

## **Data Preparation with Claim Fraud Dataset/ SAS Enterprise Miner**

Didem B. Aykurt

Colorado State University Global

MIS530; Predictive Analytics

Dr.Jennifer Catalano

January 29, 2023

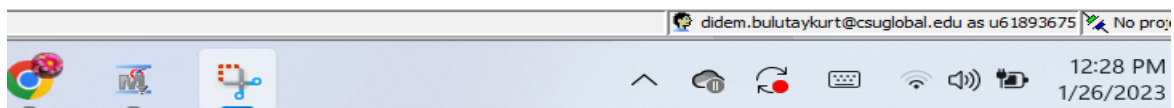
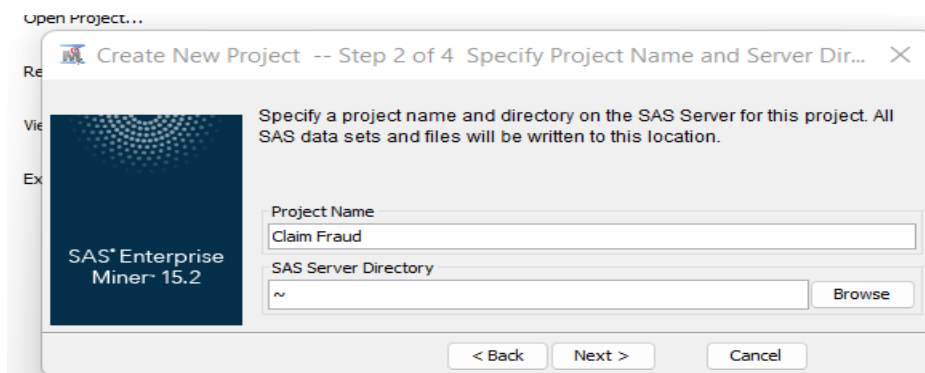
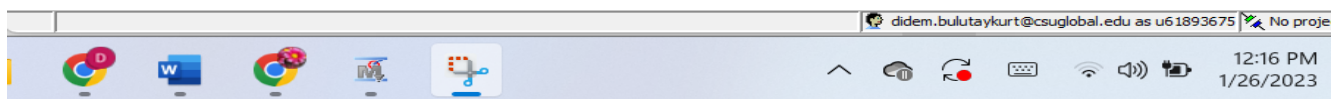
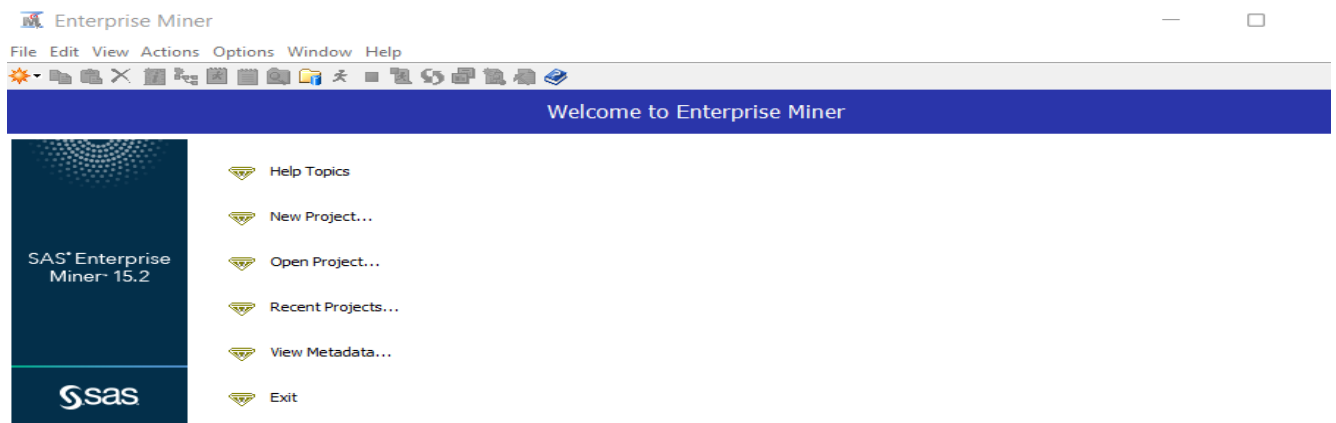
## **Data Preparation Using the SAS Enterprise Miner**

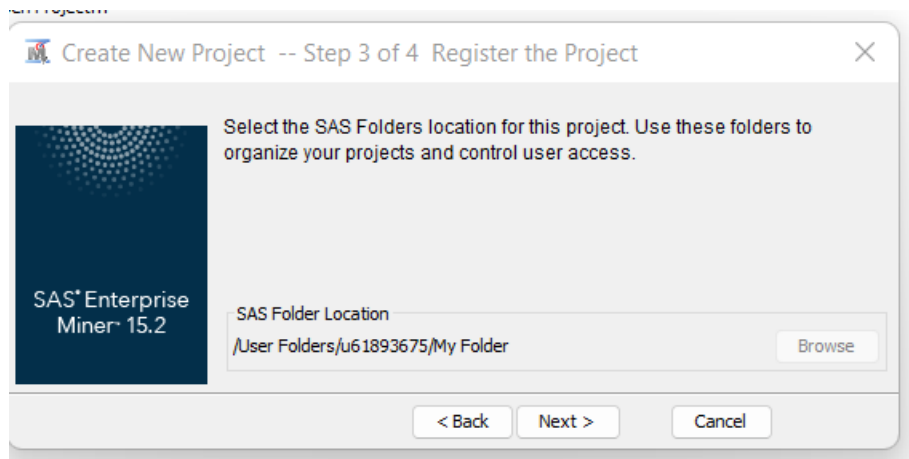
When I start to use a new programming language or software, I am always searching for why I need it or the pros and cons of this program. The first time I faced it, it was not fast and easy to use. They should update the program. Also, most learning tools were updated a few years ago; video or other devices do not support the new version. It is so complicated. The program has limited learning sources online, which tells me a lower percentage of people use this software program.

Searching for the SAS Enterprise Miner is a solution to creating accurate predictive and descriptive models on large volumes of data across different sources in the organization. (Pat Research,2021) Some business applications are for detecting fraud, minimizing risk and resource demands, reducing asset downtime and campaigns, and reducing customer attrition—the most popular category used for predictive analytics. Let's dive deep into SAS Enterprise Miner data processing step by step. There are four components for the predictive creation model in SAS Enterprise Miner:

### **Creating the Project File**

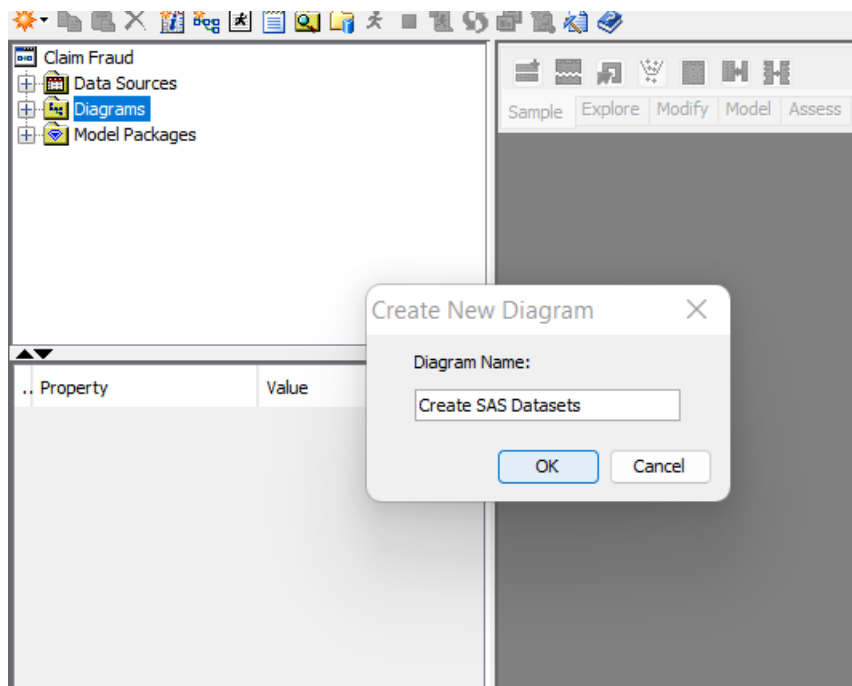
First launched, the Welcome screen appears. Create a new project, select New Project, and open new windows. I entered the project name "Claim Fraud," and the SAS Server Directory showed where the file was saved. I left the automatic field location and finished.

