

Project 2 – Predictive Analytics

DSBA-6201 - Business Intelligence and Analytics

Group members:

Yawo Eklou

Marcie Matejka

May 10, 2021

PART I: Predictive Analytics (EDA & RFM)

Completion of the initial tutorial with Catalog Data.

There are several variables that should be modified.

Variables - CATALOG2010

(none) ▾ not Equal to

Columns: Label Mining

Name	Role	Level	Report	Order	Drop
ACTBUY	Input	Interval	No		No
BOTHPAYM	Input	Binary	No		No
BUYPROP	Input	Interval	No		No
CATALOGCNT	Input	Interval	No		No
CCPAYM	Input	Binary	No		No
COUNTY	Rejected	Interval	No		No
CUST_ID	ID	Interval	No		No
DAYLAST	Input	Interval	No		No
DEPT01	Input	Interval	No		No
DEPT02	Input	Interval	No		No
DEPT03	Input	Interval	No		No
DEPT04	Input	Interval	No		No
DEPT05	Input	Interval	No		No
DEPT06	Input	Interval	No		No

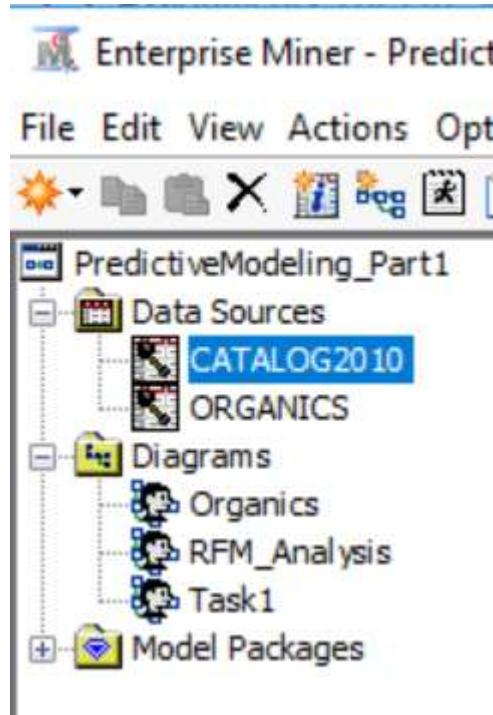
Variables - CATALOG2010

(none) ▾ not Equal to

Columns: Label Mining

Name	Role	Level	Report	Order	Drop
DEPT22	Input	Interval	No		No
DEPT23	Input	Interval	No		No
DEPT20	Input	Interval	No		No
DEPT21	Input	Interval	No		No
DEPT25	Input	Interval	No		No
DEPT24	Input	Interval	No		No
DEPT26	Input	Interval	No		No
STATE	Rejected	Nominal	No		No
DTBUYORG	Rejected	Interval	No		No
ZIP	Rejected	Nominal	No		No
DTBUYLST	Rejected	Interval	No		No
COUNTY	Rejected	Interval	No		No
ORDERSIZE	Target	Interval	No		No
RESPOND	Target	Binary	No		No

The CATALOG2010 data source is added to the Data Sources entry in the Project panel.

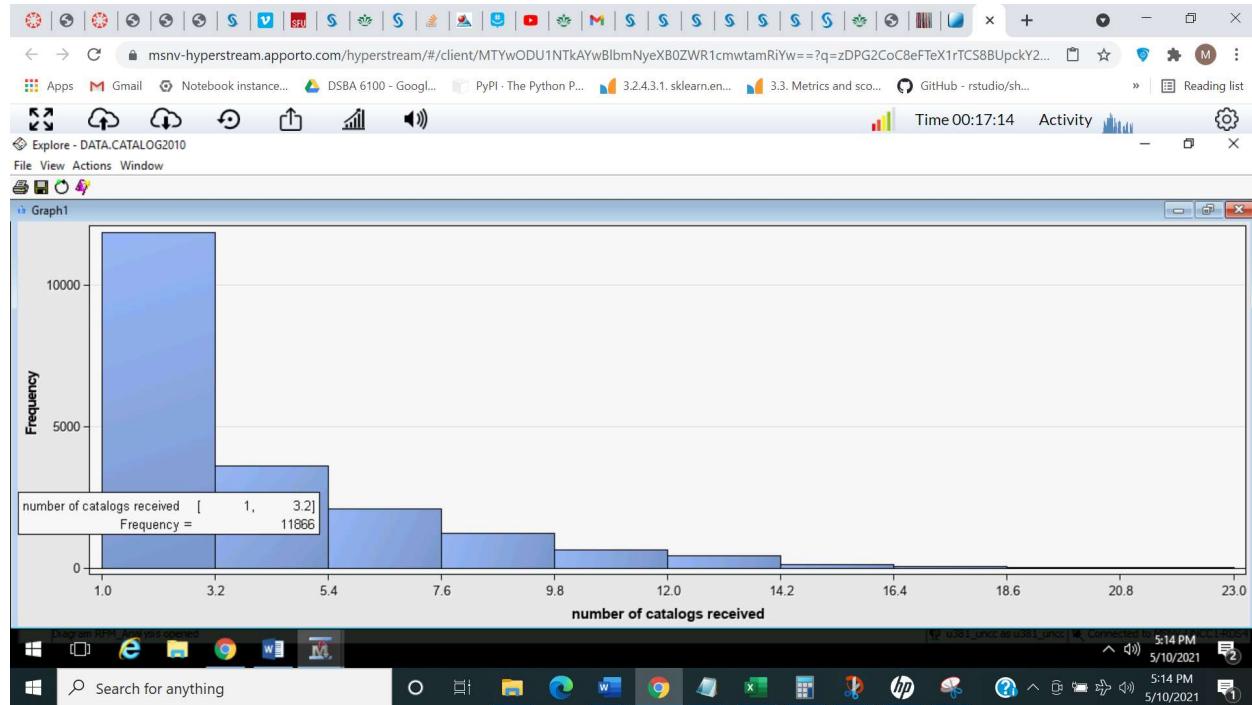


Follow the steps below to change the preference settings of SAS Enterprise Miner to use a random sample or all of the data source data in the Explore window

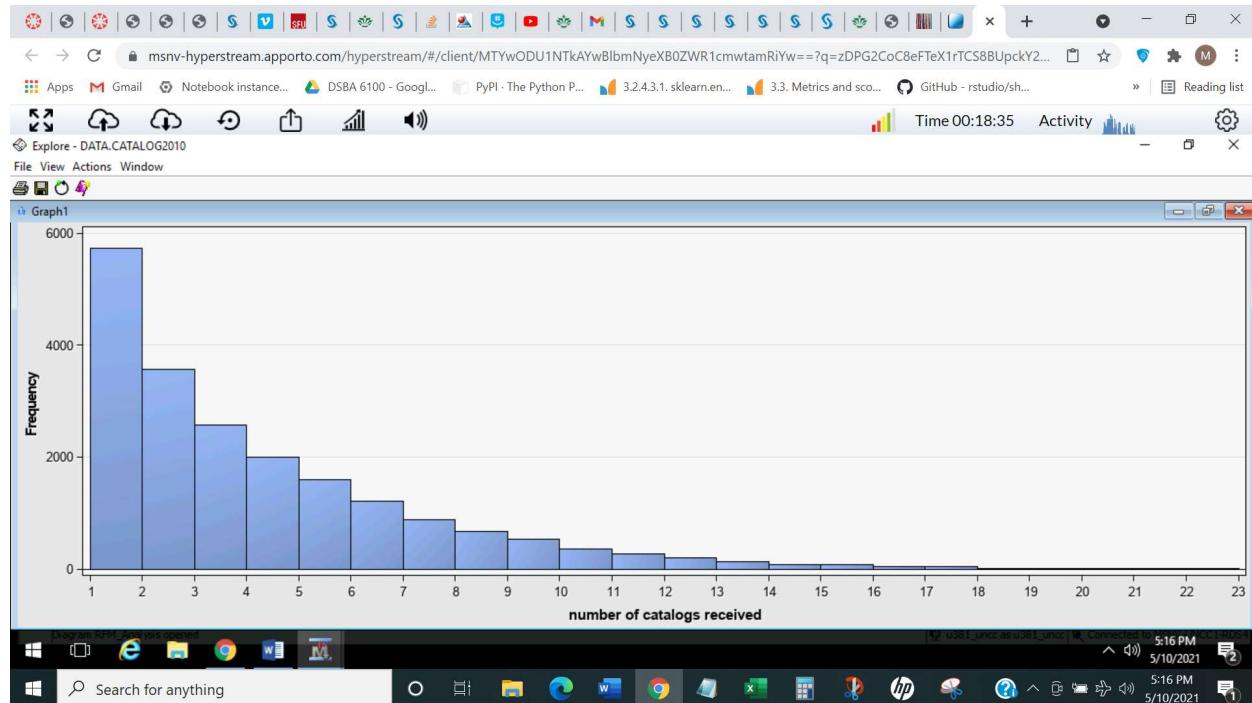
A screenshot of the "Preferences" dialog box. It has a header bar with "Preferences" and a close button. The main area is a table with two columns: "Property" and "Value".

Property	Value
User Interface	
--Property Sheet Tooltips	On
--Tools Palette Tool tips	Display tool name and description
--Open Last Opened Project Automatically	No
--Open Last Viewed Diagram Automatically	No
--Number of Recent Projects	5
Interactive Sampling	
--Sample Method	Random
--Fetch Size	Max
--Random Seed	12345
Model Package Options	

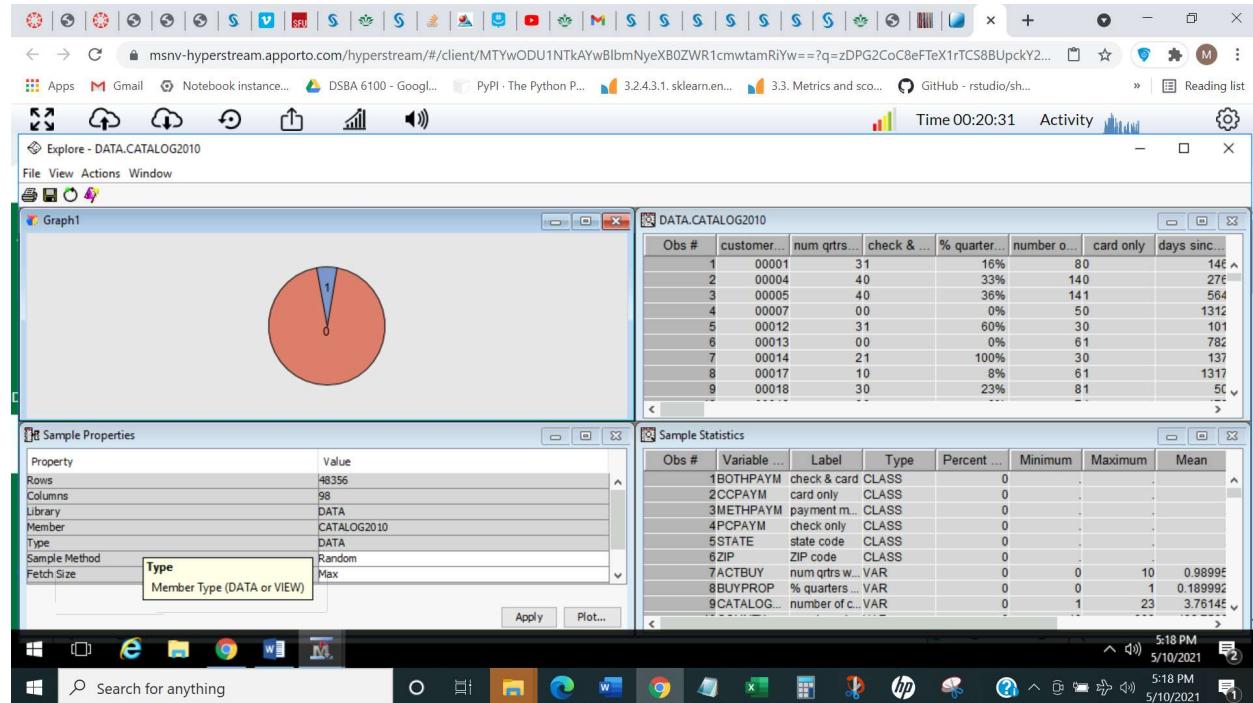
Creating a Histogram for a Single Variable



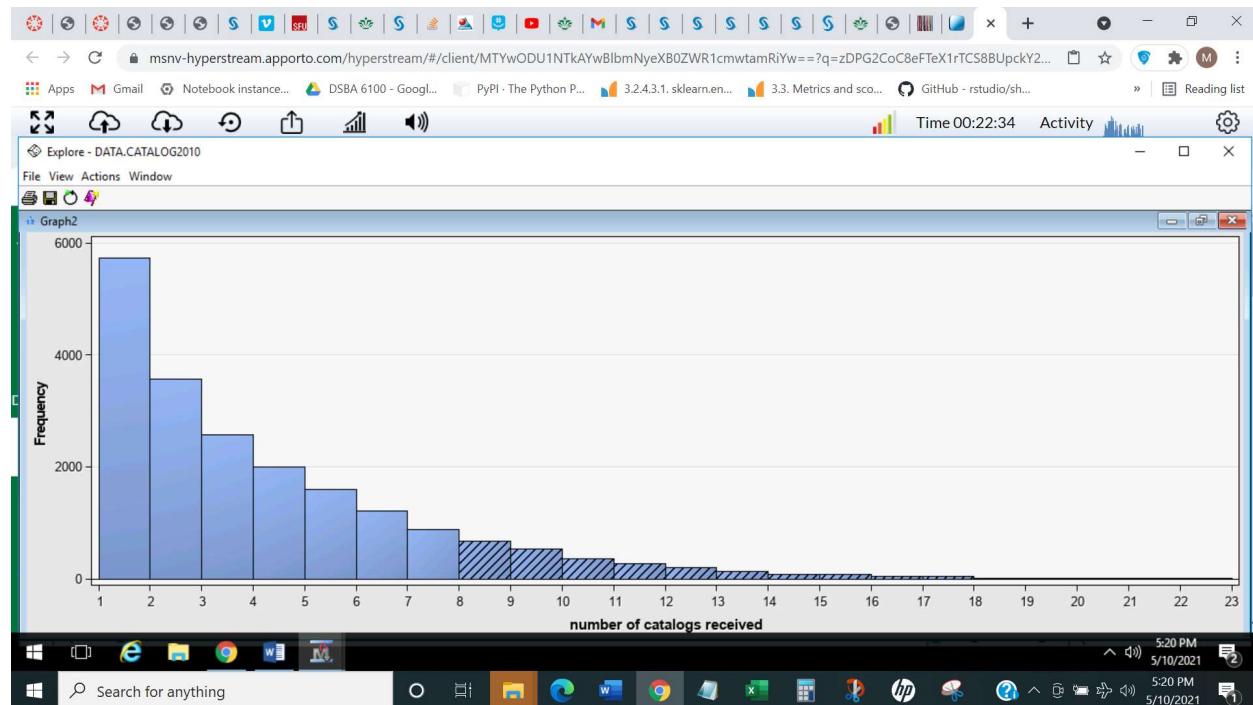
Changing the Graph Properties for a Histogram

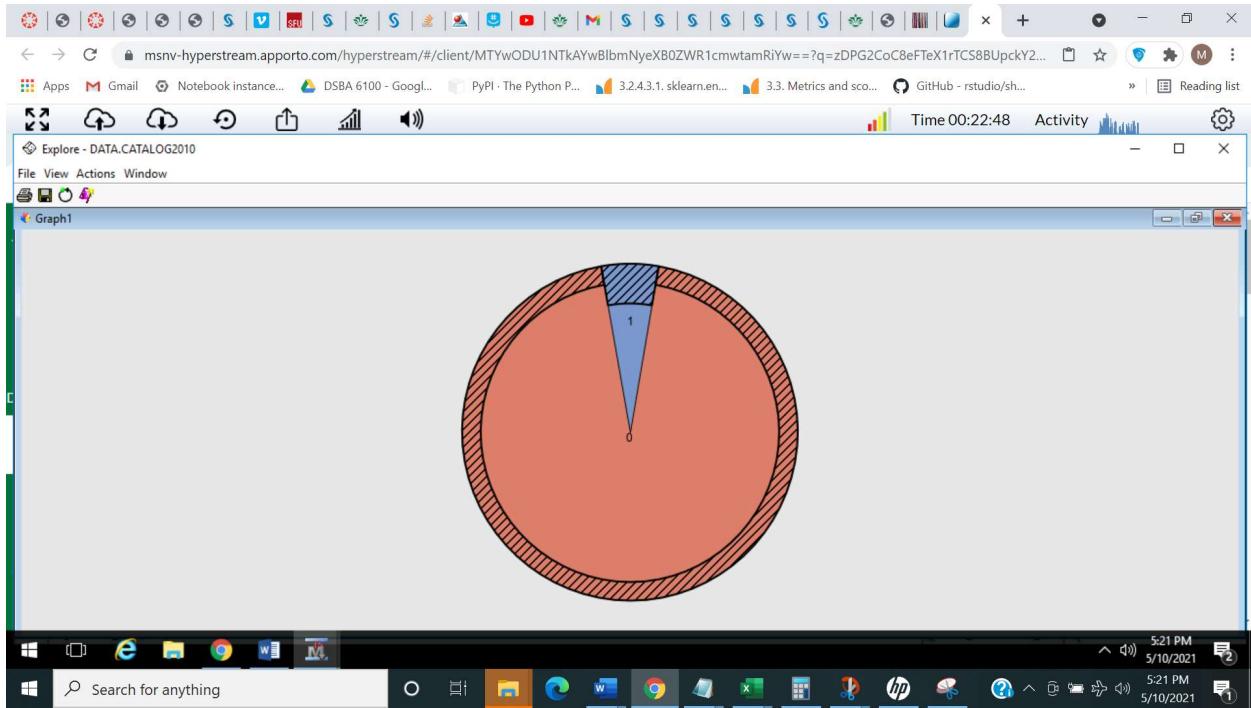


Adding Other Graphs

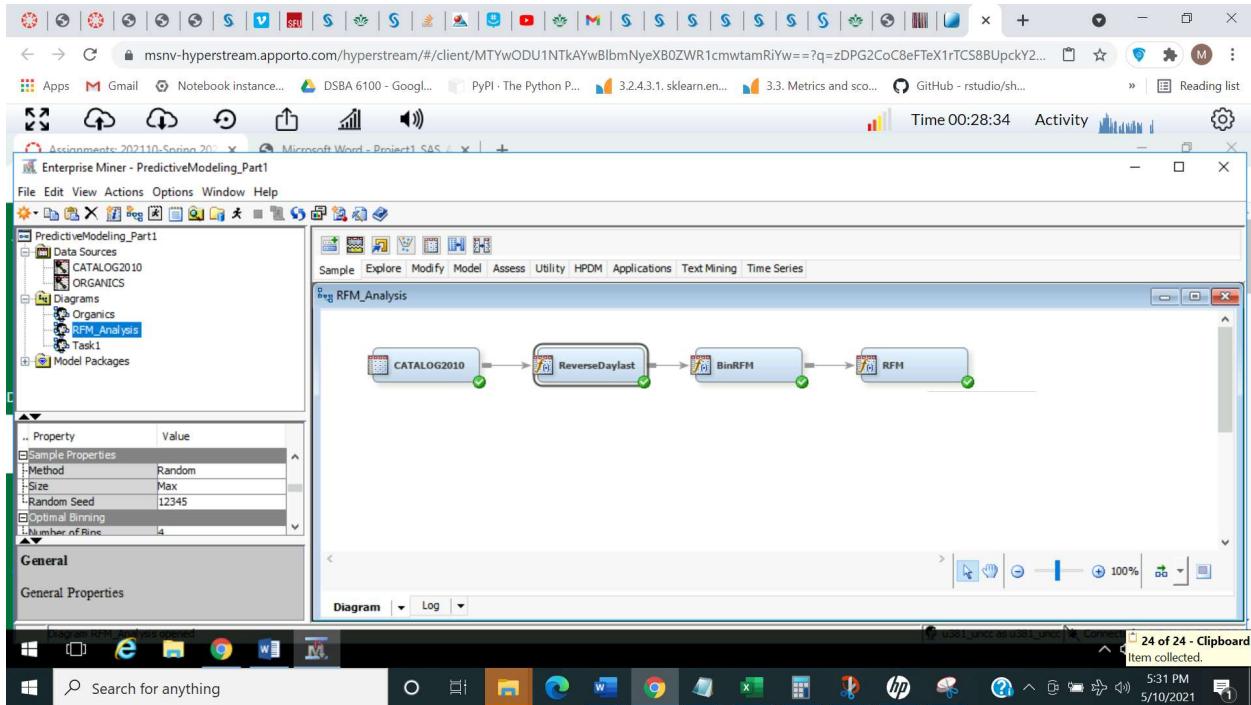


Exploring Variable Associations





Catalog Case Study: Performing RFM Analysis of the CATALOG Data



Formulas

Columns: Label Mining Basic Statistics

Name	Role	Level	Method	Number of Bins
DAYLAST	Ternary	Ternary	DEFAULT	10

Inputs

Name	Type	Length	Format	Level	Formula	Label	Role	Report
DAYLAST_REV	Numeric	8		Interval	-1 *DAYLAST		Input	No

Outputs

Sample Log

Preview OK Cancel

```

graph LR
    C1[CATALOG2010] --> R1[ReverseDaylast]
    R1 --> B1[BinRFM]
    B1 --> RFM[RFM]
  
```

Variables - Trans2

(none) ▾ not Equal to ...

Columns: Label Mining Basic Statistics

Name	Method	Number of Bins	Role	Level
DEPT23	Default	4	Input	Interval
DEPT22	Default	4	Input	Interval
DEPT24	Default	4	Input	Interval
DEPT20	Default	4	Input	Interval
DEPT19	Default	4	Input	Interval
DEPT21	Default	4	Input	Interval
DEPT27	Default	4	Input	Interval
DOLLARQ11	Default	4	Input	Interval
DOLINDEA	Default	4	Input	Interval
DEPT25	Default	4	Input	Interval
DEPT26	Default	4	Input	Interval
DOLL24	Quantile	5	Input	Interval
FREQPRCH	Quantile	5	Input	Interval
DAYLAST_REV	Quantile	5	Input	Interval

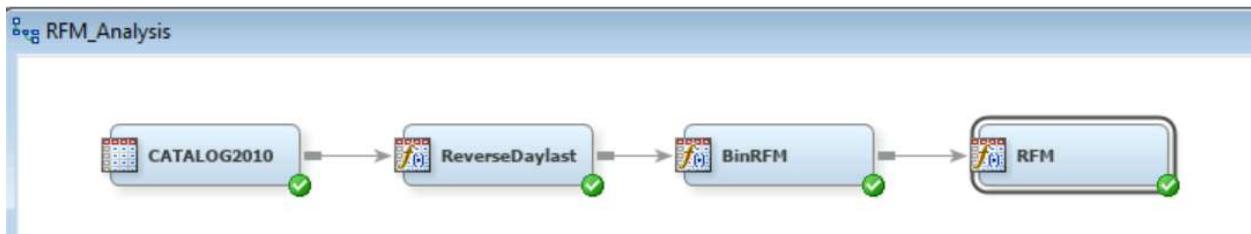
Explore... Update Path OK Cancel

The screenshot shows an RStudio interface with three main windows:

- Sample Properties** window: Displays properties like Rows (Unknown), Columns (102), Library (EMWS1), Member (TRANS2_TRAIN), Type (VIEW), Sample Method (Random), and Fetch Size (Max). Buttons for Apply and Plot... are at the bottom.
- Sample Statistics** window: A table showing descriptive statistics for each variable. The columns are Obs #, Variable ..., Label, Type, Percent ..., Minimum, Maximum, and Mean. The data includes:

Obs #	Variable ...	Label	Type	Percent ...	Minimum	Maximum	Mean
1	BOTHPAYM	check & card	CLASS	0	.	.	.
2	CCCPAYM	card only	CLASS	0	.	.	.
3	METHPAYM	payment m...	CLASS	0	.	.	.
4	PCPPAYM	check only	CLASS	0	.	.	.
5	PCCTL_DAY...	Transforme...	CLASS	0	.	.	.
6	PCCTL_DOL...	Transforme...	CLASS	0	.	.	.
7	PCCTL_FRE...	Transforme...	CLASS	0	.	.	.
8	STATE	state code	CLASS	0	.	.	.
9	ZIP	ZIP code	CLASS	0	.	.	.
- E-MWS1.Trans2_TRAIN** window: A data grid showing transactional data across 102 columns. The first few rows are:

rs...	tot orders...	DAYLAS...	Transfor...	Transfor...	Transfor...											
0	0	0	0	0	0	0	0	0	0	0	0	0	-146.05	-256-high	04.24.75-7...	05.6-high
0	0	0	0	0	0	0	0	0	0	0	0	0	-276.04	-573-256	04.24.75-7...	05.6-high
0	0	1	0	0	0	0	1	0	0	0	0	0	-564.04	-573-256	05.71.6-high	05.6-high
0	1	0	0	0	2	0	0	0	1	0	0	1	-1312.02	-1826-1...	03.0-24.75	04.3-6
1	0	0	0	0	1	0	0	0	0	0	0	0	-7823.03	-1008-5...	03.0-24.75	02.1-2
0	0	0	0	0	0	2	0	0	0	0	1	0	-1370.05	-256-high	05.71.6-high	04.3-6
0	0	0	0	0	0	0	0	0	0	0	0	0	-1317.02	-1826-1...	03.0-24.75	04.3-6
2	0	0	0	0	0	0	0	0	0	0	0	1	-500.05	-256-high	04.24.75-7...	05.6-high



Formulas

Add Transformation

Property	Value
Name	RFM
Type	Numeric
Length	8
Format	
Level	Interval
Label	
Role	Input
Report	No

Formula:

```
RFM =  
substr(pctl_daylast_rev,1,2)||substr(pctl_freqprch,1,2)||substr(pctl_doll24,1,2)
```

Build... OK Cancel

Inputs | Outputs | Sample | Preview

Catalog Case Study: Performing Graphical RFM Analysis

A Grouped Pie Chart

Explore - EMWS1.Trans3_TRAIN

View Actions Window

Sample Properties

Property	Value
Rows	101
Columns	101
Binary	EM
Number	TR
Date	VIE
Sample Method	Random
Batch Size	Max

EMWS1.Trans3_TRAIN

Obs #	customer...	num qtrns...	check & ...	%
1	00001	31		
2	00004	40		
3	00005	40		
4	00007	00		
5	00012	31		
6	00013	00		
7	00014	21		
8	00017	10		

