

# **Babu Banarasi Das University**

BBD City,Faizabad Road,Lucknow UttarPradesh



## **PROJECT REPORT – Insurance Fraud Prediction Using IBM SPSS Modeler**

**SUBMITTED TO:**  
**Mr. AYUSHMAN**  
**BHADAURIA**

**SUBMITTED BY:**  
**VISHWAJIT**  
**VISHWAS**

# Insurance Fraud Detection Using C&R Tree Algorithm

## Agenda / Definition

The project aims to detect fraudulent insurance claims using the C&R Tree (Classification and Regression Tree) method in IBM SPSS Modeler.

By analyzing claim data (such as vehicle details, claim amount, and customer info), the model identifies patterns and predicts whether a claim is fraudulent (Y) or non-fraudulent (N).

## Outcomes / Learning

- Import and explore a dataset in IBM SPSS Modeler
- Perform data cleaning (remove irrelevant columns, handle missing values)
  - Partition data into training and testing samples
- Build and evaluate a C&R Tree classification model
- Generate and interpret prediction results and graphs

**This project demonstrates the full data mining workflow – from preparation to model evaluation.**

## Required Tools

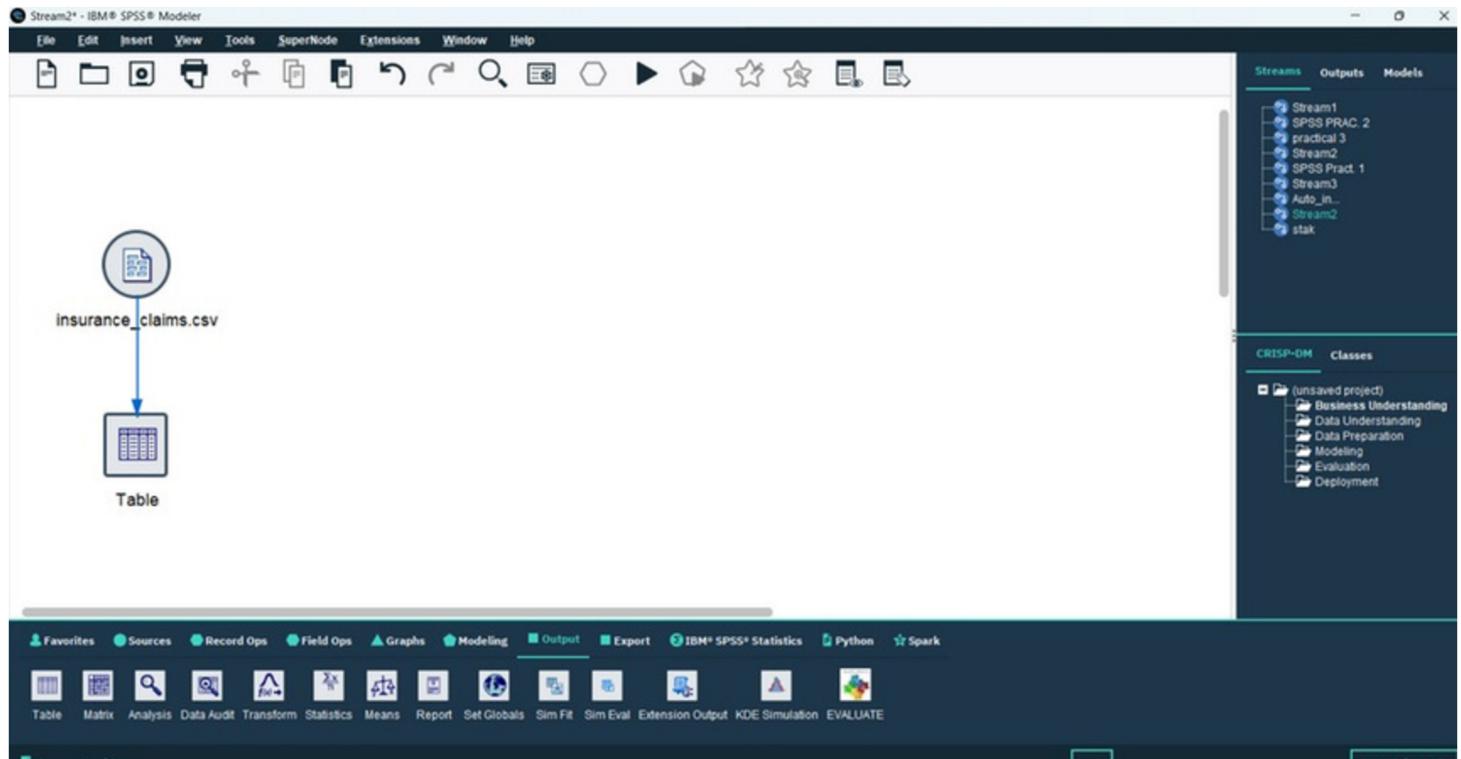
The tool used for this project is IBM SPSS Modeler.