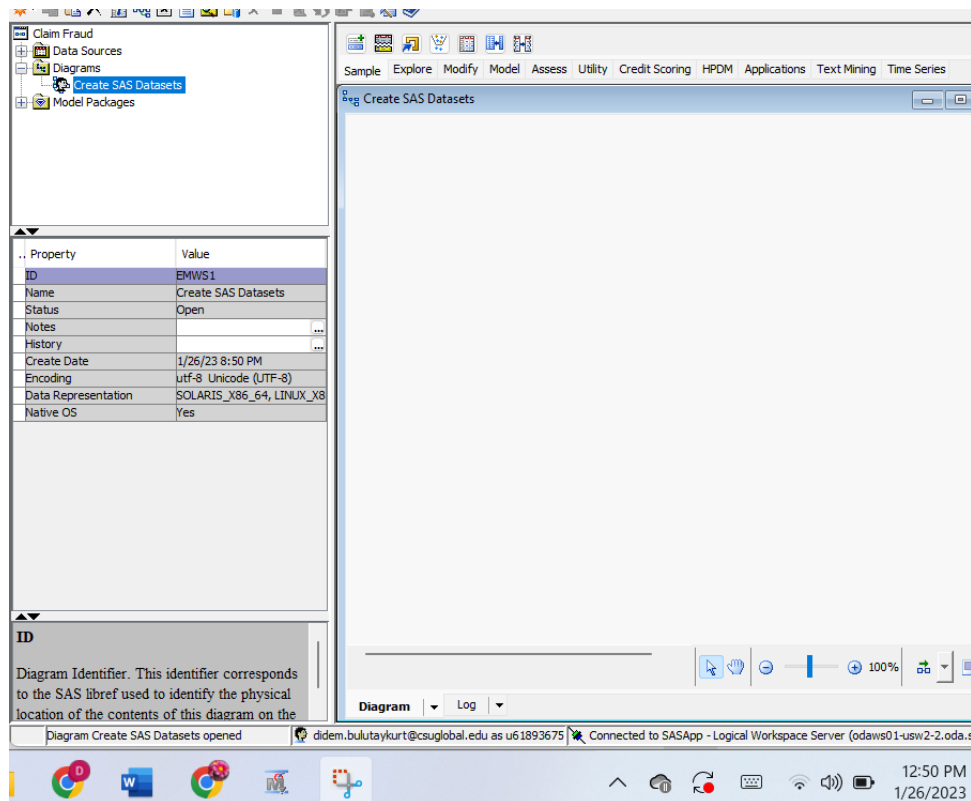




## Create the Diagram

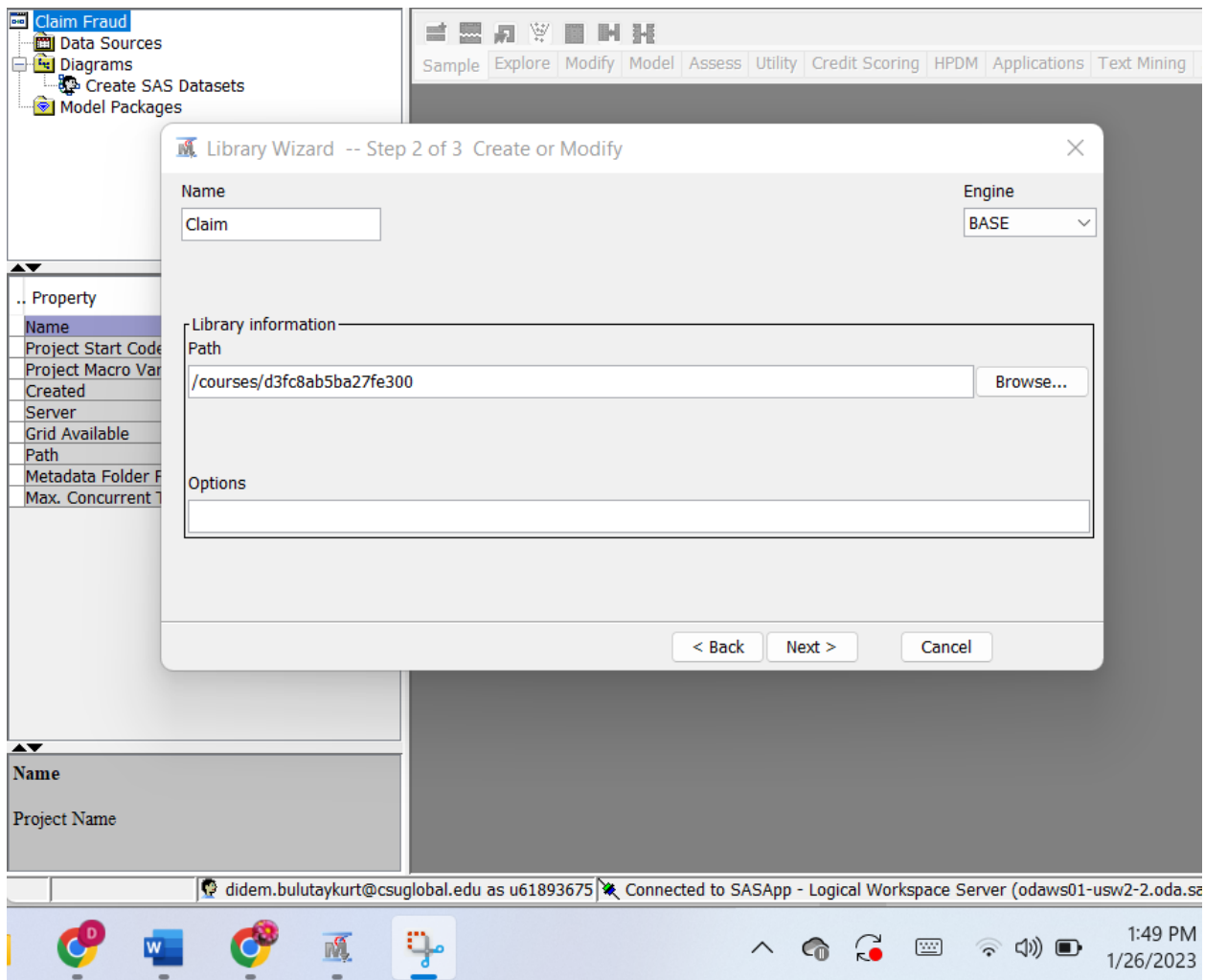
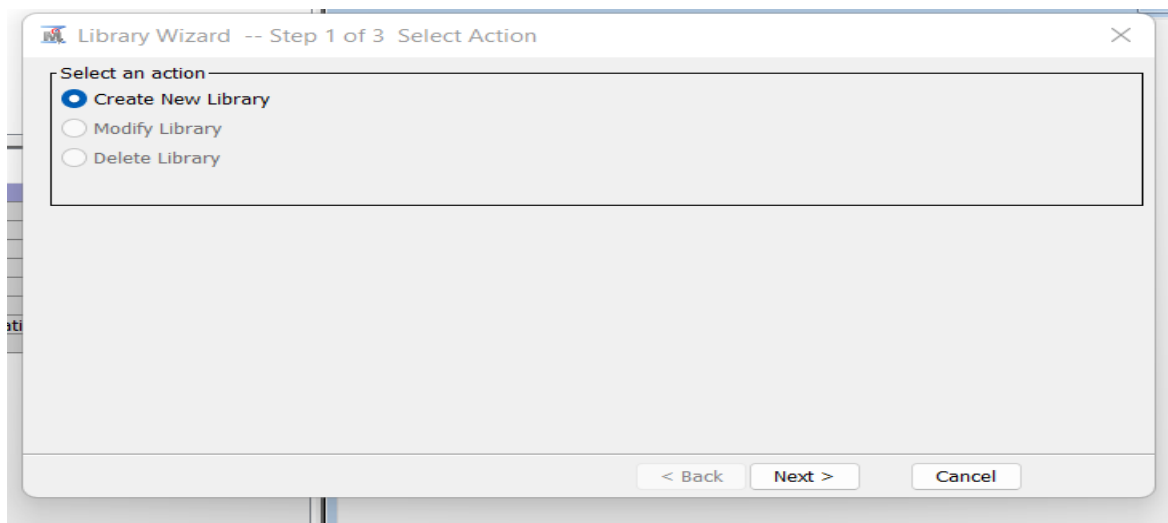
When the project window opens, the left side has a project panel list; right-click Diagrams and select Create Diagram, then pop up the new window named the diagram "Create SAS Datasets."





