results - Node: Impute Diagram: Bank Loan

File Edit View Window

Imputation Summary

Variable Name	Impute Method	Imputed Variable	Indicator Variable	Impute Value	Role	Measurement Level	Label	Number of Missing for TRAIN
REP_Age	MEAN	IMP REP_Age	M REP_Age	40.92738 INPUT		INTERVAL	Replacement age	0

Output

```

1 *-----
2 User: 301196569
3 Date: December 14, 2022
4 Time: 10:26:43
5 *-----
6 * Training Output
7 *-----
8
9
10
11
12 Variable Summary
13
14      Measurement Frequency
15      Role   Level   Count
16
17  INPUT    INTERVAL   6
18  INPUT    NOMINAL   9
19  REJECTED INTERVAL   2
20  TARGET   BINARY    1
21
22 *-----
23 *-----
```

Type here to search SAS IT_ 2022-12-06 - 30119... Balance for 301196... Enterprise Miner - ... Results - Node Im... ENG 10:13 PM US 14/12/2022

Cap and Floor:

In order to fix our skewed variables, we introduced cap and floor. However, we noticed that even after running the Cap and Floor, we still have some skewed variables. We, therefore, introduced StatExplore to ascertain the skewed variables.

Results - Node: Cap & Floor Diagram: Bank Loan

File Edit View Window

Total Replacement Counts

Variable	Label	Role	Train	Validation
IMP REP_Age	Imputed Replacement age	INPUT	192	192
balance	balance	INPUT	387	356
campaign	campaign	INPUT	408	432
day	day	INPUT	0	0
duration	duration	INPUT	489	493
previous	previous	INPUT	397	390

Output

```

1 *-----
2 User: 301196569
3 Date: December 14, 2022
4 Time: 10:34:33
5 *-----
6 * Training Output
7 *-----
8
9
10
11
12 Variable Summary
13
14      Measurement Frequency
15      Role   Level   Count
16
17  INPUT    BINARY    1
18  INPUT    INTERVAL   6
19  INPUT    NOMINAL   9
20  REJECTED INTERVAL   2
21  TARGET   BINARY    1
22
23 *-----
```

Interval Variables

Variable	Replace Variable	Lower limit	Upper limit	Label	Limit Method	Replace Method	Lower Replacement Value	Upper Replacement Value
IMP REP_Age	REP IMP REP_Age	9.15084	72.739	Imputed Replacement age	STDDEV	COMPUTED	9.15084	72.739
balance	REP balance	-7996.45	10760.51	balance	STDDEV	COMPUTED	-7996.45	10760.51
campaign	REP campaign	-6.54091	12.06988	campaign	STDDEV	COMPUTED	-6.54091	12.06988
day	REP day	-9.19575	40.83467	day	STDDEV	COMPUTED	-9.19575	40.83467
duration	REP duration	-507.407	1024.442	duration	STDDEV	COMPUTED	-507.407	1024.442
previous	REP previous	-5.10793	6.33935	previous	STDDEV	COMPUTED	-5.10793	6.33935

20 of 24 - Clipboard
Item not Collected: Format not supported by Office Clipboard

StatExplore:

After Cap and Floor, we reintroduced StatExplore to check if the skewed variables have been fixed. However, the result as can be seen below shows that some of the variables are still skewed.

Data Role	Target	Target Level	Variable	Median	Missing	Non Missing	Minimum	Maximum	Mean	Standard Deviation	Skewness	Kurtosis	Role	Label	Scaled Mean Deviation	Maximum Deviation	Level Id
TRAIN	y	no	REP_prevlo	0	0	19981	0	6.339353	0.434539	1.216911	3.312005	11.00323INPUT	Replace...	-0.13305	1.004068	1	
TRAIN	y	yes	REP_prevlo	0	0	2645	0	6.339353	1.004492	1.740355	1.832598	2.367785INPUT	Replace...	1.004068	1.004068	2	
TRAIN	y	no	REP_durall	164	0	19981	0	1024.442	218.1987	184.9291	1.922713	1.444605INPUT	Replace...	-0.13194	0.995731	1	
TRAIN	y	yes	REP_durall	42	0	2645	0	1024.442	218.1987	184.9291	1.922713	1.444605INPUT	Replace...	0.995731	0.995731	2	
TRAIN	y	no	REP_balan	428	0	19981	-4057	10760.51	5035.209	2077.153	2.731474	8.083561INPUT	Replace...	-0.04285	0.323386	1	
TRAIN	y	yes	REP_balan	733	0	2645	-3058	10760.51	1666.363	2354.205	2.238078	5.113574INPUT	Replace...	0.323386	0.323386	2	
TRAIN	y	no	REP_camp	2	0	19981	1	12.09886	2.704855	2.386508	2.202775	5.111164INPUT	Replace...	0.233386	0.233386	1	
TRAIN	y	yes	REP_camp	2	0	2645	1	12.09886	2.704855	2.386508	2.202775	5.111164INPUT	Replace...	0.233386	0.233386	1	
TRAIN	y	no	REP_camp_2	2	0	2645	1	12.09886	2.704855	2.386508	2.202775	5.111164INPUT	Replace...	-0.205049	0.205049	1	
TRAIN	y	yes	REP_camp_2	2	0	2645	1	12.09886	2.704855	2.386508	2.202775	5.111164INPUT	Replace...	-0.205049	0.205049	1	
TRAIN	y	no	REP_day	16	0	19981	1	31	15.89454	8.39745	0.084814	-1.05458INPUT	Replace...	0.004889	0.035389	1	
TRAIN	y	yes	REP_day	15	0	2645	1	31	15.26049	8.534914	0.160508	-1.07217INPUT	Replace...	-0.03539	0.035389	1	
TRAIN	y	no	REP_IMP	39	0	19981	18	72.7039	40.8144	10.04689	0.495013	-0.38607INPUT	Replace...	-0.00168	0.012496	1	
TRAIN	y	yes	REP_IMP	38	0	2645	18	72.7039	41.39296	12.99808	0.686814	-0.34782INPUT	Replace...	0.012496	0.012496	2	

Log_Transform Variables:

We further introduced Log_Transform Variables to fix variables that are skewed after Cap and Floor. However, we still have four of the variables that are still skewed after Log_Transform Variables as can be seen below in the screenshot. We also tried the standardized transform variables method to fix the skewed variables. However, the result got worst and so, we have to keep Log_Transform Variables.

See below the results of standardized transform variables.

Results - Node: Log_Transform Variables Diagram: Bank Loan

File Edit View Window

Transformations Statistics

Source	Method	Variable Name	Formula	Number of Levels	Non Missing	Missing	Minimum	Maximum	Mean	Standard Deviation	Skewness	Kurtosis	Label	
Input	Original	REP_balance		-	22699	0	-4057	10730.51	1259.195	2116.589	2.668792	7.687442	Replacement balance	
Input	Original	REP_campaign		-	22698	0	1	12.05988	2.633396	2.305485	2.270795	5.561302	Replacement campai	
Input	Original	REP_day		-	22698	0	1	31	15.82036	8.338703	0.09288	-1.06631	Replacement day	
Input	Original	REP_previous		-	22698	0	0	6.339353	0.501228	1.302088	3.029753	8.933759	Replacement previous	
Input	Original	baseage		-	22698	0	-4957	102477	1092.293	8377.599	8.159275	155.3927	Replacement baseage	
Input	Original	campaign		-	22698	0	1	63	2.764497	3.101798	5.05654	42.79048	Replacement campaign	
Input	Original	day		-	22698	0	1	31	15.82036	8.338703	0.09288	-1.06631	Replacement day	
Input	Original	duration		-	22698	0	0	4918	258.5178	255.3069	3.014064	17.08493	Replacement duration	
Input	Original	pdays		-	22698	0	-1	971	48.53917	2.031213	2.012126	0.093203	17.08493	Replacement pdays
Input	Original	previous		-	22698	0	0	58	0.585982	1.91793	7.500899	101.4746	Replacement previous	
Output	Computed	LOG REP_balance	log(REP_balance + 4...)	-	22698	0	0	9.603632	8.524738	0.905463	0.64604	29.3886	Transformed balance	
Output	Computed	LOG REP_campaign	log(REP_campaign + ...)	-	22698	0	0	0.693147	2.57031	1.15253	0.485773	1.07898	0.439898	
Output	Computed	LOG REP_day	log(REP_day + 1)	-	22698	0	0	3.44869	2.688885	2.089544	0.89777	0.004747	0.004747	
Output	Computed	LOG REP_previous	log(REP_previous + 1)	-	22698	0	0	1.993251	0.222907	0.507023	2.187237	3.574311	Transformed Replace	
Output	Computed	LOG_balance	log(balance + 4058)	-	22698	0	0	11.57294	8.530209	0.330039	1.443702	27.0098	Transformed balance	
Output	Computed	LOG_campaign	log(REP_campaign + 1)	-	22698	0	0	0.693147	4.6307	1.00000	1.37407	2.00000	0.47514	
Output	Computed	LOG_day	log(day + 1)	-	22698	0	0	0.693147	3.495736	2.858885	0.829164	-0.8777	0.043881	
Output	Computed	LOG_duration	log(duration + 1)	-	22698	0	0	8.500891	5.174497	0.920817	-0.45131	0.853881	Transformed duration	
Output	Computed	LOG_pdays	log(pdays + 2)	-	22698	0	0	6.771936	0.987707	2.057108	1.593326	0.948182	Transformed pdays	
Output	Computed	LOG_previous	log(previous + 1)	-	22698	0	0	4.077537	0.230141	0.537759	2.469031	5.823371	Transformed previous	

Output

```

1 -----
2 User:          301196569
3 Date:        December 14, 2022
4 Time:        12:10:16
5 -----
6 * Training Output
7 -----
8
9
10
11
12 Variable Summary
13
14      Measurement Frequency
15      Role    Level   Count
16
17 INPUT    BINARY      1
18 INPUT    INTERVAL     6
19 INPUT    NOMINAL     9
20 REJECTED INTERVAL     8
21 TARGET   BINARY      1
22
23

```

Type here to search

File Edit View Window

Results - Node: Log...