## Working

- **The project involves:**
- Importing the insurance claim dataset
- Cleaning and preparing the data
- Setting variable roles and partitioning data
- Configuring and running the C&R Tree model
- Viewing prediction results in a table and histogram

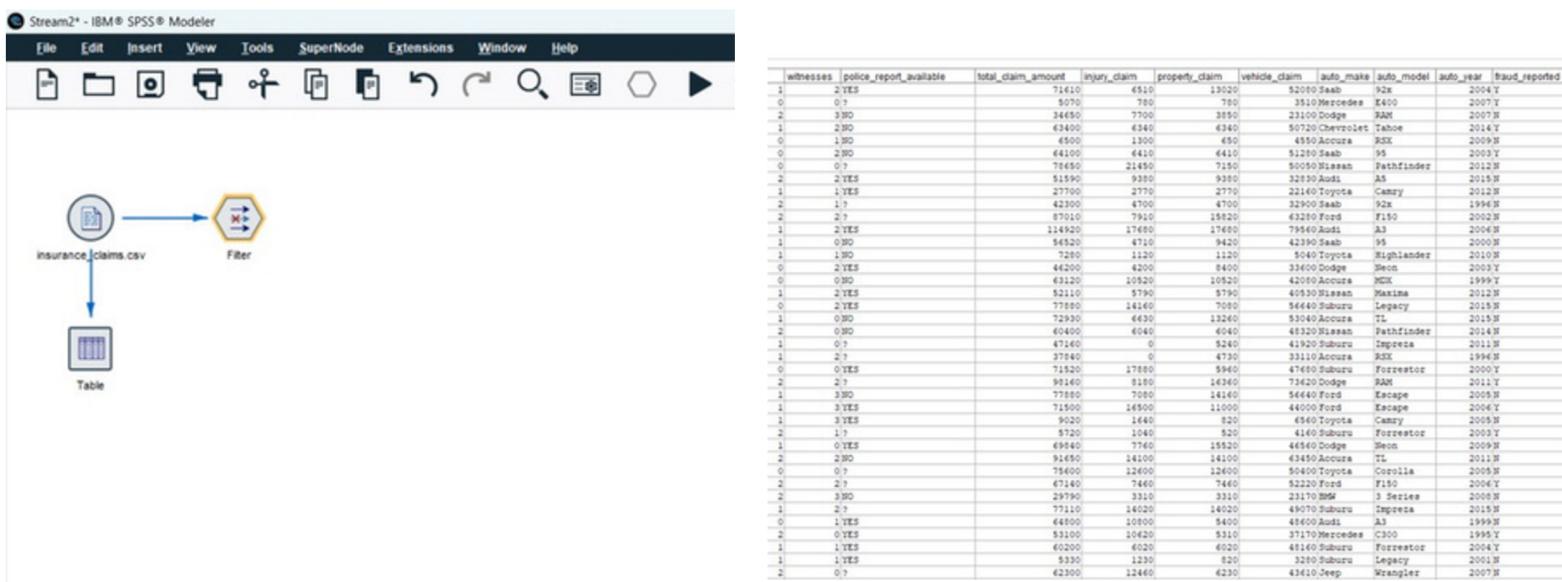
## Step 1: Import Data

Loaded the dataset (insurance\_claims.csv) into SPSS Modeler using the Var.File Node under Sources palette After reading metadata, all fields were correctly recognized.



## Step 2: Remove unnecessary Data

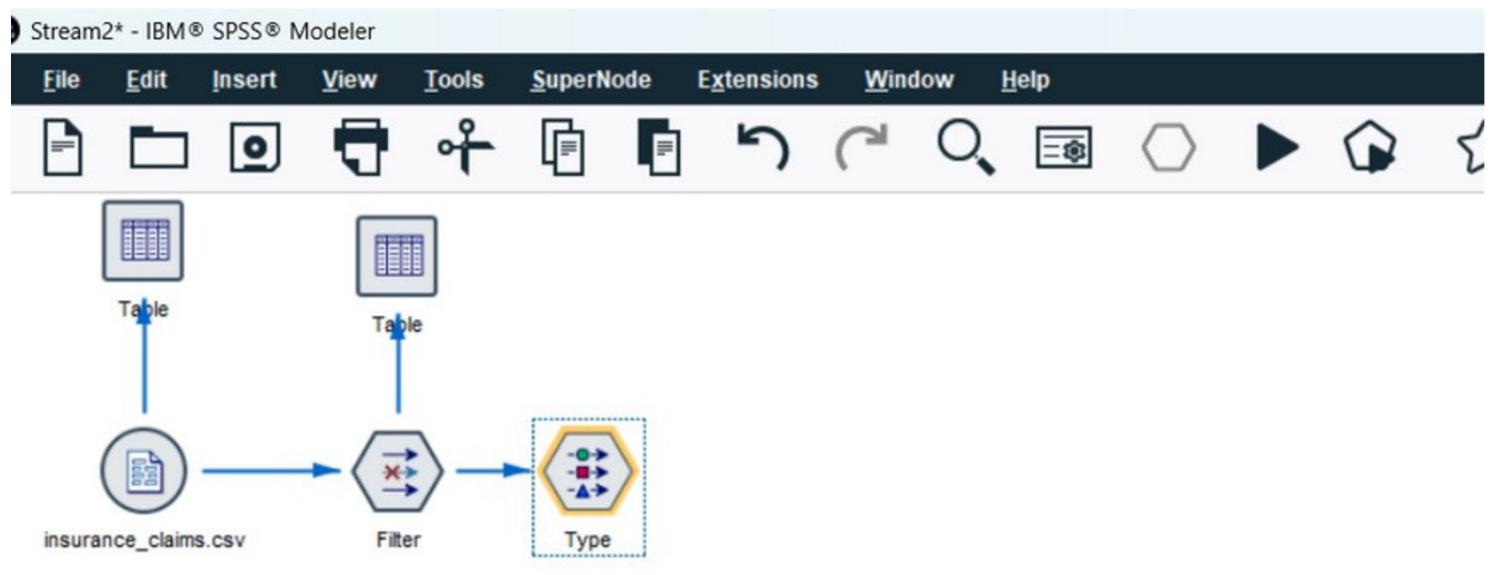
The Filter Node was used to exclude the irrelevant column \_c39 from the dataset. This column contained empty or meaningless values that could interfere with model accuracy. By filtering it out, we ensured that only useful fields (such as claim amount, vehicle claim, auto make, model, year, and fraud status) were retained for analysis.