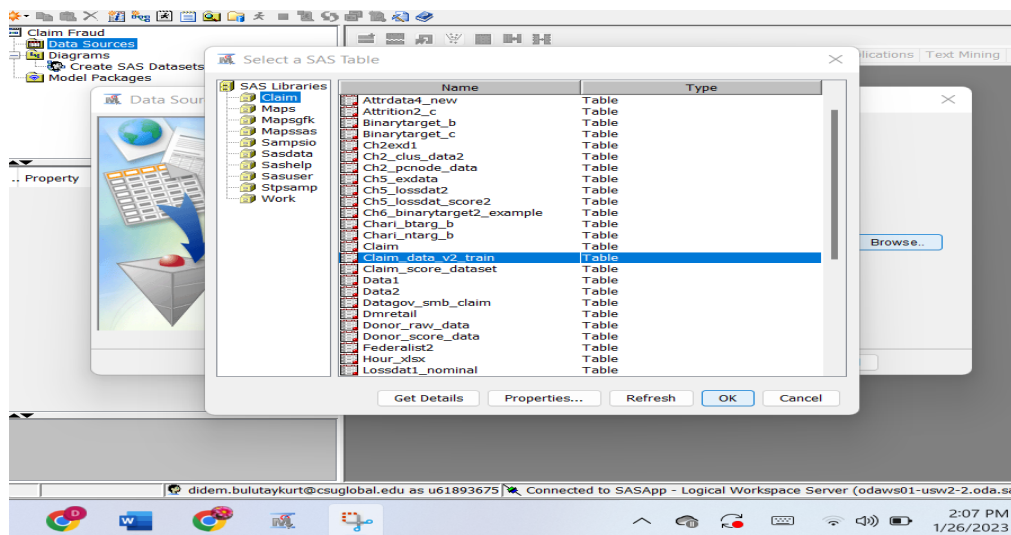
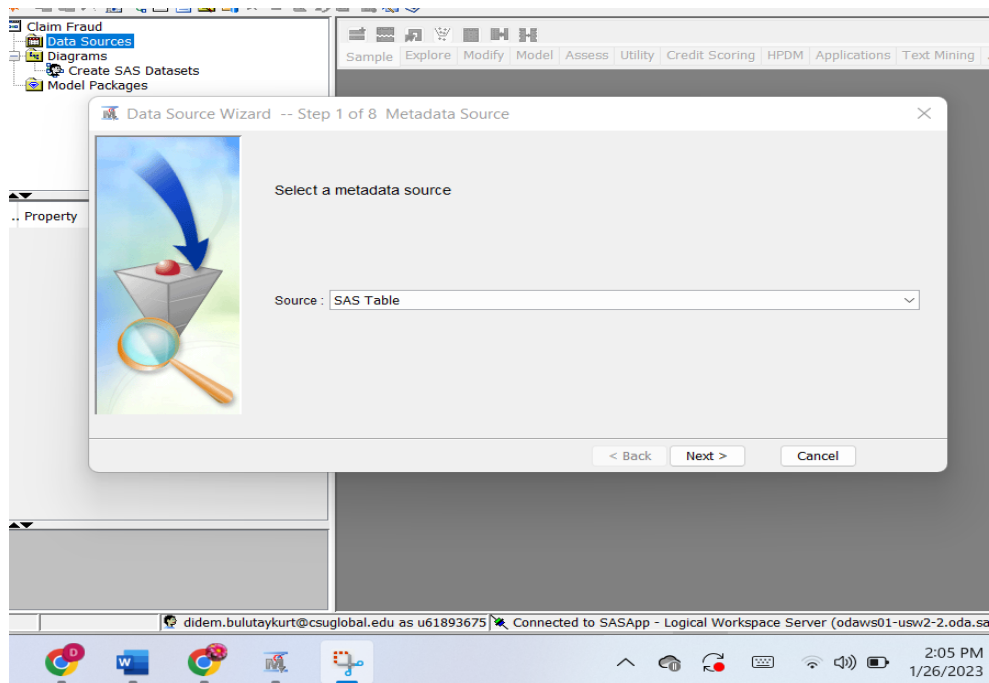
## Create the Library

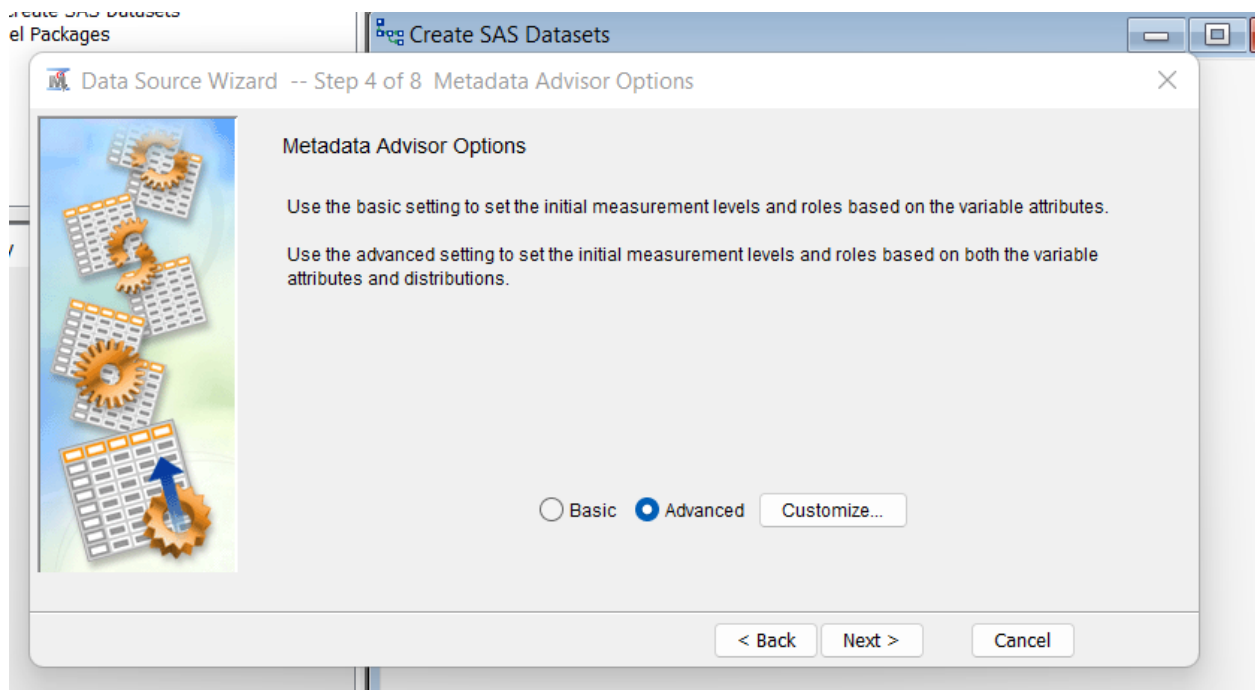
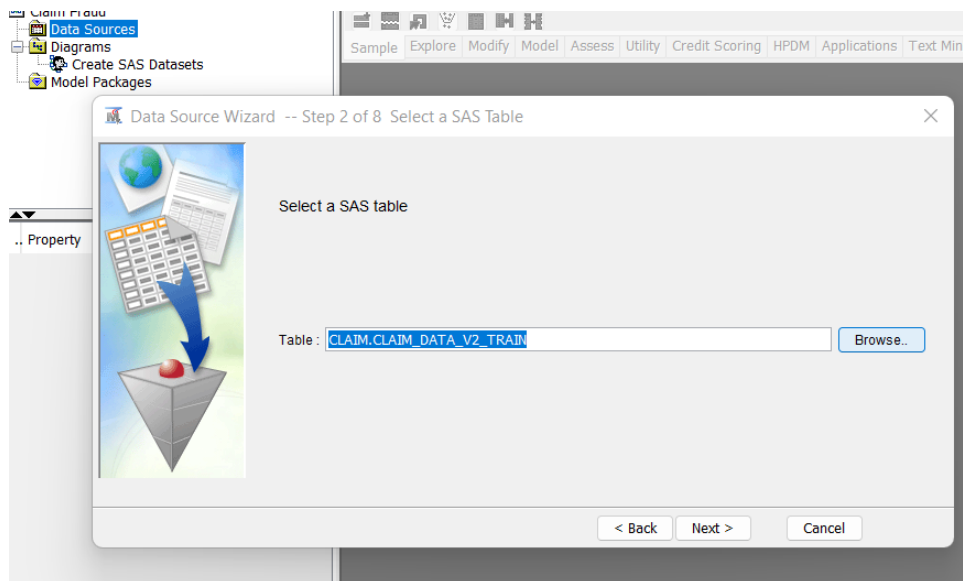
The SAS library creates the directory location of the stored SAS data files. To build a SAS library, click on the File Menu-New-Library, then the library wizard opens, and I named "Claim" the used library path as /courses/d3fc8ab5ba27fe300 next and finished. This library path already loaded the claim dataset.