File Edit View Window

Interval Variables

Data Role	Target	Target Level	Variable	Median	Missing	Non Missing	Minimum	Maximum	Mean	Standard Deviation	Skewness	Kurtosis	Role	Label	Scaled Mean	Maximum Deviation	Level Id	
TRAIN	y	no	REP_duration	164	0	19981	0	1024.462	248.937	184.9231	1.022713	-0.13194	-0.95713	1	Replacement: d.	-0.13194	-0.95713	1
TRAIN	y	yes	REP_campaign	430	0	199846	8	1024.462	591.858	296.0095	0.486852	1.038269	INPUT	Replacement: d.	0.98971	0.98971	2	
TRAIN	y	yes	LOG REP_previous	0	0	19981	0	1.993251	0.194117	0.477368	2.415317	4.714008	INPUT	Transformed: d.	0.12798	0.985548	1	
TRAIN	y	yes	LOG REP_previous	0	0	2545	0	1.993251	0.437612	0.520504	1.143841	-0.14539	INPUT	Transformed: R.	0.985848	0.985848	2	
TRAIN	y	no	LOG REP_campaign	1.089112	0	19981	0.693147	2.57031	1.1688	0.492956	0.981014	0.334991	INPUT	Transformed: R.	0.019344	0.105228	1	
TRAIN	y	yes	LOG REP_campaign	1.089112	0	2545	0.693147	2.57031	1.1688	0.492956	0.981014	0.334991	INPUT	Transformed: R.	-0.0522	0.105228	2	
TRAIN	y	no	LOG REP_balance	2.833213	0	19981	0.693147	3.465736	2.666903	0.228887	-0.86602	0.042253	INPUT	Transformed: R.	0.002753	0.020774	1	
TRAIN	y	yes	LOG REP_balance	2.772589	0	2545	0.693147	3.465736	2.603454	0.672209	-0.84626	-0.03041	INPUT	Transformed: R.	-0.02077	0.020774	2	
TRAIN	y	no	REP_IMP REP_dpd	39	0	19981	18	72.7029	40.8144	10.04689	0.496503	-0.38607	INPUT	Replacement: I.	-0.00168	0.012499	1	
TRAIN	y	yes	REP_IMP REP_dpd	39	0	19981	18	72.7029	40.8144	10.04689	0.496503	-0.38607	INPUT	Replacement: I.	0.01049	0.01049	2	
TRAIN	y	no	LOG REP_balance	8.408717	0	19981	0	6.603632	8.515853	0.30181	0.528255	34.47457	INPUT	Transformed: R.	0.00104	0.007865	1	
TRAIN	y	yes	LOG REP_balance	8.474494	0	2545	6.907755	6.603632	8.581789	0.324031	1.333237	1.746478	INPUT	Transformed: R.	0.007865	0.007865	2	

Type here to search

File Edit View Window

Results - Node: Stat...

File Edit View Window

File Edit View Window

Results of Standardized Transform Variables.

Results - Node: Log_Transform Variables Diagram: Bank Loan

File Edit View Window

Transformations Statistics

Source	Method	Variable Name	Formula	Number of Levels	Non Missing	Missing	Minimum	Maximum	Mean	Standard Deviation	Skewness	Kurtosis	Label
Input	Original	REP_balance		22606	0	-4057	10760.51	1201.166	2116.589	2.69792	7.697442Replacement balance		
Input	Original	REP_campaign		22606	0	1	12.06989	2.31336	2.317485	2.270795	5.561302Replacement campai...		
Input	Original	REP_day		22606	0	1	31	15.82036	3.38703	0.09289	1.06511Replacement day		
Input	Original	REP_previous		22606	0	0	6.339353	0.501228	1.302088	3.029753	8.83759Replacement previous		
Input	Original	balance		22606	0	-4057	102127	1382.029	3128.159	8.970275	158.207		
Input	Original	campaign		22606	0	1	63	42.00007	1.00000	0.00000	42.00000		
Input	Original	day		22606	0	1	31	15.82036	3.38703	0.09289	1.06511		
Input	Original	duration		22606	0	0	4918	258.5178	255.3082	3.014064	17.05493		
Input	Original	pdays		22606	0	-1	871	40.53017	100.413	2.612216	6.993526		
Input	Original	previous		22606	0	0	68	0.98662	1.91783	7.500899	101.4748		
Output	Computed	STD_REP_balance	(REP_balance - 1259.	22606	0	-2.51167	4.498987	-5.7E-12	1	2.658792	7.697442Transformed Replace...		
Output	Computed	STD_REP_campaign	(REP_campaign - 2.6.	22606	0	-0.70844	4.093098	-3.6E-12	1	2.270795	5.561302Transformed Replace...		
Output	Computed	STD_REP_day	(REP_day - 15.820357.	22606	0	-1.7773	1.820384	2.78E-11	1	0.09288	1.06511Transformed Replace...		
Output	Computed	STD_balance	(balance - 102127.01)	22606	0	-40380.0	4.498987	-3.4E-11	1	3.029753	8.83759Transformed previous		
Output	Computed	STD_balance	(balance - 1382.0288)	22606	0	-1.73984	32.22644	-1.3E-11	1	8.970275	158.207Transformed balance		
Output	Computed	STD_campaign	(campaign - 2764487.	22606	0	-0.56888	19.41955	-1.4E-11	1	5.06544	42.79048Transformed campaign		
Output	Computed	STD_duration	(duration - 4918.)	22606	0	-1.7773	1.820384	2.78E-11	1	0.09288	1.06511Transformed duration		
Output	Computed	STD_pdays	(pdays - 40.5301688)	22606	0	-0.41359	8.270543	-3.2E-12	1	3.014064	17.05493Transformed pdays		
Output	Computed	STD_previous	(previous - 0.588821.	22606	0	-0.30548	29.93703	-6E-12	1	7.500899	101.4748Transformed previous		

Output

```

1 *-----*
2 User: 301196559
3 Date: December 14, 2022
4 Time: 22:39:08
5 *-----*
6 * Training Output*
7 *-----*
8
9
10
11
12 Variable Summary
13
14      Measurement   Frequency
15      Role       Level   Count
16
17 INPUT    BINARY      1
18 INPUT    INTERVAL     6
19 INPUT    NORMAL      9
20 PREDICTED INTERVAL     6
21 TARGET   BINARY      1
22
23

```

Results - Node: StatExplore (3) Diagram: Bank Loan

File Edit View Window

Interval Variables

Data Role	Target	Target Level	Variable	Median	Missing	Non Missing	Minimum	Maximum	Mean	Standard Deviation	Skewness	Kurtosis	Role	Label	Scaled Mean Deviation	Maximum Deviation	Level Id
TRAIN	y	no	STD REP_campaign	-0.2747	0	19951	-0.70844	4.093098	0.031034	1.026468	2.202775	5.111164INPUT	Transformed R..	-8.639E-9	6.519E10	1	
TRAIN	y	yes	STD REP_campaign	-0.2747	0	2645	-0.70844	4.093098	0.031034	0.730295	2.765691	10.28157INPUT	Transformed R..	6.519E10	6.519E10	2	
TRAIN	y	no	STD REP_balance	-0.39269	0	19951	-2.51167	4.488995	0.031034	0.915988	2.731174	8.00011INPUT	Transformed R..	4.494E10	3.297E10	1	
TRAIN	y	yes	STD REP_balance	-0.39269	0	2645	-2.51167	4.488995	0.031034	0.915988	2.731174	5.113574INPUT	Transformed R..	4.38E10	3.357E10	2	
TRAIN	y	no	STD REP_previous	-0.38494	0	19951	-0.38494	4.483973	-0.05122	0.934589	3.312905	11.00325INPUT	Transformed R..	1.5013E9	1.133E10	1	
TRAIN	y	yes	STD REP_previous	-0.38494	0	2645	-0.38494	4.483973	0.05122	0.934589	3.312905	2.367785INPUT	Transformed R..	-1.13E10	1.133E10	2	
TRAIN	y	no	STD REP_day	0.021543	0	19951	-1.7773	4.820344	0.060746	0.982297	0.984814	-1.074565INPUT	Transformed R..	3.209E9	2.4164E9	1	
TRAIN	y	yes	STD REP_day	0.021543	0	2645	-1.7773	4.820344	0.060746	0.982297	0.984814	1.074565INPUT	Transformed R..	-2.4164E9	2.4164E9	2	
TRAIN	y	no	REP_duration	164	0	19951	0	1024.442	218.1987	184.9291	1.922713	4.444990INPUT	Replacement d..	-0.13194	0.995731	1	
TRAIN	y	yes	REP_duration	430	0	2645	8	1024.442	501.6558	296.8005	0.488356	-1.03828INPUT	Replacement d..	0.995731	0.995731	2	
TRAIN	y	no	REP_IMP REP_age	39	0	19951	18	72.7039	40.8144	10.04889	0.495013	-0.38607INPUT	Replacement I..	-0.00165	0.012498	1	
TRAIN	y	yes	REP_IMP REP_age	38	0	2645	18	72.7039	41.39298	12.99808	0.086814	-0.34782INPUT	Replacement I..	0.012498	0.012498	2	

Type here to search    

^      EN5 10:42 PM
US 14/12/2022

Dummies Recode:

We introduced dummies recode in our model to avoid overfitting. However, the results of our models' assessment measures became worst with dummies recode (replacement values). The results below show the effect of dummies recodes on the models.

