

## **Start SAS Enterprise Miner with CLAIM\_DATA\_V2\_TRAIN**

Didem B. Aykurt

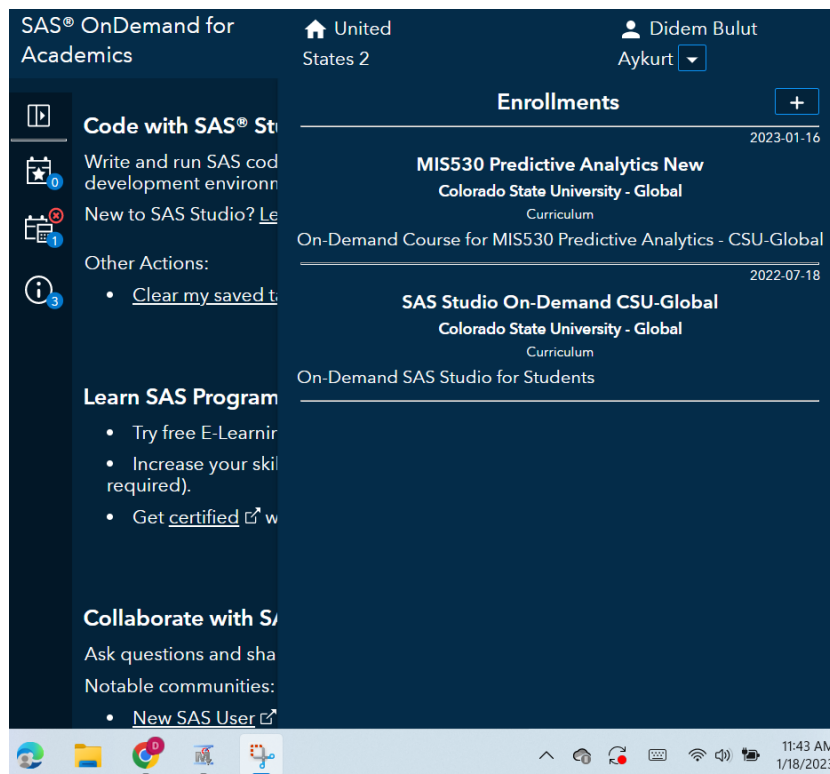
Colorado State University Global

MIS530; Predictive Analytics

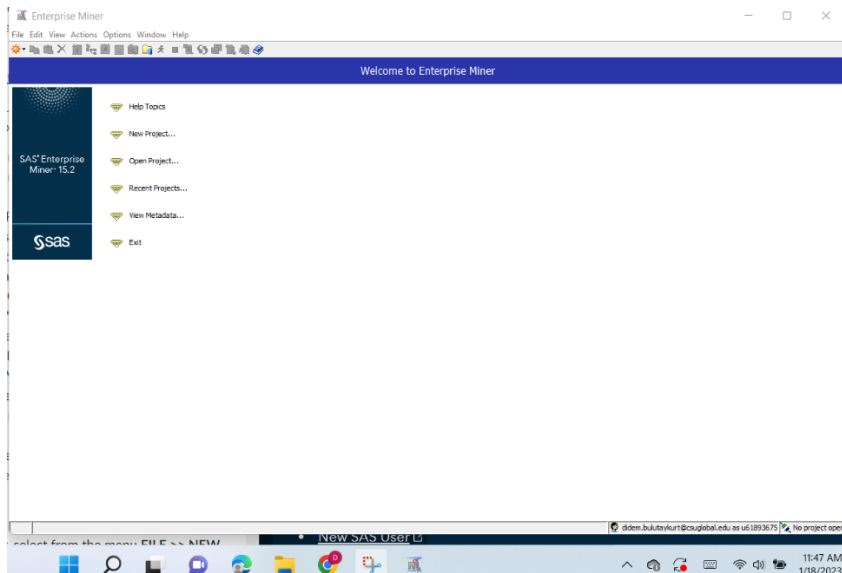
Dr.Jennifer Catalano

January 22, 2023

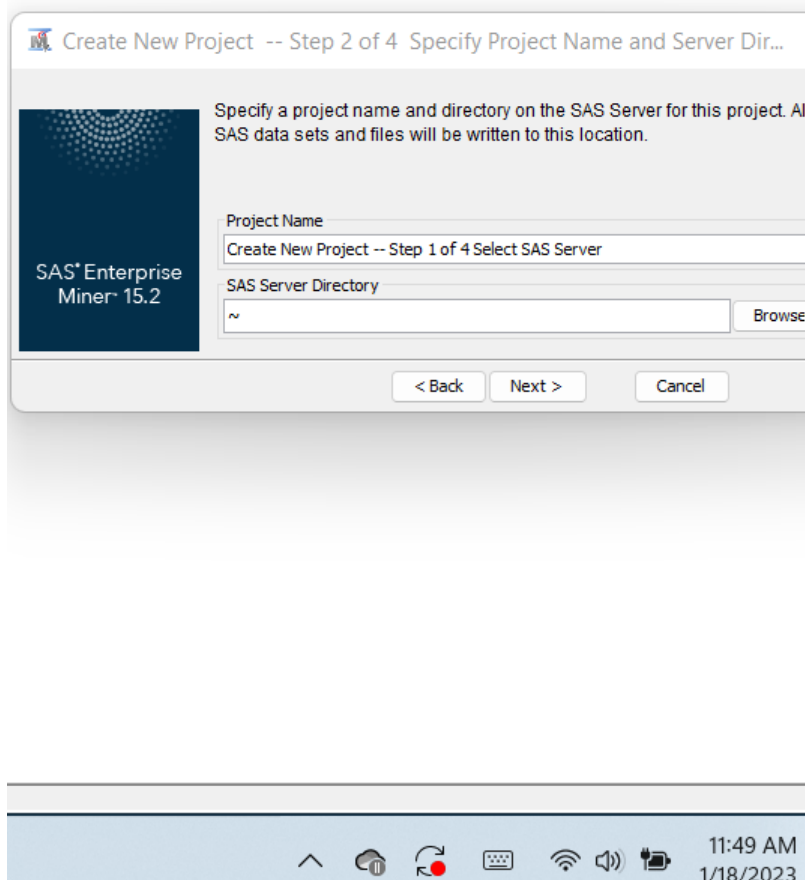
1-Course enrollment: I used the SAS OnDemand for Academic website to add a new course, left the site with the enrolment icon, and then clicked the plus icon to pop up the new tab to give system information.



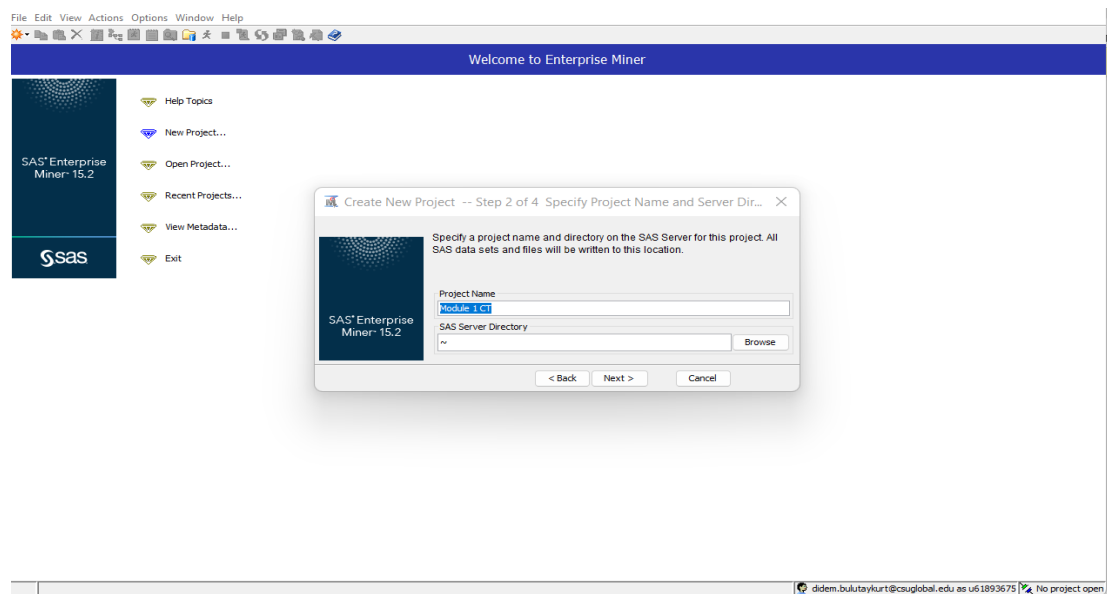
2-Running SAS Enterprise Miner: The left side of the SAS OnDemand for Academics page has a hammer icon that shows how to download and the URL to connect to SAS Enterprise Miner. This part was a little strange because when I gave an email and password to log on, it said, “Public access denied.” How I solved this: I uninstalled the program, cleared data on Google’s history page, and deleted all SAS links. Then I opened the SAS OnDemand Academic page to do all the steps without any open page and web. I also didn’t close the SAS Academic page; I got a URL link from that page and solved the problem.