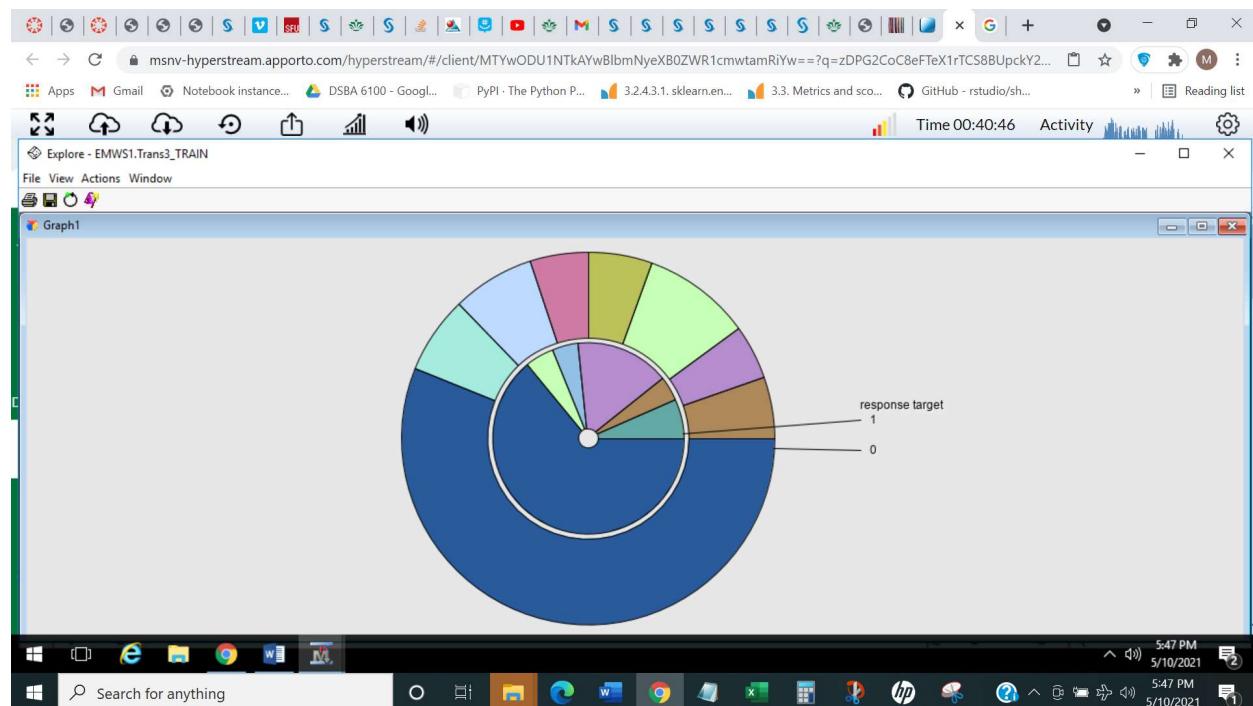
Select Chart Roles

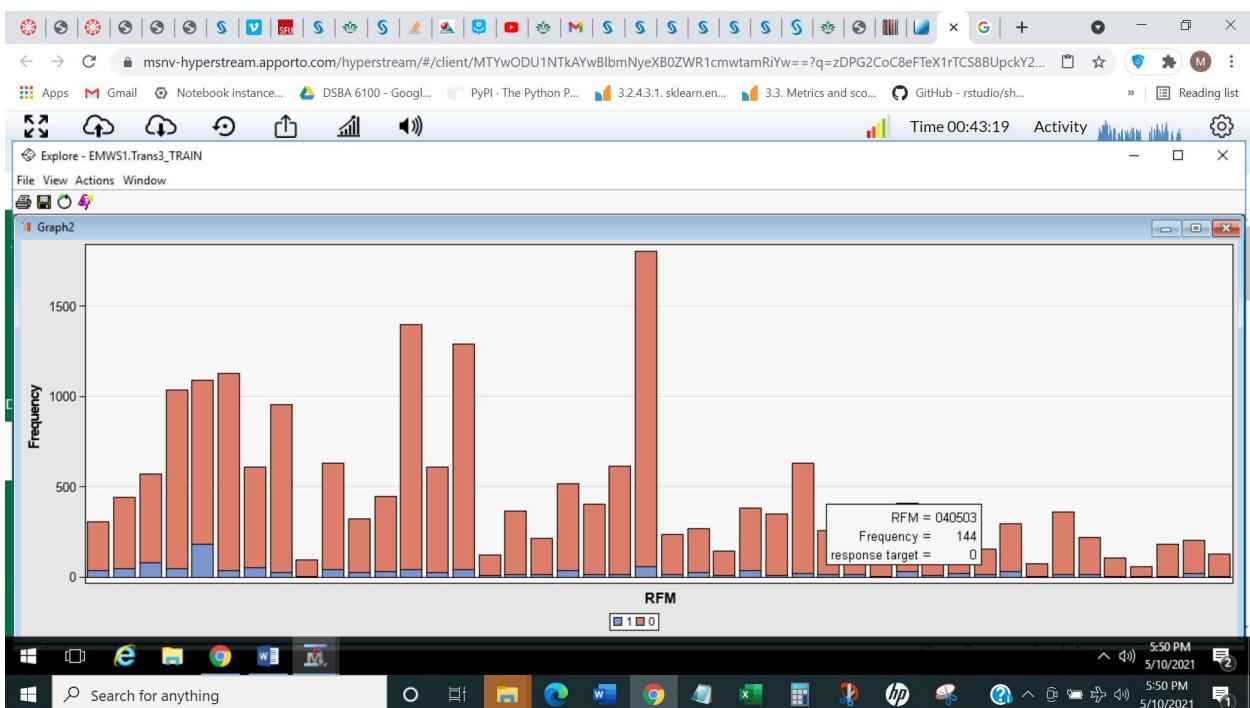
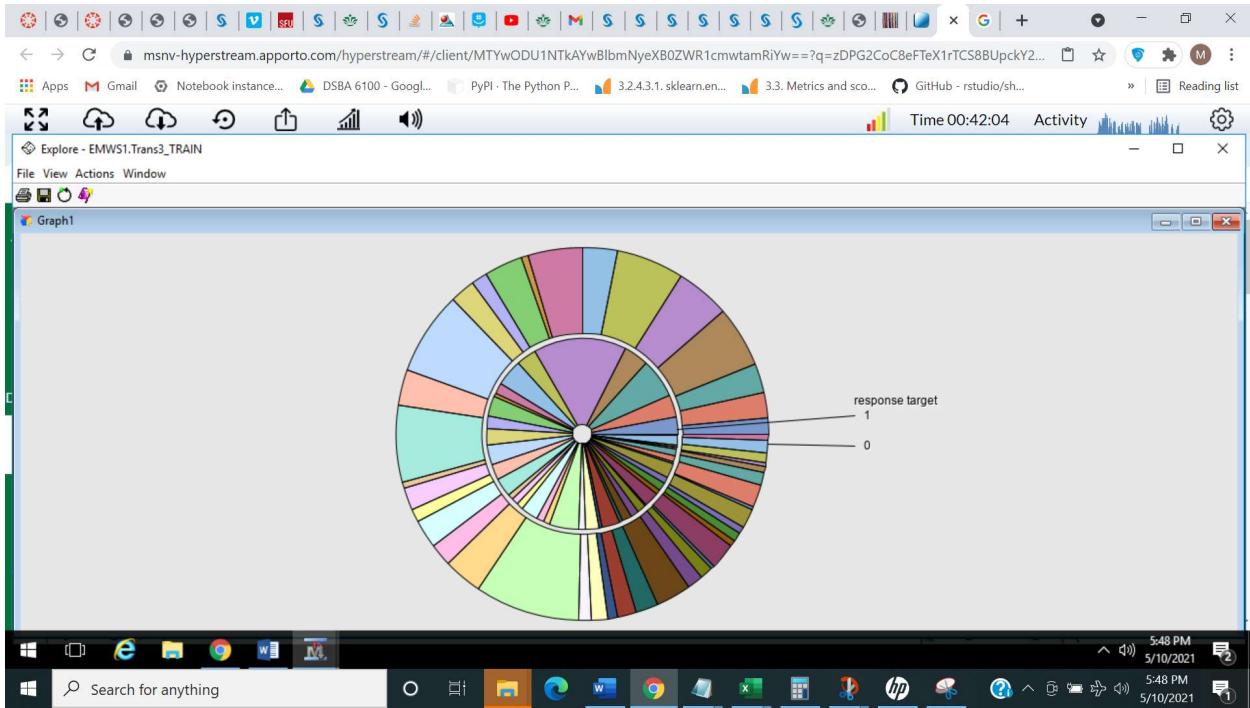
Use default assignments

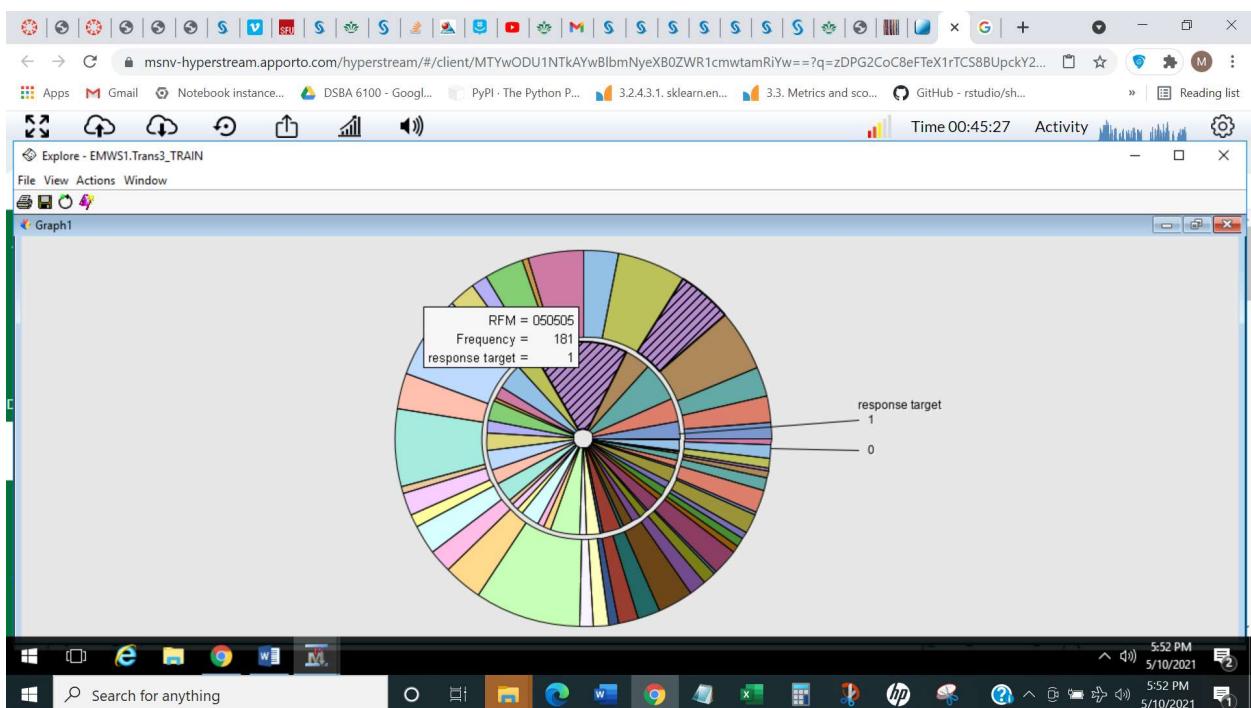
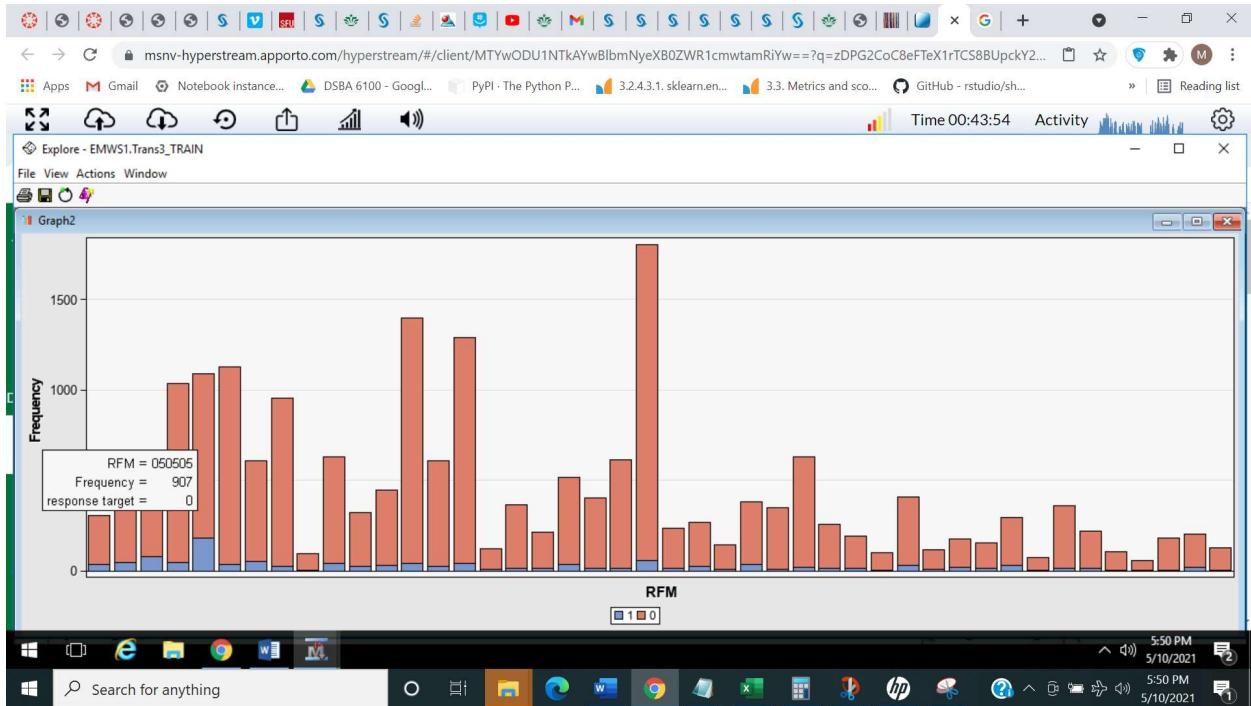
Variable	Role	Type	Description	Format
PCTL_DAYLAST_REV	Character	Transformed DAYLA...		
PCTL_DOLL24	Character	Transformed: \$ last 2...		
PCTL_FREQPRCH	Character	Transformed: lifetime...		
RESPOND	Group	Numeric	response target	BEST12.
RFM	Category	Character	RFM	
STATE	Character	state code	\$2.	
TENURE	Numeric	months since 1st	BEST12.	
TOTORDQ01	Numeric	tot orders 93Q1	BEST12.	
TOTORDQ02	Numeric	tot orders 93Q2	BEST12.	
TOTORDQ03	Numeric	tot orders 93Q3	BEST12.	
TOTORDQ04	Numeric	tot orders 93Q4	BEST12.	
TOTORDQ05	Numeric	tot orders 94Q1	BEST12.	
TOTORDQ06	Numeric	tot orders 94Q2	BEST12.	
TOTORDQ07	Numeric	tot orders 94Q3	BEST12.	
TOTORDQ08	Numeric	tot orders 94Q4	BEST12.	

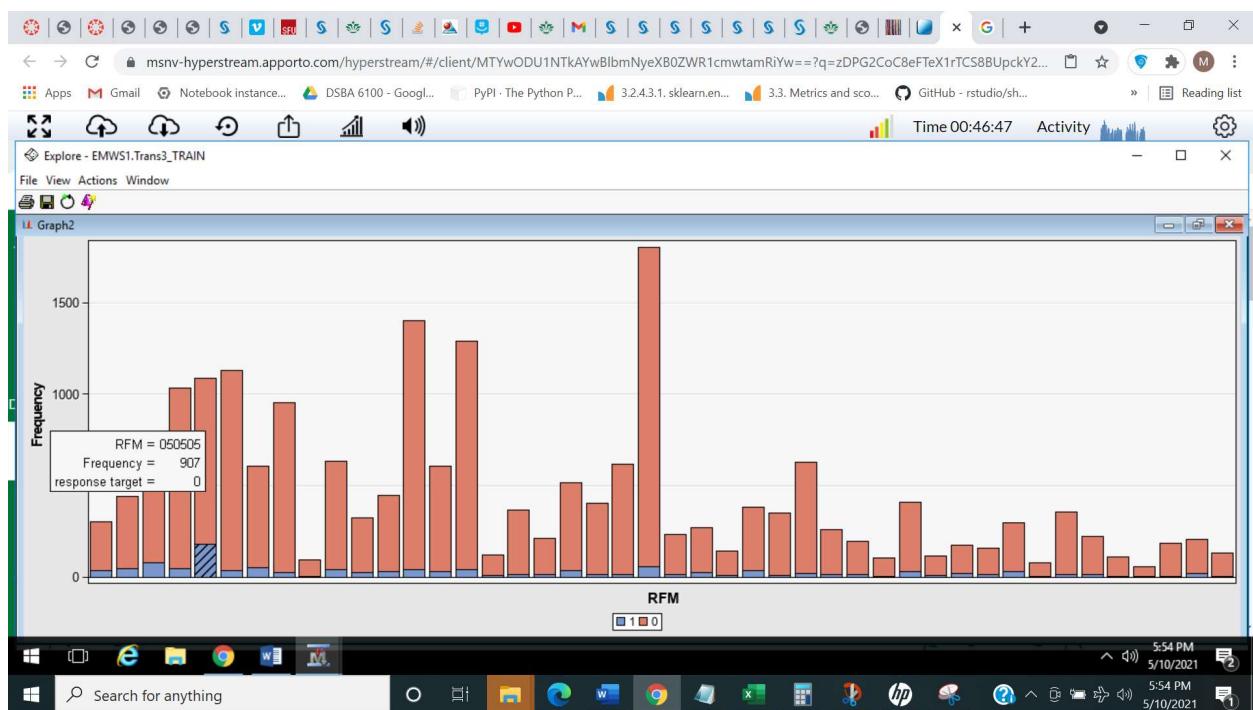
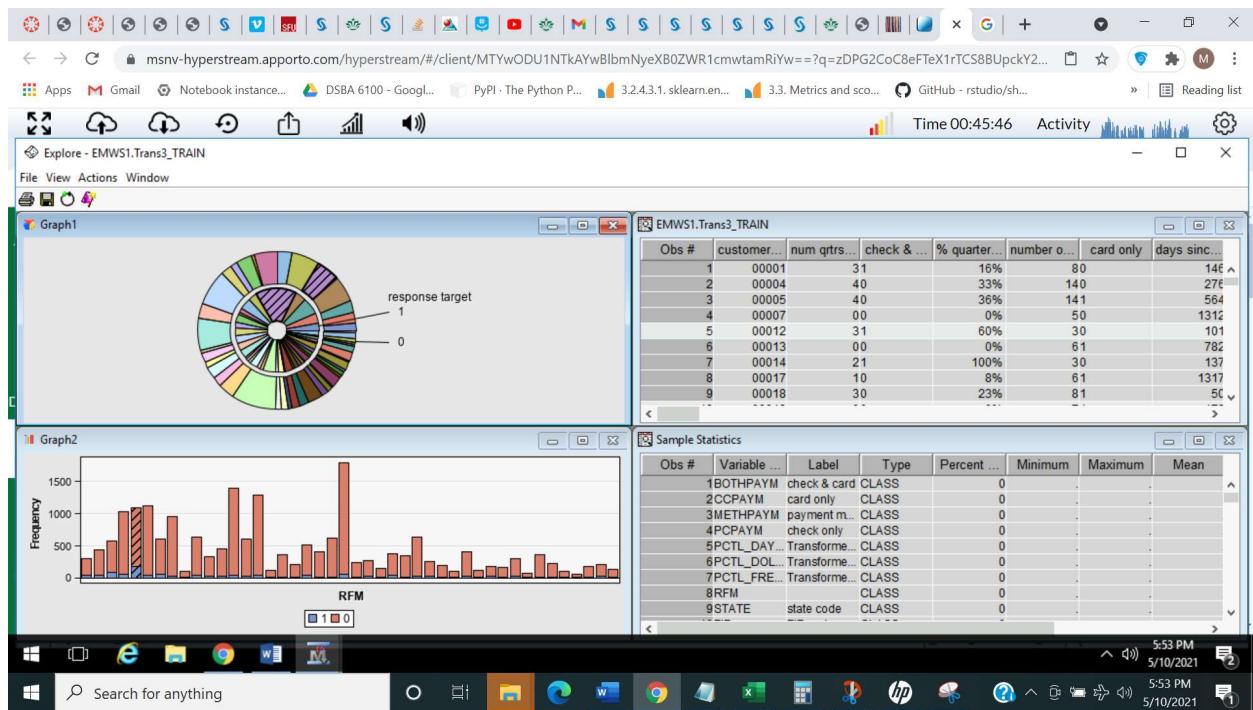
Allow multiple role assignments

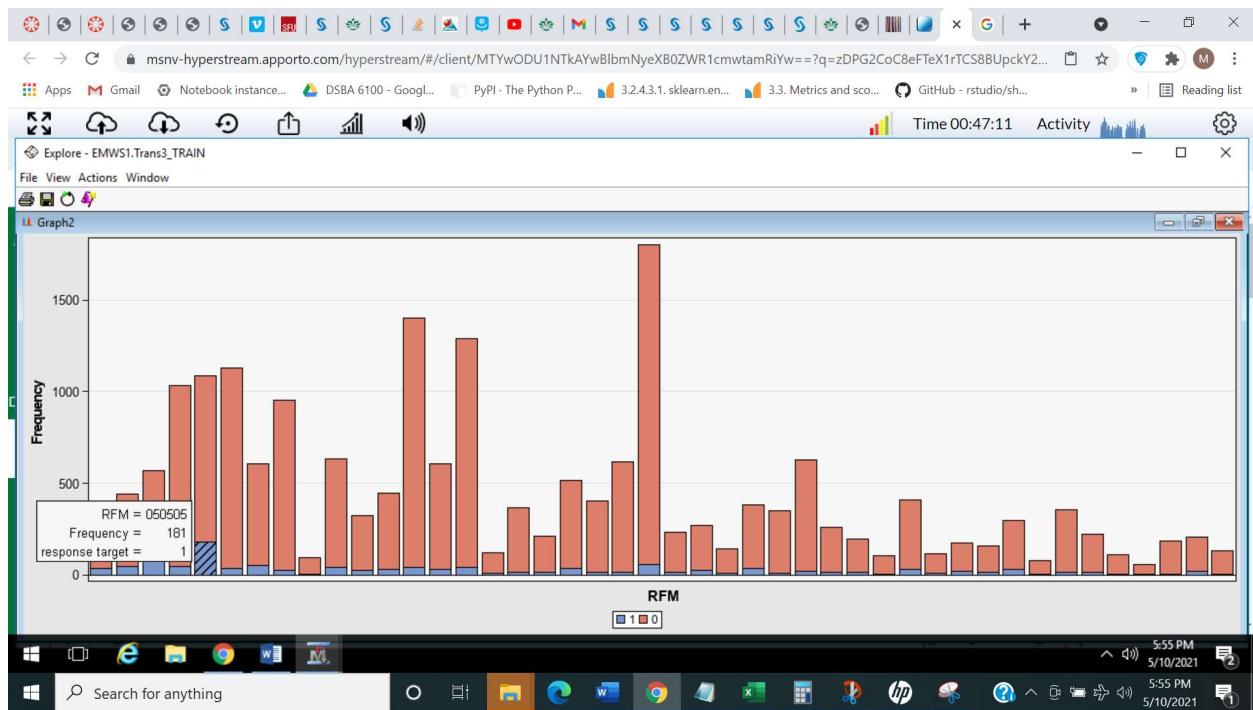
Cancel < Back Next > Finish





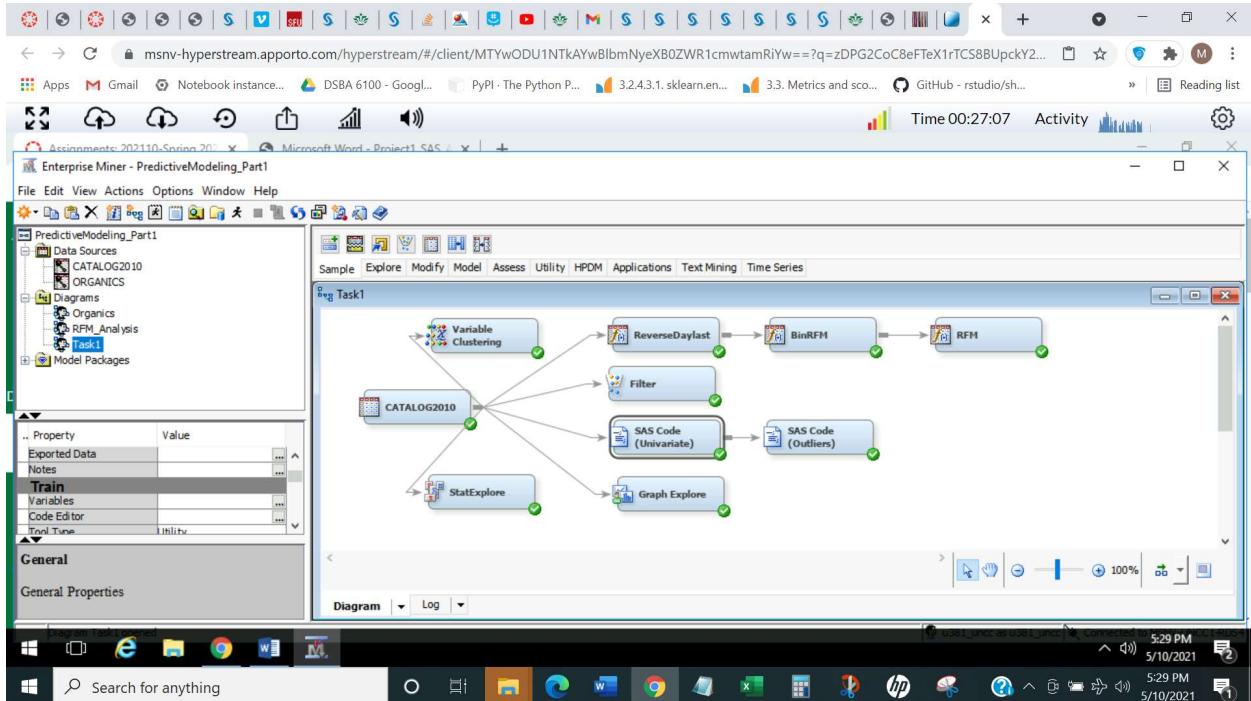






Task 1:**SAS Code (5points)**

- **Univariate Analysis**
- **Outlier**
- **Plotting**



You need to perform additional analysis on various variables and make a report

- 1) You might want to study which variables are highly correlated. If you find such variables you can suggest dimension reduction by dropping one of the variables.
- 6points**

The following variables are highly correlated based a correlation of greater than 0.8.
[\(<http://www.sfu.ca/~dsignori/buec333/lecture%2016.pdf>\)](http://www.sfu.ca/~dsignori/buec333/lecture%2016.pdf)

Variable 1	Variable 2	Correlation
days since last	months since last	0.999975251
months since last	days since last	0.999975251
avg \$ net demand	total \$ demand	0.993953254
total \$ demand	avg \$ net demand	0.993953254
avg \$ demand	tot \$ net demand	0.953178097
tot \$ net demand	avg \$ demand	0.953178097
tot units demand	total \$ demand	0.881179133

total \$ demand	tot units demand	0.881179133
avg \$ net demand	tot units demand	0.877361971
tot units demand	avg \$ net demand	0.877361971
lifetime orders	total \$ demand	0.815401727
total \$ demand	lifetime orders	0.815401727
avg \$ net demand	lifetime orders	0.812388543
lifetime orders	avg \$ net demand	0.812388543
lifetime orders	tot units demand	0.804471837
tot units demand	lifetime orders	0.804471837

Highly correlated variables may overly skew data analysis in specific directions. One remedy is to drop one of each pair of highly correlated variables. In the CATALOG2010 data, we would drop “months since last” and keep “days since last” because days is a more precise measure. Between variables “avg \$ net demand” and “total \$ demand”, we would drop “avg \$ net demand” because it is a comparison value and keep “total \$ demand” because it is a final demand value rather than a value for comparison. For the same reason, we would keep “avg \$ demand” instead of the comparison value of “tot \$ net demand”. When reviewing “tot unit demand” and “total \$ demand”, we have already indicated we will keep “total \$ demand” so “tot unit demand” can be dropped. The next pair is “avg \$ net demand” and “tot units demand”. We have already indicated that we will drop both. If we must retain one, it would be “tot units demand” because it represents a final concrete value. When comparing variables “lifetime orders” and “total \$ demand”, we have already indicated that we will keep “total \$ demand”. Between “avg \$net demand” and “lifetime orders”, we can indicated that we will drop both of these. However, if we must retain one it will be “lifetime orders”. In comparing “lifetime orders” and “tot units demand”, we will drop “tot units demand” because it is more highly correlated with another variable – “total \$ demand”.