We will retain the dummies recode node, however, we will not replace the dummy variables because our assessment measures will be optimized without dummies recode.

Replacement Editor-WORK.OUTCLASS

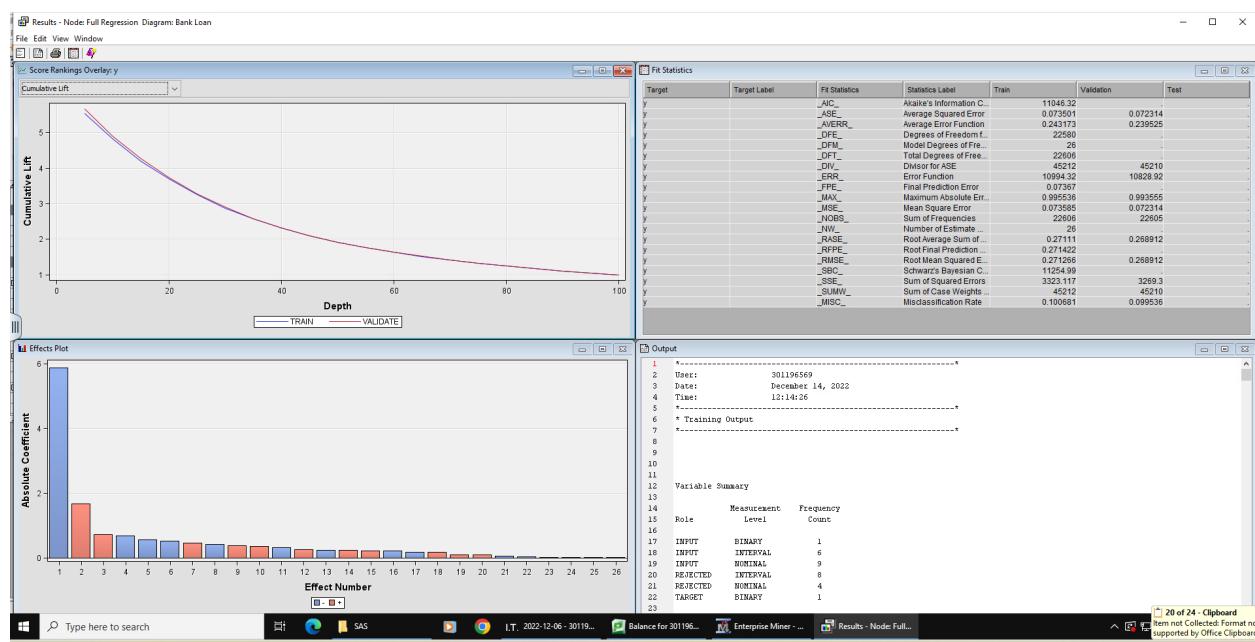
Variable	Formatted Value	Replacement Value	Frequency Count	Type	Character Unformatted Value	Numeric Value
1 REP_age	0		22597N			0
1 REP_age	1		9N			1
1 REP_age	_UNKNOWN_	_DEFAULT_	N		.	.
contact	cellula	phone	14752C	cellula	.	.
contact	unknown	unknown	6366C	unknown	.	.
contact	telepho	phone	1488C	telepho	.	.
contact	_UNKNOWN_	_DEFAULT_	C		.	.
default	no		22177C	no	.	.
default	yes		429C	yes	.	.
default	_UNKNOWN_	_DEFAULT_	C		.	.
education	secondary	learning	11528C	secondary	.	.
education	tertiary	learning	6693C	tertiary	.	.
education	primary	learning	3462C	primary	.	.
education	unknown		923C	unknown	.	.
education	_UNKNOWN_	_DEFAULT_	C		.	.
housing	yes		12619C	yes	.	.
housing	no		9987C	no	.	.
housing	_UNKNOWN_	_DEFAULT_	C		.	.
ob	blue-collar	trades	4957C	blue-collar	.	.
ob	management	office	4762C	management	.	.
ob	technician	trades	3767C	technician	.	.
ob	admin.	office	2621C	admin.	.	.
ob	services	office	1988C	services	.	.
ob	retired	retired	1087C	retired	.	.
ob	self-employed	self	811C	self-employed	.	.
ob	entrepreneur	self	721C	entrepreneur	.	.
ob	unemployed	nomoney	675C	unemployed	.	.
ob	housemaid	trades	604C	housemaid	.	.
ob	student	nomoney	467C	student	.	.
ob	unknown	unknown	146C	unknown	.	.
ob	_UNKNOWN_	_DEFAULT_	C		.	.
pan	no		19029C	no	.	.
pan	yes		3577C	yes	.	.
pan	_UNKNOWN_	_DEFAULT_	C		.	.
marital	married		13618C	married	.	.
marital	single		6409C	single	.	.

OK Cancel

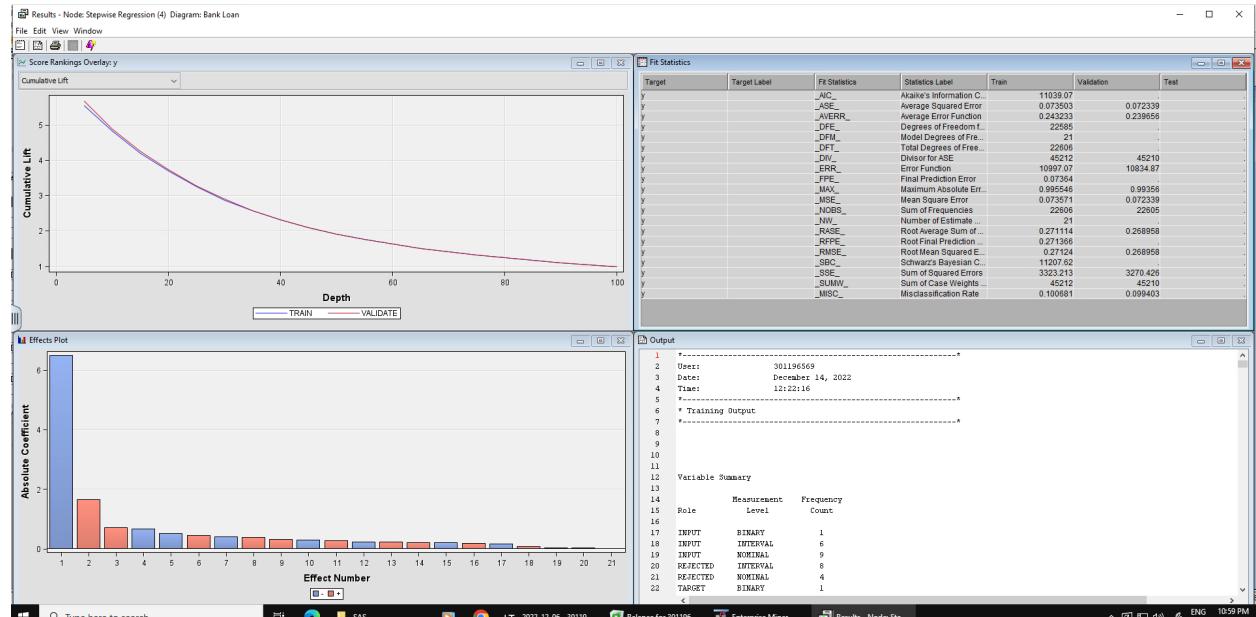
n	yes		3577C	yes	.
n	_UNKNOWN_	_DEFAULT_	C	.	.
rital	married		13618C	married	.
rital	single		6409C	single	.
rital	divorced		2579C	divorced	.
rital	_UNKNOWN_	_DEFAULT_	C	.	.
nth	may	q2	6918C	may	.
nth	jul	q3	3452C	jul	.
nth	aug	q3	3171C	aug	.
nth	jun	q2	2591C	jun	.
nth	nov	q4	1928C	nov	.
nth	apr	q2	1460C	apr	.
nth	feb	q1	1357C	feb	.
nth	jan	q1	718C	jan	.
nth	oct	q4	372C	oct	.
nth	sep	q3	294C	sep	.
nth	mar	q1	237C	mar	.
nth	dec	q4	108C	dec	.
nth	_UNKNOWN_	_DEFAULT_	C	.	.
utcome	unknown		18440C	unknown	.
utcome	failure		2463C	failure	.
utcome	other		948C	other	.
utcome	success		755C	success	.
utcome	_UNKNOWN_	_DEFAULT_	C	.	.
	no		19961C	no	.
	yes		2645C	yes	.
	UNKNOWN	_DEFAULT_	C	.	.