## Create Data Source

Now, I have four components to start creating the data sources. I used the auto insurance claim data set that had already loaded. To make the data sources, click the File menu-New-Data Source, or the left side Project panel has Data Sources; Right click, then pop up Data Sources Wizard and select Claim.





To change roles and levels, the next Metadata Advisor Option. The `Fraudulent_Claim` variable's column Role change Input to Target as a dependent variable; it is binary (yes/no). The `Claimant_Number`'s Role column changed Input to ID variable. The state variable will not be



used in the model as changed Role Input to Reject.

Enterprise Miner - Claim Fraud

File Edit View Actions Options Window Help

Claim Fraud  
Data Sources  
Diagrams  
Create SAS Datasets  
Model Packages

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

Create SAS Datasets

Data Source Wizard -- Step 5 of 8 Column Metadata

(none) ☐ not Equal to ☐ Apply Reset


Columns: ☐ Label ☐ Mining ☐ Basic ☐ Statistics


Name	Role	Level	Report	Order	Drop	Lower Limit
Annual_Premiu	Input	Interval	No		No	.
Claim_Amount	Input	Interval	No		No	.
Claim_Cause	Input	Nominal	No		No	.
Claim_Date	Input	Nominal	No		No	.
Claim_Report	Input	Nominal	No		No	.
Claimant_Nur	ID	Interval	No		No	.
Education	Input	Nominal	No		No	.
Employment_S	Input	Nominal	No		No	.
Fraudulent_Cla	Target	Binary	No		No	.
Gender	Input	Binary	No		No	.
Income	Input	Interval	No		No	.
Location	Input	Nominal	No		No	.
Marital_Status	Input	Nominal	No		No	.

Show code Explore Refresh Summary < Back Next > Cancel

Diagram Log

Diagram Create ... didem.bulutaykurt@csuglobal.edu as u61893675 Connected to SASApp - Logical Workspace Server (odaws01-usw2-2.oda.

 Data Source Wizard -- Step 6 of 10 Decision Configuration



### Decision Processing

Do you want to build models based on the values of the decisions ?

If you answer yes, you may enter information about the cost or profit of each possible decision, prior probability and cost function. The data will be scanned for the distributions of the target variables.

☒ No ☐ Yes

< Back Next > Cancel

Claim Fraud


**Data Sources**

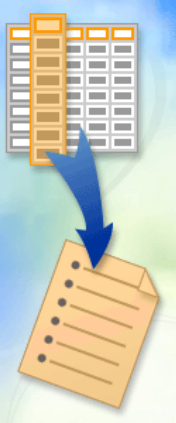
Diagrams

Create SAS Datasets

Model Packages

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

 Data Source Wizard -- Step 6 of 8 Create Sample



Do you wish to create a sample data set?

☒ No ☐ Yes

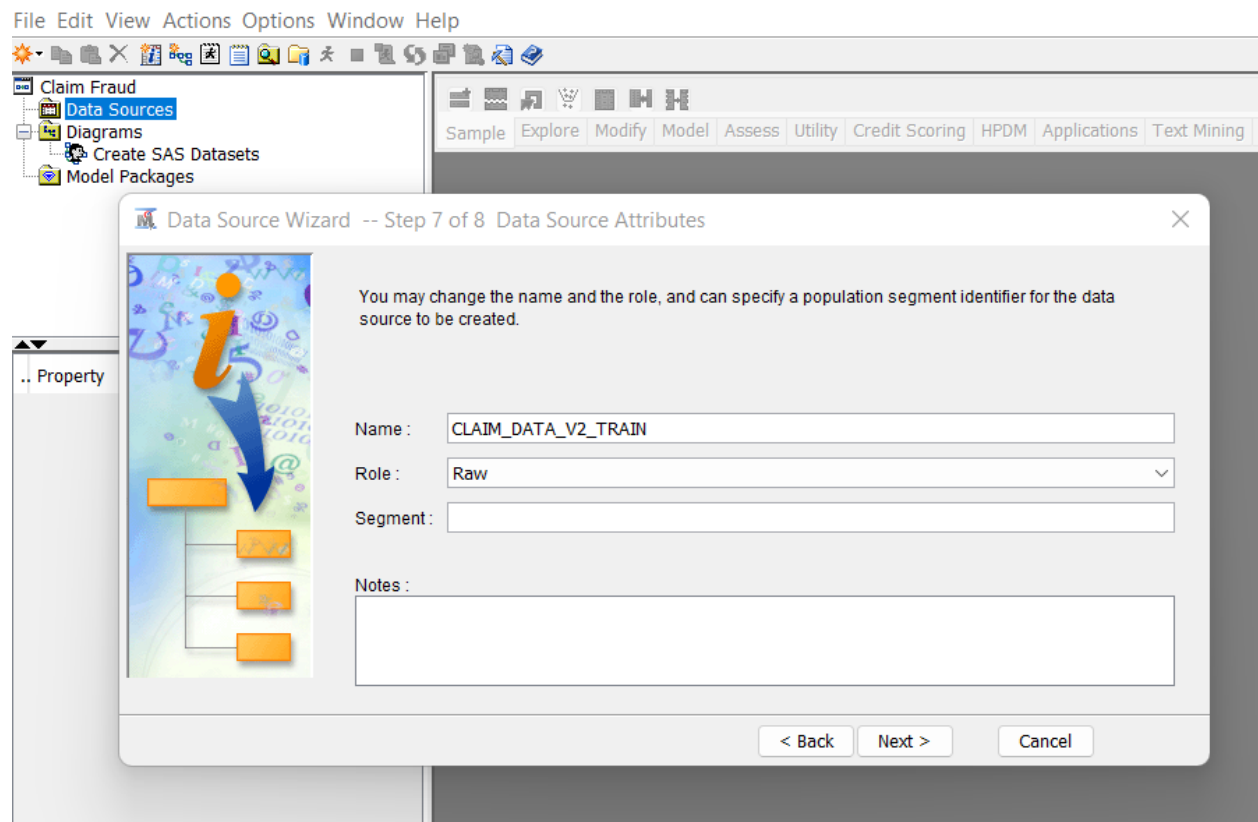
**Table Info**

Columns	22
Rows	5001

**Sample Size**

Type	Percent
Percent	20
Rows	

< Back Next > Cancel



The next step is I drag the CLAIM\_DATA\_V2\_TRAIN from the Project Panel to the Diagram  
Create SAS Datasets in the workspace. Then right-click on the node and RUN. Two ways to look  
at results. One is that the end of the run window has a result option or right-click on the dataset  
and Result.