2) You can study in if there are outliers in your variables 6points

After adding and running the filter node to the project, we have the following results:

Number Of Observations			
Data			
Role	Filtered	Excluded	DATA
TRAIN	20607	27749	48356

With the StatExplore Node, we generated the Interval Variables table which displays a table of summary statistics for interval variables including mean and standard deviation.

Data Role	Target	Target Level	Variable	Median	Missing	Non Missing	Minimum	Maximum	Mean	Standard Deviation	Skewness	Kurtosis	Role	Label	Scaled Mean Deviation	Max Dev
TRAIN	RESPOND	0	ORDERSIZE	0	0	45617	0	0	0	0	0	0	TARGET	dollar value...	-1	
TRAIN	RESPOND	1	ORDERSIZE	43.85	0	2739	5.4	510.72	62.77013	59.94363	2.771339	10.801821TARGET	dollar value...	16.65462		
TRAIN	RESPOND	0	DOLLARQ22	0	0	45617	0	670.2	4.49098	22.46453	9.046402	128.1413INPUT	tot \$ 98Q2	-0.08609		
TRAIN	RESPOND	1	DOLLARQ22	0	0	2739	0	1128.75	11.95974	46.40408	9.822622	161.0246INPUT	tot \$ 98Q2	1.433795		
TRAIN	RESPOND	0	TOTORQD22	0	0	45617	0	4	0.081154	0.296624	3.982015	18.54265INPUT	tot orders 9...	-0.07837		
TRAIN	RESPOND	1	TOTORQD22	0	0	2739	0	6	0.202994	0.522376	3.47717	17.47316INPUT	tot orders 9...	1.3053		
TRAIN	RESPOND	0	DEPT25	0	0	45617	0	186	1.634829	4.356396	8.898405	177.32281INPUT	food	-0.07335		
TRAIN	RESPOND	1	DEPT25	0	0	2739	0	126	3.919314	9.162664	5.516103	44.46609INPUT	food	1.221546		
TRAIN	RESPOND	0	DOLLARQ20	0	0	45617	0	768.85	8.064814	30.50179	7.515434	92.45299INPUT	tot \$ 97Q4	-0.07293		
TRAIN	RESPOND	1	DOLLARQ20	0	0	2739	0	908.5	19.25051	51.52793	6.036659	62.18658INPUT	tot \$ 97Q4	1.214567		
TRAIN	RESPOND	0	DOLLARQ18	0	0	45617	0	752.1	4.612252	22.39056	9.339538	144.8153INPUT	tot \$ 97Q2	-0.07228		
TRAIN	RESPOND	1	DOLLARQ18	0	0	2739	0	486.15	10.95618	36.27364	5.370204	39.47621INPUT	tot \$ 97Q2	1.20376		
TRAIN	RESPOND	0	TOTORQD20	0	0	45617	0	22	0.148782	0.423859	6.257399	174.5595INPUT	tot orders 9...	-0.06903		
TRAIN	RESPOND	1	TOTORQD20	0	0	2739	0	27	0.343556	0.849894	13.10508	362.1374INPUT	tot orders 9...	1.149715		
TRAIN	RESPOND	0	DOLL24	0	0	45617	0	2140.1	42.60035	86.8664	5.45256	57.27983INPUT	\$ last 24 m...	-0.06764		
TRAIN	RESPOND	1	DOLL24	43.9	0	2739	0	2433.5	97.16399	168.497	5.069334	42.82712INPUT	\$ last 24 m...	1.126547		
TRAIN	RESPOND	0	DOLLARQ21	0	0	45617	0	1084.3	4.432568	23.1125	12.7998	328.14881INPUT	tot \$ 98Q1	-0.06547		
TRAIN	RESPOND	1	DOLLARQ21	0	0	2739	0	542.05	9.914816	34.62575	6.405912	58.5225INPUT	tot \$ 98Q1	1.090368		
TRAIN	RESPOND	0	DOLLARQ19	0	0	45617	0	559.2	4.171569	21.59825	9.664066	140.6688INPUT	tot \$ 97Q3	-0.06279		
TRAIN	RESPOND	1	DOLLARQ19	0	0	2739	0	677.88	9.105451	33.8163	7.024176	81.5345INPUT	tot \$ 97Q3	1.045693		
TRAIN	RESPOND	0	TOTORQD18	0	0	45617	0	5	0.089046	0.319699	4.142542	21.77453INPUT	tot orders 9...	-0.0623		

We moved this table into Excel to calculate the Q1 & Q3, IQR and outlier values for each variable.

1	Data Role	Target	Target Lev	Variable	Median	Missing	Non Miss	Minimum	Maximum	Mean	Standard Deviation	IQ1	Q3	IRQ	Q1 - 1.5(IQR)	Q3 + 1.5(IQR)	Skewness	Kurtosis
2	TRAIN	RESPOND	0	ACTBUY	1	0	45617	0	10	0.951597	1.1202	-0.1686	2.071797	2.2404	-3.529202145	5.43239613	1.864435	5
3	TRAIN	RESPOND	1	ACTBUY	1	0	2739	0	10	1.52647	1.611549	-0.08508	3.138018	3.223098	-4.919725601	7.97266463	1.559801	2
4	TRAIN	RESPOND	0	BUYPROP	0.1	0	45617	0	1	0.18395	0.253858	-0.06991	0.437808	0.507716	-0.831481204	1.199381951	1.91111	3
5	TRAIN	RESPOND	1	BUYPROP	0.181818	0	2739	0	1	0.265676	0.293457	-0.02778	0.559133	0.586913	-0.908149761	1.439502674	1.278582	0
6	TRAIN	RESPOND	0	CATALOGCNT	3	0	45617	1	24	3.678453	3.052072	0.626381	6.730525	6.104144	-8.529834304	15.88674072	1.591797	2
7	TRAIN	RESPOND	1	CATALOGCNT	4	0	2739	1	27	5.220884	3.99266	1.228224	9.213543	7.985319	-10.74975524	21.19152231	1.260401	1
8	TRAIN	RESPOND	0	DAYLAST	782	0	45617	0	8265	1204.939	1236.789	-31.8496	2441.728	2473.578	-374.2216042	6152.094552	1.741862	3
9	TRAIN	RESPOND	1	DAYLAST	362	0	2739	0	7859	759.7401	981.9968	-222.25	1741.737	1963.994	-3168.24702	4687.727122	2.517462	
10	TRAIN	RESPOND	0	DEPT01	0	0	45617	0	59	0.469386	1.686947	-1.21756	2.156333	3.373893	-6.278399852	7.217172678	7.020023	0
11	TRAIN	RESPOND	1	DEPT01	0	0	2739	0	49	0.917853	3.032121	-2.11427	3.949975	6.064243	-11.21063275	13.04633921	6.733124	€
12	TRAIN	RESPOND	0	DEPT02	0	0	45617	0	24	0.275928	1.101996	-0.82607	1.377924	2.203993	-4.132057982	4.68391365	5.496759	4
13	TRAIN	RESPOND	1	DEPT02	0	0	2739	0	19	0.564074	1.76376	-1.19969	2.327834	3.52752	-6.490965289	7.619114249	4.44541	2
14	TRAIN	RESPOND	0	DEPT03	0	0	45617	0	60	1.025692	2.662534	-1.63684	3.688227	5.325069	-9.624445181	11.67582953	4.713144	3
15	TRAIN	RESPOND	1	DEPT03	0	0	2739	0	54	2.085433	4.697419	-2.61199	6.782852	9.394838	-16.70424309	20.87510837	4.338569	2
16	TRAIN	RESPOND	0	DEPT04	0	0	45617	0	45	0.645921	2.007705	-1.55485	2.846696	4.401549	-8.157176645	9.449019599	6.159116	€
17	TRAIN	RESPOND	1	DEPT04	0	0	2739	0	47	1.396495	3.791373	-2.39488	5.187864	7.582746	-13.76899695	16.56198709	4.582803	2
18	TRAIN	RESPOND	0	DEPT05	0	0	45617	0	28	0.520552	1.479894	-0.95934	2.000446	2.959788	-5.399024703	6.4401278	3.882806	
19	TRAIN	RESPOND	1	DEPT05	0	0	2739	0	19	0.874407	2.014991	-1.14058	2.889398	4.029982	-7.185556981	8.934370416	3.080609	1
20	TRAIN	RESPOND	0	DEPT06	0	0	45617	0	32	0.825591	1.917107	-1.09152	2.742698	3.834213	-6.842835614	8.494018287	3.328319	1
21	TRAIN	RESPOND	1	DEPT06	0	0	2739	0	29	1.241329	2.520296	-1.27897	3.761625	5.040592	-8.839855472	11.32251338	3.423343	2

We compared this info to the Univariate Histogram for each variable generated from SAS Node Code.

File Edit Run View

Macro Tree:

- .. Macro
- Train**
 - Utility
 - EM_REGISTER
 - EM_REPORT
 - EM_DATA2CODE
 - EM_DECDATA
 - EM_CHECKMACRO
 - EM_CHECKSETINIT
 - EM_ODSLISTON

Report Code:

```

proc univariate data=&em_import_data noplay;
  class &em_dec_target;
  histogram %em_interval_input;
run;

```

Report Graphs:

The histogram displays the percentage distribution of 'days since last' across various intervals. The x-axis represents 'days since last' from 0 to 8100, and the y-axis represents 'Percent' from 0 to 10.0. The distribution is skewed right, with the highest frequency occurring between 0 and 450 days.

Days Since Last (Bin)	Percent
0 - 150	~6.0%
150 - 300	~12.0%
300 - 450	~11.0%
450 - 600	~7.0%
600 - 750	~6.0%
750 - 900	~5.5%
900 - 1050	~5.5%
1050 - 1200	~5.5%
1200 - 1350	~4.5%
1350 - 1500	~3.5%
1500 - 1650	~3.0%
1650 - 1800	~3.0%
1800 - 1950	~2.5%
1950 - 2100	~2.0%
2100 - 2250	~2.0%
2250 - 2400	~1.5%
2400 - 2550	~1.5%
2550 - 2700	~2.0%
2700 - 2850	~1.5%
2850 - 3000	~1.5%
3000 - 3150	~1.5%
3150 - 3300	~1.5%
3300 - 3450	~1.0%
3450 - 3600	~1.0%
3600 - 3750	~1.0%
3750 - 3900	~1.0%
3900 - 4050	~1.0%
4050 - 4200	~1.0%
4200 - 4350	~1.0%
4350 - 4500	~1.0%
4500 - 4650	~1.0%
4650 - 4800	~1.0%
4800 - 4950	~1.0%
4950 - 5100	~1.0%
5100 - 5250	~1.0%
5250 - 5400	~1.0%
5400 - 5550	~1.0%
5550 - 5700	~1.0%
5700 - 5850	~1.0%
5850 - 6000	~1.0%
6000 - 6150	~1.0%
6150 - 6300	~1.0%
6300 - 6450	~1.0%
6450 - 6600	~1.0%
6600 - 6750	~1.0%
6750 - 6900	~1.0%
6900 - 7050	~1.0%
7050 - 7200	~1.0%
7200 - 7350	~1.0%
7350 - 7500	~1.0%
7500 - 7650	~1.0%
7650 - 7800	~1.0%
7800 - 7950	~1.0%
7950 - 8100	~1.0%

Using the table and histograms for each variable, we were able to determine which variables have outliers that need to be removed.

```
IF ACTBUY > 7.97266463 then delete;  
IF DAYLAST > 6152.094552 then delete;  
IF DEPT05 > 8.934370416 then delete;  
IF DEPT07 > 1.999186007 then delete;  
IF DEPT10 > 11.19664554 then delete;  
IF DEPT11 > 4.347513557 then delete;  
IF DEPT15 > 6.712963144 then delete;  
IF DEPT16 > 6.355999274 then delete;  
IF DEPT18 > 3.775612457 then delete;  
IF DEPT19 > 3.556790332 then delete;  
IF DEPT20 > 1.63582488 then delete;  
IF DEPT21 > 1.21526434 then delete;  
IF DEPT27 > 8.91633293 then delete;  
IF DOLINDEA > 199.4817483 then delete;  
IF DOLLARQ01 > 104.7142525 then delete;  
IF DOLLARQ04 > 112.6889834 then delete;  
IF DOLLARQ05 > 119.9368083 then delete;  
IF DOLLARQ06 > 108.8219948 then delete;  
IF DOLNETDA > 191.7379656 then delete;  
IF MONLAST > 202.0299063 then delete;  
IF TOTORDQ01 > 1.847747185 then delete;  
IF TOTORDQ03 > 1.373660894 then delete;  
IF TOTORDQ05 > 1.958658868 then delete;  
IF TOTORDQ09 > 1.813416952 then delete;  
IF TOTORDQ10 > 1.86231867 then delete;  
IF TOTORDQ12 > 2.48326094 then delete;  
IF TOTORDQ13 > 1.985993095 then delete;  
IF TOTORDQ14 > 1.784962061 then delete;  
IF TOTORDQ16 > 2.658483902 then delete;  
IF TOTORDQ17 > 2.261503116 then delete;  
IF TOTORDQ19 > 1.964440261 then delete;  
IF UNITSLAP > 105.2575375 then delete;  
IF UNTLANPO > 12.00516903 then delete;
```

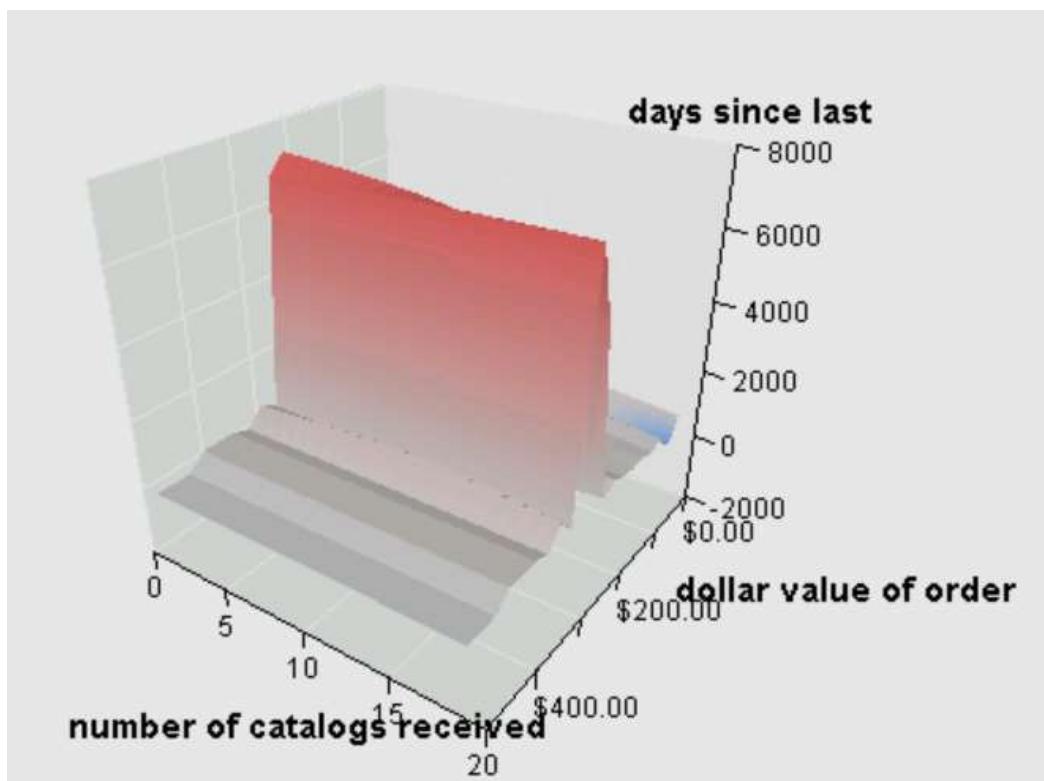
Report Code

```

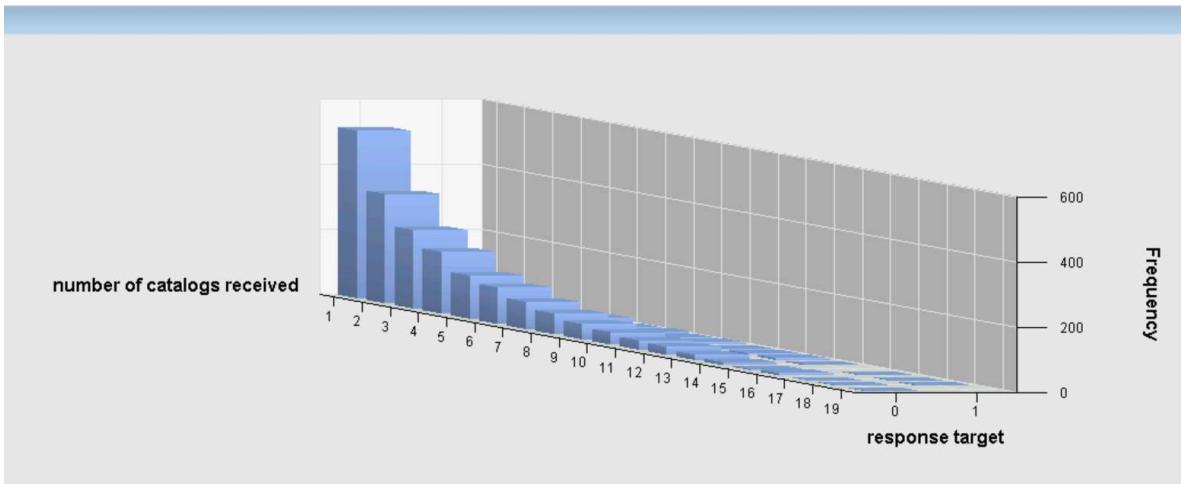
data ExportedCatalog2010; set &EM_IMPORT_DATA;
IF ACTBUY > 7.97266463 then delete;
IF DAYLAST > 6152.094552 then delete;
IF DEPT05 > 8.934370416 then delete;
IF DEPT07 > 1.999186007 then delete;
IF DEPT10 > 11.19664554 then delete;
IF DEPT11 > 4.347513557 then delete;
IF DEPT15 > 6.712963144 then delete;

```

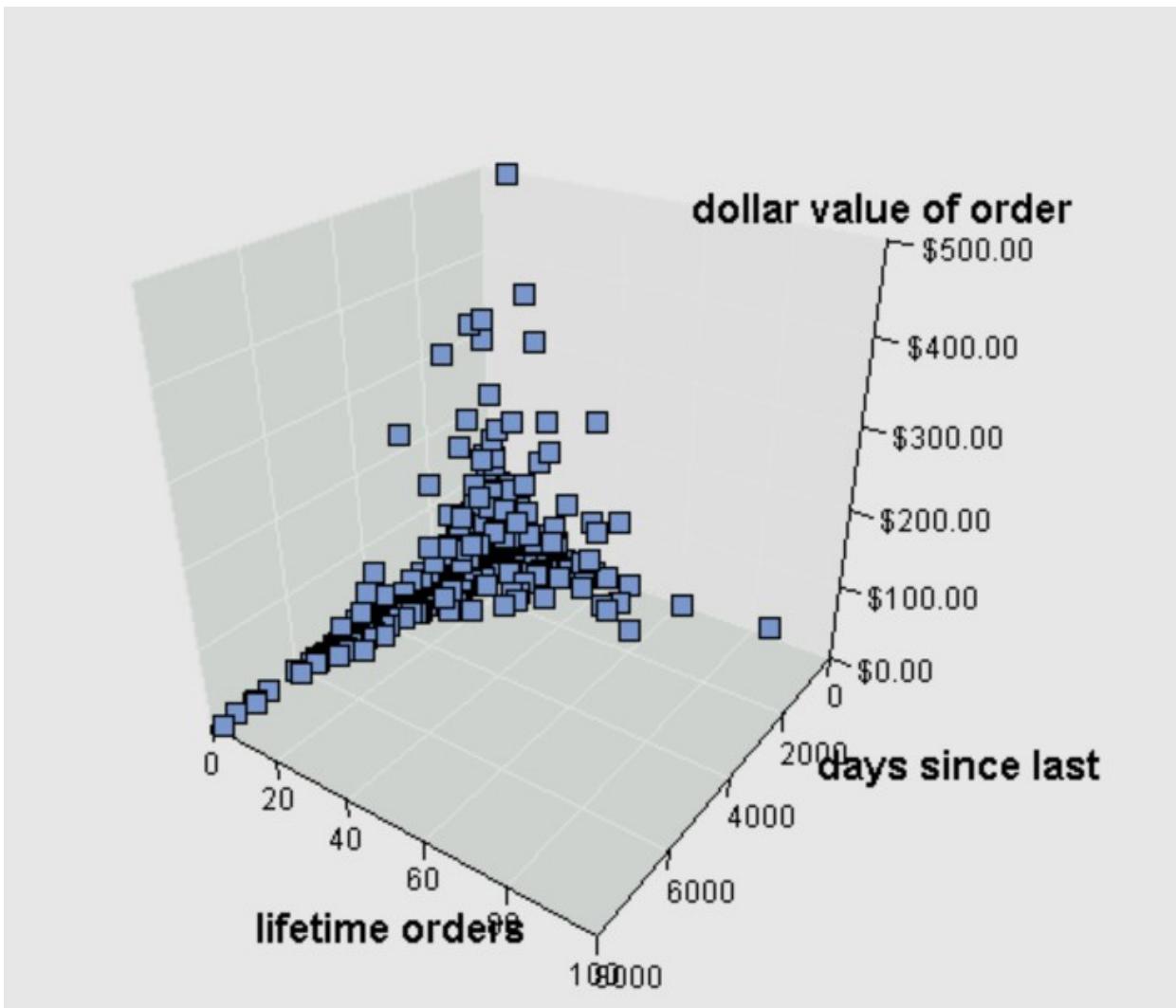
- 3) You can make 3-D plots to get a better sense of how independent variables affect the dependent variables 6points



This 3D surface area graph shows the relationship between Number of catalogs received, days since last purchase and the target variable of dollar value of order. When the days since last order increases part 6000, then dollar value of the order approaches \$200.00.



This bar chart with both a category and series variable indicates that most customers received 1 catalog. The number of catalogs received does not increase the response target.



This scatter plot in 3 dimensions indicates that the dollar value of the order does not increase

with the number of lifetime order or the days since last order. Customers with lifetime orders > 30 tend to have ordered in the last 2000 days. The highest dollar value orders are from customers in the lower half of lifetime number of order and the lower half of days since last order. This informs us that order dollar value does not increase with the number of lifetime orders or the number of days since last order.

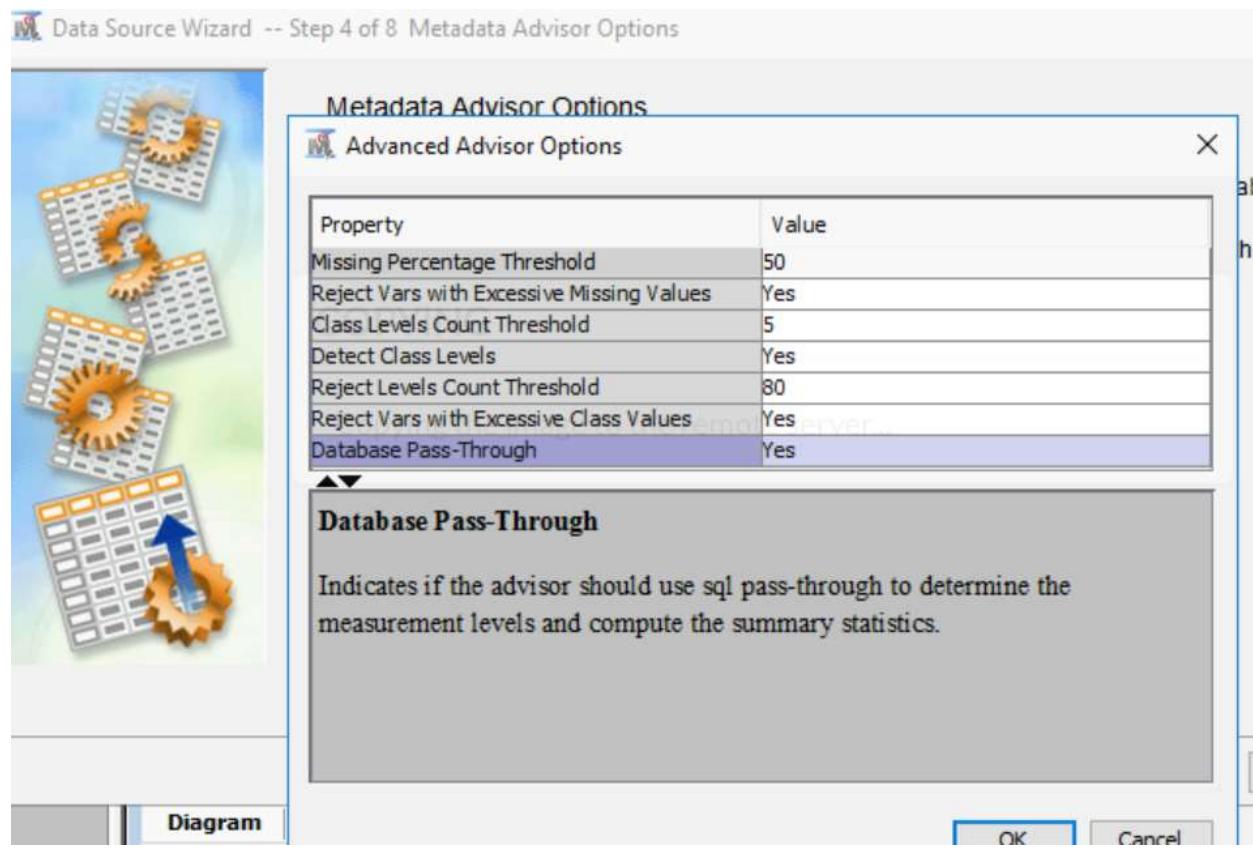
Task 2: 10 points

A national veterans' organization seeks to better target its solicitations for donation. By soliciting only the most likely donors, less money is spent on solicitation efforts and more money is available for charitable concerns. Solicitations involve sending a small gift to an individual and including a request for a donation. Gifts to donors include mailing labels and greeting cards.

The organization has more than 3.5 million individuals in its mailing database. These individuals are classified by their response behaviors to previous solicitation efforts. Of particular interest is the class of individuals identified as *lapsing donors*. These individuals made their most recent donation between 12 and 24 months ago. The organization seeks to rank its lapsing donors based on their responses to a greeting card mailing sent in June of 1997. (The charity calls this the 97NK Campaign.) With this ranking, a decision can be made to either solicit or ignore a lapsing individual in the June 1998 campaign.

1. RFM Analysis of Charity Direct Mail Data

- a) Define PVA97NK as a data source in SAS Enterprise Miner. Use the Advanced Metadata Advisor options to customize the following:
- Change the Class Levels Count Threshold from 20 to 5.
 - Change the Reject Levels Count Threshold from 20 to 80.
 - Reject the variable TargetD



Data Source Wizard -- Step 5 of 8: Column Metadata

Name	Role	Level	Report	Order	Drop
americancards0	Input	Interval	No		No
GiftCntCardAll	Input	Interval	No		No
GiftTimeFirst	Input	Interval	No		No
GiftTimeLast	Input	Interval	No		No
ID	ID	Nominal	No		No
PromCnt12	Input	Interval	No		No
PromCnt36	Input	Interval	No		No
PromCntAll	Input	Interval	No		No
PromCntCard12	Input	Interval	No		No
PromCntCard36	Input	Interval	No		No
PromCntCardAll	Input	Interval	No		No
StatusCat96NK	Input	Nominal	No		No
StatusCatStarAll	Input	Binary	No		No
TargetB	Target	Binary	No		No
TargetD	Rejected	Interval	No		No

Show code Explore Refresh Summary < Back Next >

- b) Create a new diagram and transform the R, F, and M variables as described previously to create four bins of each variable. Concatenate them to create an RFM variable.