## Step 4 : Type Node

Defined roles for each field:

- Input Fields: Claim amount, incident severity, age, etc.
- Target Field: fraud\_reported



The 'Type' dialog box displays field settings for various variables. The 'Types' tab is selected, showing the following table:

Field	Measurement	Values	Missing	Check	Role
total_claim...	Continuous	<Read>		None	Input
injury_claim	Continuous	<Read>		None	Input
property_clai...	Continuous	<Read>		None	Input
vehicle_claim	Continuous	<Read>		None	Input
auto_make	Categorical	<Read>		None	Input
auto_model	Categorical	<Read>		None	Input
auto_year	Continuous	<Read>		None	Input
fraud_report...	Categorical	<Read>		None	Target

At the bottom, there are two radio buttons: 'View current fields' (selected) and 'View unused field settings'. Below the table are 'OK', 'Cancel', 'Apply', and 'Reset' buttons.

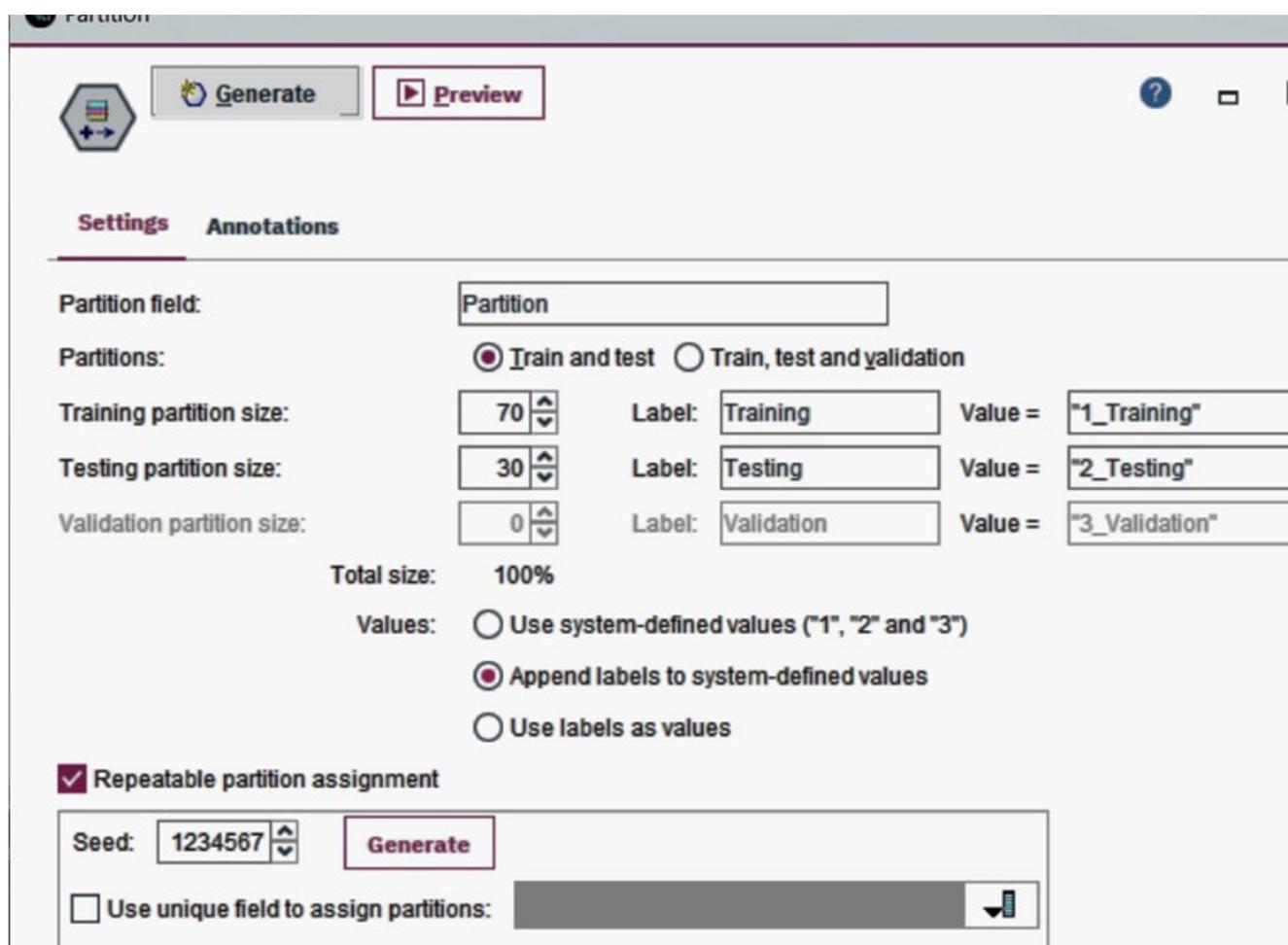
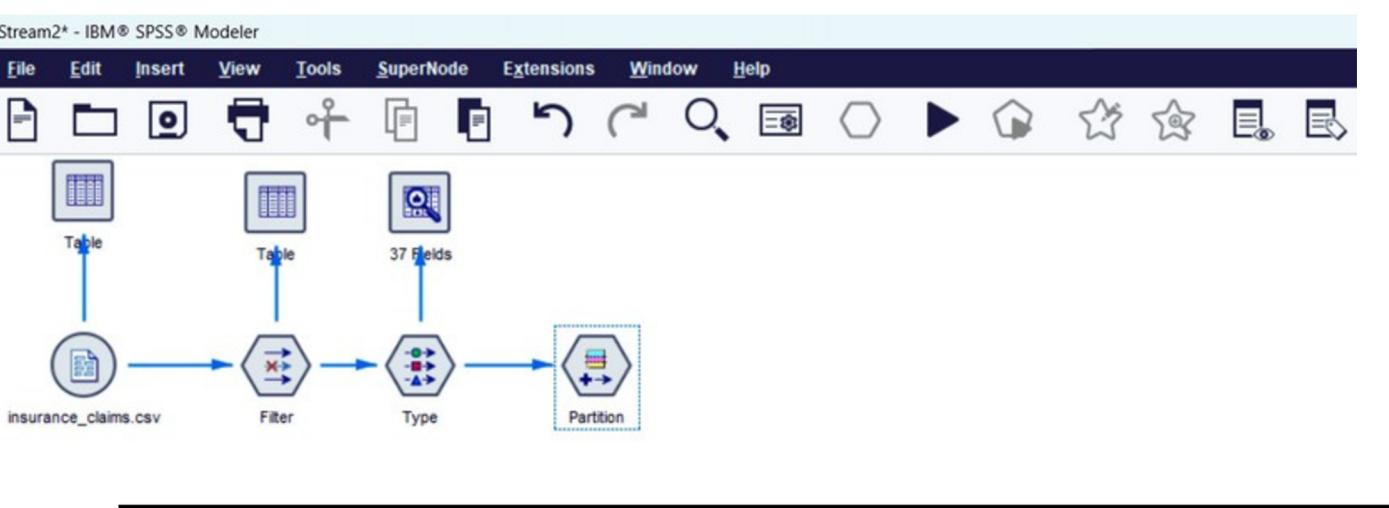
## Step 4 : Partition Data