Claim Fraud

- Data Sources
  - CLAIM\_DATA\_V2\_TRAIN
- Diagrams
  - Create SAS Datasets
- Model Packages

Property Value

General	
Node ID	Ids
Imported Data	...
Exported Data	...
Notes	
Train	
Output Type	View
Role	Raw
Re-run	No
Summarize	No
Drop Map Variables	Yes
Columns	
Variables	...
Decisions	...
Refresh Metadata	...
Advisor	Basic
Advanced Options	...

Sample Explore Modify Model Assess Utility Credit Scoring HPDM Applications Text

Create SAS Datasets

CLAIM\_DATA\_V2\_TRAIN

Run Status

Run completed  
Diagram: Create SAS Datasets  
Path: CLAIM\_DATA\_V2\_TRAIN

OK Results...

Output

```

1 *-----*
2 User:          u61893675
3 Date:          26 January 2023
4 Time:          23:39:07
5 *-----*
6 * Training Output
7 *-----*
8
9
10
11
12 Variable Summary
13
14      Measurement   Frequency
15 Role      Level      Count
16
17 ID          INTERVAL      1
18 INPUT       BINARY        1
19 INPUT       INTERVAL      7

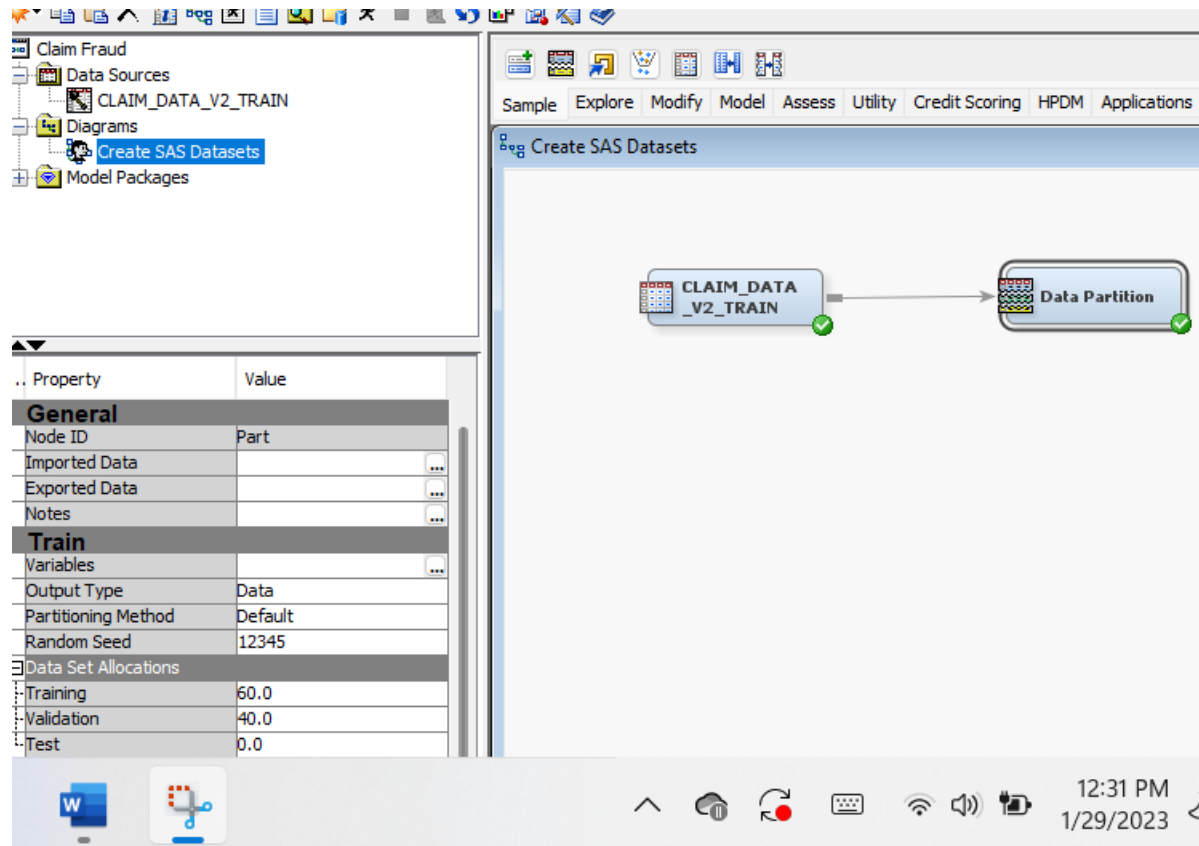
```

Variables

Variable Name	Role	Measurement Level	Order
Annual Premium	Input	Interval	
Claim Amount	Input	Interval	
Claim Cause	Input	Nominal	
Claim Date	Input	Nominal	
Claim Report Type	Input	Nominal	
Claimant Number	ID	Interval	
Education	Input	Nominal	
Employment Status	Input	Nominal	
Fraudulent Claim	Target	Binary	
Gender	Input	Binary	
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	
Outstanding Balance	Input	Interval	
State	Rejected	Nominal	