Add Transformation

Property	Value
Name	Recency
Type	Numeric
Length	8
Format	
Level	Interval
Label	
Role	Input
Report	No

Formula:
TRANS_0 =
-1*GiftTimeLast

Build... OK Cancel

Results - Node: Recency Diagram: Charity Direct Mail

File Edit View Window

Transformations Statistics

Source	Method	Variable Name	Formula	Number of Levels	Non Missing	Missing	Minimum	Maximum	Mean	Standard Deviation	Skewness	Ku
Input	Original	GiftTimeLast			9686	0	4	27	18.00217	4.073549	-0.77805	
Output	Formula	Recency	-1*GiftTime...		9686	0	-27	-4	-18.0022	4.073549	0.778047	

Output

Enterprise Miner - PredictiveModeling_Part1

File Edit View Actions Options Window Help

PredictiveModeling_Part1

- Data Sources
 - CATALOG2010
 - ORGANICS
 - PVA97NK
- Diagrams
 - Charity Direct Mail
 - Organics
 - RFM_Analysis
 - Task1
- Model Packages

.. Property Value

Train

Variables	...
Formulas	...
Interactions	...
SAS Code	...
Default Methods	...

General

ORGANICS
PVA97NK

Diagrams

- Charity Direct Mail
- Organics
- RFM_Analysis
- Task1

Model Packages

.. Property Value

Sample Properties

Method	Random
Size	Max
Random Seed	12345

Frequency

Age

Columns: Label Mining Basic Statistics

Name	Role	Level	Method	Number of Bins
DemAge	Input	Interval	DEFALKT	4.0

Inputs

Preview plot for the selected variable "Monetary".

Name	Type	Length	Format	Level	Formula	Label	Role	Report
Recency	Numeric	8		Interval	$-1*GiftTimeLast$		Input	No
Monetary	Numeric	8		Interval	$GiftAvgAll*GiftCount$		Input	No

Charity Direct Mail

```

graph LR
    PVA97NK[PVA97NK] --> Transform[Transform Variables]
  
```

Enterprise Miner - PredictiveModeling_Part1

File Edit View Actions Options Window Help

PredictiveModeling_Part1

- Data Sources
 - CATALOG2010
 - ORGANICS
 - PVA97NK
- Diagrams
 - Charity Direct Mail
 - Organics
 - RFM_Analysis
 - Task1
- Model Packages

Property

Exported Data

Notes

Train

Variables

Formulas

Interactions

Variables - Trans2

(none) ▾ not Equal to

Columns: Label Mining

Name	Method	Number of Bins	Role	Level
DemHomeOwner	Default	4	Input	Binary
DemPctVeterans	Default	4	Input	Interval
DemMedIncome	Default	4	Input	Interval
DemCluster	Default	4	Input	Nominal
DemAge	Default	4	Input	Interval
DemGender	Default	4	Input	Nominal
GiftCnt36	Default	4	Input	Interval
GiftAvgCard36	Default	4	Input	Interval
GiftAvgLast	Default	4	Input	Interval
GiftAvgAll	Default	4	Input	Interval
GiftAvg36	Default	4	Input	Interval
GiftCntAll	Quantile	5	Input	Interval
Recency	Quantile	5	Input	Interval
Monetary	Quantile	5	Input	Interval

msnv-hyperstream.apporto.com/hyperstream/#/client/MTYwODU1NTkAYwBlbmNyeXB0ZWR1cmwtamRiYw==?q=zDPG2CoC8eFTeX1rTCS8BUpkcY2...

Time 01:19:55 Activity

Explore - EMWS4.Trans2_TRAIN

File View Actions Window

Sample Properties

Property	Value
Rows	Unknown
Columns	33
Library	EMWS4
Member	TRANS2_TRAIN
Type	VIEW
Sample Method	Random
Fetch Size	Max

Sample Statistics

Obs #	Variable ...	Label	Type	Percent ...	Minimum	Maximum	Mean
1	DemCluster	Demographic CLASS	CLASS	0	.	.	.
2	DemGender	Gender CLASS	CLASS	0	.	.	.
3	DemHome...	Home Owner CLASS	CLASS	0	.	.	.
4	ID	Control Num... CLASS	CLASS	0	.	.	.
5	PCTL_Gift...	Transformed CLASS	CLASS	0	.	.	.
6	PCTL_Mon...	Transformed CLASS	CLASS	0	.	.	.
7	PCTL_Rec...	Transformed CLASS	CLASS	0	.	.	.
8	StatusCat9...	Status Category CLASS	CLASS	0	.	.	.
9	DemAge	Age VAR	VAR	24.8503	0	87	59.15084

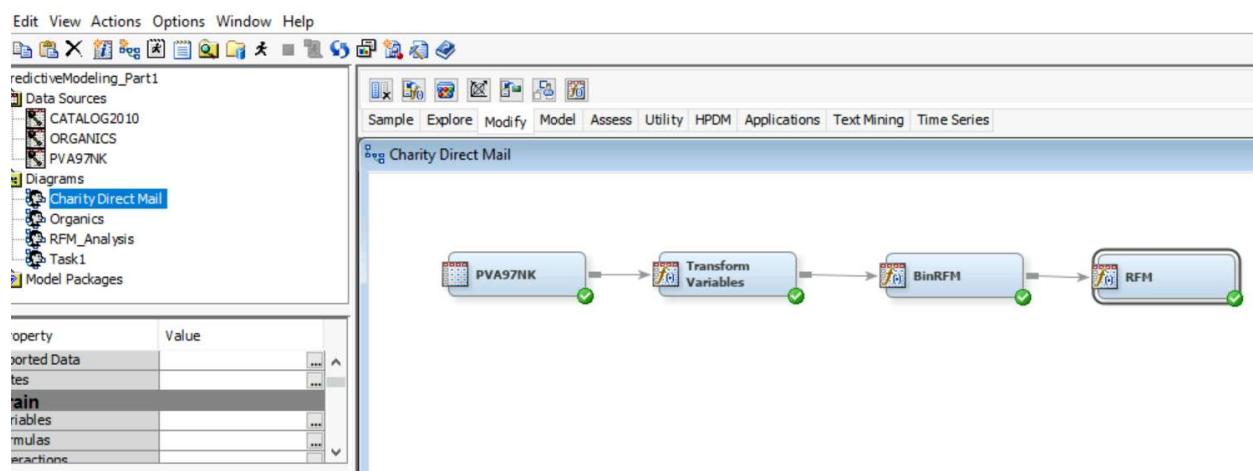
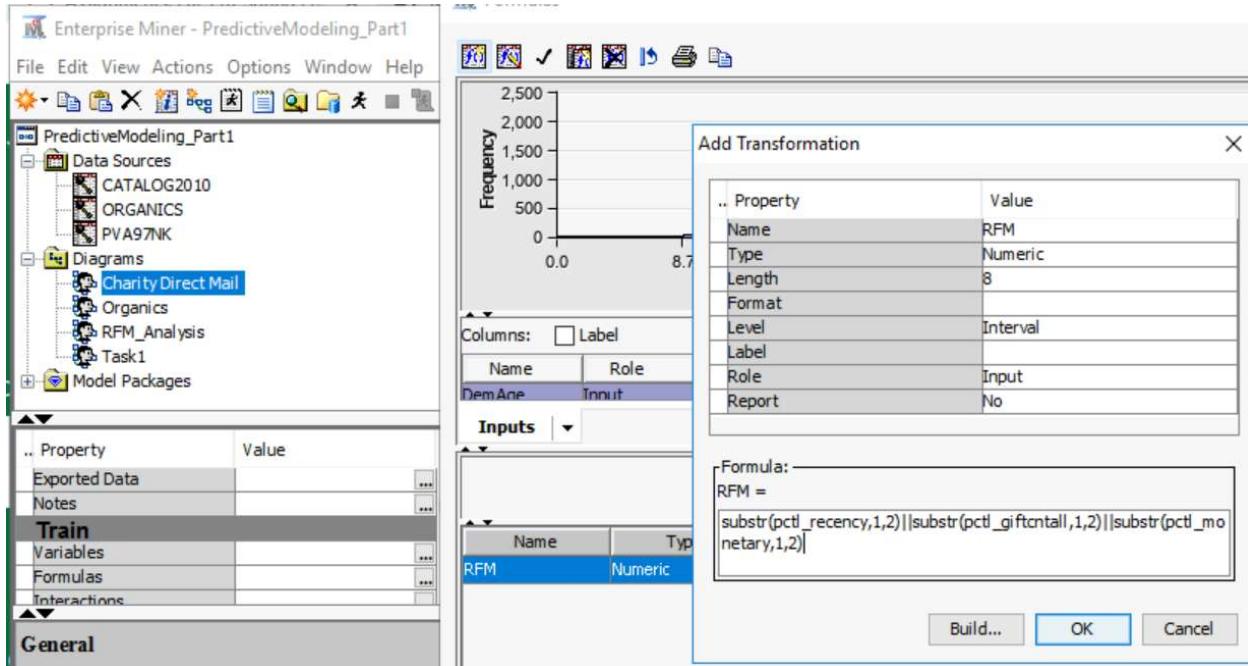
EMWS4.Trans2_TRAIN

io...	Promotio...	Promotio...	Status C...	Status C...	Demogra...	Age	Gender	Home Ow...	Median H...	Percent ...	Median In...	Recency	Monetary	Transfor...	Transfor...	Transfor...
3	8	13A	000	023	.F	U	\$0	0	\$0	-21	37023.6	02:36:55.01	02-21-18			
5	5	24A	023	023	.F	U	\$186.800	85	\$0	-26	127.04036-10	04:59:56.1...	01:low-21			
5	11	22S	100	100	.M	U	\$87.600	36	\$38.750	-18	152.930516-high	05:150.96...	03-18-17			
2	6	16E	100	100	.M	U	\$139.200	27	\$38.942	-9	10204:10-16	04:59:56.1...	05-16-high			
4	7	6F	035	035	.53M	U	\$168.100	37	\$71.509	-21	2001low-3	01:low-36	02-21-18			
5	10	22S	100	100	.47M	H	\$253.100	0	\$92.514	-22	90.9704:10-16	03:65.01-9...	01:low-21			
6	16	18A	035	035	.58M	H	\$234.700	22	\$72.868	-17	52023.6	02:36:55.01	04-17-16			
4	4	15A	008	008	.U	U	\$207.000	44	\$0	-18	46023.6	02:36:55.01	03-18-17			
6	6	22A	035	035	.F	U	\$137.300	32	\$0	-17	84.99023.6	03:65.01-9...	04-17-16			

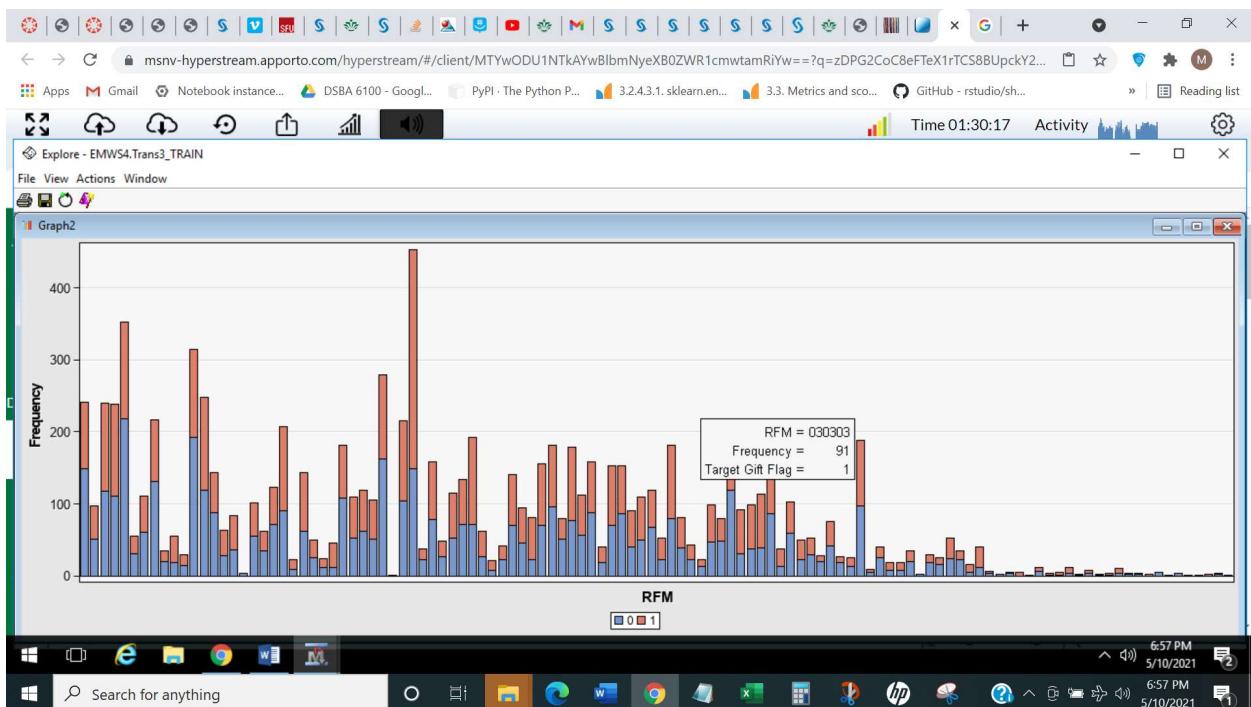
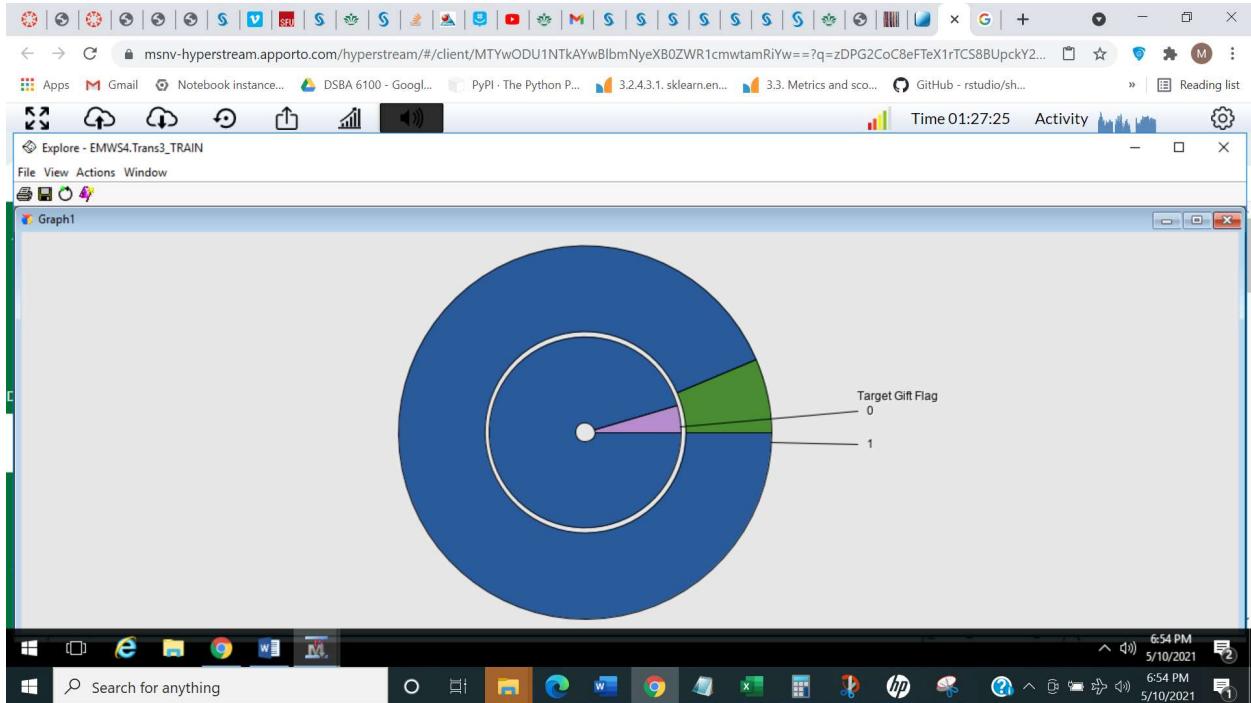
6:44 PM 5/10/2021

Search for anything

6:44 PM 5/10/2021



c) Explore the data and perform graphical RFM analysis using a grouped pie chart and a stacked bar chart.



d) Calculate response rate for 040404 and 030303 group?

$$040404 = 54/51+54 = 51.4\%$$

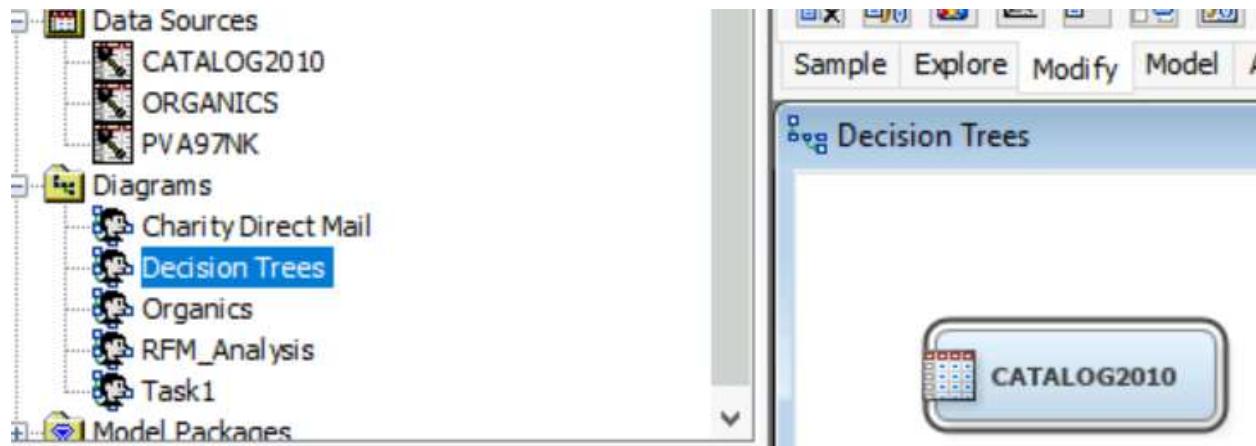
$$030303 = 91/97+91 = 48.4\%$$

e) Each promotional mailing (request for a gift) costs \$2.3, and the average donation is about \$21. What is the break-even response rate for this promotion? Do any RFM cells exceed this response rate? Remember to account for the fact that in the population, 95% of mailings are not responded to, while this sample is oversampled to 50% responders and 50% non-responders.

Break-even is 10.95%. Many RFM cells exceed that rate including 040404 & 030303.

Part 2: Predictive Modeling: Decision Trees

CATALOG CASE STUDY: CONSTRUCTING A DECISION TREE PREDICTIVE MODEL



The screenshot shows the SAS Enterprise Miner interface with the "Variables - Ids" tab selected in the palette on the right.

The "Variables" palette includes the following sections:

- Filter: (none) dropdown, "not" checkbox, "Equal to" dropdown, and a search input field.
- Columns: "Label" checkbox, "Mining" checkbox, and "Basic" checkbox.

The main area displays a table of variables:

Name	Role	Level	Report	Order	Drop	Lower Limit
DOLLARQ19	Input	Interval	No		No	.
DOLLARQ20	Input	Interval	No		No	.
DOLLARQ21	Input	Interval	No		No	.
DOLLARQ22	Input	Interval	No		No	.
DOLNETDA	Input	Interval	No		No	.
DOLNETDT	Input	Interval	No		No	.
DTBUYLST	Rejected	Interval	No		No	.
DTBUYORG	Rejected	Interval	No		No	.
FREQPRCH	Input	Interval	No		No	.
METHPAYM	Input	Nominal	No		No	.
MONLAST	Input	Interval	No		No	.
ORDERSIZE	Rejected	Interval	No		No	.
PCPAYM	Input	Binary	No		No	.
RESPOND	Target	Binary	No		No	.

PVA97NK

- Diagrams
 - Charity Direct Mail
 - Decision Trees**
 - Organics
 - RFM_Analysis
 - Task1
- Model Packages

Property	Value
Data Set Allocations	
Training	2.0
Validation	1.0
Test	0.0

Decision Trees

```

graph LR
    Catalog["CATALOG2010"] --> DataPartition["Data Partition"]
    DataPartition --> DecisionTree["Decision Tree"]
  
```

Results - Node: Data Partition Diagram: Decision Trees

File Edit View Window

Output

```

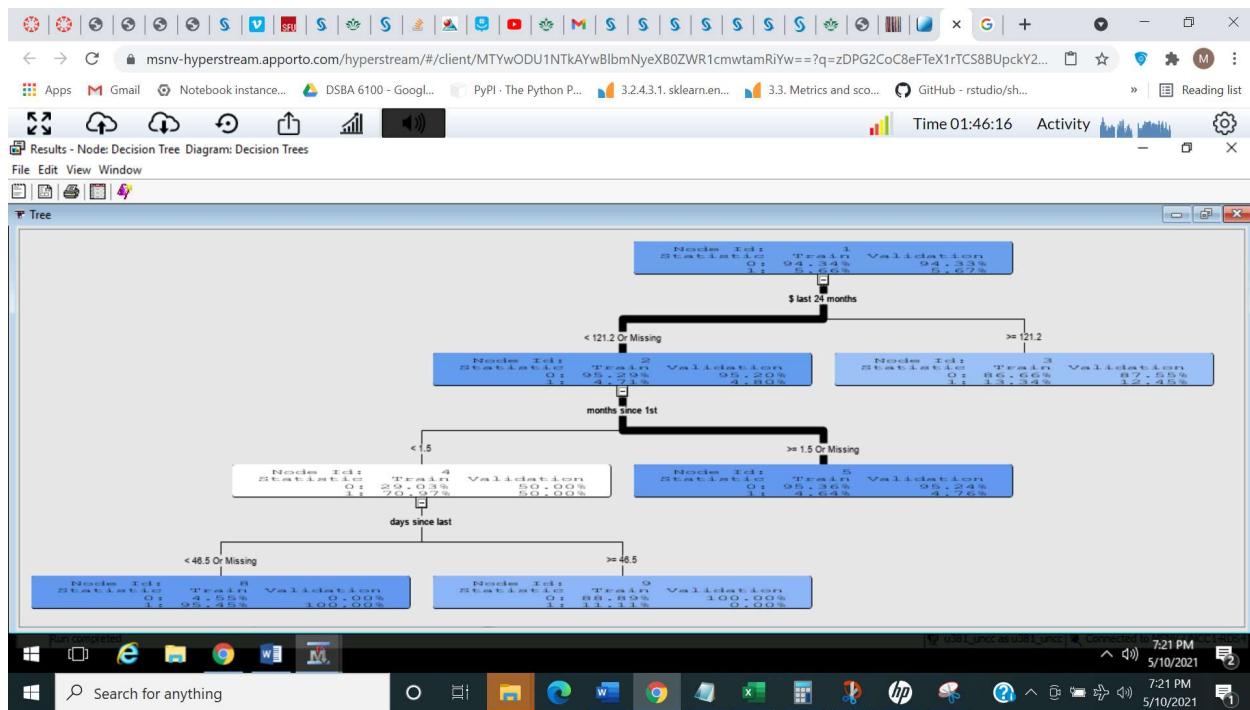
50  Summary Statistics for Class Targets
51
52  Data=DATA
53
54      Numeric   Formatted   Frequency
55  Variable   Value       Value     Count   Percent   Label
56
57  RESPOND    0          0         45617  94.3358  response target
58  RESPOND    1          1         2739   5.6642   response target
59
60
61  Data=TRAIN
62
63      Numeric   Formatted   Frequency
64  Variable   Value       Value     Count   Percent   Label
65
66  RESPOND    0          0         30410  94.3385  response target
67  RESPOND    1          1         1825   5.6615   response target
68
69
70  Data=VALIDATE
71
72      Numeric   Formatted   Frequency
73  Variable   Value       Value     Count   Percent   Label
  
```

7:19 PM 5/10/2021 7:19 PM 5/10/2021

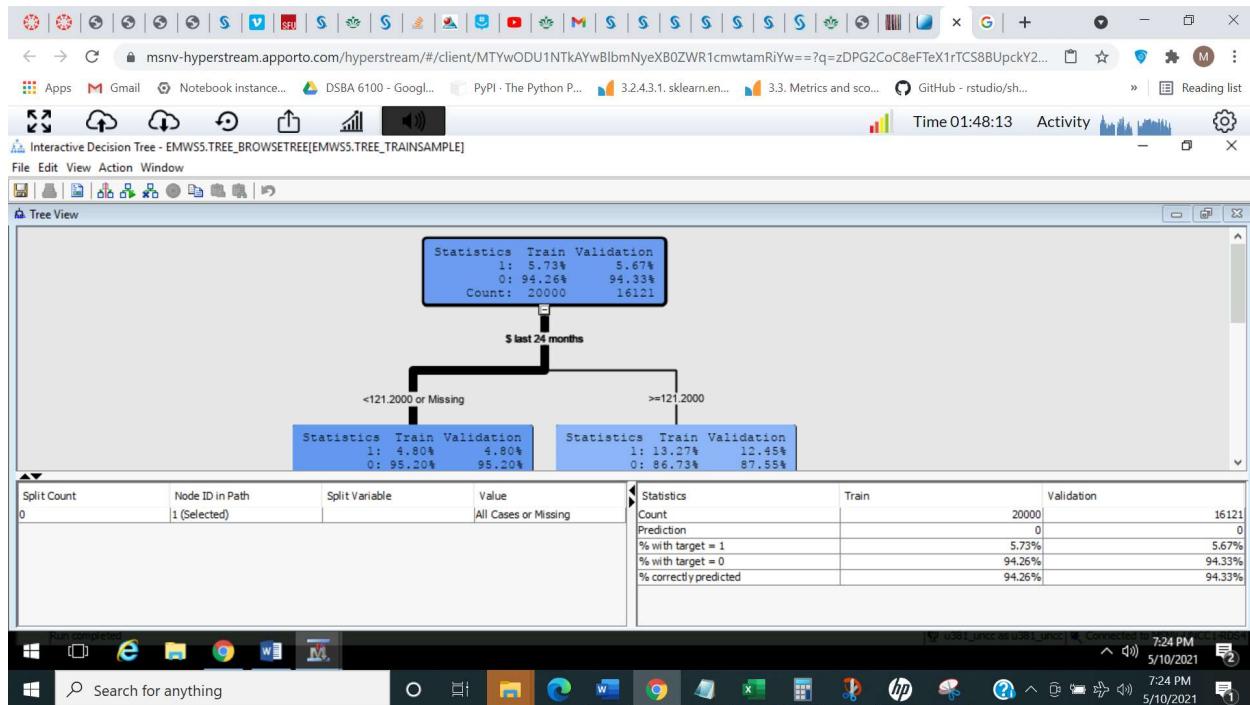
Decision Trees

```

graph LR
    Catalog["CATALOG2010"] --> DataPartition["Data Partition"]
    DataPartition --> DecisionTree["Decision Tree"]
  
```



Interactive Decision Tree (Self-Study)