OK Cancel

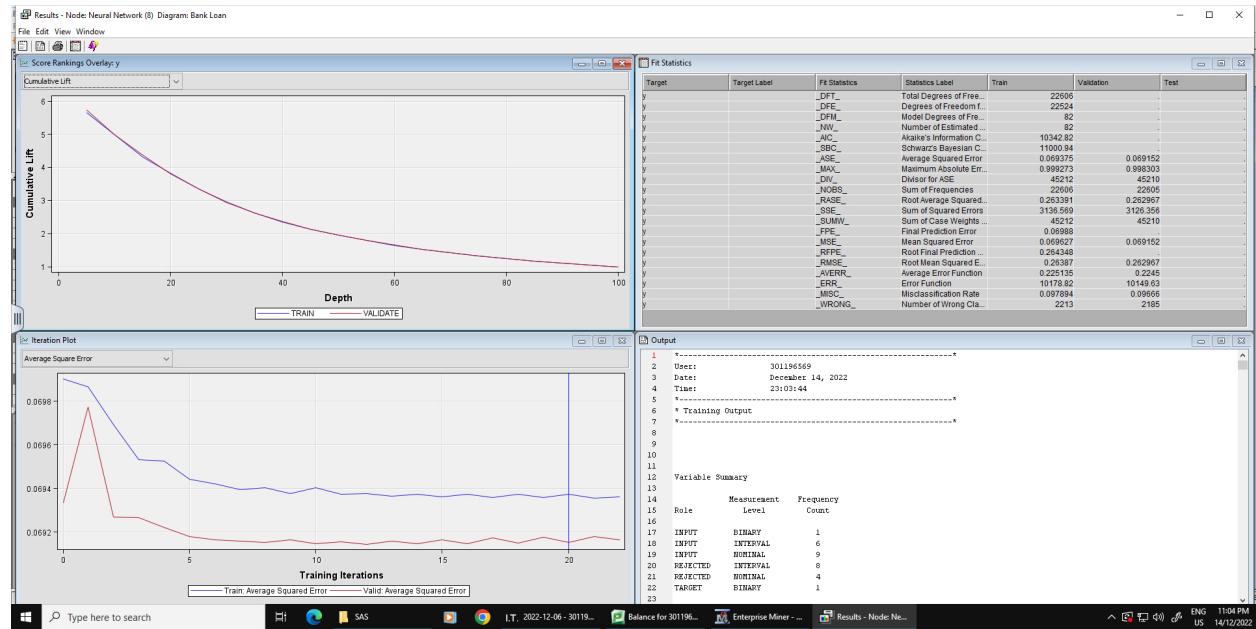
Full Regression with Dummies Recode (with replacement values). The ASE is 0.072314.



Stepwise Regression with Dummies Recode (with replacement values). The ASE is 0.072329.

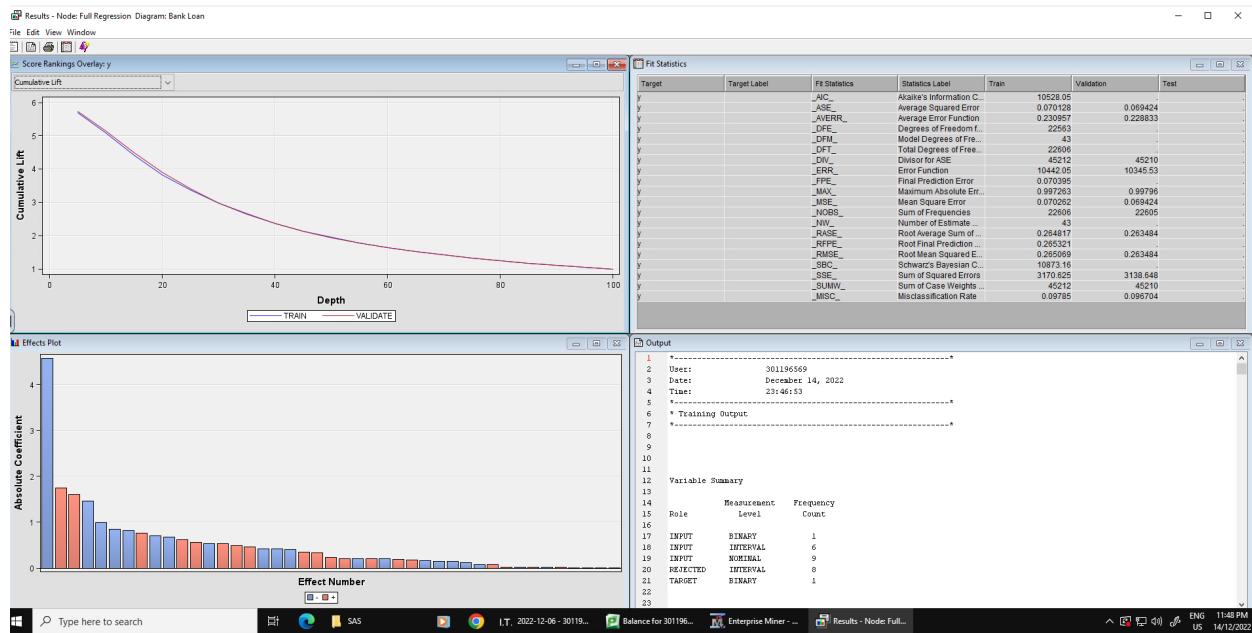


Neural Network with Dummies Recode (replacement values). The ASE is 0.069152.