3- New project; click Create a new project on the list, named “Create New Project – Step 1 of 4 Select SAS Server”.



4- New project named “Module 1 CT. “



5- A new project named "Create New Project – Step 3 of 4 Register the Project."

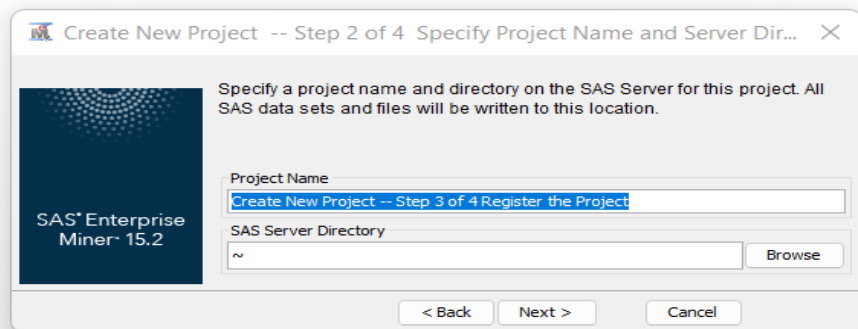
New Project...

Open Project...

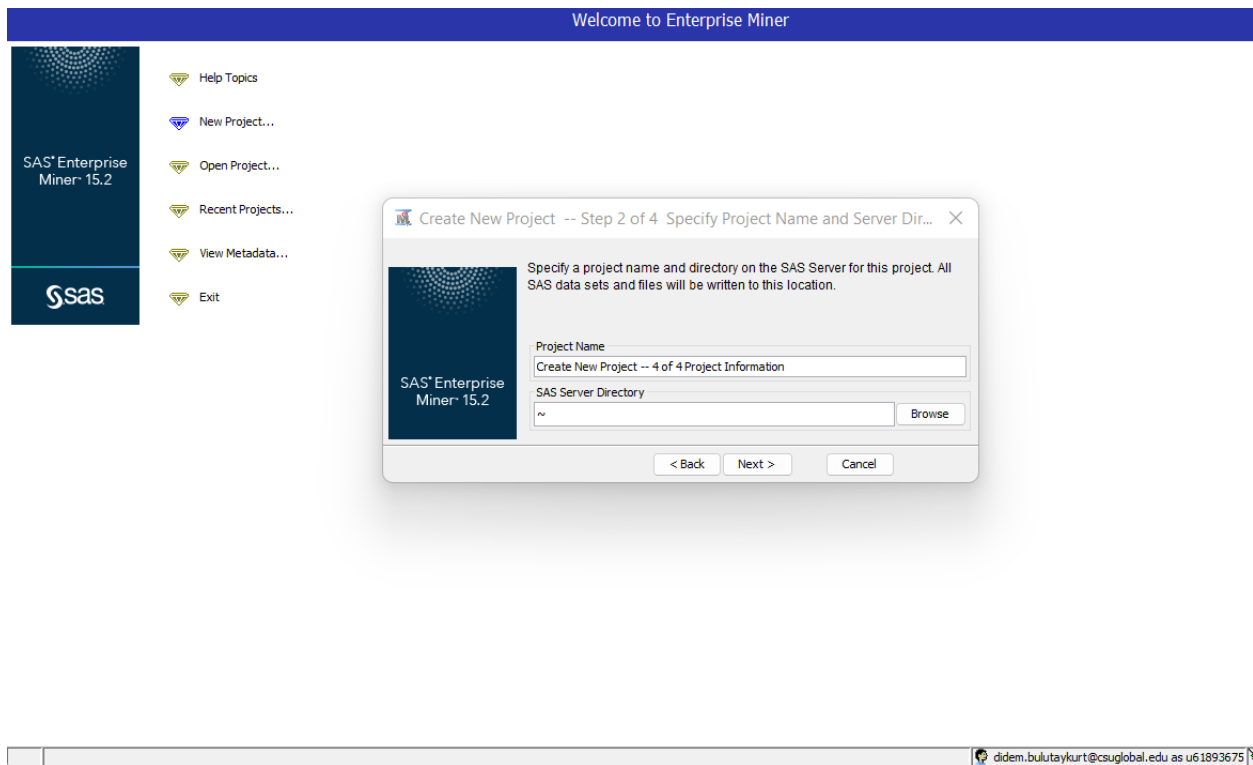
Recent Projects...