Interactive Decision Tree - EMWS5.TREE_BROWSETREE[EMWS5.TREE_TRAINSAMPLE]

File Edit View Action Window

Tree View

Split Node 1

Target Variable: RESPOND

Variable	Variable Description	-Log(p)	Branches
DOLL24	\$ last 24 months	59.6643	2
DAYLAST	days since last	52.6652	2
MONLAST	months since last	51.9173	2
UNITSIDD	tot units demand	48.3687	2
FREQPRCH	lifetime orders	45.6174	2

DOLL24 - Interval Split Rule

Target Variable: RESPOND

Assign missing values to:

- A specific branch: 1
- A separate missing values branch
- All branches

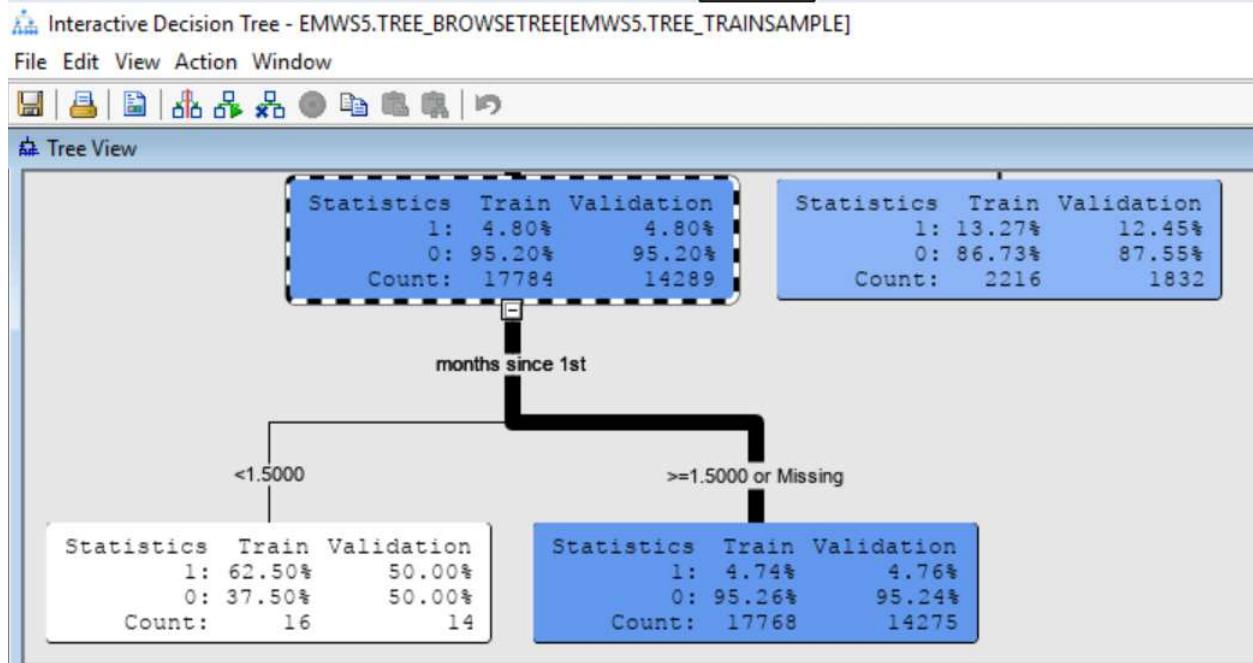
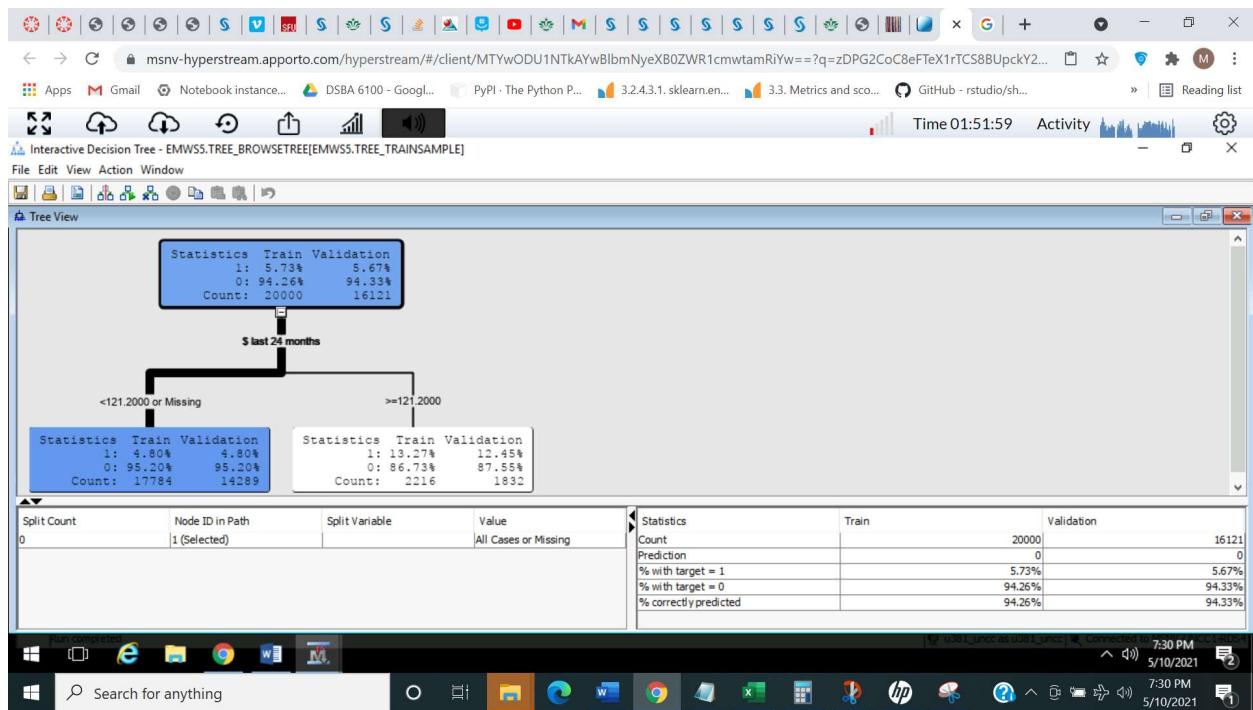
Branches:

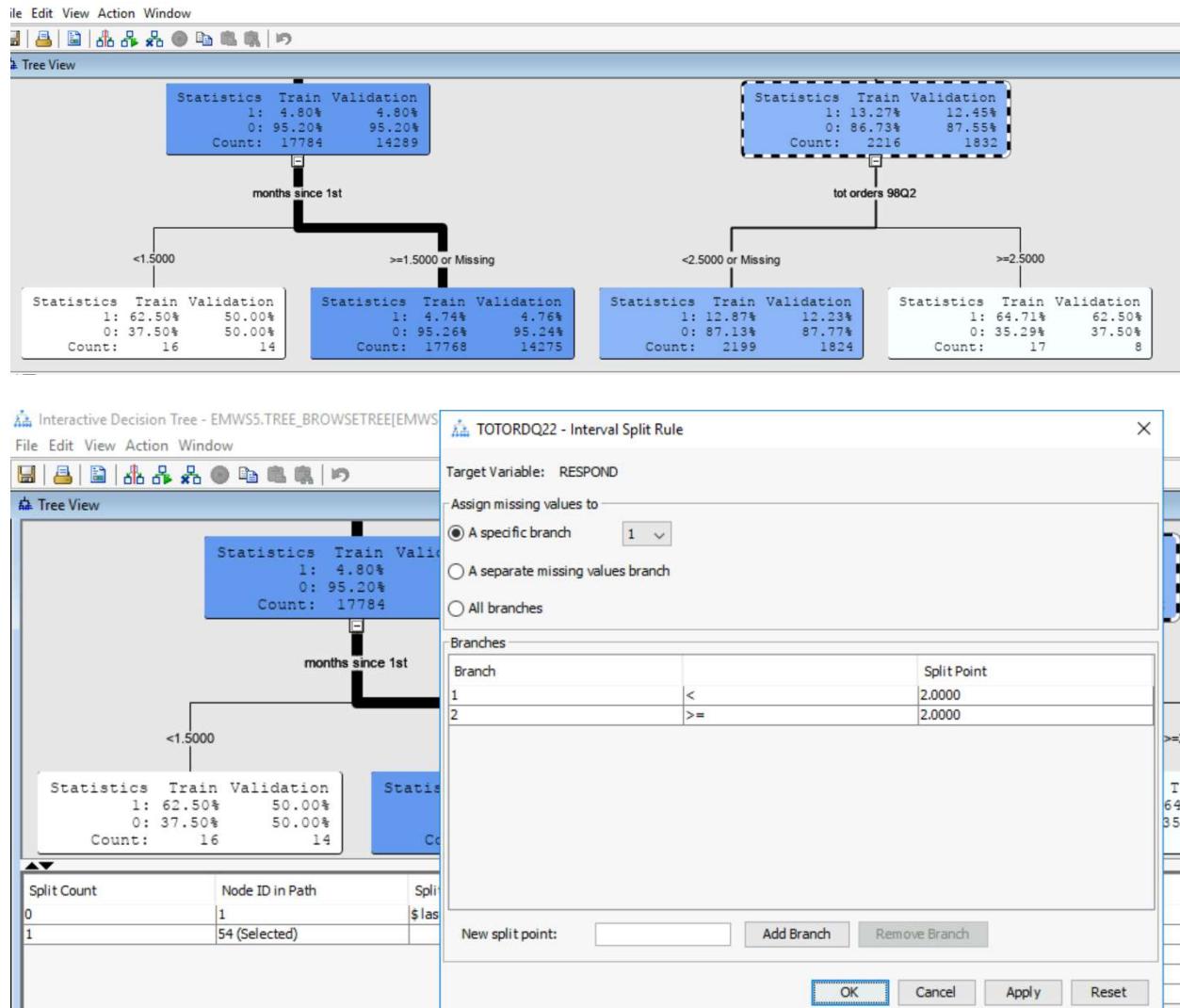
Branch	Split Point
1	< 121.2000
2	>= 121.2000

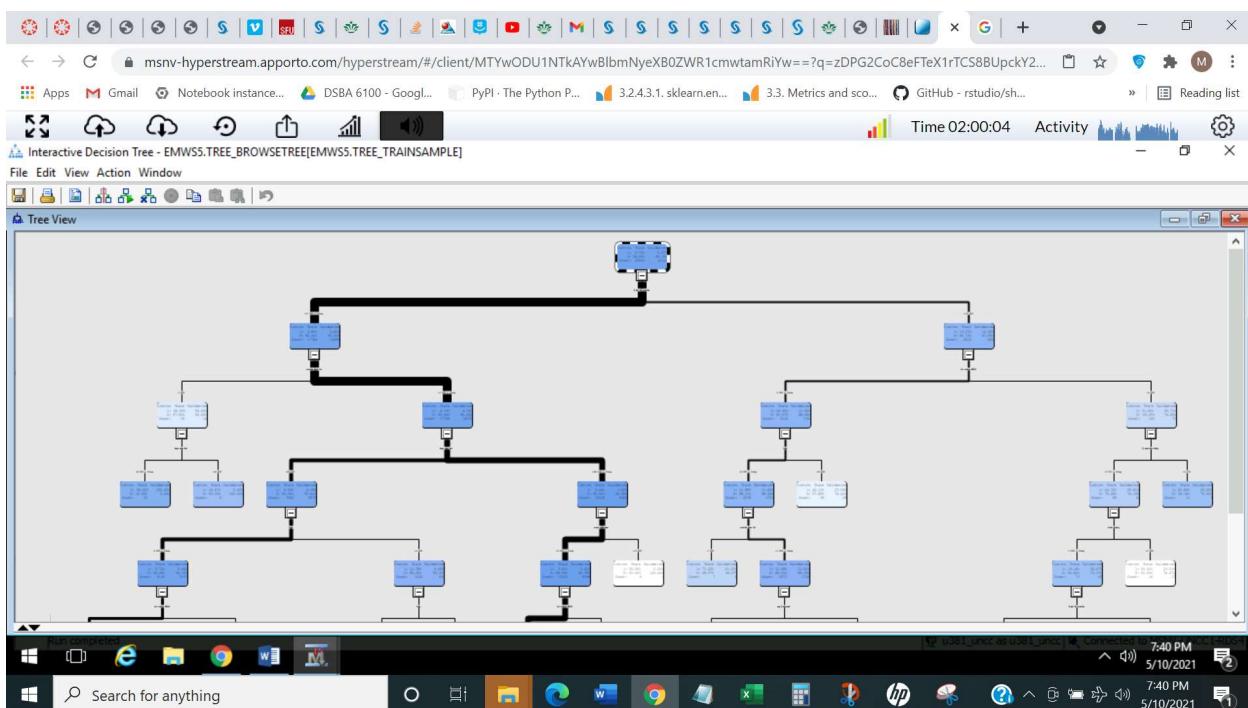
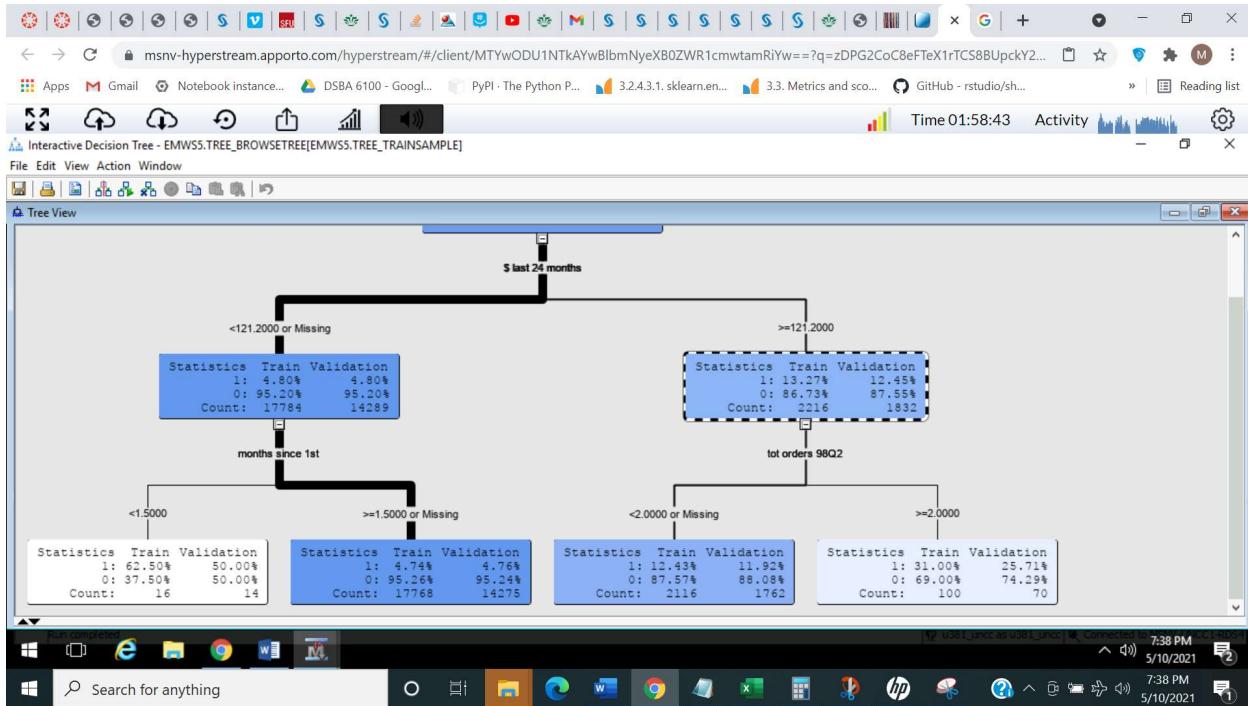
New split point: Add Branch Remove Branch

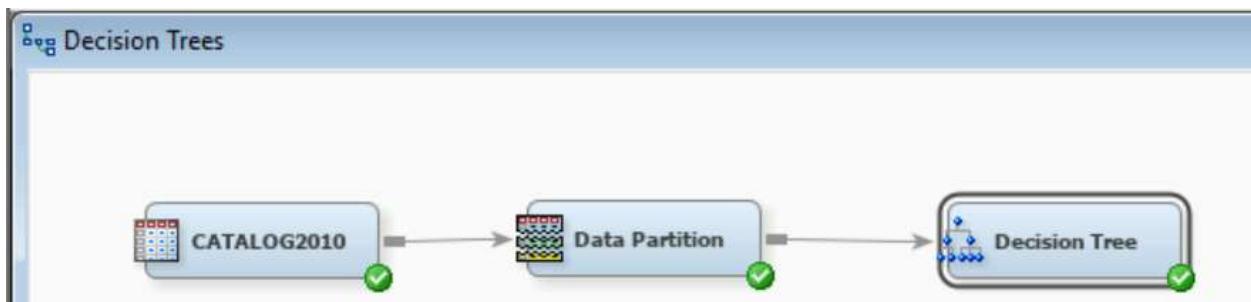
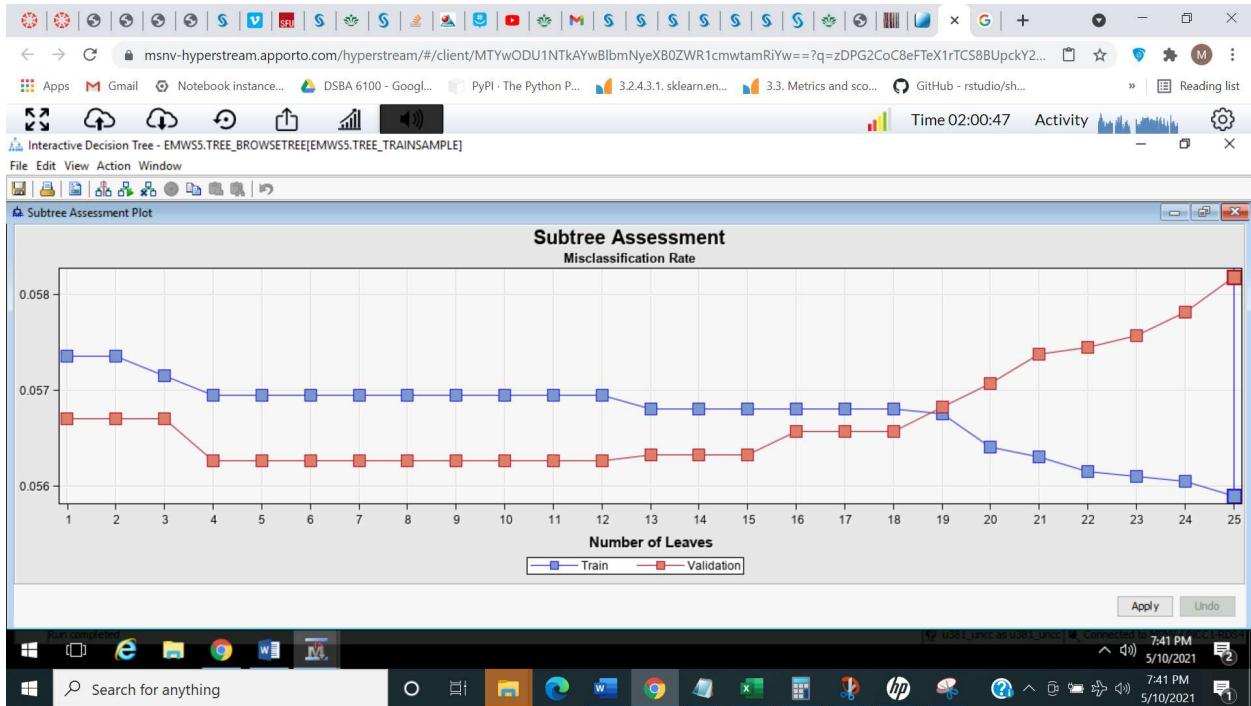
OK Cancel Apply Reset

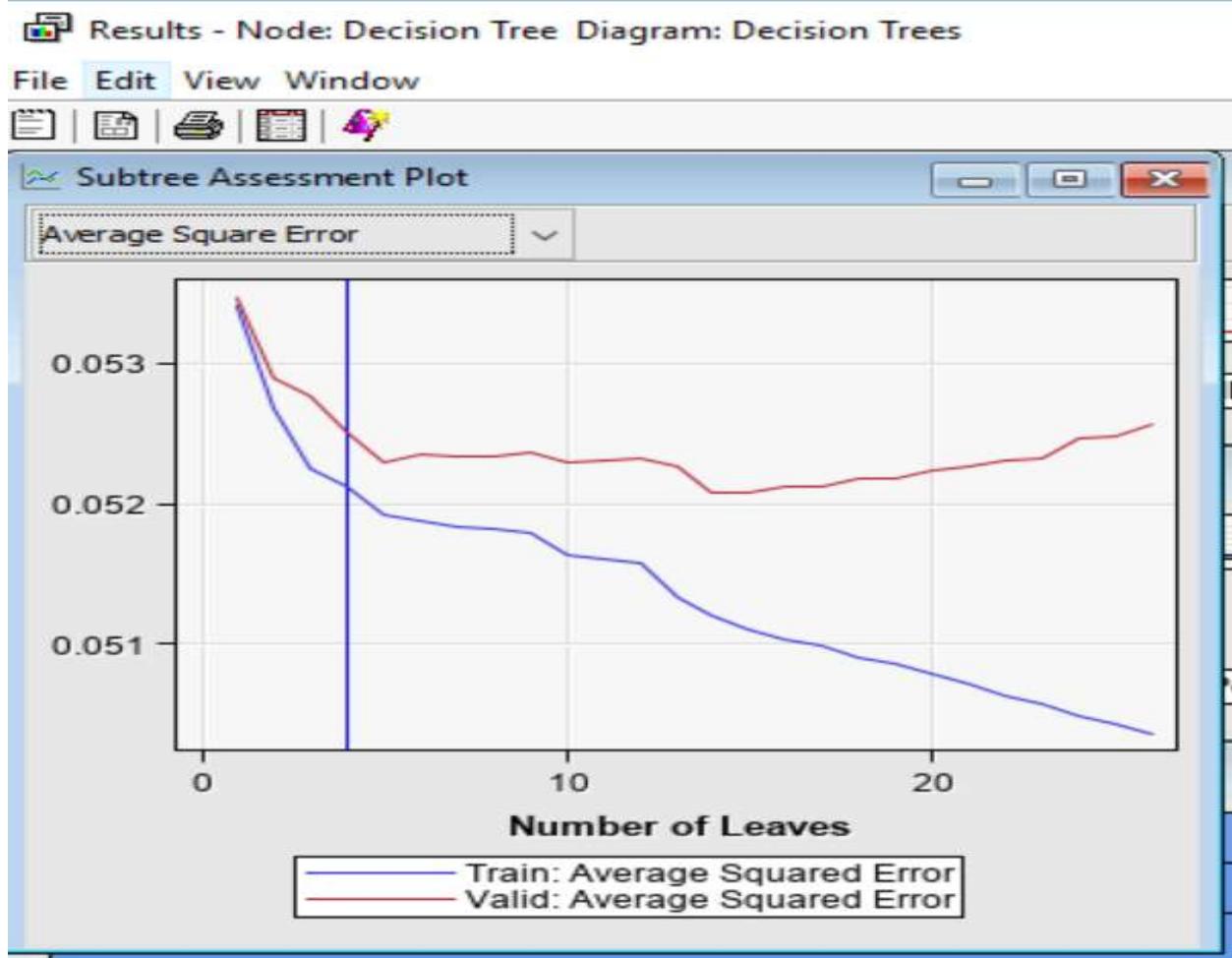
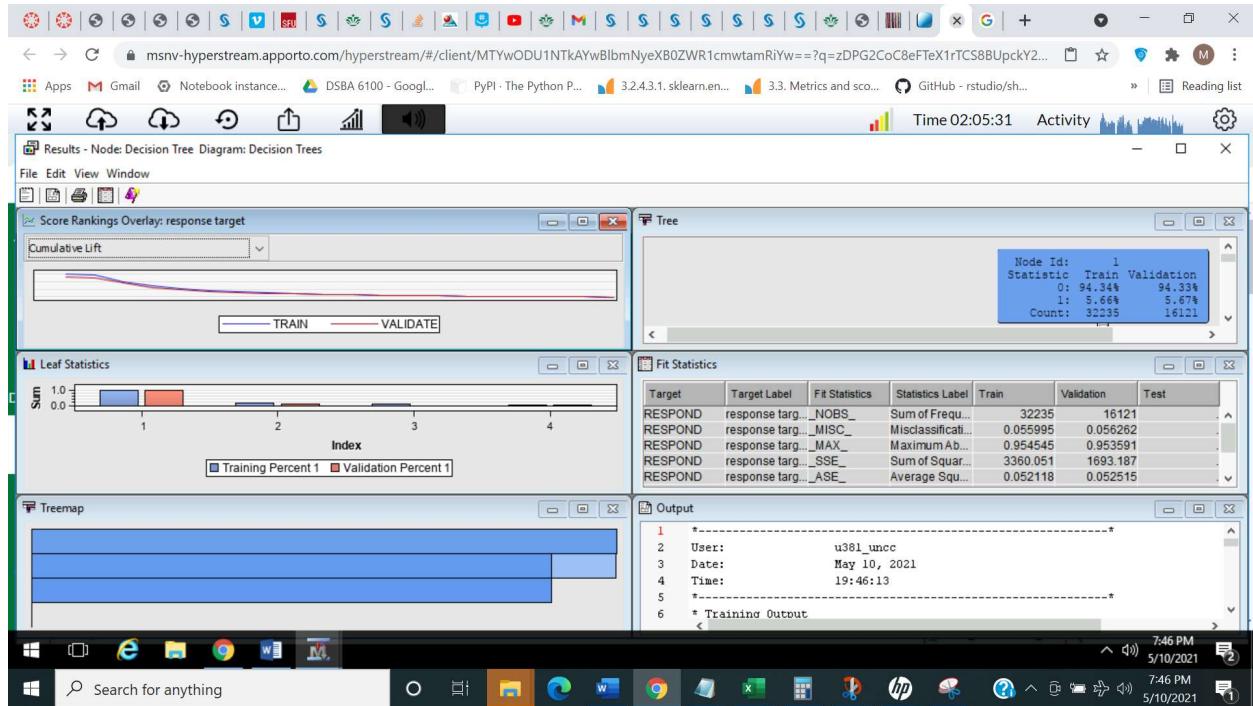
Split Count: 0 Node ID in Path: 1 (Selected)

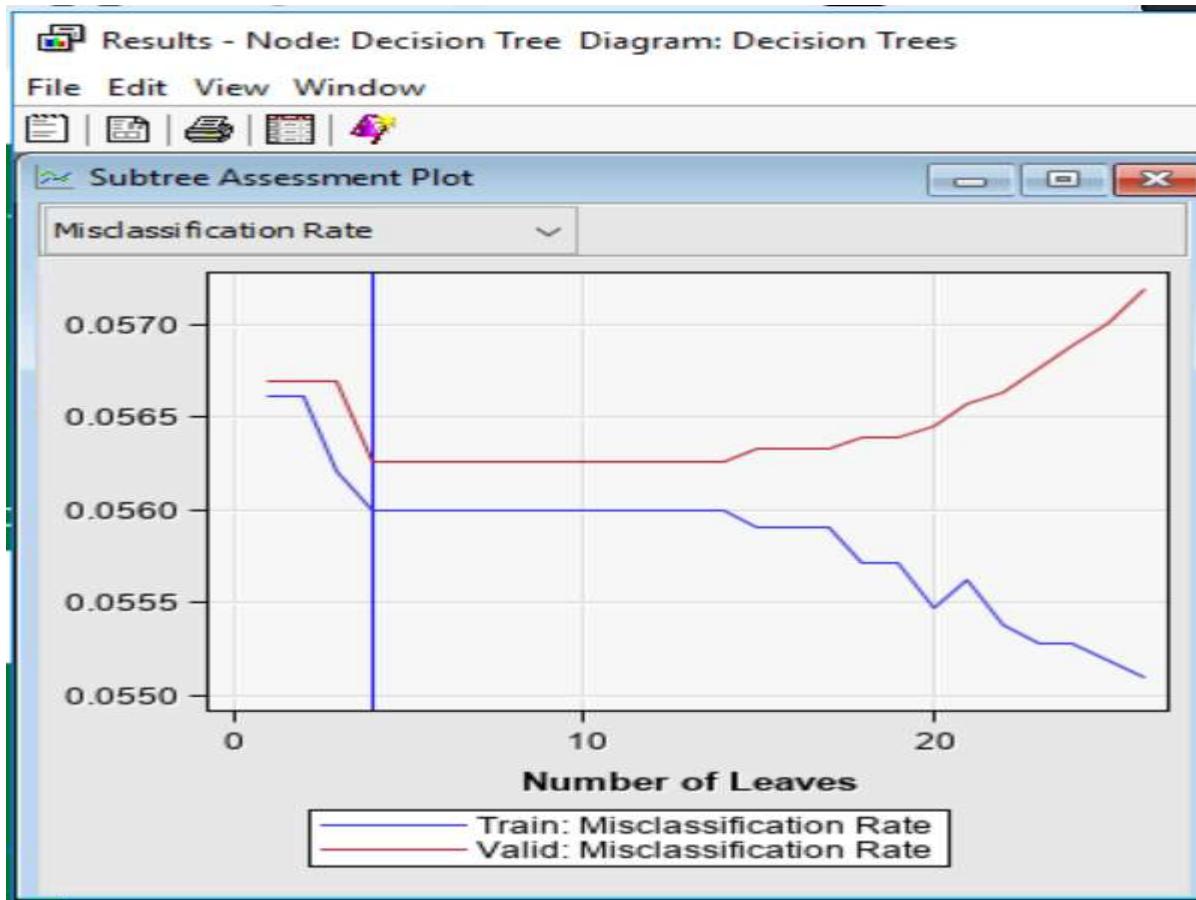












msnv-hyperstream.apporto.com/hyperstream/#/client/MTYwODU1NTkAYwBbmNyeXB0ZWR1cmwtamRiYw==?q=zDPG2CoC8eFteX1rTCS8BUpckY2...