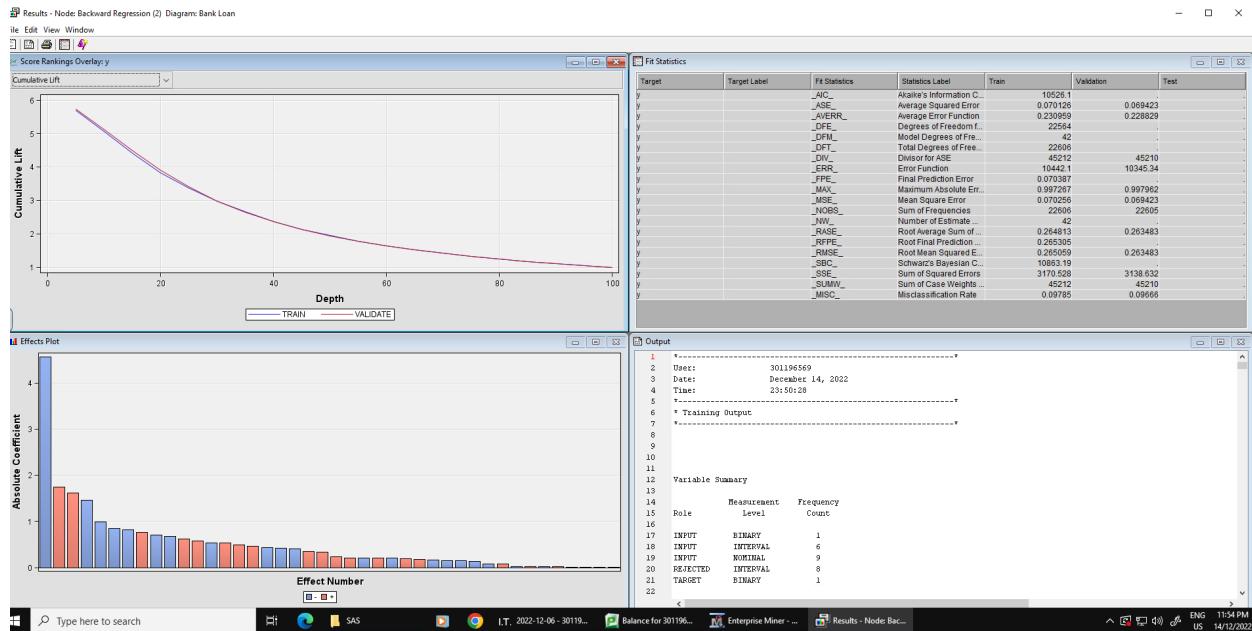


Regression without Dummies Recode

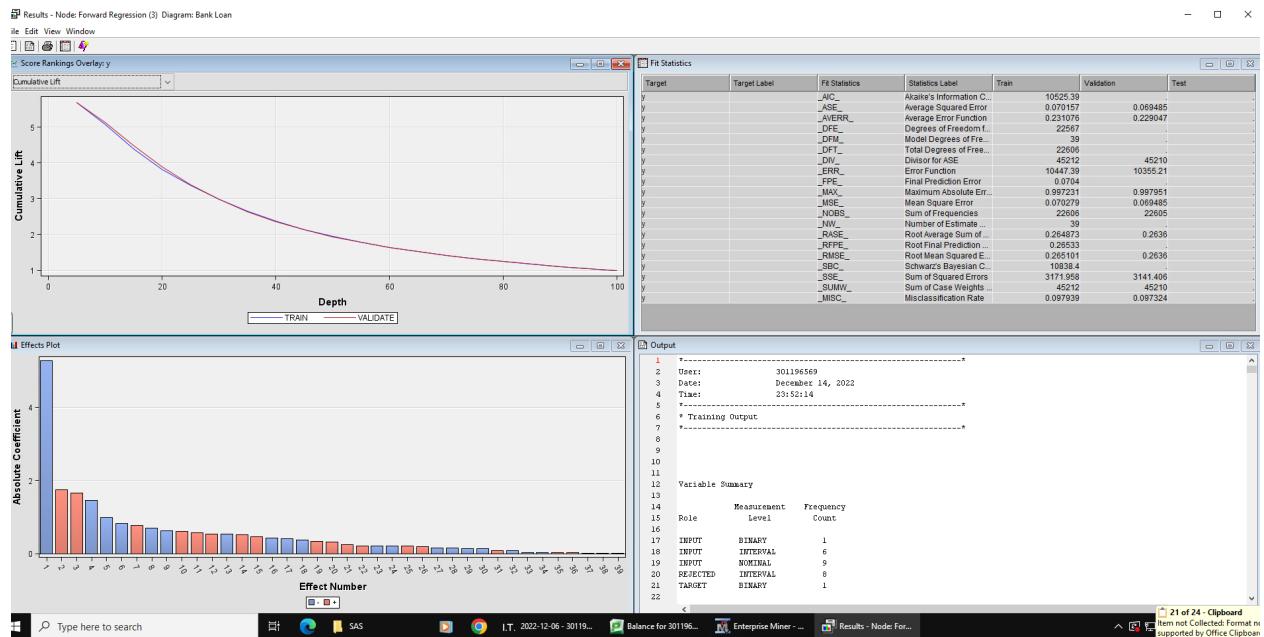
Full Regression: In our full regression, we did not change anything in the property value, including the selection model. Our ASE for full regression is **0.069424**. The result of the full without dummies recode is better than with dummies recode.



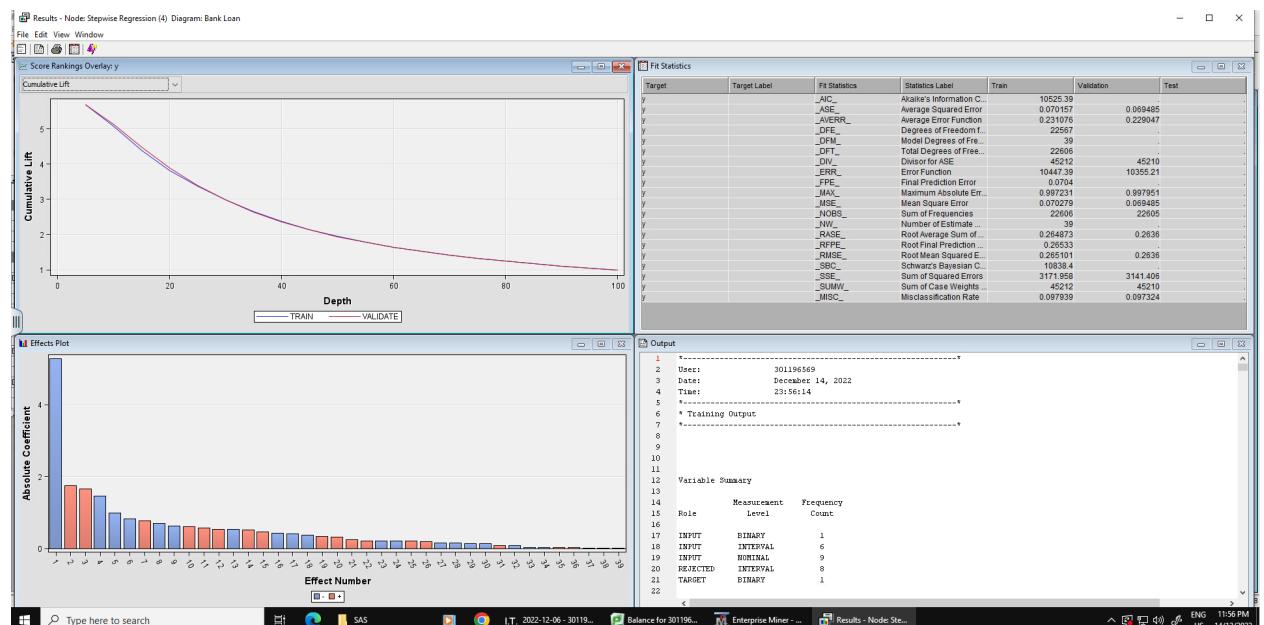
Backward Regression: In our regression, we made some changes. We changed the selection model to Backward and the selection criterion to validation error. Our ASE for backward regression is **0.069423**.



Forward Regression: For forward regression, we changed our selection model to forward and the selection criterion to validation error. Our ASE for forward regression is **0.069485**.



Stepwise Regression: In stepwise regression, we changed the selection model to stepwise while the selection criterion was validation error. The ASE for stepwise regression is **0.069485**.



Odds Ratio and interpretation.

213	Odds Ratio Estimates	
214		Point Estimate
215	Effect	
216	LOG_REP_balance	1.268
217	LOG_REP_campaign	0.582
218	LOG_REP_DAY	1.211
219	LOG_REP_previous	1.193
220	M REP age	0 vs 1
221	M REP age	0.182
222	REP_IMP_REP_age	0 vs 1
223	REP_IMP_REP_age	0.999
224	REP_duration	1.005
225	contact cellula vs unknown	1.178
226	contact telephone vs unknown	5.359
227	default no vs yes	1.055
228	education primary vs unknown	0.883
229	education secondary vs unknown	1.084
230	education tertiary vs unknown	1.275
231	housing no vs yes	1.999
232	job admin. Vs unknown	0.939
233	job blue-collar vs unknown	0.705
234	job entrepreneur vs unknown	0.712
235	job housemaid vs unknown	0.572
236	job management vs unknown	0.879
237	job retired vs unknown	1.210
238	job self-employed vs unknown	0.804
239	job services vs unknown	0.750
240	job student vs unknown	0.824
241	job technician vs unknown	0.862
242	job unemployed vs unknown	0.893
243	loan no vs yes	1.496
244	marital divorced vs single	0.896
245	marital married vs single	0.747
246	month apr vs sep	0.932
247	month aug vs sep	0.192
248	month dec vs sep	0.727
249	month feb vs sep	0.344
250	month jan vs sep	0.091
251	month jul vs sep	1.172
252	month jun vs sep	0.624
253	month may vs sep	2.247
254	month may vs sep	0.260
255	month nov vs sep	0.144
256	month oct vs sep	0.642
257	poutcome failure vs unknown	0.846
258	poutcome other vs unknown	1.084
259	poutcome success vs unknown	0.338
260		

24 of 24 - Clipboard
Item not Collected. Format supported by Office Clipbo...

Odds Ratio

Variables	Ratios	% of Change
LOG_REP_balance	1.268	27%
LOG_REP_camp	0.582	-42%
LOG_REP_REP_day	1.211	21%
LOG_REP_previous	1.193	19%
M REP_age 0 v 1	0.182	-82%
REP_IMP_REP_age	0.999	0%
REP_duration	1.005	0%
contact cellula vs unknown	5.178	418%
contact telephone vs unknown	5.359	436%
default no vs yes	1.055	5%
education primary vs unknown	0.883	-12%
education secondary vs unknown	1.004	0%
education tertiary vs unknown	1.275	28%
housing no vs yes	1.999	100%
job admin. Vs unknown	0.939	-6%
job blue-collar vs unknown	0.705	-30%
job entrepreneur vs unknown	0.712	-29%
job housemaid vs unknown	0.572	-43%
job management vs unknown	0.879	-12%
job retired vs unknown	1.21	21%
job self-employed vs unknown	0.804	-20%
job services vs unknown	0.75	-25%

job student vs unknown	1.424	42%
job technician vs unknown	0.862	-14%
job unemployed vs unknown	0.893	-11%
loan no vs yes	1.496	50%
divorced vs single	0.896	-10%
married vs single	0.747	-25%
month apr vs sep	0.382	-62%
month aug vs sep	0.192	-81%
month dec vs sep	0.727	-27%
month feb vs sep	0.344	-66%
month jan vs sep	0.091	-91%
month jul vs sep	0.172	-83%
month jun vs sep	0.624	-38%
month mar vs sep	0.624	-38%
month may vs sep	0.26	-74%
month nov vs sep	0.144	-86%
month oct vs sep	0.842	-16%
poutcome failure vs unknown	0.846	-15%
poutcome other vs unknown	1.084	8%
poutcome success vs unknown	8.338	734%

From the odds ratio, we have the following variables to interpret.