View Metadata...

Exit



6- A new project named "Create New Project – Step 4 of 4 Project Information."



7- Created Diagram; right-click on Diagrams and name "MODULE 1". That opens a page to work on it with the dataset. Then, I created a Library named "FRSTLIB" and chose the table CLAIM\_DATA\_V2\_TRAIN. Drag the data source from the project panel to the Diagram located in the workspace. Right-click on the data and run. Click the Sample tab, drag the Data Partition node onto the process flow, and then connect the Data Partition node to the CLAIM\_DATA\_TRAIN node. The sample data set is divided into 60% train, and 40% will be validation.

Data Exploration is the data preparation process. One of the tabs is the StatExplore node, which allows descriptive statistics to summarize the raw data like sums, average, and percent changes and describe the past. Also, prepare the data for future use in prescriptive analytics.

Enterprise Miner - Create New Project -- Step 1 of 4 Select SAS Server

File Edit View Actions Options Window Help

Sample Explore Modify Model Assess Utility Credit Scoring HPDM Applications Text Mining Time Series

Create New Project -- Step 1 of 4 Select SAS Server

Data Sources  
Diagrams  
Model Packages

Data Source Wizard -- Step 2 of 8

Select a SAS table

Table :

SAS Libraries

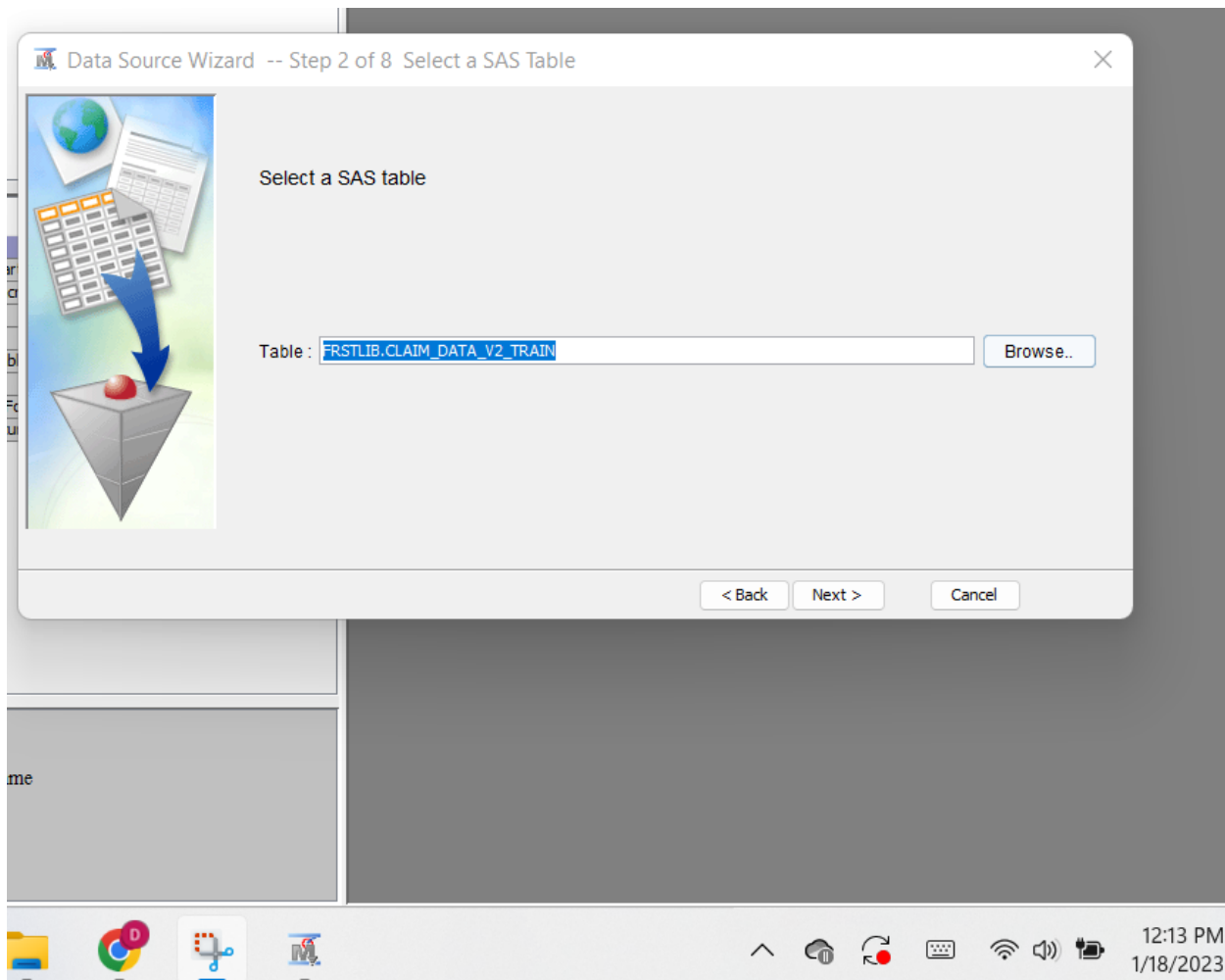
Name	Type
Ch5_lossdat_score2	Table
Ch6_binarytarget2_example	Table
Chari_btarg_b	Table
Chari_ntarg_b	Table
Claim	Table
Claim_data_v2_train	Table
Claim_score_dataset	Table
Data1	Table
Data2	Table
Datagov_smb_claim	Table
Dmretail	Table
Donor_raw_data	Table
Donor_score_data	Table
Federalist2	Table
Hour_xlsx	Table
Lossdat1_nominal	Table
Lossdat1_ordinal	Table
Lossdata6	Table
Lossfrequency	Table
Nn_resp_data	Table
Nn_resp_data2	Table
Nn_resp_score2	Table
Numri_btarg_c	Table
Numri_ntarg_b	Table
Pricetest_b	Table