Enterprise Miner - PredictiveModeling_Part1

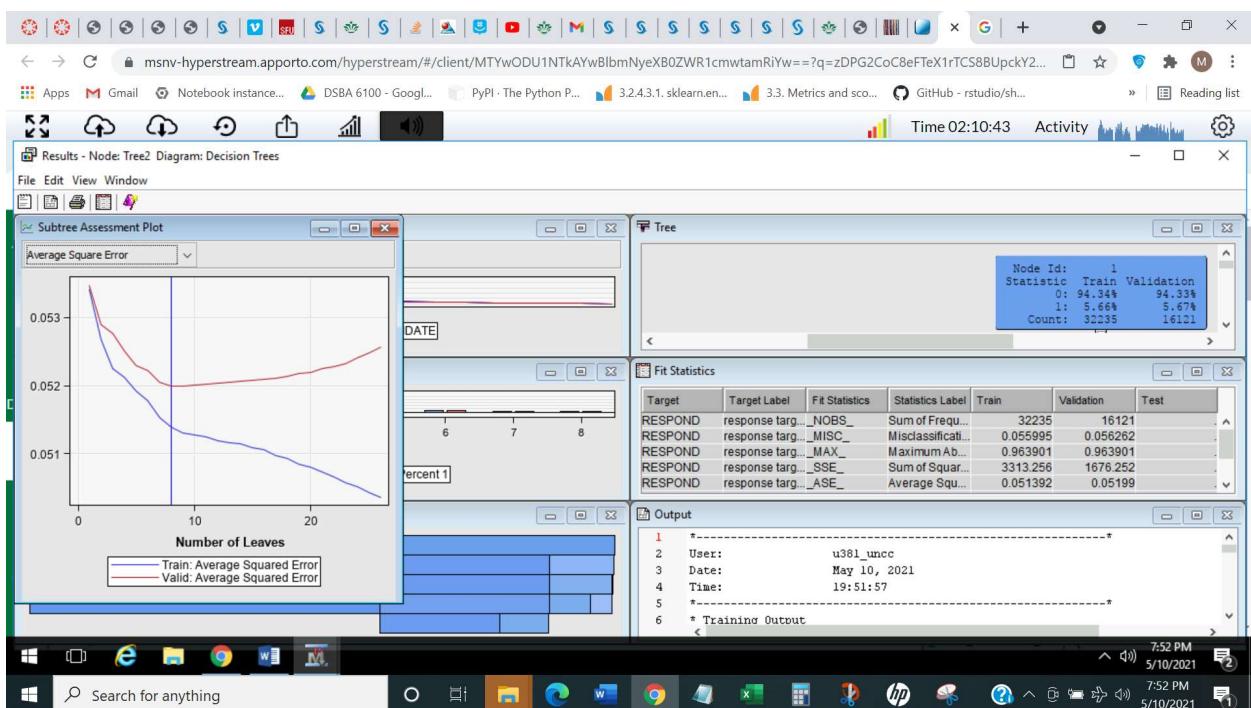
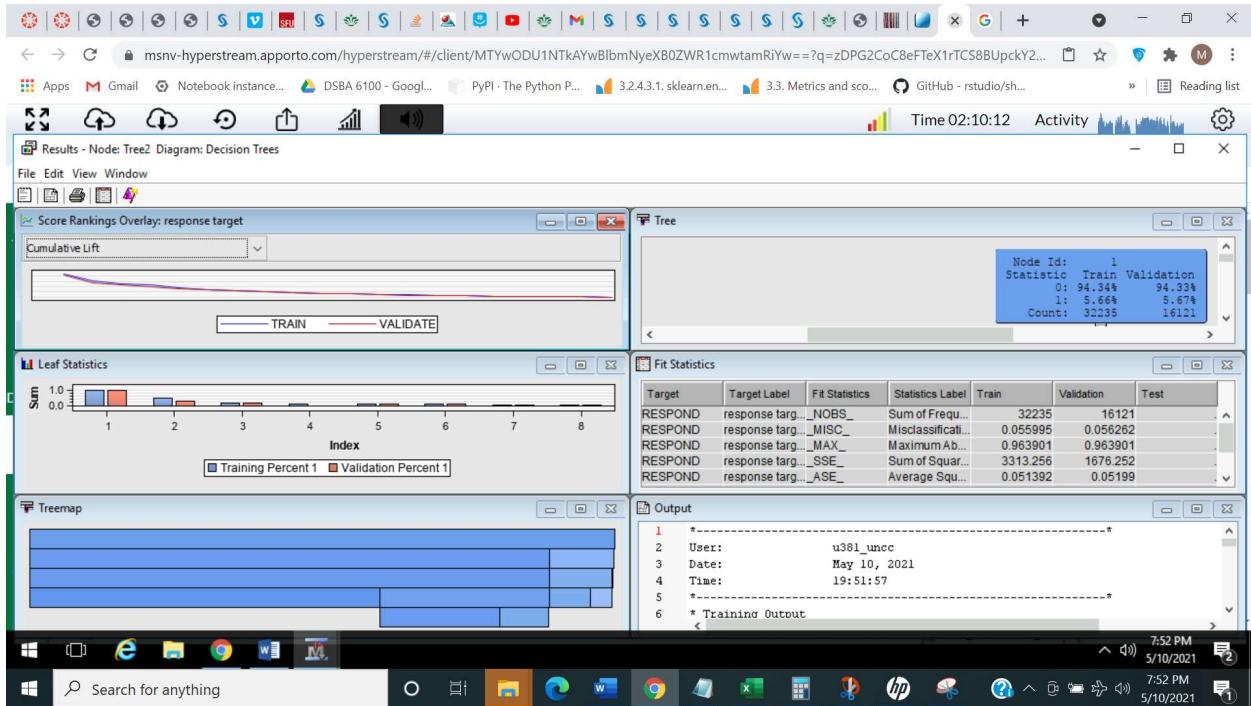
File Edit View Actions Options Window Help

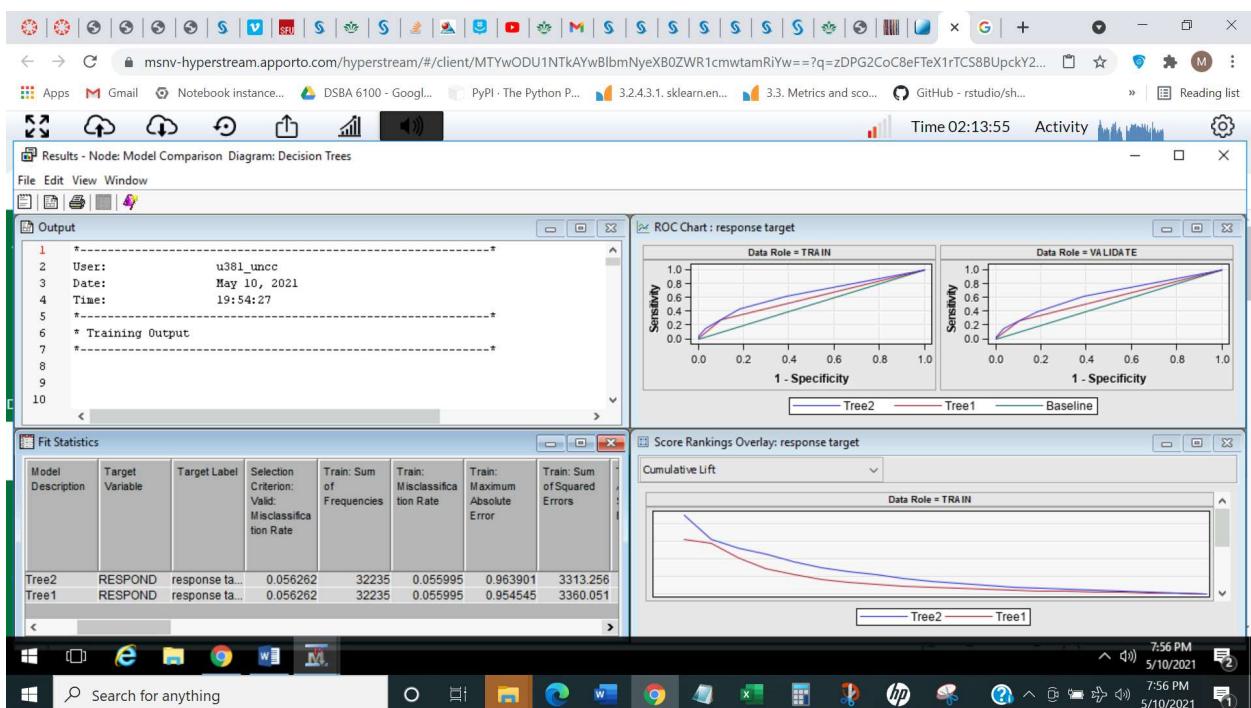
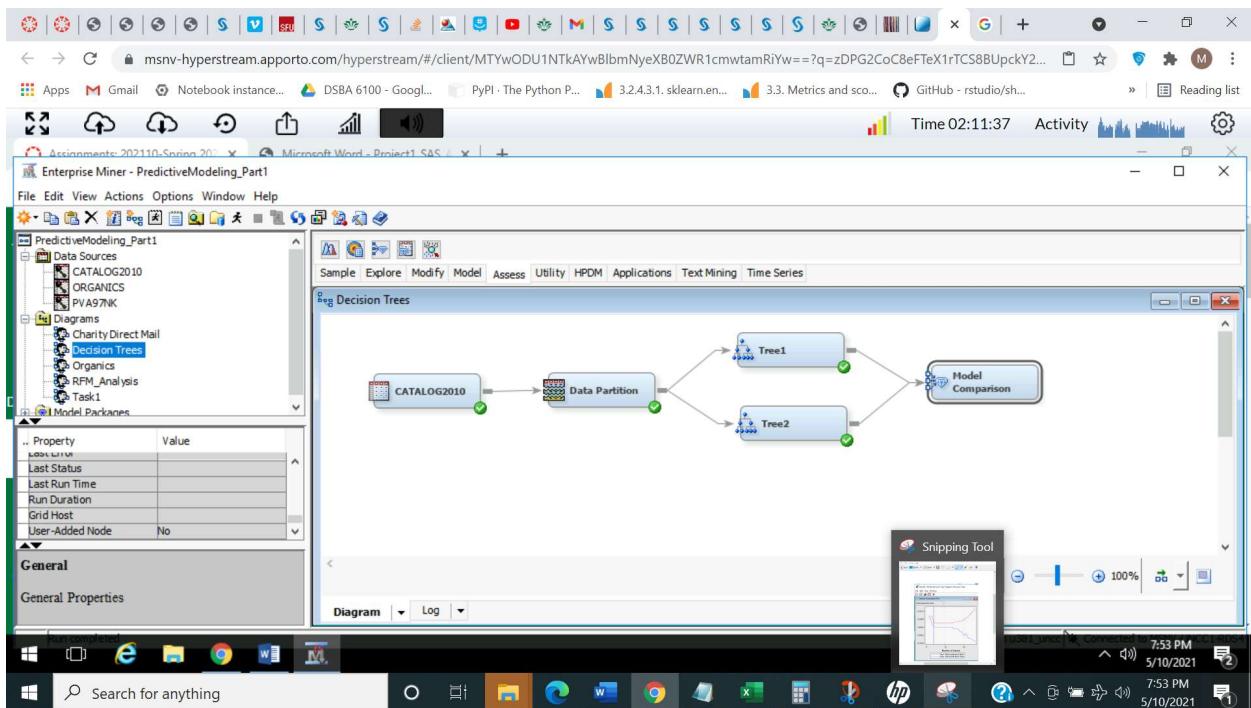
CATALOG2010 → Data Partition → Tree1, Tree2

Assessment Measure: Average Square Error

Diagram Log 100%

7:50 PM 5/10/2021 7:50 PM 5/10/2021





Task 3 Organic Dataset Decision Tree

- Set the roles for the analysis variables as shown above.

Name	Role	Level	Report	Order	Drop	Lower Limit	Upper Limit
DemAffl	Input	Interval	No	No	No	.	.
DemAge	Input	Interval	No	No	No	.	.
DemCluster	Rejected	Nominal	No	No	No	.	.
DemClusterGroup	Input	Nominal	No	No	No	.	.
DemGender	Input	Nominal	No	No	No	.	.
DemReg	Input	Nominal	No	No	No	.	.
DemTVReg	Input	Nominal	No	No	No	.	.
ID	ID	Nominal	No	No	No	.	.
PromClass	Input	Nominal	No	No	No	.	.
PromSpend	Input	Interval	No	No	No	.	.
PromTime	Input	Interval	No	No	No	.	.
RejectAmt	Rejected	Interval	No	No	No	.	.
TargetBuy	Target	Binary	No	No	No	.	.

- Examine the distribution of the target variable. What is the proportion of individuals who purchased organic products?

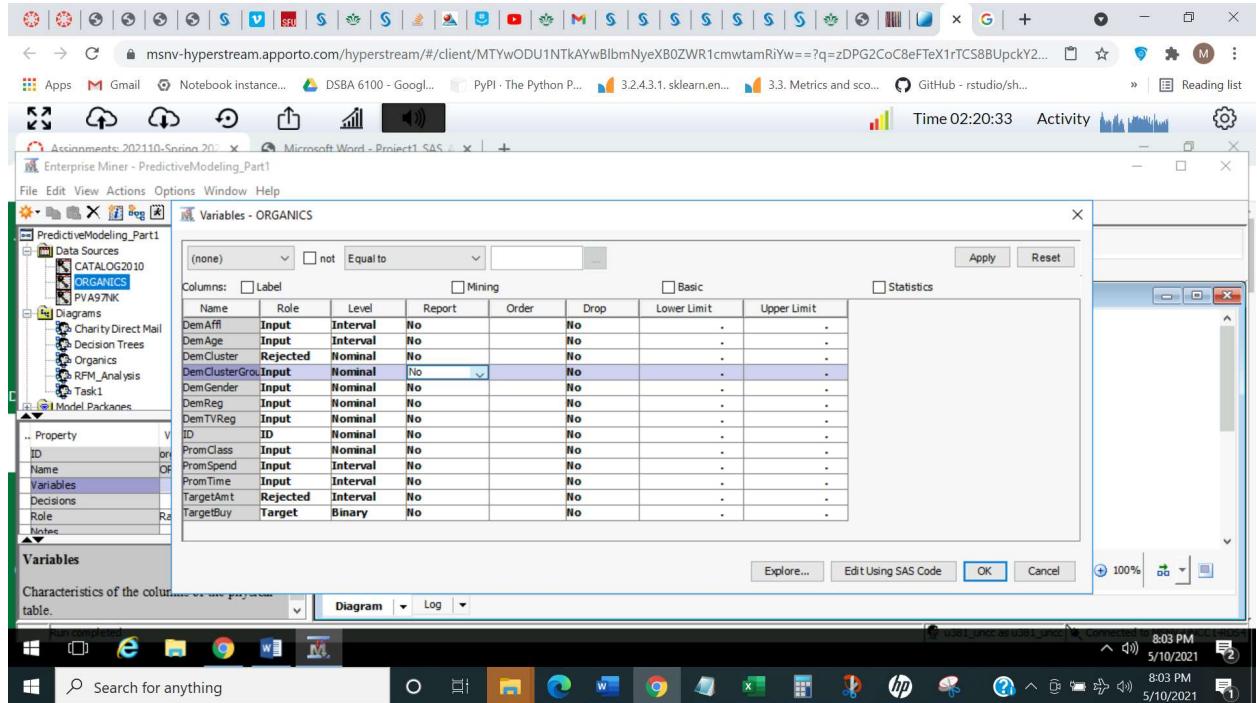
```

54
55 Distribution of Class Target and Segment Variables
56 (maximum 500 observations printed)
57
58 Data Role=TRAIN
59
60 Data Variable Frequency
61 Role Name Role Level Count Percent
62
63 TRAIN TargetBuy TARGET 0 16718 75.2284
64 TRAIN TargetBuy TARGET 1 5505 24.7716
65
66
67
68 Interval Variable Summary Statistics
69 (maximum 500 observations printed)
70
71 Data Role=TRAIN
72
73 Variable Standard Non
74 Role Mean Deviation Missing Missing Minimum Median Maximum Skewness Kurtosis
75
76 DemAffl INPUT 8.711893 3.420125 21138 1085 0 8 34 0.891684 2.09686
77 DemAge INPUT 53.79715 13.20605 20715 1508 18 54 79 -0.07983 -0.84389
78 PromSpend INPUT 4420.59 7559.048 22223 0 0.01 2000 296313.9 8.037186 184.8715
79 PromTime INPUT 6.56467 4.657113 21942 281 0 5 39 2.28279 8.077622
80
81
82

```

The proportion of individuals who purchased organic products is 24.7716 %.

3. The variable **DemClusterGroup** contains collapsed levels of the variable **DemCluster**. Presume that, based on previous experience, you believe that **DemClusterGroup** is sufficient for this type of modeling effort. Set the model role for **DemCluster** to **Rejected**.

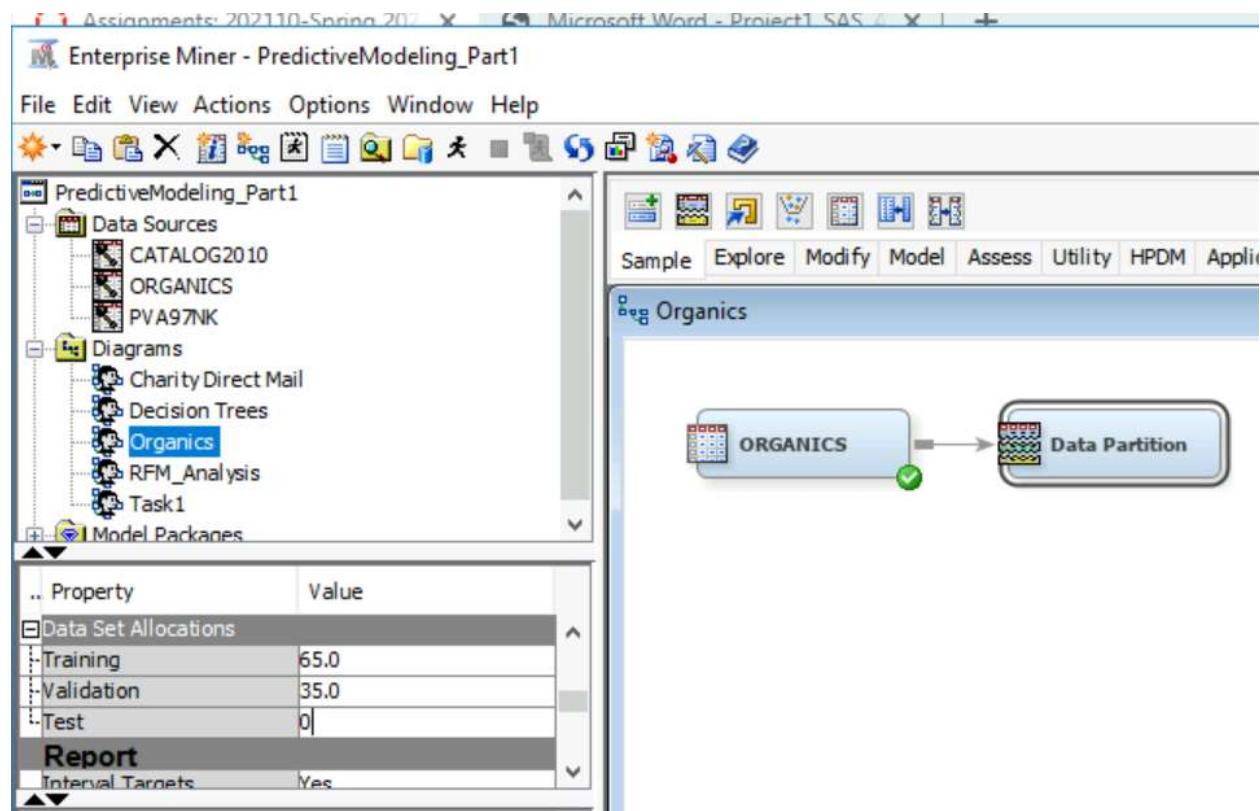


4. As noted above, only **TargetBuy** is used for this analysis, and it should have a role of **Target**. Can **TargetAmt** be used as an input for a model that is used to predict **TargetBuy**? Why or why not?

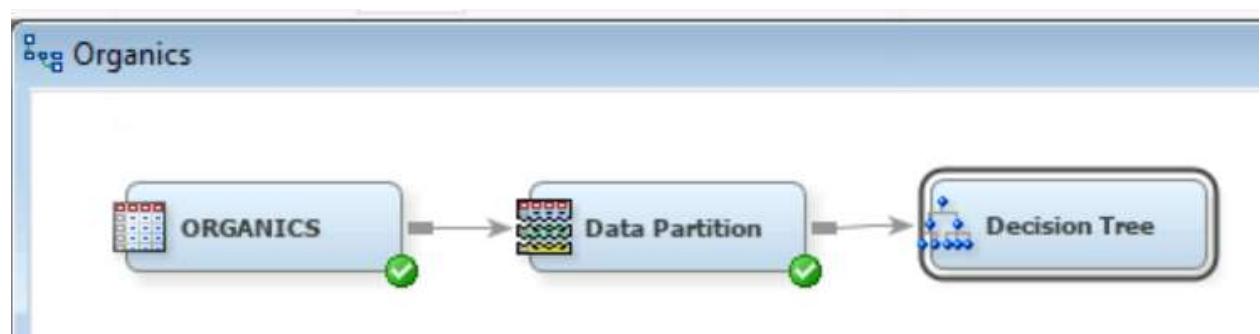
No, **TargetAmt** can not be used to predict **TargetBuy** because any amount purchased indicates a buy so it is a redundant variable. Adding **TargetAmt** to our existing variables will skew data incorrectly because of its 100% correlation with the **TargetBuy Target** role variable.

5. Finish the **ORGANICS** data source definition.

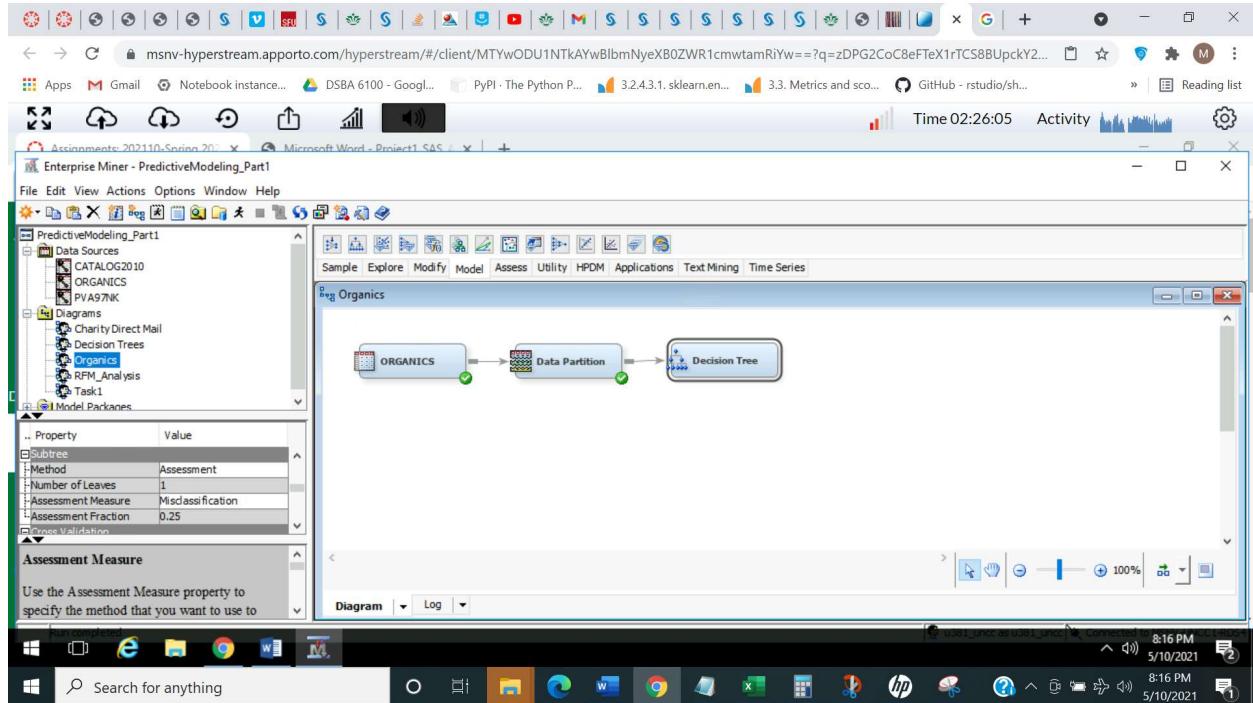
- c) Add the **ORGANICS** data source to the Organics diagram workspace.
- d) Add a **Data Partition** node to the diagram and connect it to the **Data Source** node. Assign 65% of the data for training and 35% for validation.