## Data Partition

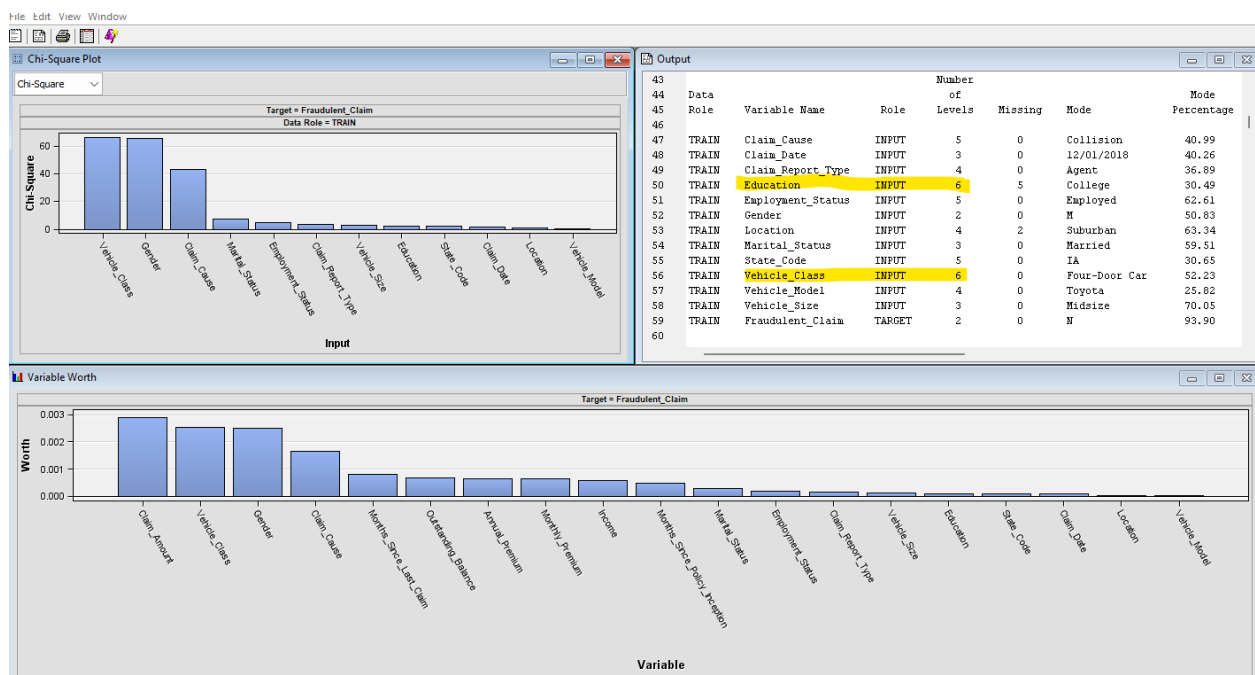
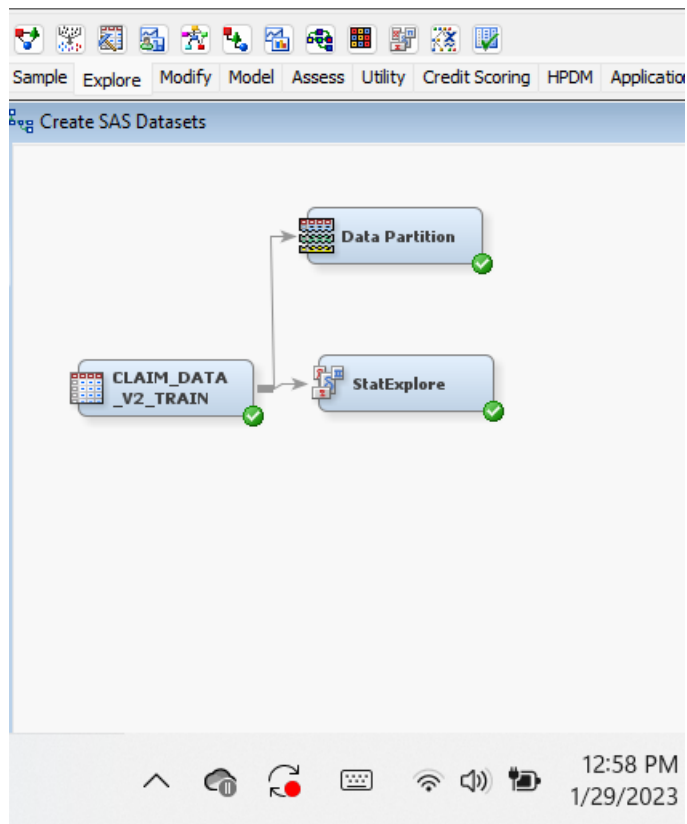
The data set divides the dataset into two or three parts: training is used to build the model, validation checks the model's accuracy, and testing partition tests the model. To create the section, click the Sample tab, then drag the Data Partition node onto the process flow. Finally, connect the Data Partition node to the CLAIM\_DATA\_TRAIN node. I divided the training partition into 60% and the validation partition into 40%.



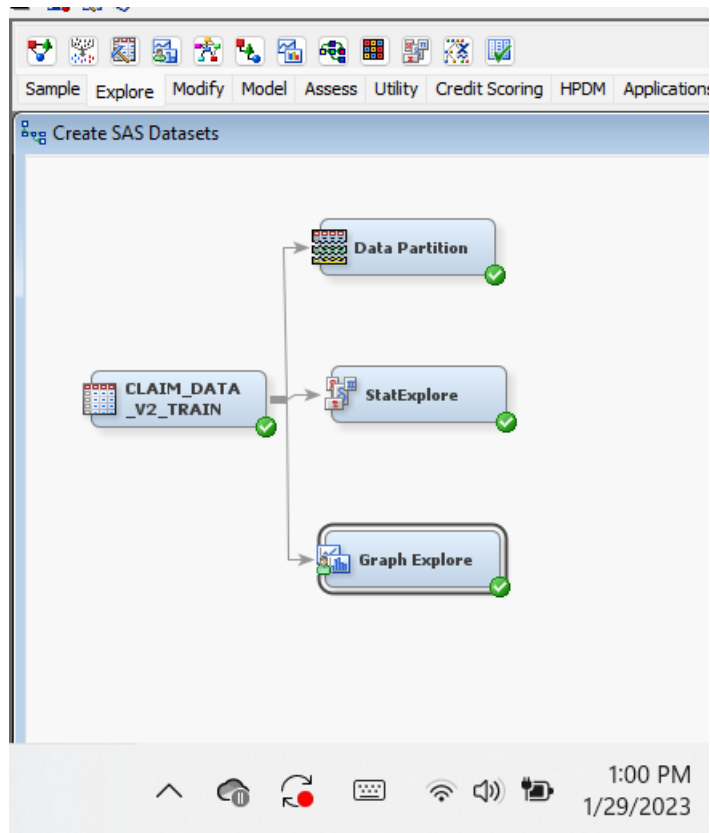
## Data Exploration

The StatExplore node shows the data summary and identifies the missing values. The Graph Explore node helps to see data point behavior by creating a histogram, stem-and-leaf plots, and box plot; both nodes are significant for descriptive statistics.

Click the Explore tab, then drag the StatExplore node onto the process flow. Finally, connect the StatExplore node to the CLAIM\_DATA\_TRAIN node and right-click on the node-run.



The Graph Explore node; Click the Explore tab, then drag the Graph Explore node onto the process flow. Connect the Graph Explore node to the CLAIM\_DATA\_TRAIN node and right-click on the node-run.



After running the Graph Explore node, right-click and select Edit Variable. For sorting by column, click on the Role column title to sort by Role. For example, to determine all the input variables you want, click the Marital\_Status variable name, hold the shift down, and click on Education.