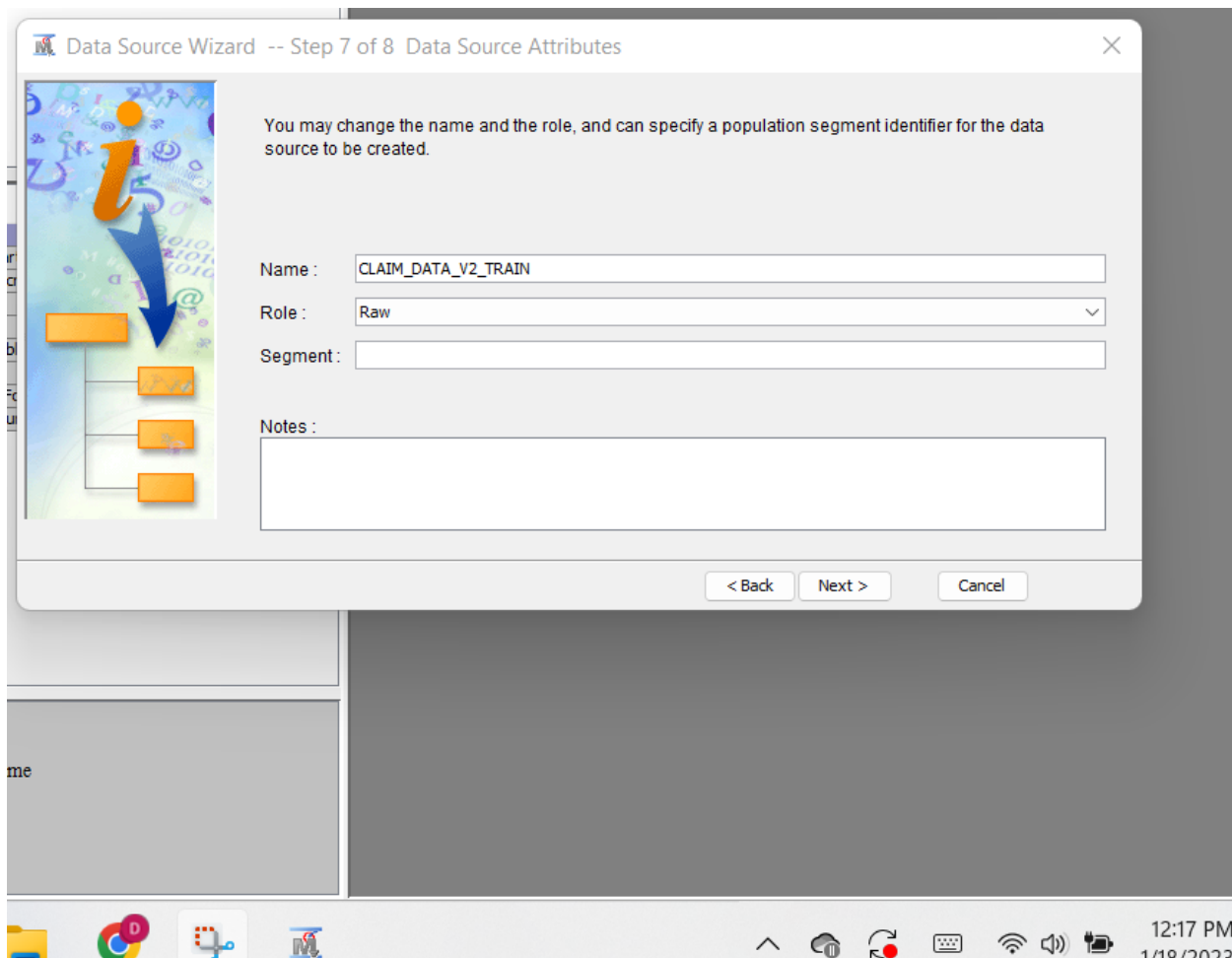
Get Details Properties... Refresh OK Cancel

Name

Project Name

12:07 PM  
1/18/2023







# Enterprise Miner - Create New Project -- Step 1 of 4 Select SAS Server

File Edit View Actions Options Window Help



Create New Project -- Step 1 of 4 Select SAS Server

- Data Sources
  - CLAIM\_DATA\_V2\_TRAIN
- Diagrams
  - Create SAS Dataset
  - MODULE 1
- Model Packages

Property	Value
<b>General</b>	
Node ID	Ids
Imported Data	...
Exported Data	...
Notes	...
<b>Train</b>	
Output Type	View
Role	Raw
Rerun	No
Summarize	No
Drop Map Variables	Yes
<b>Columns</b>	
Variables	...
Decisions	...
Refresh Metadata	...
Advisor	Basic
Advanced Options	...
<b>Data</b>	
Data Selection	Data Source
Sample	Default