**Added Partition Node to split data:**

70% for Training

30% for Testing

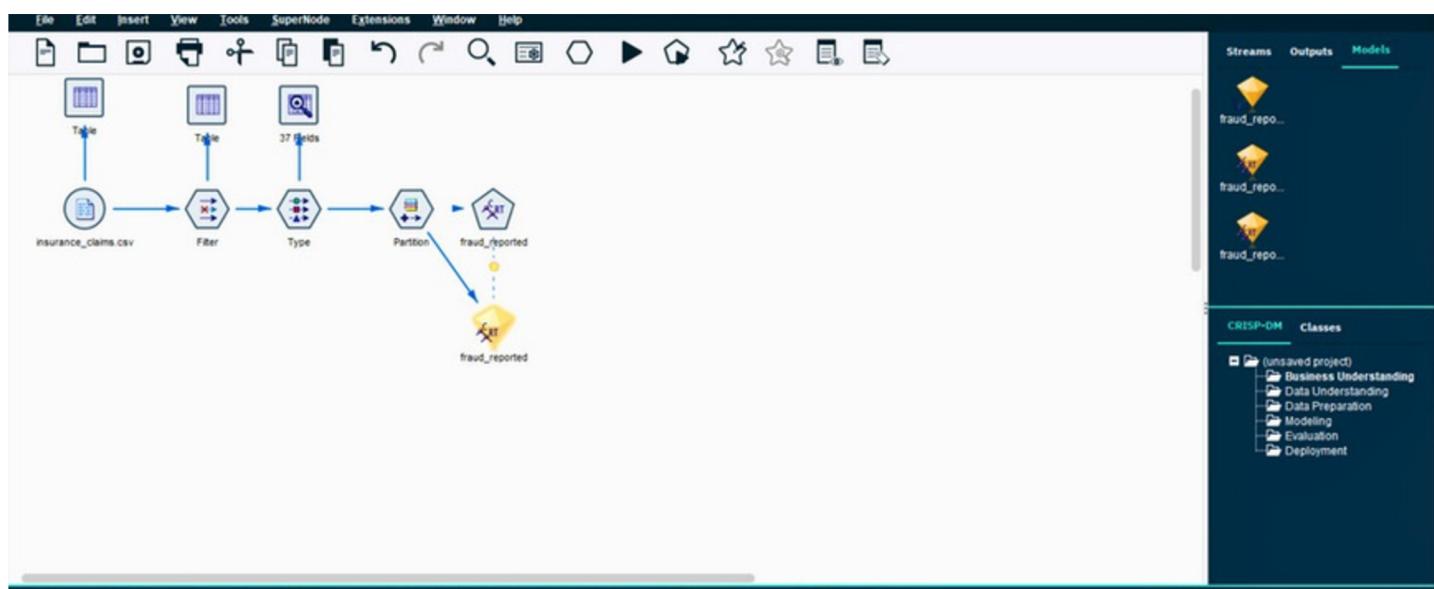
This allows model evaluation on unseen data