- LOG_REP_balance: This means that for every log(euro) increase in the customers' log_balances, they are 27% more likely to accept the loan campaign of the bank. Alternatively, for every increase in euro of 2.74, the probability of accepting the loan campaign goes up 27%.
- LOG_REP_camp: For every log(campaign), the number of customers contacted is likely to reduce by 42%. Alternatively, for every decrease in campaign of 2.74, the probability of rejecting the loan campaign goes up 42%.
- LOG_REP_REP_day: Depending on the log(day) of the month the customers are contacted, there is a 21% chance that the customer will accept the loan campaign. Alternatively, for every increase in day of 2.74, the probability of accepting the loan campaign goes up 21%.
- LOG_REP_previous: For every log(previous) marketing of the bank loan to the customers there is a 19% chance of influencing the customers' decision to accept the bank loan. Alternatively, for every increase in previous marketing campaign of 2.74, the probability of accepting the loan campaign goes up 19%.
- M REP_age 0 v 1: For every addition in the known age, there is an 82% chance of not being likely to accept the loan campaign.

- REP_IMP_REP_age: For every increase in the imputed age, there are unlikely that the loan campaign will influence the customers' decision to accept a bank loan.
- REP_duration: For every increase in the number of seconds spent explaining bank loans to customers, there is a 0.5% of chance of the customers accepting the loan campaign.
- Contact - cellula vs unknown: Customers contacted via cellula are 5 times more likely to accept the bank loan campaign than those contacted via unknown means of communication.
- Contact - telephone vs unknown: Customers contacted through telephone are also 5 times more likely to accept loan campaigns than the unknown.
- default - no vs yes: For every customer contacted for the bank loan, there is a 5% chance that they have no credit default.
- education - primary vs unknown: For every customer, there is a 12% chance that they do not have primary education than the unknown.
- education - secondary vs unknown: For every customer, there is a 0.4% chance that they do have secondary education than the unknown.
- education - tertiary vs unknown: Every customer contacted is 28% likely to have tertiary education.
- housing - no vs yes: Every customer is 100% likely to not have a housing loan.
- job - admin. Vs unknown: Every customer contacted is 6% more unlikely to be in an admin job than the unknown.
- Job - blue-collar vs unknown: 30% of the customers are more unlikely to have blue-collar than the unknown.
- job - entrepreneur vs unknown: There is a 29% chance that the customers are unlikely to be entrepreneurs.
- job - housemaid vs unknown: There is a 43% chance that the customers are more unlikely to be housemaids than the unknown.
- job - management vs unknown: The customers are 12% unlikely to be in a management cadre in their jobs.
- job - retired vs unknown: Retired customers are 21% more likely to accept loan campaigns.
- job - self-employed vs unknown: Self-employed customers are 20% more unlikely to accept loan campaigns than the unknown.
- job - services vs unknown: Customers in services jobs are 25% more likely not to accept loan campaigns than the unknown.
- job - student vs unknown: Students are 45% most likely to accept loan campaigns than the unknown.
- job - technician vs unknown: Customers who are technicians are 14% more likely not to accept loan campaigns than the unknown.
- job - unemployed vs unknown: Customers who are unemployed are 11% unlikely to accept loan campaigns.

- loan - no vs yes: Customers who have existing personal loans are 50% more likely to say no loan campaigns than yes.
- divorced - vs single: Divorced customers are 10% more unlikely to accept loan campaigns than single.
- married - vs single: Bank customers who are single are 25% more likely to accept loan campaigns than married.
- month - apr vs sep: Customers contacted in April are 62% more unlikely to accept loan campaigns than in September.
- month - aug vs sep: Customers contacted in August are 81% more unlikely to accept loan campaigns than in September.
- month - dec vs sep: Customers contacted in December are 27% more unlikely to accept loan campaigns than in September.
- month - feb vs sep: Customers contacted in February are 66% more unlikely to accept loan campaigns than in September.
- month - jan vs sep: Customers contacted in January are 91% more unlikely to accept loan campaigns than in September.
- month - jul vs sep: Customers contacted in July are 83% more unlikely to accept loan campaigns than in September.
- month - jun vs sep: Customers contacted in June are 38% more unlikely to accept loan campaigns than in September.
- month mar - vs sep: Customers contacted in March are 38% more unlikely to accept loan campaigns than in September.
- month - may vs sep: Customers contacted in May are 74% more unlikely to accept loan campaigns than in September.
- month - nov vs sep: Customers contacted in November are 86% more unlikely to accept loan campaigns than in September.
- month - oct vs sep: Customers contacted in October are 16% more unlikely to accept loan campaigns than in September.
- poutcome - failure vs unknown: There is a likelihood of 15% failure of the outcome of the previous campaign than the unknown.

- poutcome - other vs unknown: There is an 8% chance of other outcomes of the previous loan campaigns than the unknown.
- poutcome - success vs unknown: The outcome of the previous loan campaign is 8 times more likely to be successful than the unknown.

Summary of our regression models

S/N	Regression	ASE
1	Full regression	0.069424
2	Backward regression	0.069423
3	Forward and stepwise regression	0.069485

The most important variables in the regression are those variables with the least Ch-square of 0.0001 as below.

```

969
970
971      Type 3 Analysis of Effects
972
973      Wald
974   Effect      DF   Chi-Square   Pr > ChiSq
975
976 LOG_REP_balance    1     8.3401    0.0039
977 LOG_REP_campaign    1    67.2839   <.0001
978 LOG_REP_day         1    17.1160   <.0001
979 LOG_REP_previous    1     2.6306    0.1048
980 M_REP_age          1     2.7811    0.0954
981 REP_IMP_REP_age    1     0.1083    0.7421
982 REP_duration       1   2572.8082   <.0001
983 contact            2    254.3552   <.0001
984 education          3     14.7029    0.0021
985 housing             1    122.3450   <.0001
986 job                 11    41.7323   <.0001
987 loan                1    22.6922   <.0001
988 marital             2     21.2814   <.0001
989 month               11    534.6409   <.0001
990 poutcome            3    437.2939   <.0001
991

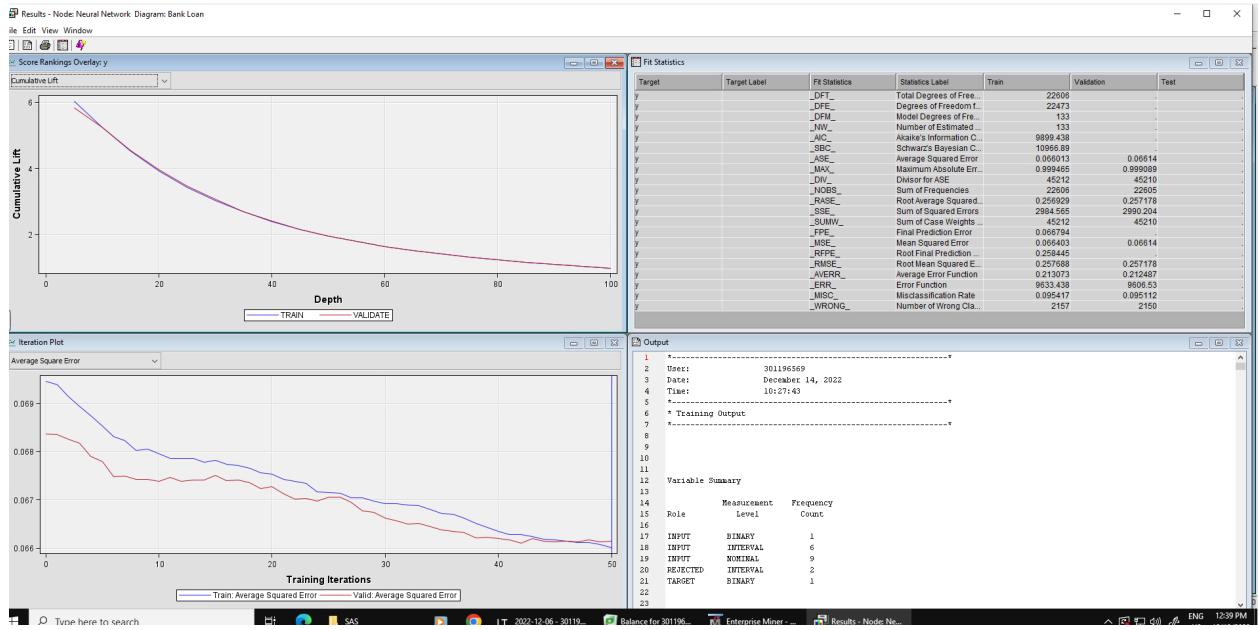
```

Conclusion: From the above, our best regression model is backward regression with an ASE of **0.069423**.