## Variables - GrfExpl

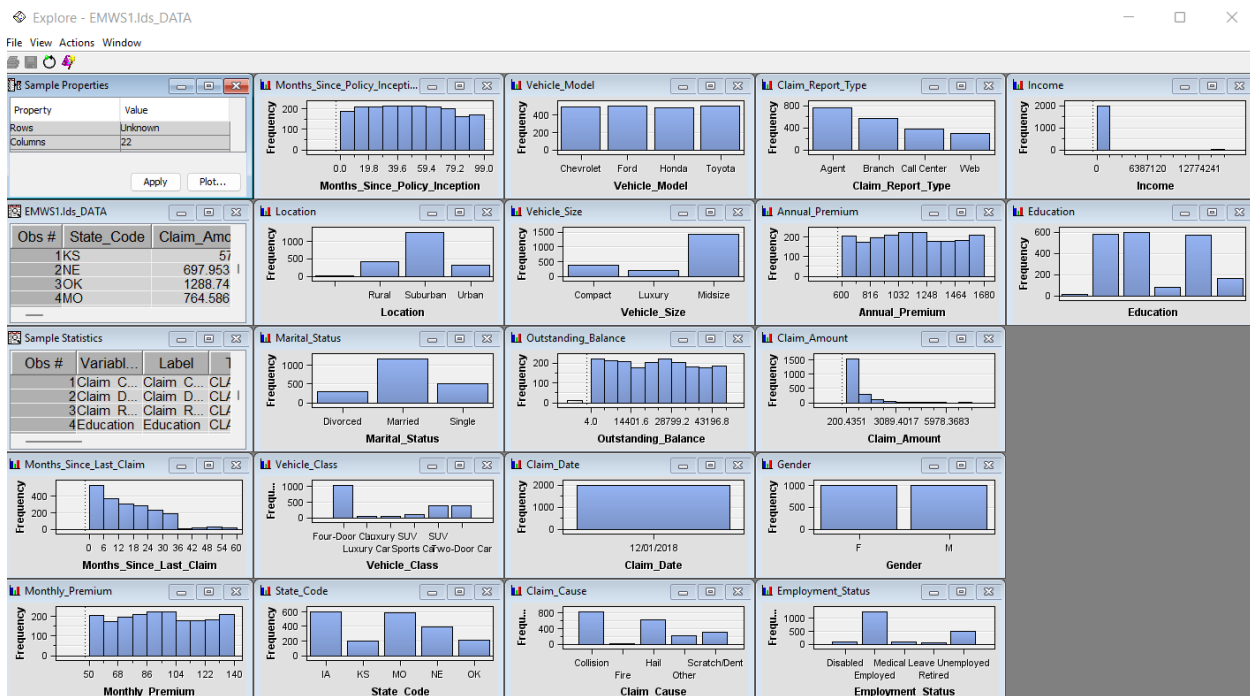
(none) ☐ not Equal to

Columns: ☐ Label ☐ Mining

Name	Use	Report	Sample Role	Role	Level
Claimant_Numb	Default	No	Default	ID	Interval
Months_Since_L	Default	No	Default	Input	Interval
Monthly_Premiur	Default	No	Default	Input	Interval
Months_Since_P	Default	No	Default	Input	Interval
Location	Default	No	Default	Input	Nominal
Marital_Status	Default	No	Default	Input	Nominal
Vehicle_Class	Default	No	Default	Input	Nominal
State_Code	Default	No	Default	Input	Nominal
Vehicle_Model	Default	No	Default	Input	Nominal
Vehicle_Size	Default	No	Default	Input	Nominal
Outstanding_Ba	Default	No	Default	Input	Interval
Claim_Date	Default	No	Default	Input	Nominal
Claim_Cause	Default	No	Default	Input	Nominal
Claim_Report_T	Default	No	Default	Input	Nominal
Annual_Premium	Default	No	Default	Input	Interval
Claim_Amount	Default	No	Default	Input	Interval
Gender	Default	No	Default	Input	Binary
Employment_Sta	Default	No	Default	Input	Nominal
Income	Default	No	Default	Input	Interval
Education	Default	No	Default	Input	Nominal
State	Default	No	Default	Rejected	Nominal
Fraudulent_Clair	Default	No	Default	Target	Binary

1:43 PM  
1/29/2023

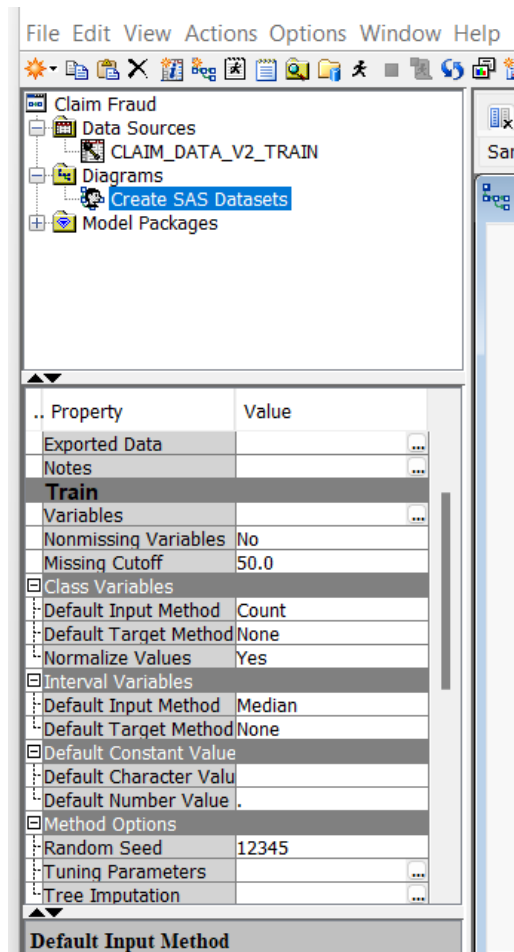
Click the Explore button, and a new window displays a histogram of all the input variables. And the Income variable has an outlier.

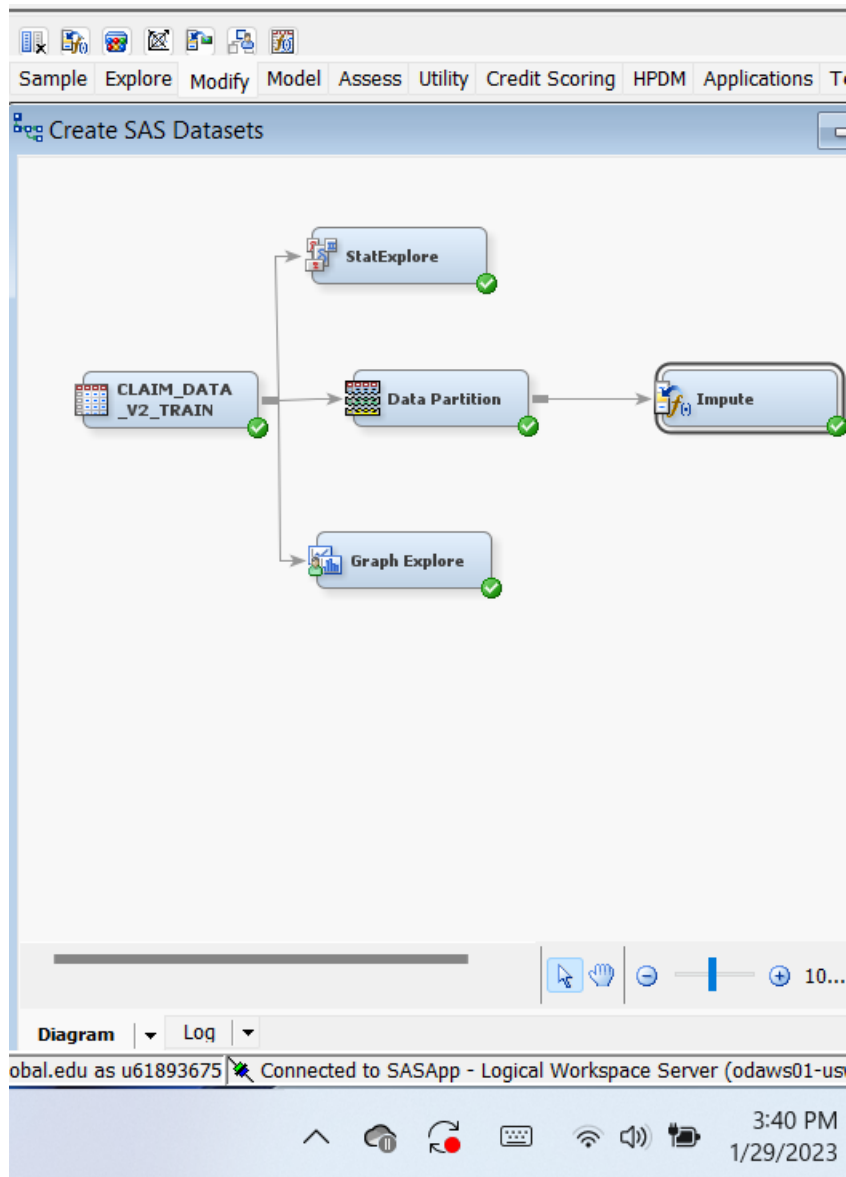




## Missing Data

Missing values can cause the model result, which is why to improve data. Two methods to handle missing values; are listwise deletion and imputation, replacing the missing values with substitute values. The Impute node is used to deal with missing values. Click the Modify tab, then drag the Impute node to the process flow diagram. The Default Input Method for the Interval Variables changed from Mean to Median.