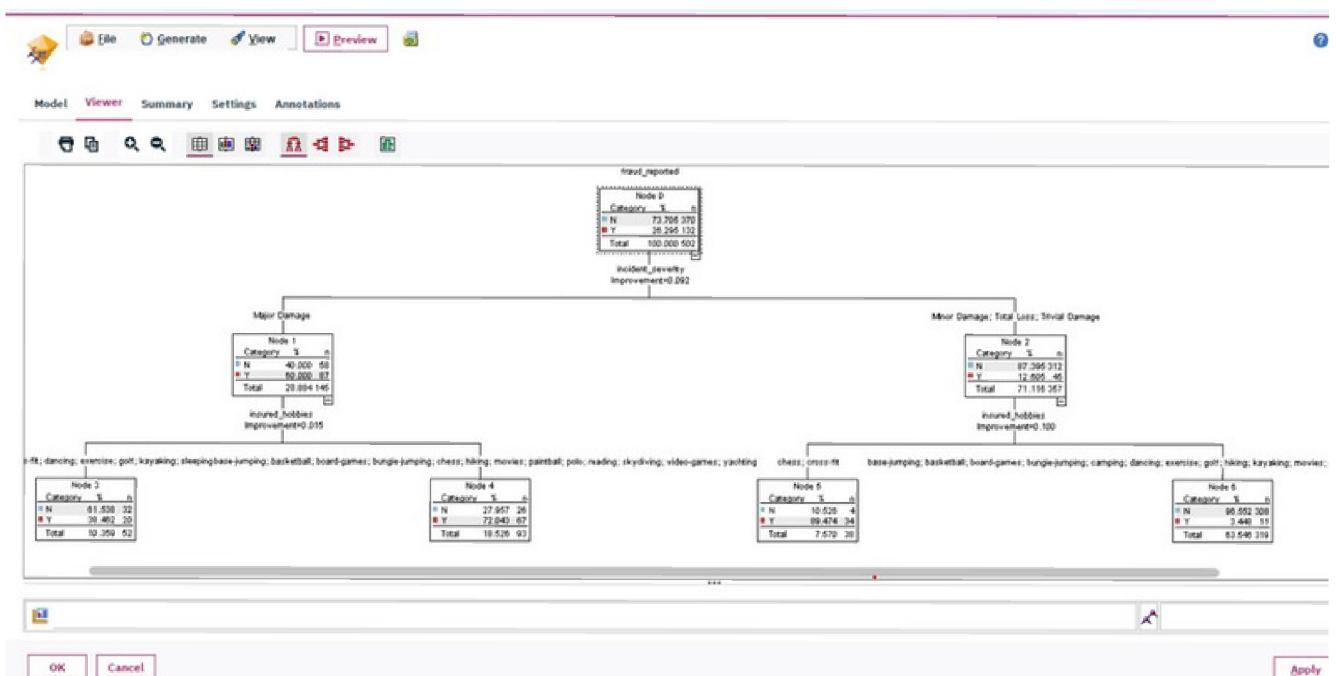
# Step 5 : Build Models

## Model 1 - C&R Tree

- From the Modeling palette, drag a C&R Tree Node.
- Connect it to the Partition Node.
- Open it → confirm:
- Target: fraud\_reported
- Inputs: Auto-selected.
- Click Run → view the decision tree output (splits, accuracy, etc.).