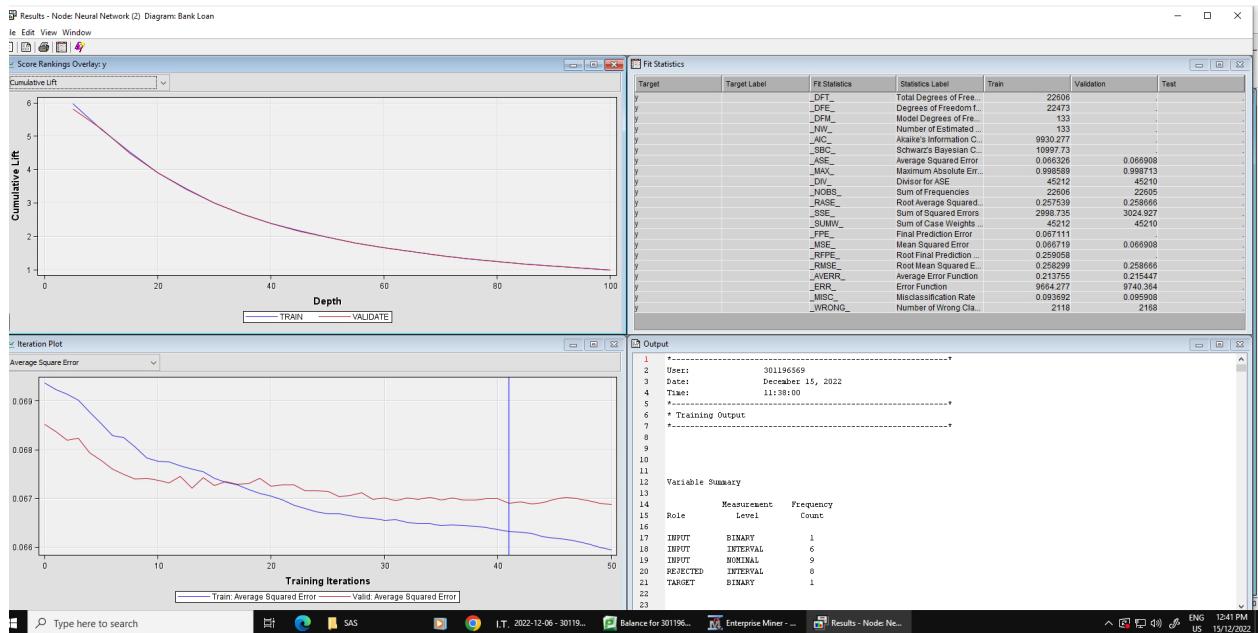
Neural Network:

We built neural network regression models using profit and loss as the selection criterion. The following are the neural network models and the results.

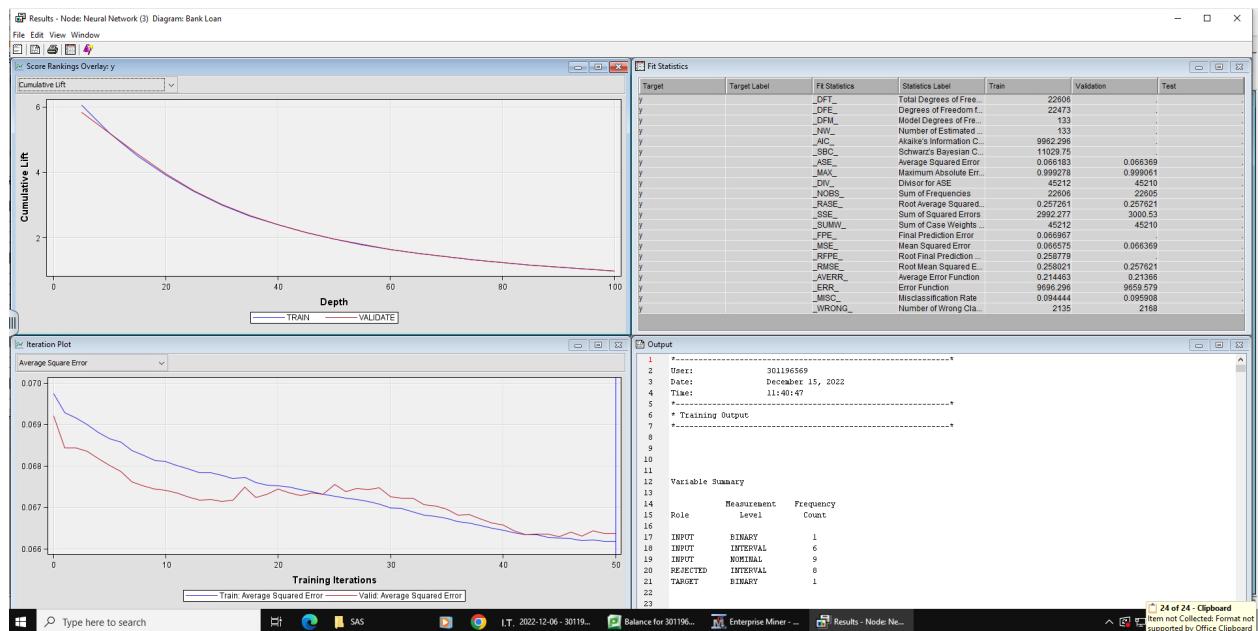
1. Neural network: The first network is connected to the impute without changing anything in the model optimization. The ASE for the first neural network is **0.06614**.



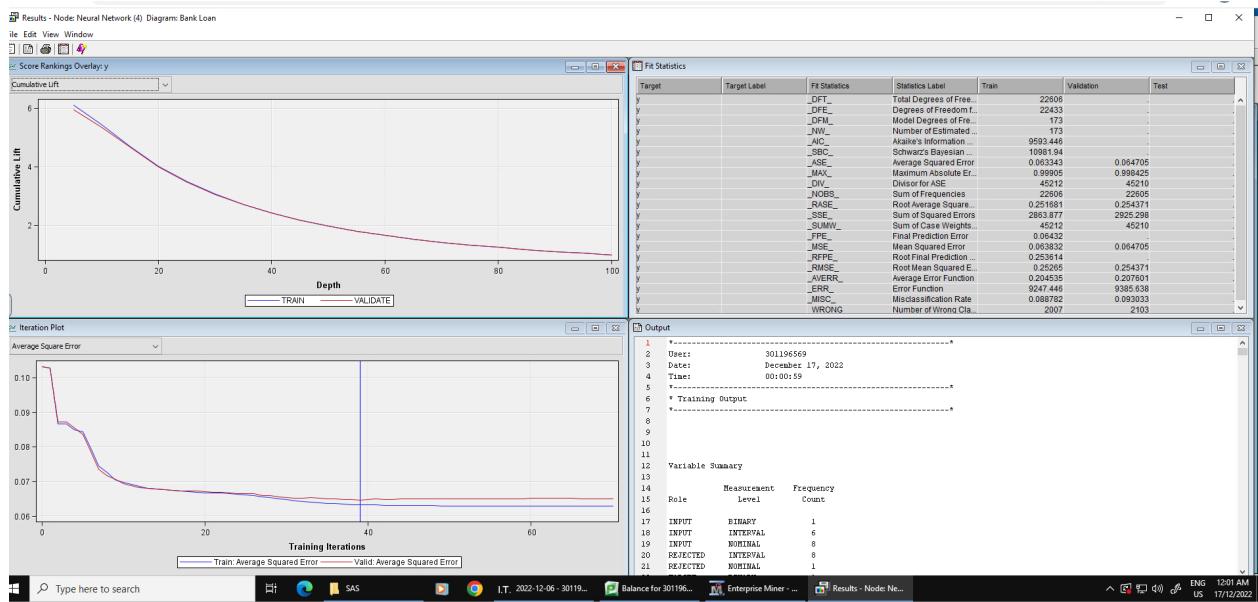
2. Neural network (2): The second neural network is connected Cap and Floor. The cap and floor are intended to treat variables that are skewed. Hence, the difference between the ASE of the first neural network and the second. The ASE of the second neural network is **0.066908**.



3. Neural network (3): The third to connected to the Log Tranform Variables which is also intended to treat the skewed variables after cap and floor. However, the ASE of the third neural network is 0.066369.



4. Neural network (4): The neural network (4) is connected to the backward regression because it is our best regression model. However, we made some changes to optimize the model by changing the maximum iteration to 100 and preliminary training to no. We also changed the hidden unit number of the network from 3 to 4 and this is the optimal hidden unit number. The ASE for the 4th neural work is 0.064705.



Event Classification:

Our event classification refers to the classification of the impact of a bank loan campaign. The classification is below.

- **True negative:** This represents customers who were not targeted (contacted) by the bank for the loan campaign and did not subscribe to the bank loan. They represent a total of 19, 265 customers, or 85.2245% of the customers.
- **False negative:** This represents bank customers who were targeted (contacted) for the bank loan but did not subscribe to the loan. They are 1, 407 representing 6.2243% of the customers.
- **False positive:** These are the customers not targeted but subscribed to the loan. They are 696 representing 3.079% of the customers.
- **True positive:** These are bank customers targeted and who subscribed to the bank loan. They are 1, 237 representing 5.4722% of the customers.