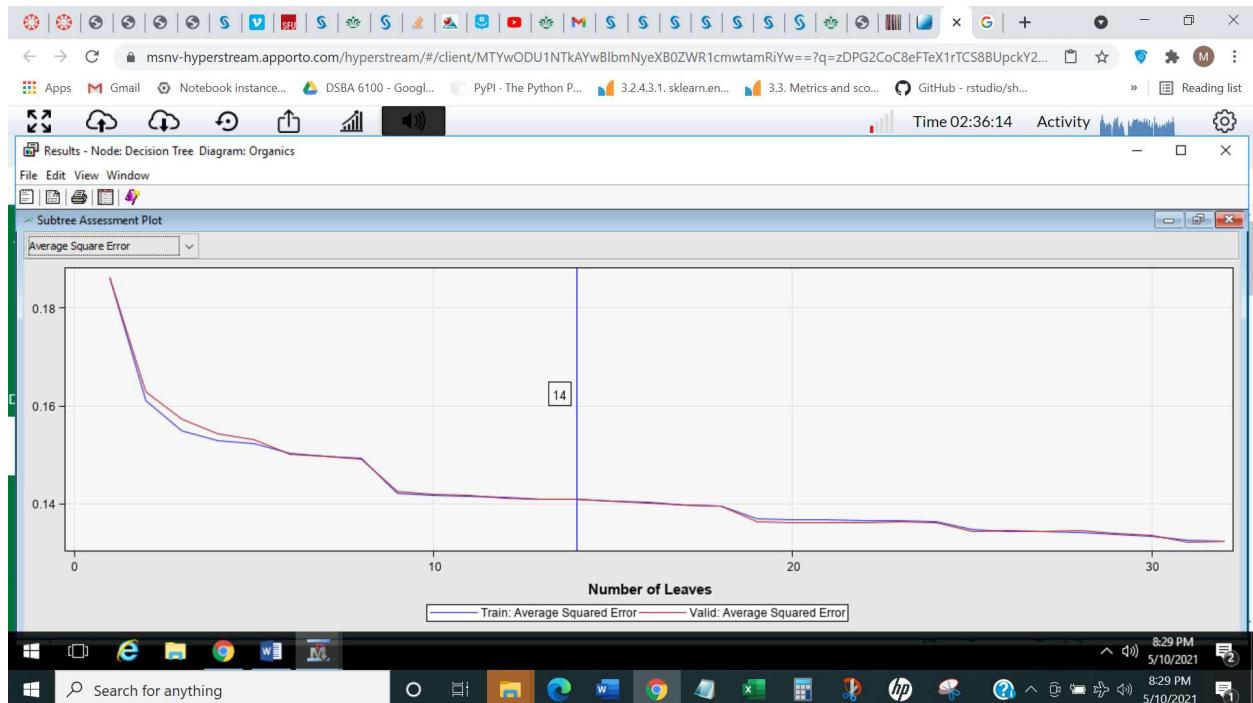
- e) Add a **Decision Tree** node to the workspace and connect it to the **Data Partition** node.



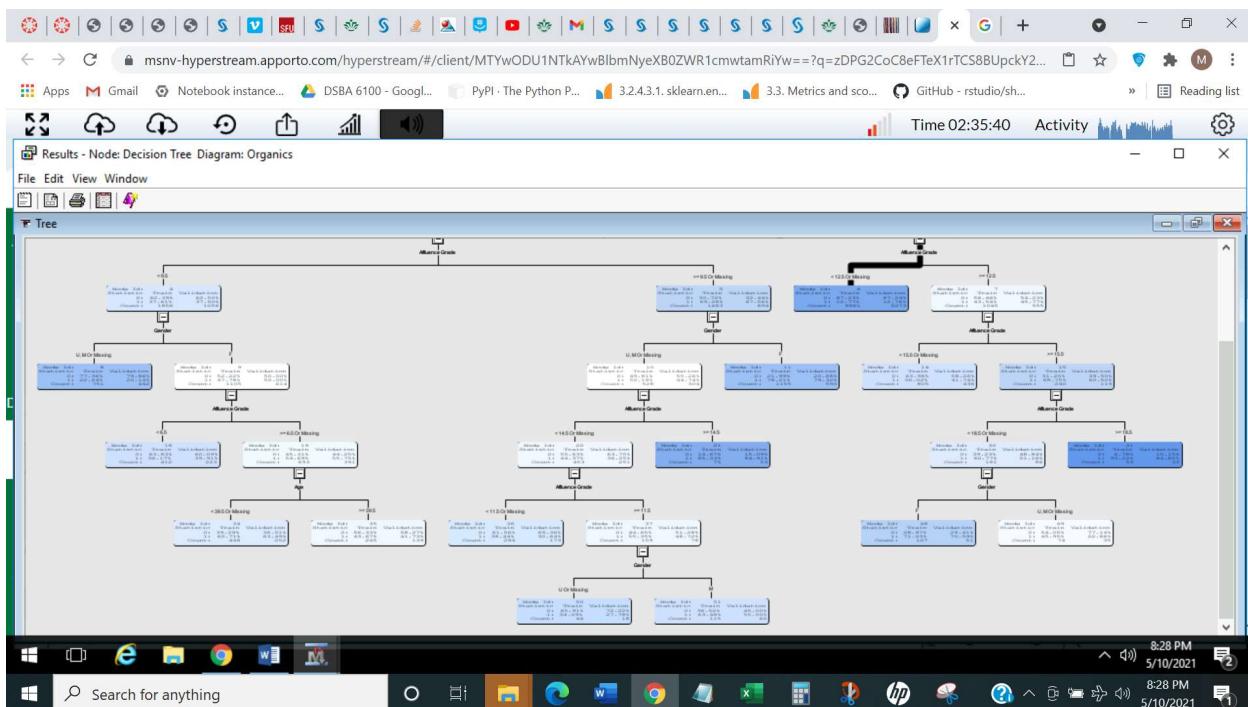
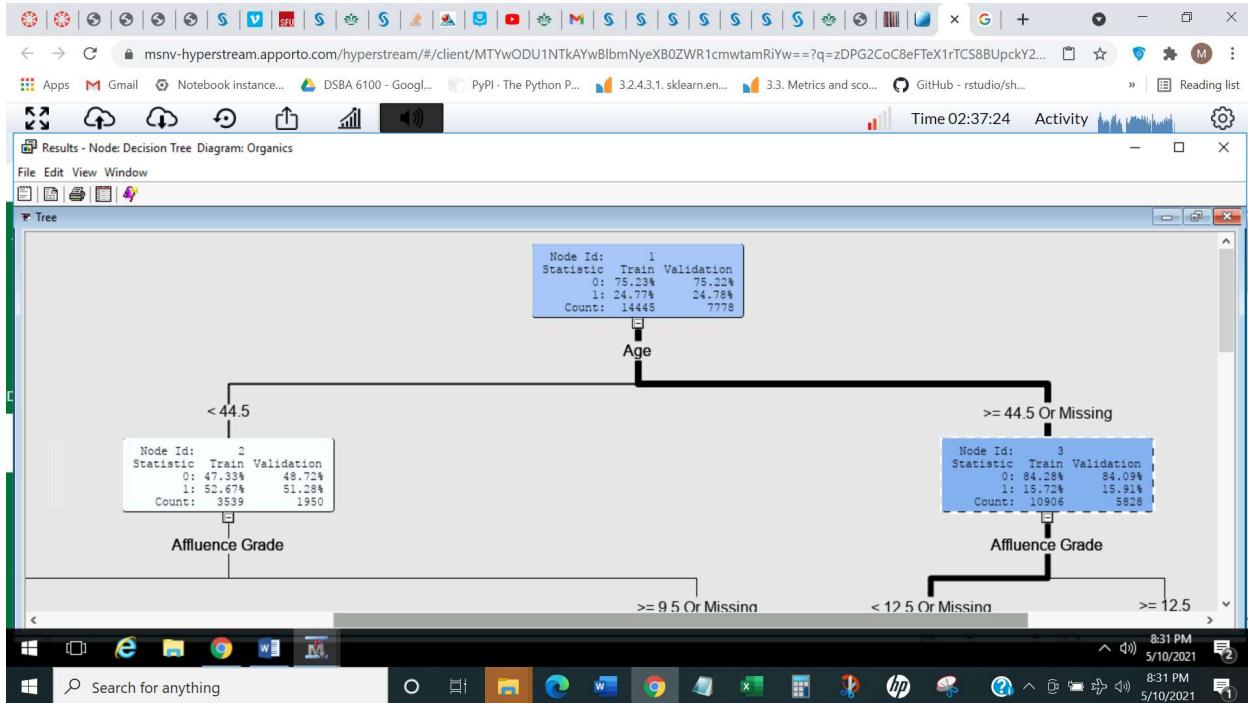
f) Create a **decision tree** model autonomously. Use **Misclassification** as the model assessment statistic.



1) How many leaves are in the optimal tree? 14 leaves



- 2) Which variables were used for the first split? Age < 44.5 and Age >= 44.5 or Age is missing



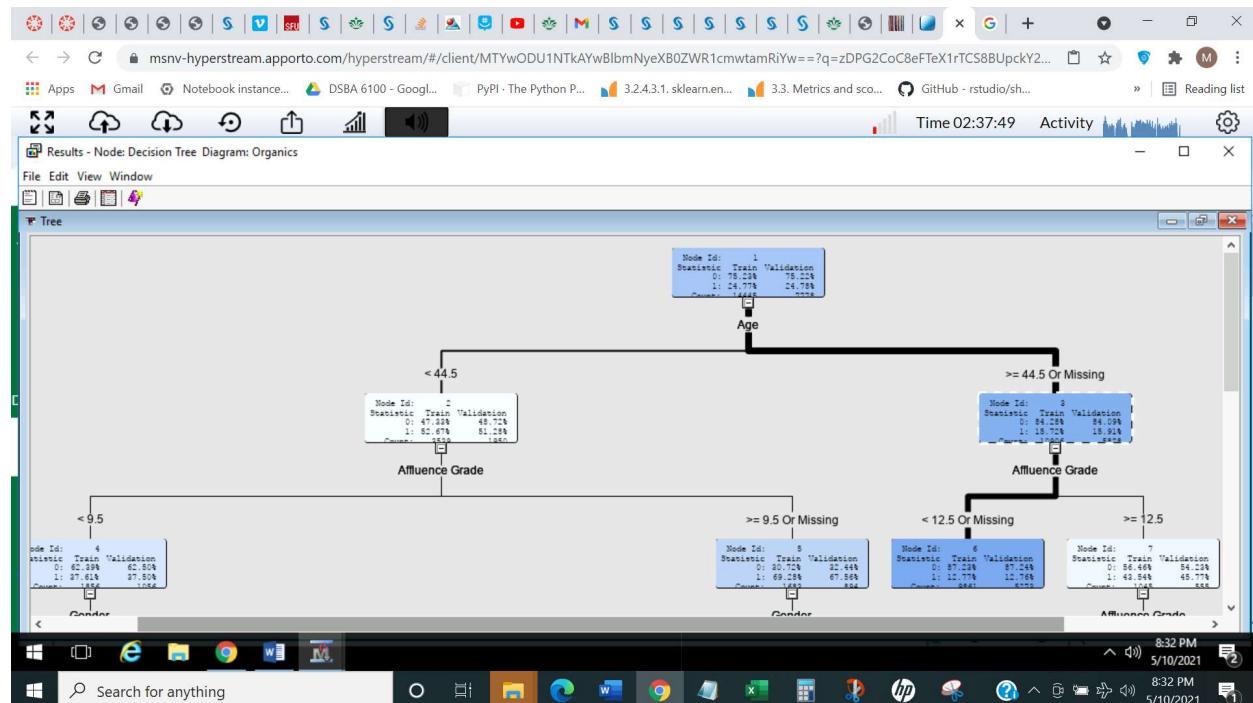
3) What were the competing splits for this first split?

Affluence Grade (Importance=0.7806) and Gender (Importance=0.4090) were also competing for the first split.

Results - Node: Decision Tree Diagram: Organics					
File	Edit	View	Window	Output	
80					
81				Number of	Ratio of
82	Variable			Splitting	Validation
83	Name	Label		Rules	to Training
84					Importance
85	DemAge	Age		1.0000	1.0000
86	DemAffl	Affluence Grade	7	0.7806	0.7929
87	DemGender	Gender	4	0.4090	0.5184
88					1.0158
					1.2674

4) Which variables were used for the second split for all branches from first split?

Affluence Grade < 9.5, Affluence Grade >=9.5 or Missing, Affluence Grade < 12.5 or Missing, Affluence Grade >= 12.5

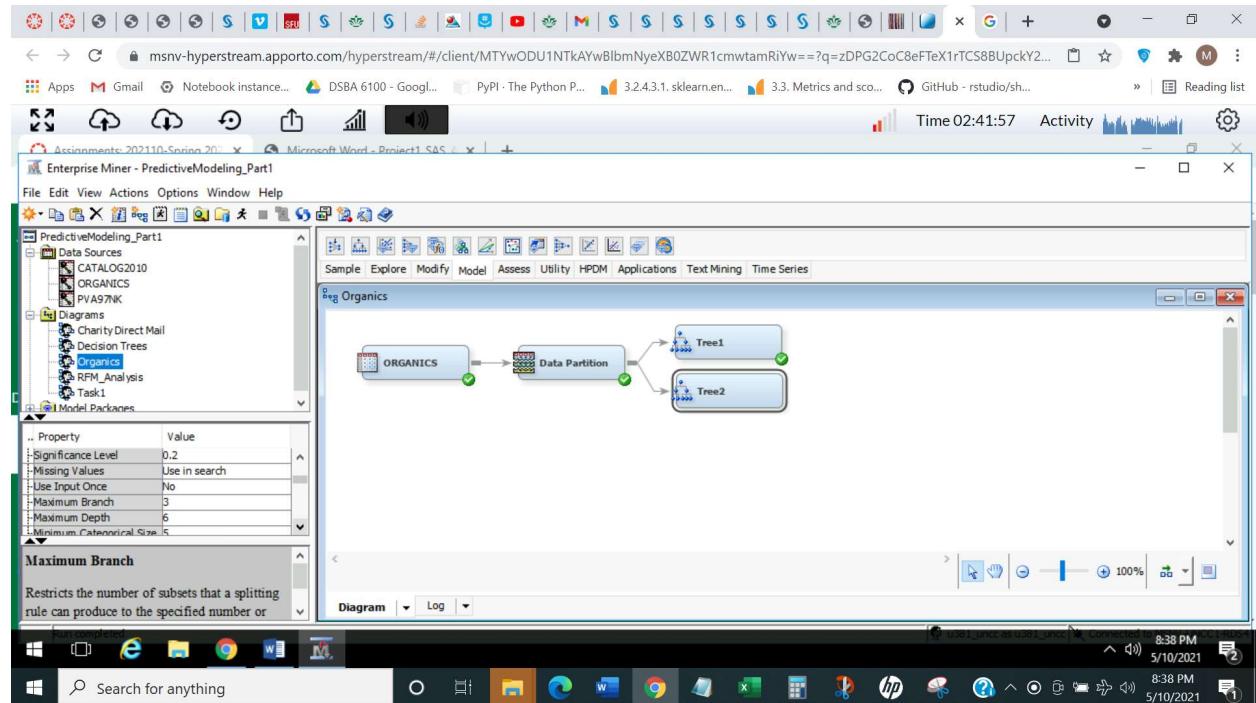


5) Discuss the results and provide your insights?

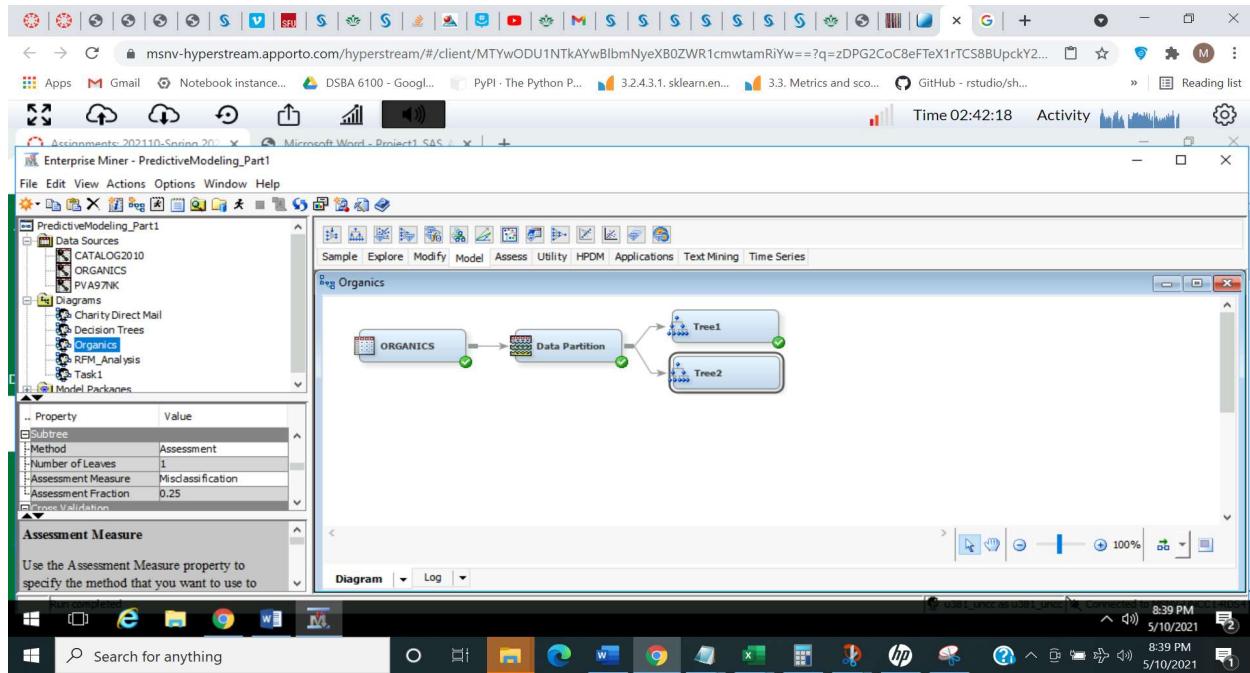
I see that Affluence Grade is split multiple times with different values for example in Node 7, Node 15 and Node 20. This is a complex model with 51 total nodes and nodes are reused to split as needed.

Add a second **Decision Tree** node to the diagram and connect it to the **Data Partition** node.

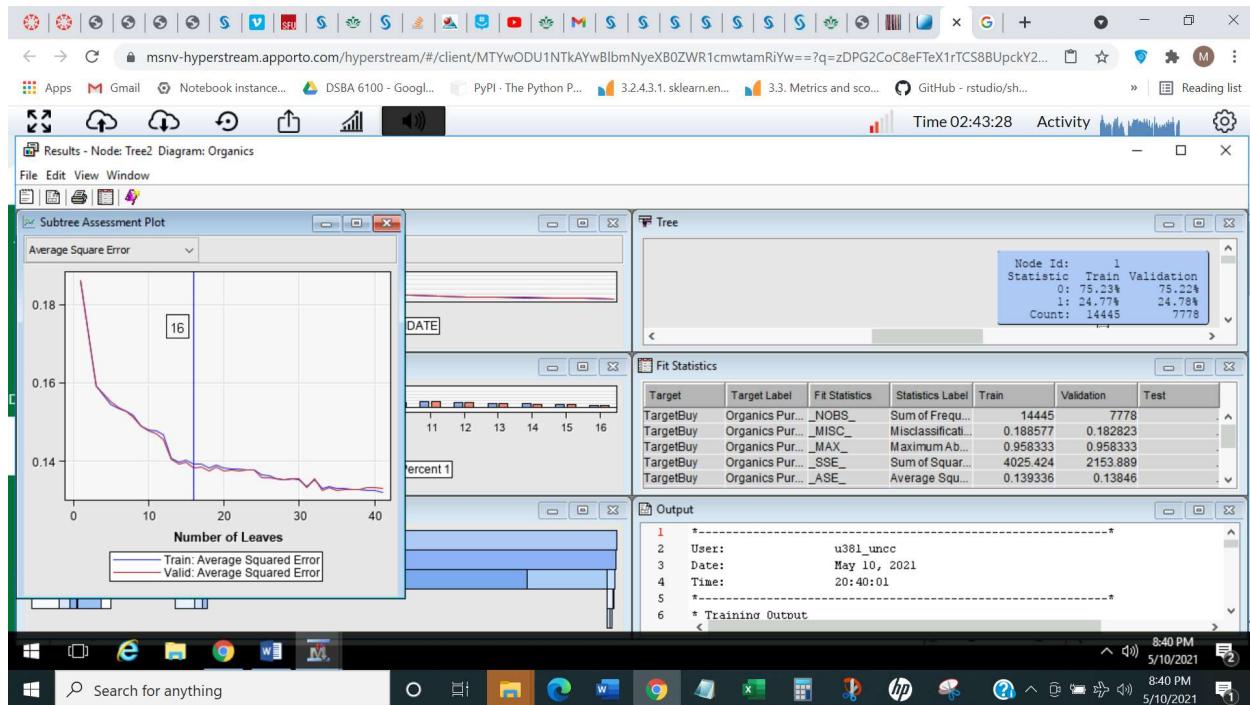
- 1) In the Properties panel of the new Decision Tree node, change the maximum number of branches from a node to 3 to allow for three-way splits.



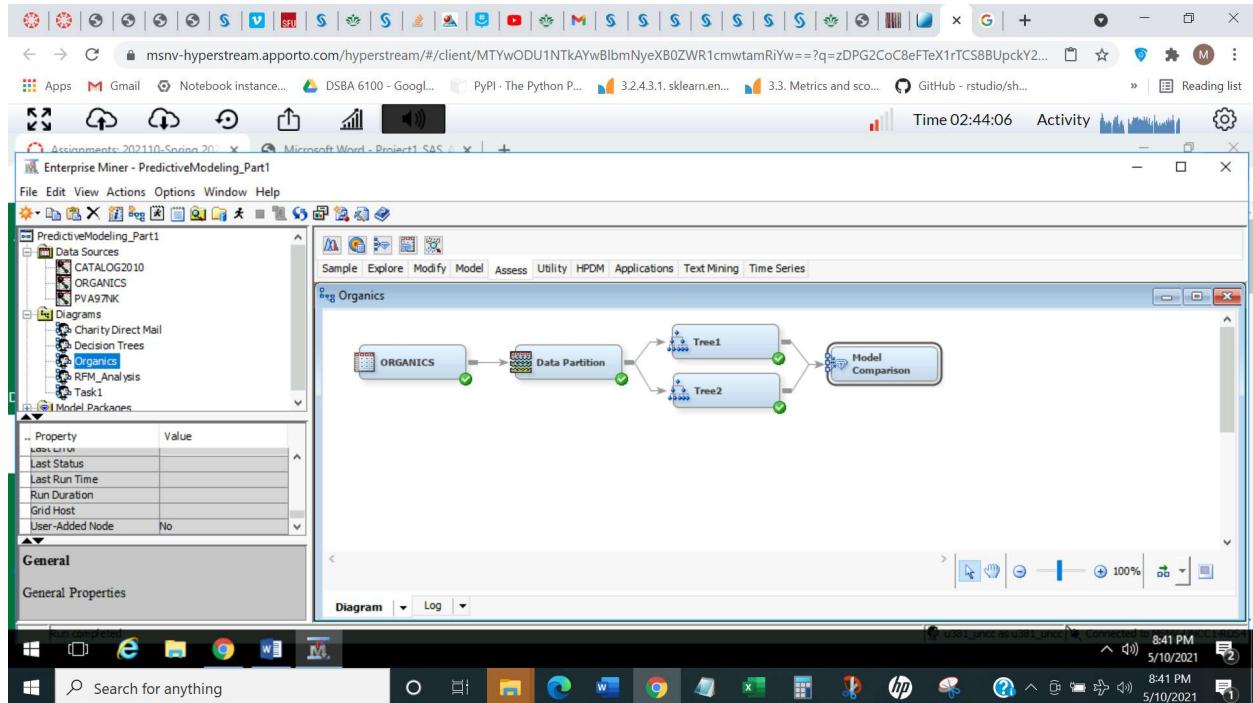
- 2) Create a decision tree model using **Misclassification** as the model assessment statistic.



- 3) How many leaves are in the optimal tree? 16 in Tree2



h) Based on **Misclassification rate**, which of the decision tree models appears to be better?



Tree2 has a slightly better misclassification rate of 0.188577 versus Tree1 0.183801.