## General

General Properties

Sample Explore Modify Model Assess Utility Credit Scoring HPDM Applications Text Mining Time Series

MODULE 1

CLAIM\_DATA\_V2\_TRAIN

Diagram Log



12:31 PM  
1/18/2023

File Edit View Actions Options Window Help

Sample Explore Modify Model Assess Utility Credit Scoring HPDM Applications Text Mining Time Series

MODULE 1

CLAIM\_DATA\_V2\_TRAIN

Data Partition

StatExplore

Property Value

**General**

Node ID Stat

Imported Data

Exported Data

Notes

**Train**

Variables

**Data**

Number of Observations 100000

Validation No

Test No

**Standard Reports**

Interval Distributions Yes

Class Distributions Yes

Level Summary Yes

Use Segment Variables No

Cross-Tabulation

**Variable Selection**

Hide Rejected Variables Yes

Number of Selected Variables 1000

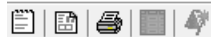
**General**

General Properties

12:39 PM  
1/18/2023

Results - Node: Data Partition Diagram: MODULE 1

File Edit View Window



Output

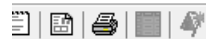
```
4 Time: 20:37:40
5 *-----*
6 * Training Output
7 *-----*
8
9
10
11
12 Variable Summary
13
14      Measurement      Frequency
15 Role      Level      Count
16
17 INPUT      INTERVAL      8
18 INPUT      NOMINAL      14
19
20
21
22
23 Partition Summary
24
25      Number of
26 Type      Data Set      Observations
27
28 DATA      EMWS2.Ids_DATA      5001
29 TRAIN      EMWS2.Part_TRAIN      2308
30 VALIDATE    EMWS2.Part_VALIDATE      1539
31 TEST       EMWS2.Part_TEST      1154
32
33
34 *-----*
35 * Score Output
36 *-----*
37
38
39 *-----*
40 * Report Output
41 *-----*
42
```



2:17 PM  
1/18/2023

Results - Node: CLAIM\_DATA\_V2\_TRAIN Diagram: MODULE 1

File Edit View Window



Output

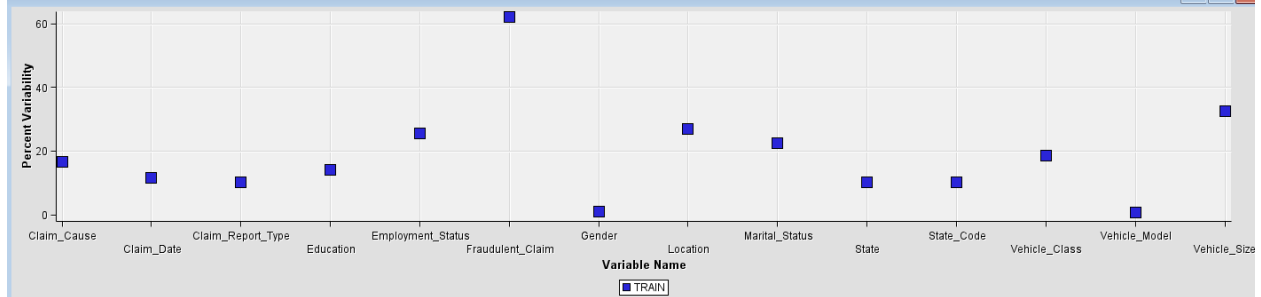
```

1  *-----*
2  User:          u61893675
3  Date:          18 January 2023
4  Time:          20:37:34
5  *-----*
6  * Training Output
7  *-----*
8
9
10
11
12 Variable Summary
13
14      Measurement   Frequency
15 Role      Level      Count
16
17 INPUT     INTERVAL      8
18 INPUT     NOMINAL      14
19