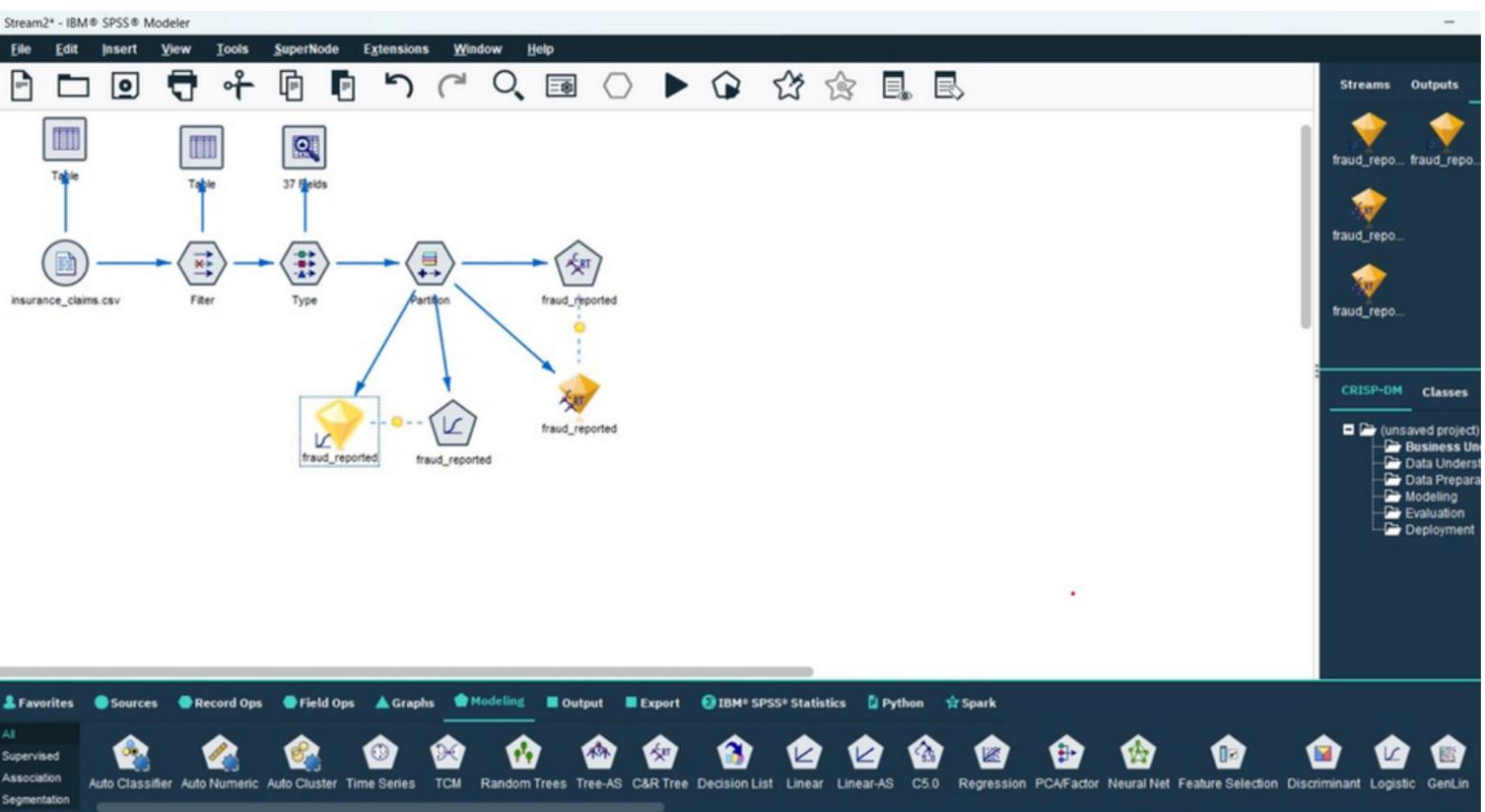


### OUTPUT:



## **Model 2 :Logistic Regression**

- 1.Drag a Logistic Regression Node from the Modeling palette.
2. Connect it to the same Partition Node.
3. Click Run → Check model summary and predictor significance.



## **Step 7: Generate Predictions**

After training, connected the Model Nugget to the dataset and added a Table Node.

**Output fields included:**

- fraud\_reported → Actual
- \$L-fraud\_reported → Predicted
- \$LP-fraud\_reported → Prediction probability

id	fraud_reported	Partition	SL-fraud_reported	SLP-fraud_reported
004	Y	1_Training	Y	0.800
007	Y	1_Training	N	0.595
007	N	1_Training	N	0.840
014	Y	2_Testing	Y	0.998
009	N	1_Training	N	1.000
003	Y	1_Training	Y	0.922
012	N	1_Training	N	0.998
015	N	1_Training	N	0.996
012	N	1_Training	N	1.000
006	N	1_Training	N	0.998
002	N	1_Training	N	0.981
006	N	1_Training	Y	0.910
000	N	2_Testing	N	0.745
010	N	1_Training	N	0.963
003	Y	2_Testing	N	0.903
099	Y	2_Testing	N	0.986
012	N	1_Training	N	0.542
015	N	1_Training	N	0.905
015	N	1_Training	N	0.961