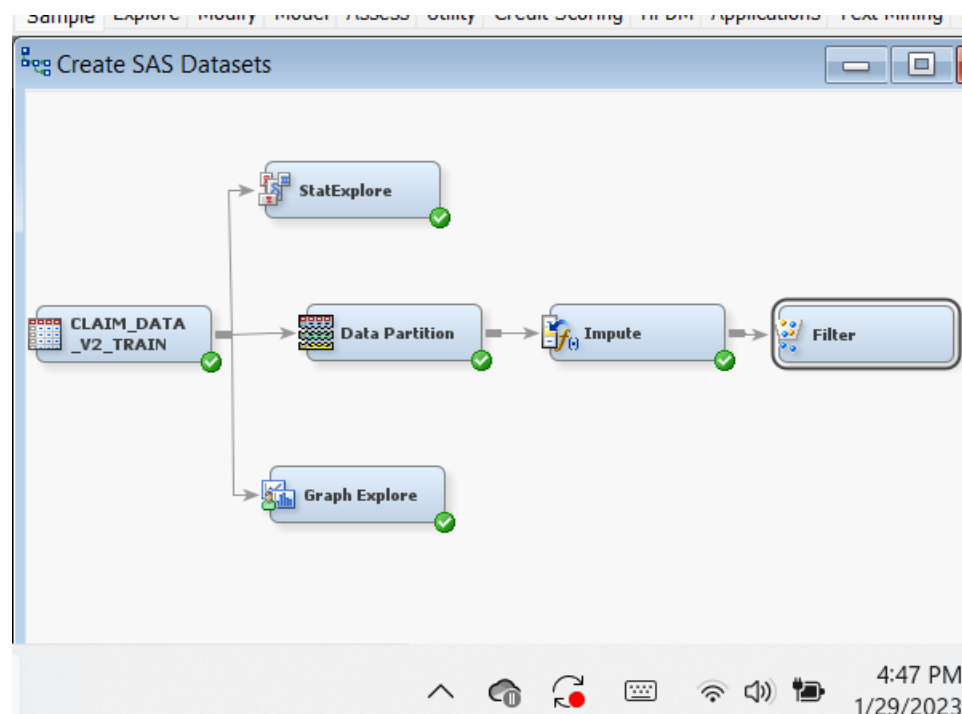


The Impute node creates a new table with a new variable replacement value for the missing data. The Impute node creates a default table as the original data set in the original dataset are not overwritten. New variables containing the impute values can identify with the prefix IMP\_.

28	Output							
29	*-----*							
30	* Report Output							
31	*-----*							
32								
33								
34								
35	Imputation Summary							
36	Number Of Observations							
37								
38								
39		Impute			Measurement		Number of	
40	Variable Name	Method	Imputed Variable	Impute Value	Role	Level	Label	Missing
41								for TRAIN
42	Education	COUNT	IMP_Education	College	INPUT	NOMINAL	Education	5
43	Location	COUNT	IMP_Location	Suburban	INPUT	NOMINAL	Location	2
44	Outstanding_Balance	MEDIAN	IMP_Outstanding_Balance	24069.5	INPUT	INTERVAL	Outstanding_Balance	4
45								
46								

## Handling Outliers

The Filter node helps to identify and eliminate outliers and filter the dataset. Filtering uses the data from the training dataset for the better result of models. Additionally, the Filter node ignores target and rejected variables. Use the Sample Tab, drag a Filter node to the process flow diagram, and connect the Impute to the Filter node.



Let's look at the Claim fraud dataset, as only Income has outliers on the Graph Explore node showed. Just the Income variable needs to be filtered; the filter setting needs to change under the Train group, Table to Filter set All Data Sets, and under the Interval Variables group, click

the Interval Variables' ellipsis and set the minimum income zero and the maximum income to 103,677.78 that value probably 99.7% of the values are within three standard deviations. The Default Filtering Method is set to User-Specified Limits.

The screenshot displays the SAS Studio interface. At the top, a workflow diagram titled "Create SAS Datasets" shows a sequence of steps: "CLAIM\_DATA\_V2\_TRAIN" (with a green checkmark), "Data Partition" (with a green checkmark), "Impute" (with a green checkmark), and "Filter" (with a green checkmark). From "CLAIM\_DATA\_V2\_TRAIN", arrows also point to "StatExplore" and "Graph Explore", both of which also have green checkmarks.

Below the workflow diagram, the "Interactive Interval Filter" dialog is open. It contains a table with columns: Name, Label, Keep Missing Values, Filter Lower Limit, Filter Upper Limit, Role, and Level. The table lists several variables, with "Income\_Default" highlighted, showing a Filter Lower Limit of 0 and a Filter Upper Limit of 103677.8.

To the right of the dialog, a "Property" pane is visible, showing various settings for the "Filter" step. The "Default Filtering Method" is set to "User-Specified Limits".

Property	Value
Export Table	Filtered
Tables to Filter	All Data Sets
Distribution Data Sets	Yes
Class Variables	
Class Variables	
Default Filtering Method	Rare Values (Percentage)
Keep Missing Values	Yes
Normalized Values	Yes
Minimum Frequency Cutoff	1
Minimum Cutoff for Percentage	0.01
Maximum Number of Levels	25
Interval Variables	
Interval Variables	
Default Filtering Method	User-Specified Limits
Keep Missing Values	Yes
Tuning Parameters	
Score	
Create Score Code	Yes
Update Measurement	No
Status	

43	Output						
44							
45	Variable	Role	Level	Train Count	Train Percent	Label	Filter Method
46							
47	Claim_Cause	INPUT	FIRE	1	0.033356	Claim_Cause	MINPCT
48							
49							
50							
51							
52	Number Of Observations						
53							
54	Data						
55	Role	Filtered	Excluded	DATA			
56							
57	TRAIN	2996	2	2998			
58	VALIDATE	2001	2	2003			
59							
60							
61							
62	Statistics for Original and FILTERED Data						
63	(maximum 500 observations printed)						
64							
65	Data Role=TRAIN Variable=Income						
66							
67	Statistics		Original	Filtered			
68							
69	Non Missing		2998.00	2996.00			
70	Missing		0.00	0.00			
71	Minimum		0.00	0.00			
72	Maximum		933288.00	99981.00			
73	Mean		38326.89	38040.96			
74	Standard Deviation		34559.26	30448.93			
75	Skewness		5.97	0.28			
76	Kurtosis		148.50	-1.11			
77							
78							
79	Data Role=VALIDATE Variable=Income						
80							
81	Statistics		Original	Filtered			
82							
83	Non Missing		2003.00	2001.00			
84	Missing		0.00	0.00			
85	Minimum		0.00	0.00			
86	Maximum		15967801.00	99960.00			
87	Mean		45776.10	37717.00			
88	Standard Deviation		357290.69	30708.64			
89	Skewness		44.25	0.29			
90	Kurtosis		1972.80	-1.09			
91							
92							
93	*-----*						

As a result, training and validation have two observations that were filtered. And the train data display that the maximum income value in the partition is \$99,981, but the original maximum value is \$933,288.

### Categorical Variables with Too Many Levels

One of the practical uses of many categories or class variables, like zip code variables, is that they can be combined at the city or state levels and can also combine their frequency. The Replacement node connects groups and establishes different levels for the group. The Filter node can be used to set minimum frequency and several levels.

Use the same direction for the filter node, set None for the Default Filtering Method under the Class Variables group, and click on the ellipsis by Class Variables. And click on each variable to

see the histogram top on the window and can update the Minimum and Maximum Frequency Cutoff.

The screenshot displays the SAS Enterprise Miner interface for a project named "Claim Fraud". The left sidebar shows a tree view with "Data Sources" containing "CLAIM\_DATA\_V2\_TRAIN", "Diagrams" containing "Create SAS Datasets", and "Model Packages". The main workspace shows a workflow diagram titled "Create SAS Datasets" with the following steps: "CLAIM\_DATA\_V2\_TRAIN" (input), "Data Partition", "Impute", and "Filter". There are also "StatExplore" and "Graph Explore" nodes connected to the workflow. The bottom status bar indicates "Run completed" and "Connected to SASApp - Logical Workspace Server (odaws01-usw2-2.oda.sas)".

**Properties Panel:**

Property	Value
Export Table	Filtered
Tables to Filter	All Data Sets
Distribution Data Sets	Yes
Class Variables	
Class Variables	...
Default Filtering Method	None
Keep Missing Values	Yes
Normalized Values	Yes
Minimum Frequency C1	
Minimum Cutoff for P	0.01
Maximum Number of U25	
Interval Variables	
Interval Variables	...
Default Filtering Method	User-Specified Limits
Keep Missing Values	Yes
Tuning Parameters	...
<b>Score</b>	
Create Score Code	Yes
Update Measurement	No
<b>Status</b>	
Default Filtering Method	
Default filtering method for class variables.	



## References

SAS Enterprise Miner by PAT RESEARCH, 2021. <https://www.predictiveanalyticstoday.com/sas-enterprise-miner/>

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