```

Variables

Variable Name	Role	Measurement Level
Annual Premium	Input	Interval
Claim Amount	Input	Interval
Claim Cause	Input	Nominal
Claim Date	Input	Nominal
Claim Report Type	Input	Nominal
Claimant Number	Input	Interval
Education	Input	Nominal
Employment Status	Input	Nominal
Fraudulent Claim	Input	Nominal
Gender	Input	Nominal
Income	Input	Interval
Location	Input	Nominal
Marital Status	Input	Nominal
Monthly Premium	Input	Interval
Months Since Last Claim	Input	Interval
Months Since Policy Inception	Input	Interval



## Output

```
1 *-----*
2 User:      u61893675
3 Date:      18 January 2023
4 Time:      20:38:17
5 *-----*
6 * Training Output
7 *-----*
8
9
10
11
12 Variable Summary
13
14 Role      Measurement Level      Frequency
15
16
17 INPUT     INTERVAL      8
18 INPUT     NOMINAL       14
19
```

## Output

```
22 Class Variable Summary Statistics
23 (maximum 500 observations printed)
24
25 Data Role=TRAIN
26
27
28 Data
29 Role      Variable Name      Role      Number
30                                     of
31                                     Levels
32 TRAIN     Claim_Cause          INPUT     5
33 TRAIN     Claim_Date            INPUT     3
34 TRAIN     Claim_Report_Type     INPUT     4
35 TRAIN     Education             INPUT     6
36 TRAIN     Employment_Status      INPUT     5
37 TRAIN     Fraudulent_Claim       INPUT     2
38 TRAIN     Gender                 INPUT     2
39 TRAIN     Location               INPUT     4
40 TRAIN     Marital_Status         INPUT     3
41 TRAIN     State                  INPUT     5
```

				Number of Levels	Missing	Mode	Mode Percentage	Mode2	Mode2 Percentage
31	TRAIN	Claim_Cause	INPUT	5	0	Collision	41.53	Hail	31.55
32	TRAIN	Claim_Date	INPUT	3	0	12/01/2018	39.99	12/15/2018	39.99
33	TRAIN	Claim_Report_Type	INPUT	4	0	Agent	37.77	Branch	28.21
34	TRAIN	Education	INPUT	6	11	College	29.83	Bachelor	29.81
35	TRAIN	Employment_Status	INPUT	5	0	Employed	62.51	Unemployed	25.39
36	TRAIN	Fraudulent_Claim	INPUT	2	0	N	93.86	Y	6.14
37	TRAIN	Gender	INPUT	2	0	M	50.61	F	49.39
38	TRAIN	Location	INPUT	4	3	Suburban	63.11	Rural	19.62
39	TRAIN	Marital_Status	INPUT	3	0	Married	58.59	Single	26.41
40	TRAIN	State	INPUT	5	0	Iowa	30.91	Missouri	28.91

Output											
49	(maximum 500 observations printed)										
50											
51	Data Role=TRAIN										
52											
53											
54	Variable	Role	Mean	Standard Deviation	Non Missing	Missing	Minimum	Median	Maximum	Skewness	Kurtosis
55											
56	Annual_Premium	INPUT	1134.368	313.4503	5001	0	600	1140	1680	0.016055	-1.17369
57	Claim_Amount	INPUT	787.7633	655.9633	5001	0	189.8684	577.3521	7422.852	2.922369	12.62942
58	Claimant_Number	INPUT	3501	1443.809	5001	0	1001	3500	6001	0	-1.2
59	Income	INPUT	41310.45	227690.4	5001	0	0	34621	15967801	68.48667	4790.542
60	Monthly_Premium	INPUT	94.53069	26.12086	5001	0	50	95	140	0.016055	-1.17369
61	Months_Since_Last_Claim	INPUT	15.0042	11.13965	5001	0	0	13	60	0.666671	0.052198
62	Months_Since_Policy_Inception	INPUT	48.23495	28.09665	5001	0	0	48	99	0.051865	-1.14433
63	Outstanding_Balance	INPUT	23728.15	13827.4	4992	9	4	23988	47996	0.005256	-1.20029
64											
65											
66	*-----*										
67	* Score Output										

## Reference

Richard V. McCarthy; Mary M. McCarthy; Wendy Ceccucci, 2022. *Applying Predictive Analytics Finding Value in Data*. Second edition.