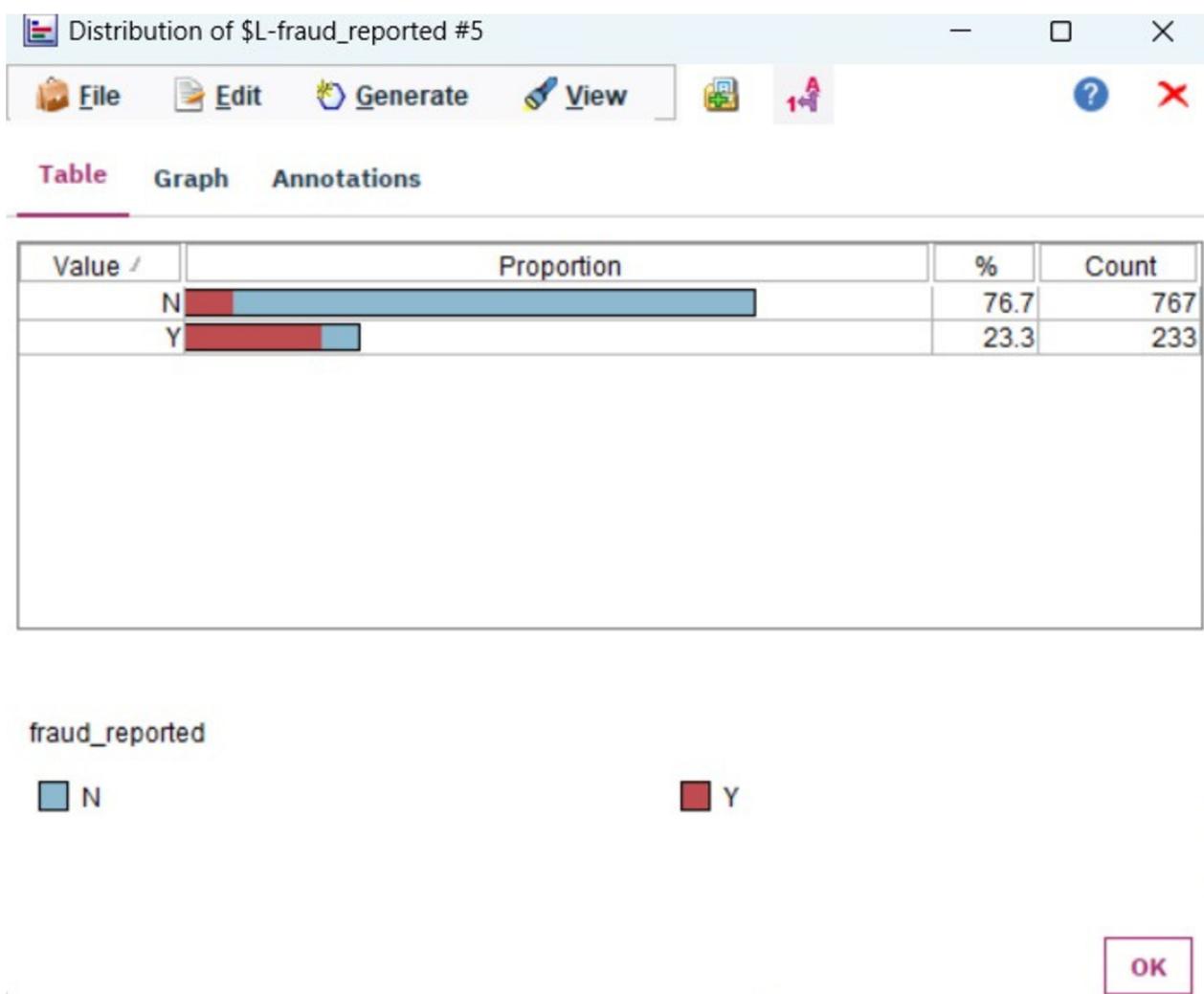
## TABLE

Table (42 fields, 1,000 records) #4														
Table		Annotations												
	bodily_injuries	witnesses	police_report_available	total_claim_amount	injury_claim	property_claim	vehicle_claim	auto_make	auto_model	auto_year	fraud_reported	Partition	\$L-fraud_reported	SLP-fraud_reported
1	1	2	YES	71610	6510	13020	52080	Saab	92x	2004	Y	1_Training	0.593	0.800
2	0	0	?	5070	780	780	3510	Mercedes	E400	2007	Y	1_Training	0.593	0.999
3	2	3	NO	34650	7700	3650	23100	Dodge	RAM	2007	N	1_Training	0.840	0.999
4	1	2	NO	63400	6340	6340	50720	Chevrolet	Tahoe	2014	Y	2_Testing	0.999	0.999
5	0	1	NO	6500	1300	650	4550	Accura	RSX	2009	N	1_Training	1.000	1.000
6	0	2	NO	44100	6410	6410	51280	Saab	95	2003	Y	1_Training	0.522	0.999
7	0	0	?	78650	21450	7150	50050	Nissan	Pathfinder	2012	N	1_Training	0.999	0.999
8	2	2	YES	51590	9380	9380	32830	Audi	A5	2015	N	1_Training	0.999	0.999
9	1	1	YES	27700	2770	2770	22160	Toyota	Camry	2012	N	1_Training	1.000	1.000
10	2	1	?	42300	4700	4700	32900	Saab	92x	1996	N	1_Training	0.599	0.999
11	2	2	?	87010	7910	15820	63280	Ford	F150	2002	N	1_Training	0.581	0.999
12	1	2	YES	114920	17680	17680	79560	Audi	A3	2006	N	1_Training	0.911	0.999
13	1	0	NO	56520	4710	9420	42390	Saab	95	2000	N	2_Testing	0.741	0.999
14	1	1	NO	7280	1120	1120	5040	Toyota	Highlander	2010	N	1_Training	0.561	0.999
15	0	2	YES	46200	4200	8400	33600	Dodge	Neon	2003	Y	2_Testing	0.901	0.999
16	0	0	NO	63120	10520	10520	42080	Accura	MDX	1999	Y	2_Testing	0.581	0.999
17	1	2	YES	52110	5790	5790	40530	Nissan	Maxima	2012	N	1_Training	0.542	0.999
18	0	2	YES	77880	14160	7080	56640	Subaru	Legacy	2015	N	1_Training	0.901	0.999
19	1	0	NO	72930	6630	13260	53040	Accura	TL	2015	N	1_Training	0.561	0.999
20	2	0	NO	60400	6040	6040	48320	Nissan	Pathfinder	2014	N	1_Training	0.970	0.999
21	1	0	?	47160	0	5240	41920	Suburu	Impreza	2011	N	1_Training	0.581	0.999
22	1	2	?	37840	0	4730	33110	Accura	RSX	1996	N	2_Testing	1.000	0.999
23	0	0	YES	71520	17880	5960	47680	Suburu	Forrestor	2000	Y	1_Training	0.781	0.999
24	2	2	?	98160	8180	16360	73620	Dodge	RAM	2011	Y	1_Training	0.796	0.999
25	1	3	NO	77880	7080	14160	56640	Ford	Escape	2005	N	1_Training	0.921	0.999
26	1	3	YES	71500	16500	11000	44000	Ford	Escape	2006	Y	1_Training	0.592	0.999
27	1	3	YES	9020	1640	820	6560	Toyota	Camry	2005	N	2_Testing	0.921	0.999
28	2	1	?	5720	1040	520	4160	Suburu	Forrestor	2003	Y	2_Testing	0.981	0.999
29	1	0	YES	69840	7760	15520	46560	Dodge	Neon	2009	N	1_Training	0.592	0.999
30	2	2	NO	91650	14100	14100	63450	Accura	TL	2011	N	1_Training	0.599	0.999
31	0	0	?	75460	12600	12600	50400	Toyota	Corolla	2005	N	1_Training	0.581	0.999
32	2	2	?	67140	7460	7460	52220	Ford	F150	2006	Y	1_Training	0.701	0.999
33	2	3	NO	29790	3310	3310	23170	BMW	3 Series	2008	N	1_Training	0.997	0.999
34	1	2	?	77110	14020	14020	49070	Suburu	Impreza	2015	N	1_Training	0.792	0.999
35	0	1	YES	64800	10800	5400	48600	Audi	A3	1999	N	2_Testing	0.999	0.999
36	2	0	YES	53100	10620	5310	37170	Mercedes	C300	1995	Y	2_Testing	0.599	0.999
37	1	1	YES	60200	6020	6020	48160	Suburu	Forrestor	2004	Y	1_Training	0.999	0.999
38	1	1	YES	5330	1230	820	3280	Suburu	Legacy	2001	N	1_Training	1.000	1.000
39	2	0	?	62300	12460	6230	43610	JEEP	Wrangler	2007	N	2_Testing	0.531	0.999
40	3	0	NO	61100	10000	10000	39999	BMW	5 Series	2011	Y	1_Training	0.999	0.999