Selected Model	Predecessor Node	Model Node	Model Description	Target Variable	Target Label	Selection Criterion: Valid: Average Profit for TargetBuy	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	Train: Average Profit for TargetBuy	Train: Profits & Target
Y	Tree2	Tree2	Tree2	TargetBuy	Organics P...	0.817177	14445	0.188577	0.958333	4025.424	0.139336	0.373278	28890	14445	0.811423	
	Tree	Tree	Tree1	TargetBuy	Organics P...	0.816662	14445	0.183801	0.932203	4070.628	0.140901	0.375368	28890	14445	0.816199	

Part 3-Predictive Modeling (LOGIT)

CATALOG CASE STUDY: FITTING A LOGISTIC REGRESSION MODEL

Variables - CATALOG2010

Name	Role	Level	Report	Order	Drop	Lower Limit	Upper Limit	Number of Levels	Percent Missing
DOLLARQ22	Input	Interval	No	No
DOLNETDA	Input	Interval	No	No
DOLNETDT	Input	Interval	No	No
DTBUYLST	Rejected	Interval	No	No
DTBUYORG	Rejected	Interval	No	No
FREPRCH	Input	Interval	No	No
METHPAYM	Input	Nominal	No	No	.	.	.	4	.
MONLAST	Input	Interval	No	No
ORDERSIZE	Target	Interval	No	No
PCPAYM	Input	Binary	No	No	.	.	.	2	.
RESPOND	Target	Binary	No	No	.	.	.	2	.
STATE	Rejected	Nominal	No	No	.	.	.	21	.
TENURE	Input	Interval	No	No

Time 02:50:56 Activity

Diagram LOGIT

```

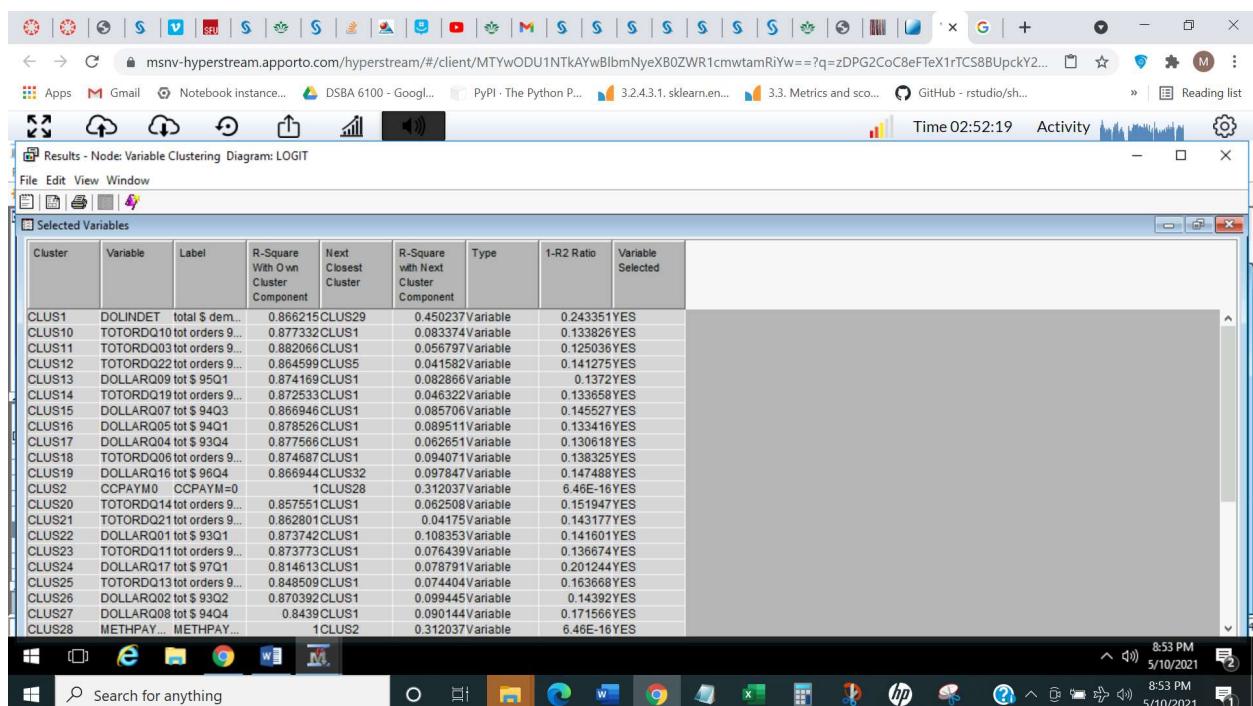
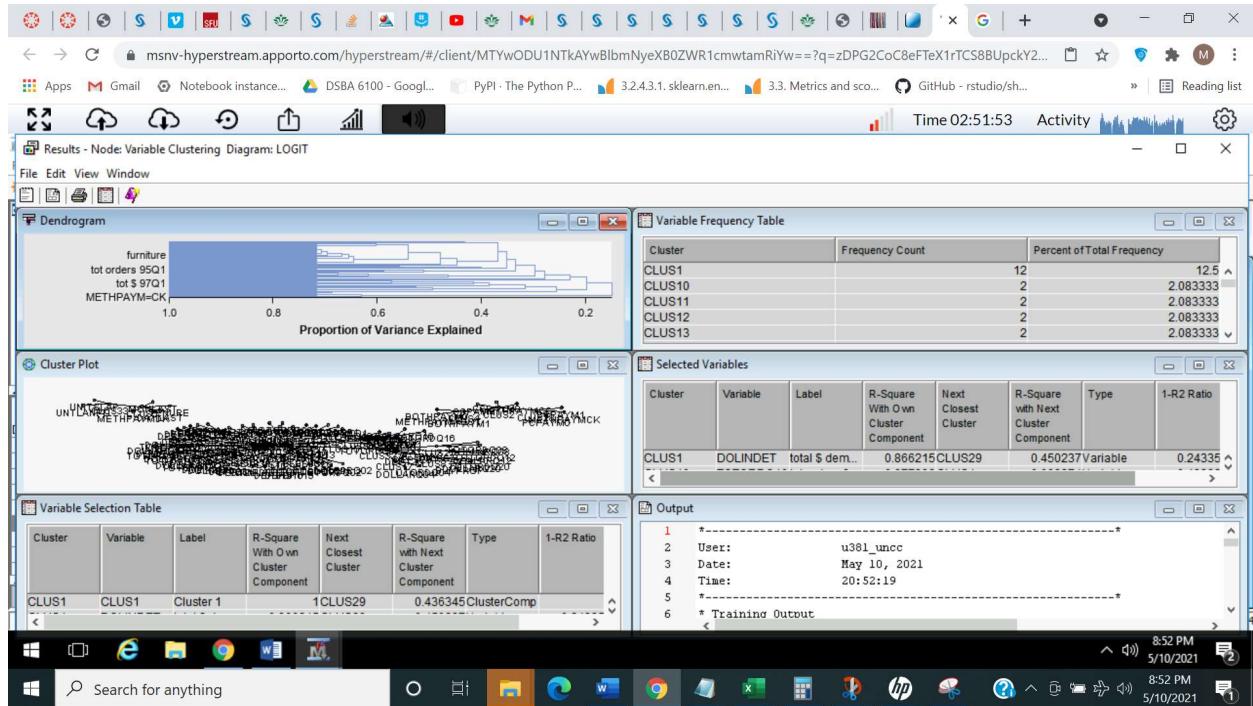
graph LR
    C1[CATALOG2010] --> V1[Variable Clustering]
    V1 --> LOGIT[LOGIT]
  
```

Score

- Variable Selection: Best Variables
- Interactive Selection: ...
- Hides Rejected Variables: Yes

Variable Selection

Diagram LOGIT opened



Results - Node: Variable Clustering Diagram: LOGIT

File Edit View Window

Output

	Cluster	Variable	Total	Proportion	Minimum	Maximum	
5285	Cluster 30	DEPT07	0.1047	0.0238	0.9171	mens apparel	
5286		DEPT08	0.2844	0.0912	0.7874	mens sleepwear	
5287		DEPT09	0.3088	0.0516	0.7288	mens underwear	
5288		DEPT10	0.3074	0.0794	0.7523	mens hosiery	
5289		DEPT11	0.2009	0.0524	0.8432	mens footwear	
5290		DEPT12	0.3714	0.1064	0.7035	mens misc	
5291							-----
5292	Cluster 31	DEPT23	0.6818	0.2844	0.4447	beauty	
5293		DEPT24	0.6818	0.1738	0.3851	health	
5294							-----
5295	Cluster 32	ACTBUY	0.6583	0.3653	0.5383	num qrtrs w/buy	
5296		DEPT25	0.5737	0.2171	0.5446	food	
5297		DEPT26	0.2855	0.0544	0.7556	gift	
5298							-----
5299	Cluster 33	UNITSLAP	0.6172	0.2465	0.5081	avg price/unit	
5300		UNTLANPO	0.6172	0.2613	0.5182	avg units/order	
5301							-----
5302	Cluster 34	DEPT20	0.5167	0.0188	0.4925	furniture	
5303		DEPT21	0.5167	0.0102	0.4883	light	
5304							-----
5305	Cluster 35	DEPT15	0.5112	0.1126	0.5509	bath	
5306		DEPT16	0.5218	0.1175	0.5419	floor	
5307		DEPT19	0.2059	0.0265	0.8157	window	
5308							

Windows Taskbar: Search for anything, 8:55 PM, 5/10/2021

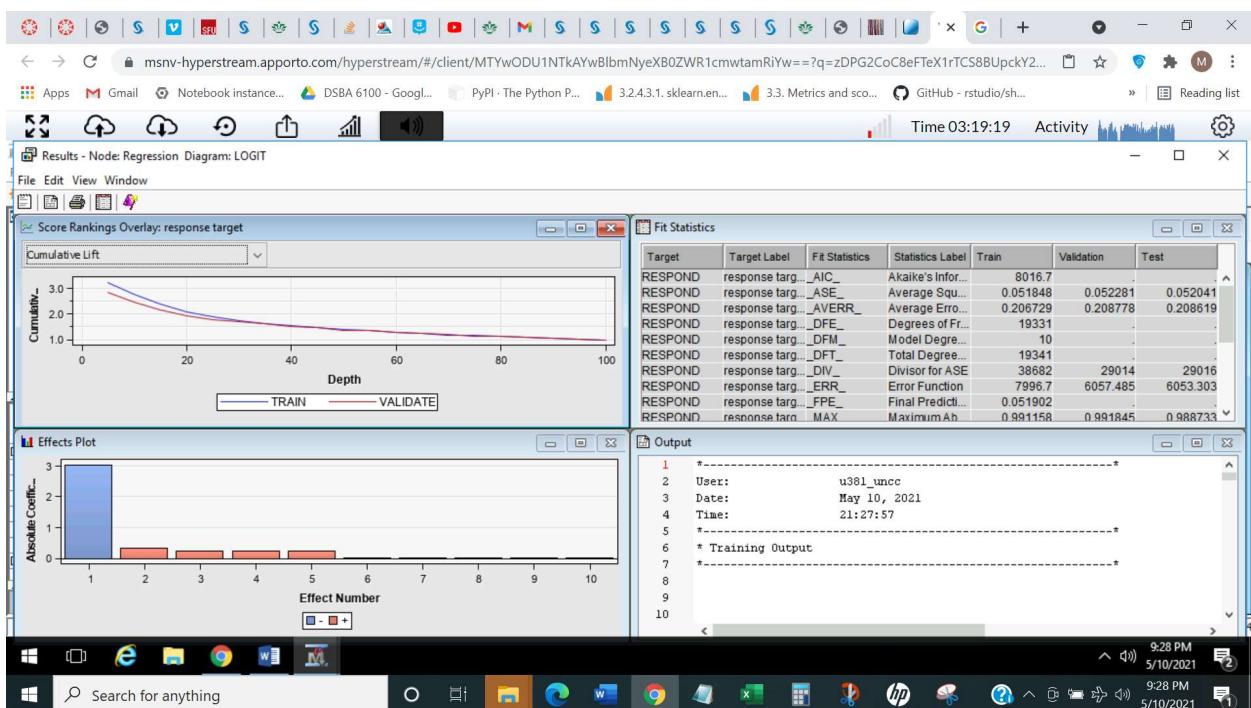
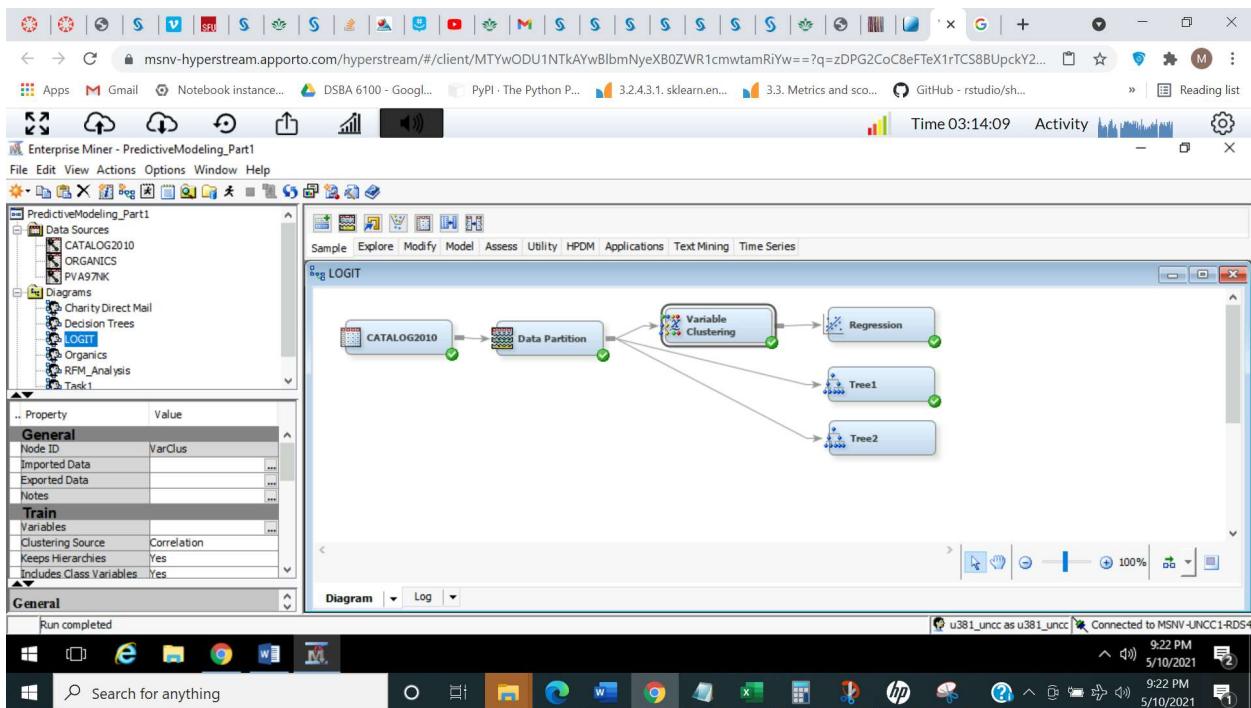
Results - Node: Variable Clustering Diagram: LOGIT

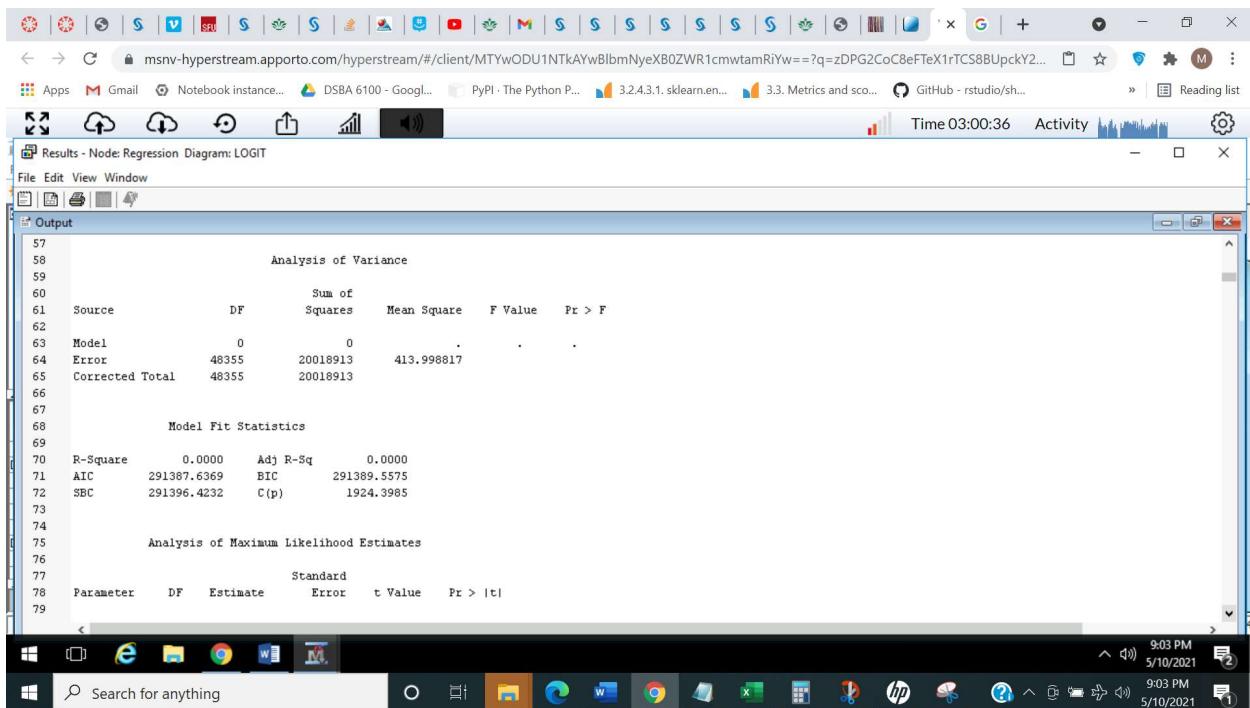
File Edit View Window

Output

	Number of Clusters	Total Variation Explained by Clusters	Proportion of Variation Explained by Clusters	Minimum Proportion Explained by Clusters	Maximum Proportion Explained by Clusters	Second Eigenvalue in a Cluster	Minimum R-squared for a Variable	Maximum R-squared for a Variable	1-R**2 Ratio
5312	1	14.409222	0.1501	0.1501	5.059839	0.0048			
5313	2	19.021411	0.1981	0.1594	3.716946	0.0097	0.9928		
5314	3	22.021680	0.2294	0.1846	2.781652	0.0110	0.9911		
5315	4	24.580286	0.2560	0.1864	2.326299	0.0112	0.9908		
5316	5	26.724262	0.2784	0.1864	2.242560	0.0112	0.9978		
5317	6	28.575376	0.2977	0.1907	2.028430	0.0110	1.1978		
5318	7	30.287189	0.3155	0.2004	1.976234	0.0110	1.1978		
5319	8	32.077210	0.3341	0.2004	1.893210	0.0110	1.1978		
5320	9	33.724665	0.3513	0.2148	1.759903	0.0113	1.1978		
5321	10	35.316554	0.3679	0.2311	1.643163	0.0118	1.1978		
5322	11	36.891238	0.3843	0.2458	1.626987	0.0118	1.1978		
5323	12	38.354461	0.3995	0.2458	1.561945	0.0118	1.1978		
5324	13	39.906380	0.4157	0.2458	1.559701	0.0118	1.1978		
5325	14	41.369877	0.4309	0.2458	1.559301	0.0118	1.1978		
5326	15	42.929140	0.4472	0.2458	1.555050	0.0118	1.1978		
5327	16	44.451080	0.4630	0.2529	1.553516	0.0120	1.1978		
5328	17	45.930600	0.4784	0.2529	1.543773	0.0120	1.1978		
5329	18	47.441137	0.4942	0.2608	1.534784	0.0121	1.1978		

Windows Taskbar: Search for anything, 8:56 PM, 5/10/2021

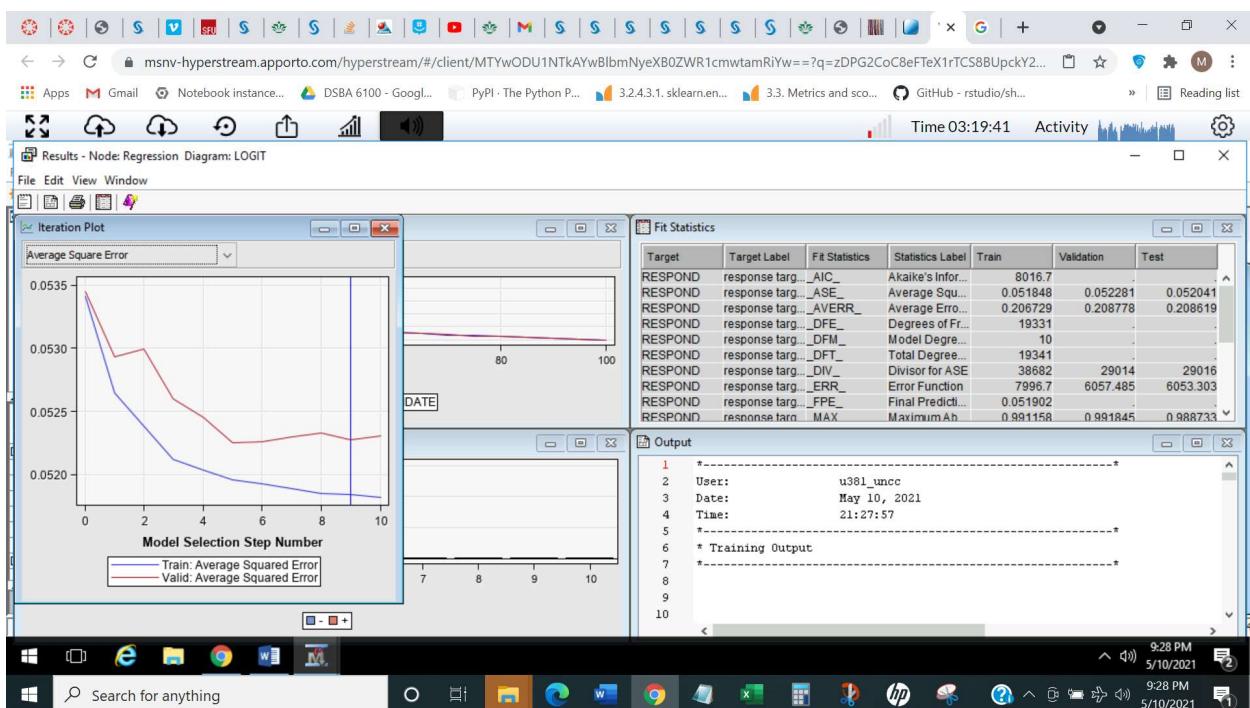


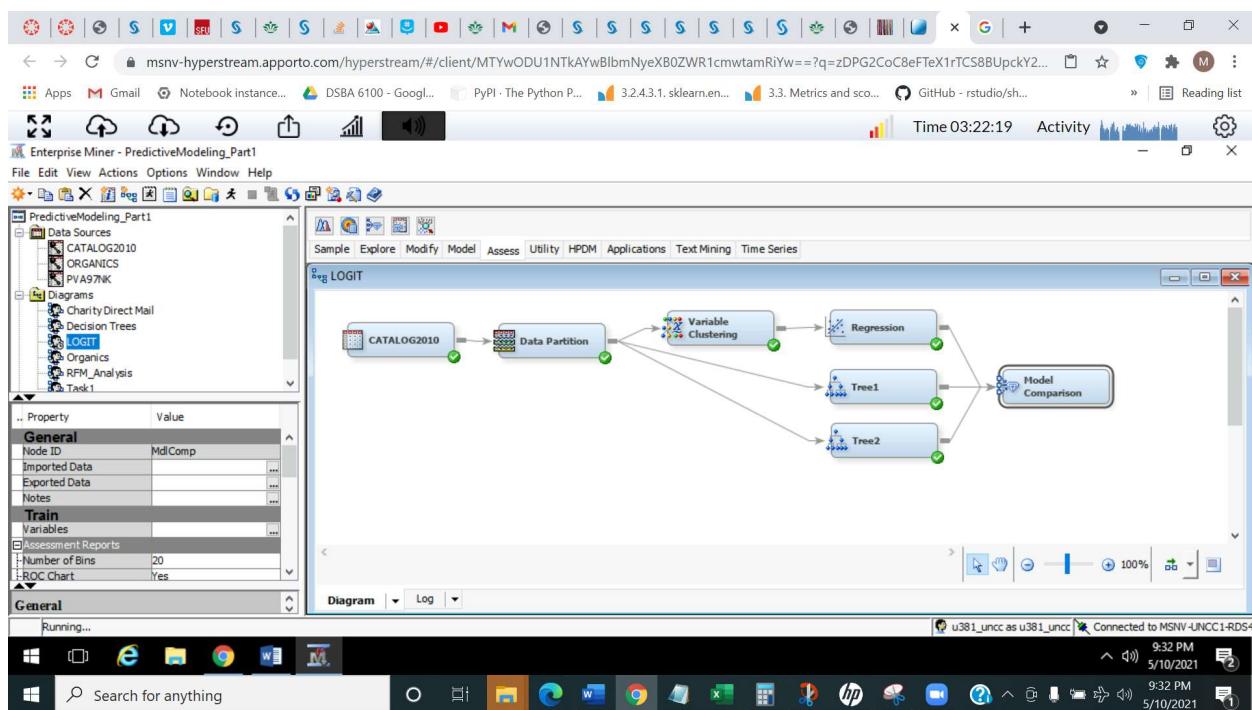
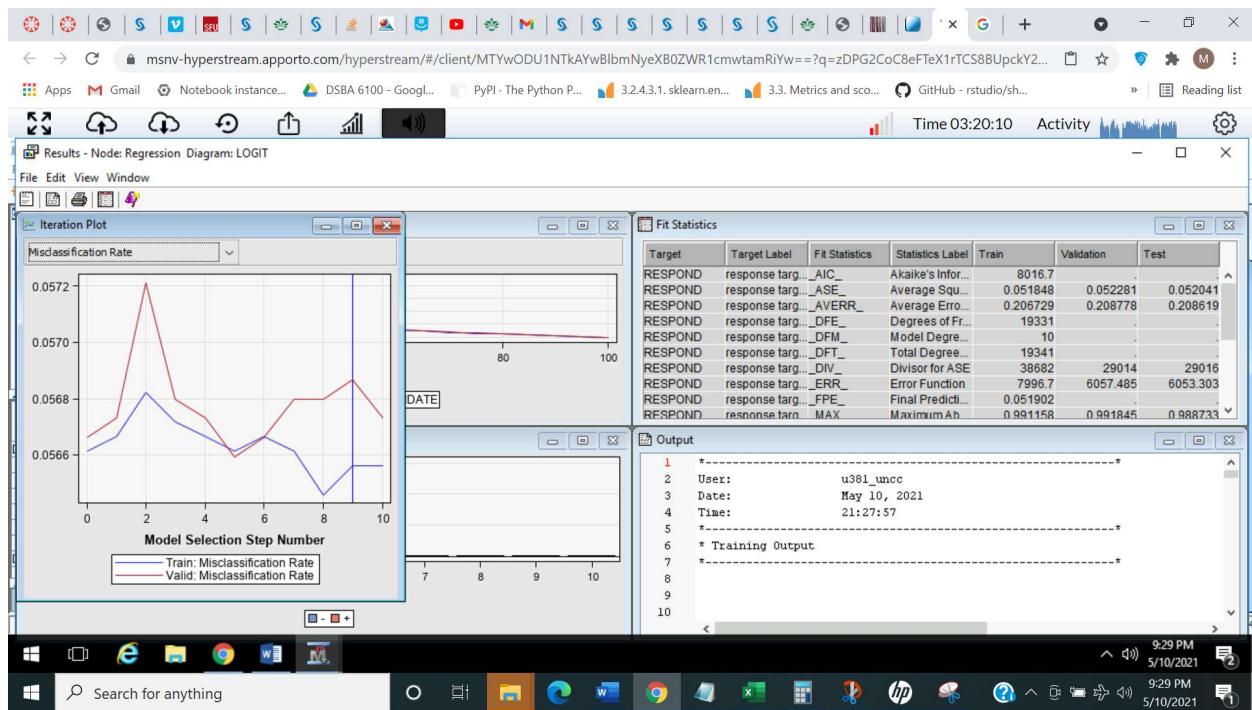


```

57
58          Analysis of Variance
59
60      Source           DF     Sum of Squares   Mean Square   F Value    Pr > F
61      Model            0         0             413.998817   .
62      Error           48355  20018913        413.998817   .
63      Corrected Total  48355  20018913
64
65
66          Model Fit Statistics
67
68      R-Square       0.0000   Adj R-Sq  0.0000
69      AIC            291387.6369   BIC       291389.5575
70      SBC            291396.4232   C(p)      1924.3985
71
72
73          Analysis of Maximum Likelihood Estimates
74
75      Standard
76      Parameter   DF   Estimate     Error   t Value   Pr > |t|
77
78

```





Results - Node: Model Comparison Diagram: LOGIT

File Edit View Window

Fit Statistics

Selected Model	Predecessor Node	Model Node	Model Description	Target Variable	Selection Criterion:	Train: Akaike's Information Criterion
Tree	Tree	Tree1	RESPOND	response ta...	Valid: Misclassifica...	0.055766
Tree2	Tree2	Tree2	RESPOND	response ta...	0.055766	
Reg	Reg	Regression	RESPOND	response ta...	0.056069	8016.7

Output

```

1 -----
2 User: u381_umcc
3 Date: May 10, 2021
4 Time: 21:32:36
5 -----
6 * Training Output
7 -----
8
9
10

```

9:32 PM 5/10/2021 9:32 PM 5/10/2021

Results - Node: Model Comparison Diagram: LOGIT

File Edit View Window

Output

```

184 Data Role=Valid
185
186 Statistics
187 Tree Tree2 Reg
188 Valid: Kolmogorov-Smirnov Statistic 0.16 0.16 0.23
189 Valid: Average Squared Error 0.05 0.05 0.05
190 Valid: Roc Index 0.59 0.59 0.65
191 Valid: Average Error Function .
192 Valid: Bin-Based Two-Way Kolmogorov-Smirnov Probability Cutoff 0.07 0.07 0.06
193 Valid: Cumulative Percent Captured Response 22.12 22.12 24.53
194 Valid: Percent Captured Response 6.70 6.70 10.30
195 Valid: Divisor for VASE 29014.00 29014.00 29014.00
196 Valid: Error Function .
197 Valid: Gain 121.16 121.16 145.29
198 Valid: Gini Coefficient 0.17 0.17 0.31
199 Valid: Bin-Based Two-Way Kolmogorov-Smirnov Statistic 0.15 0.15 0.23
200 Valid: Kolmogorov-Smirnov Probability Cutoff 0.05 0.05 0.06
201 Valid: Cumulative Lift 2.21 2.21 2.45
202 Valid: Lift 1.34 1.34 2.06
203 Valid: Maximum Absolute Error 0.96 0.96 0.99
204 Valid: Misclassification Rate 0.06 0.06 0.06
205 Valid: Mean Square Error .
206 Valid: Sum of Frequencies 14507.00 14507.00 14507.00
207 Valid: Root Average Squared Error 0.23 0.23 0.23

```

9:33 PM 5/10/2021 9:33 PM 5/10/2021