See the event classification table below.

```

596
597
598 Classification Table
599
600 Data Role=TRAIN Target Variable=y Target Label=' '
601
602             Target          Outcome      Frequency      Total
603     Target    Outcome   Percentage   Percentage   Count   Percentage
604
605     NO        NO       93.5734     96.5783     19278    85.2782
606     YES       NO       6.4266      50.0567     1324     5.8569
607     NO        YES      34.0818     3.4217      683     3.0213
608     YES       YES      65.9182     49.9433     1321    5.8436
609
610
611 Data Role=VALIDATE Target Variable=y Target Label=' '
612
613             Target          Outcome      Frequency      Total
614     Target    Outcome   Percentage   Percentage   Count   Percentage
615
616     NO        NO       93.1937     96.5132     19265    85.2245
617     YES       NO       6.8063      53.2148     1407     6.2243
618     NO        YES      36.0062     3.4868      696     3.0790
619     YES       YES      63.9938     46.7852     1237    5.4722
620
621
622
623
624 Event Classification Table
625
626 Data Role=TRAIN Target=y Target Label=' '
627
628     False      True      False      True
629 Negative  Negative  Positive  Positive
630
631     1324     19278     683     1321
632
633

```

Summary of neural network

S/N	Neural Network	ASE
1	Neural Network (connected to impute).	0.06614
2	Neural network (2) connected to cap & floor.	0.066908
3	Neural network (3) connected to log transform variable.	0.066369
4	Neural network (4) connected backward regression.	0.064705

From the above summary table, our best neural network model is the one with the lowest ASE of **0.064705** which is the neural network (4) connected to backward regression.

Model comparison:

For our model comparison, use the ROC index and ASE to compare the three models, decision tree, regression, and neural network. Both the ROC index and the ASE are in agreements on the best model. The best model has the highest ROC index and the lowest ASE. Below are the screenshots of the model comparison, ROC curve and Fit statistics.

The specificity and sensitivity are measures used to evaluate the performance of the model in a binary classification task. Sensitivity also known as the True Positive rate or recall measures the portion of actual true positive case correctly identified. Thus, our sensitivity of 90.96% implies that the model correctly identifies 90.96% of the customers who said yes to the bank loan.

On the other hand, a specificity of 21.03% correctly identified the proportion of the customers who said no to the bank loan campaign. This implies that model only identifies 21% of actual customers who said no to the bank loan campaign as can be seen from the ROC curve above.

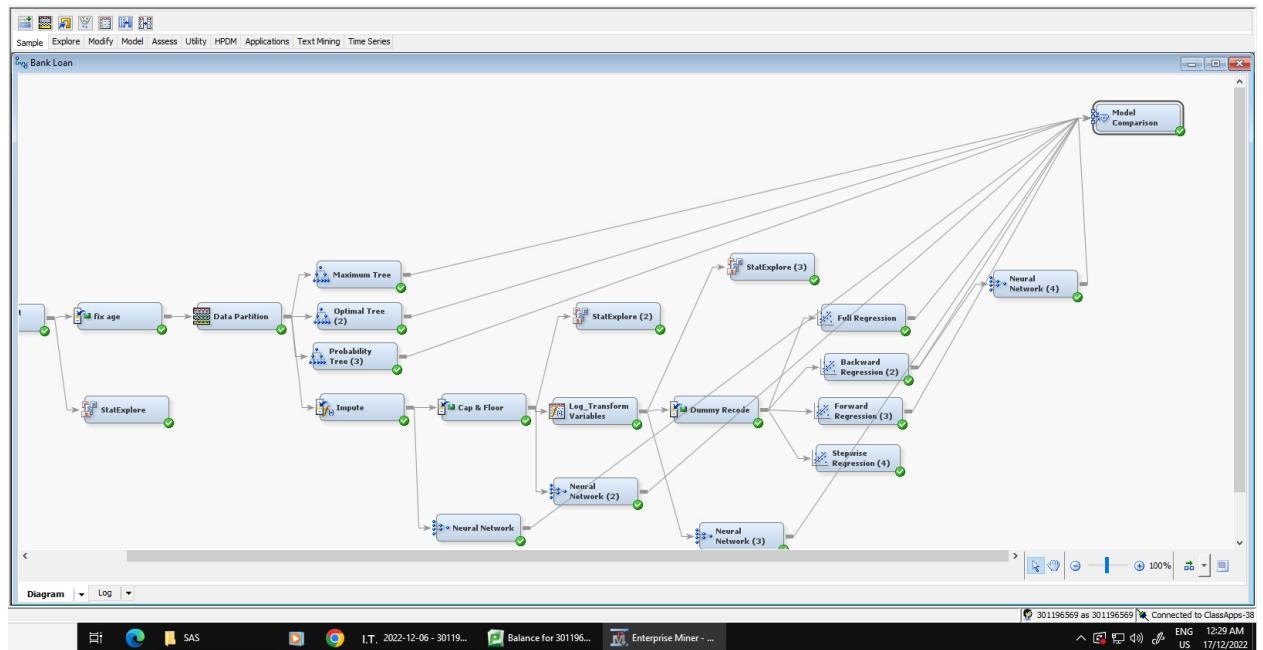
Fit Statistics																		
Selected Model	Predecessor Node	Model Node	Model Description	Target Variable	Valid: Roc Index	Valid: Average Squared Error	Target Label	Selection Criterion: Valid Average Squared Error	Train: Sum of Frequencies	Train: Misclassification Rate	Train: Maximum Absolute Error	Train: Sum of Squared Errors	Train: Average Squared Error	Train: Root Average Squared Error	Train: Divisor for ASE	Train: Total Degrees of Freedom	Valid: Sum of Frequencies	Valid: Misclassification Rate
Y	Neural4	Neural4	Neural Net..	y	0.925	0.064705	0.064705	22300	0.089782	0.99805	2863.877	0.093343	0.251681	45212	22606	22305	0.095	
	Neural	Neural	Neural Net..	y	0.921	0.06914	0.06914	22300	0.095417	0.99495	2984.566	0.096113	0.256929	45212	22606	22305	0.095	
	Neural3	Neural3	Neural Net..	y	0.918	0.06939	0.06939	22300	0.094444	0.99278	2992.277	0.095183	0.257411	45212	22606	22305	0.095	
	Neural2	Neural2	Neural Net..	y	0.918	0.069208	0.069208	22300	0.093692	0.988589	2998.735	0.095326	0.257539	45212	22606	22305	0.095	
	Reg3	Reg3	Backward...y	y	0.911	0.069423	0.069423	22300	0.09765	0.997267	3170.528	0.070126	0.264813	45212	22606	22305	0.05	
	Reg	Reg	Full Regress...y	y	0.911	0.069424	0.069424	22300	0.09765	0.997263	3170.625	0.070128	0.264817	45212	22606	22305	0.05	
	Reg2	Reg2	Forward Reg...y	y	0.911	0.069485	0.069485	22300	0.097939	0.997231	3171.956	0.070167	0.264873	45212	22606	22305	0.097	
	Reg4	Reg4	Sterwise Reg...y	y	0.911	0.069485	0.069485	22300	0.097939	0.997231	3171.958	0.070167	0.264873	45212	22606	22305	0.097	
	Tree3	Tree3	probability t...y	y	0.909	0.065942	0.065942	22300	0.02924	0.997384	2933.709	0.064898	0.254731	45212	22606	22305	0.094	
	Tree	Tree	maximal tree	y	0.909	0.066752	0.066752	22306	0.090054	0.997384	2875.951	0.063861	0.252211	45212	22606	22305	0.091	
	Tree2	Tree2	Class Tree	y	0.857	0.069907	0.069907	22306	0.094975	0.974399	3217.253	0.071159	0.266757	45212	22606	22305	0.092	

Talking:

Type here to search

File I.T. 2023-05-10 - ... SAS Balance for 3... Enterprise Mi... Results - No... ENG 2:23 PM CAFR 12/07/2023





Conclusion and recommendations

In summary, we performed model analytics on the bank dataset to predict the effectiveness of a marketing campaign by a bank using three different analytics models. The dataset consists of 45,211 observations and 17 variables. Our important variables are y (target variable), contact, housing and job. All our important variables have Chi-squares less than 0.0001.

In this project, we built decision trees, regressions, and neural networks to ascertain model optimization. Below is the summary of our best models with the model assessment measures.

S/N	MODELS	ASE	VALIDATION INDEX	ROC
1	Probability Decision Tree	0.065942	0.909	
2	Backward Regression	0.069423	0.911	
3	Neural Network 4	0.064705	0.925	

From the above table, our best model is the neural network model with the highest validation ROC index and least average squared error of 0.925 and 0.064705 respectively.

Limitation:

In the course of performing the analytics, we countered some challenges. First, we have skewed variables. We used cap and floor, standardize and finally log transformation but still have about three of the variables skewed. We were not able to transform all the variables. Secondly, we have dummy variables and we tried to use a dummy recode. However, our regression and neural network models' optimization were worse off with dummy recode replacement values. We, therefore, realized that our model without dummy recode have the optimal results, and it was difficult to interpret the odds ratio estimates of the model.

Recommendation:

We recommend that as analysts, there is a need to thoroughly review the data to ensure that irrelevant and redundant variables are rejected to ensure model fit.

From the business point of view, we will advise the client to pay attention to these variables in targeting customers to market the loan campaign. They are housing, loan, poutcome, and contact. These variables have been interpreted in the odds ratio interpretation section. Please refer to this section for better understanding of these variables and their impacts on the success of the loan campaign.

References

Camilosierra. (2020, June 30). *Bankdataset*. Kaggle.
<https://www.kaggle.com/datasets/juancamilosierra/bankdataset>