## Step 8: Visualize Results

- After model training (C&R Tree and Logistic Regression), I connected the Distribution Graph Node to visualize the target field fraud\_reported.
- The field fraud\_reported was selected as the target variable for visualization.
- The output shows a clear bar chart distribution of “Yes” (fraud) and “No” (non-fraud) claims.
  - **No (N): 76.7% (767 records)**
  - **Yes (Y): 23.3% (233 records)**

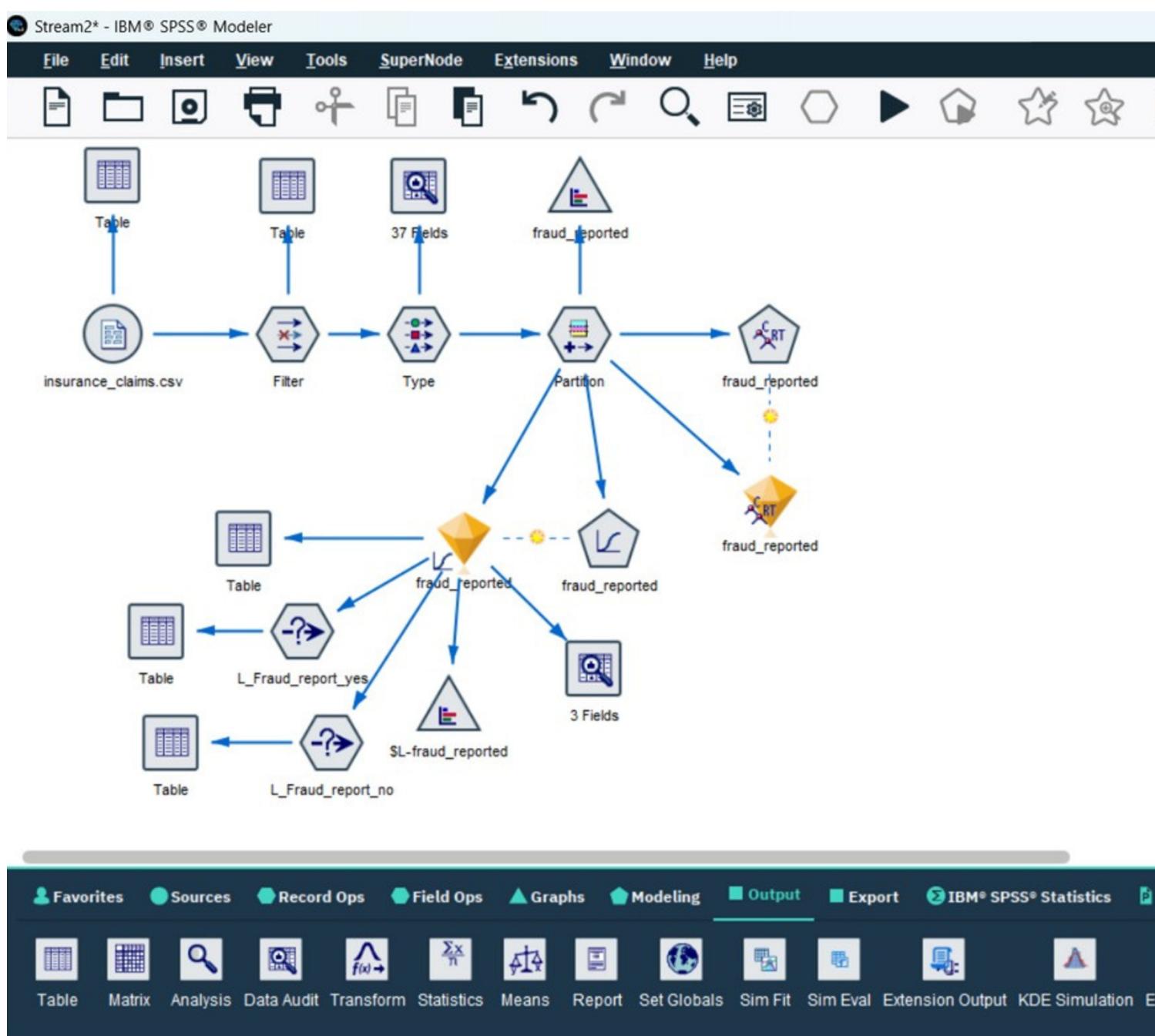
*This indicates that most insurance claims are non-fraudulent, but about one-fourth of cases (23.3%) are identified as fraud, which is still a significant number for analysis.*



# Insights :

- Major Damage and High Claim Amounts are the most common indicators of fraud.
- Customers with risky hobbies (like skydiving or motor racing) are more prone to fraudulent claims.
- Fraud detection is better when combining C&R Tree visualization and Logistic Regression probability scoring.

## FINAL VIEW



## **Model Comparison**

<b>Model Type</b>	<b>Description</b>	<b>Performance</b>	<b>Best Use</b>
C&R Tree	Decision tree showing split by features	High accuracy	High accuracy Visual, interpretative classification
Logistic Regression	Statistical model estimating fraud probability	High accuracy	Numerical fraud probability prediction

## **Conclusion :**

***The project successfully built and evaluated two models to detect insurance fraud using IBM SPSS Modeler.***

***The models help the insurance company:***

- *Identify potential fraudulent claims early.*
- *Save costs by reducing false claims.*
- *Improve the reliability of claim verification systems.*

*Overall, the C&R Tree and Logistic Regression models provided valuable insights into fraudulent behavior patterns.*