

SCHOOL OF COMPUTER APPLICATION

PREDICTIVE ANALYTICS

BCADS1501



Submitted By:

Name: Siddhant Kumar Patel

Section: BCA DS 36

University Roll No: 1230258423

Semester: 5th

Submitted To:

Mr. Ayushman Bhadauria

Babu Banarasi Das University, Lucknow

BBD City, Ayodhya Road, Lucknow Uttar Pradesh- 226028 Bharat

Case Study

“ACME, a company selling sports products, wants to promote its new product: the XL Original Orange Baseball Cap. To test customer interest, ACME sent a test mailing to 10,000 randomly selected customers and recorded their responses.”

Agenda/Definition:

- The agenda is to **build a predictive model** using the specified dataset to estimate an **unknown or future value (the target/output variable)**.
- To apply a suitable machine learning algorithm (e.g., Linear Regression, Decision Tree, Neural Network) to the ACME dataset within SPSS Modeler to establish a relationship between input variables (predictors) and a target variable, and then use that established relationship (the model) to predict the output for new, unseen data.

Outcomes/Learning:

- Proficiency in **data importation** and **pre-processing** (e.g., data auditing, missing value handling, variable type setting) within SPSS Modeler.
- Ability to **select and deploy appropriate modeling nodes** (algorithms) for the specific prediction task (e.g., classification or regression).
- Skills in **evaluating model performance** using appropriate graphs and statistics (e.g., gain charts, accuracy matrices, R^2 , RMSE).
- Understanding the end-to-end **CRISP-DM** (Cross-Industry Standard Process for Data Mining) process within a visual environment.

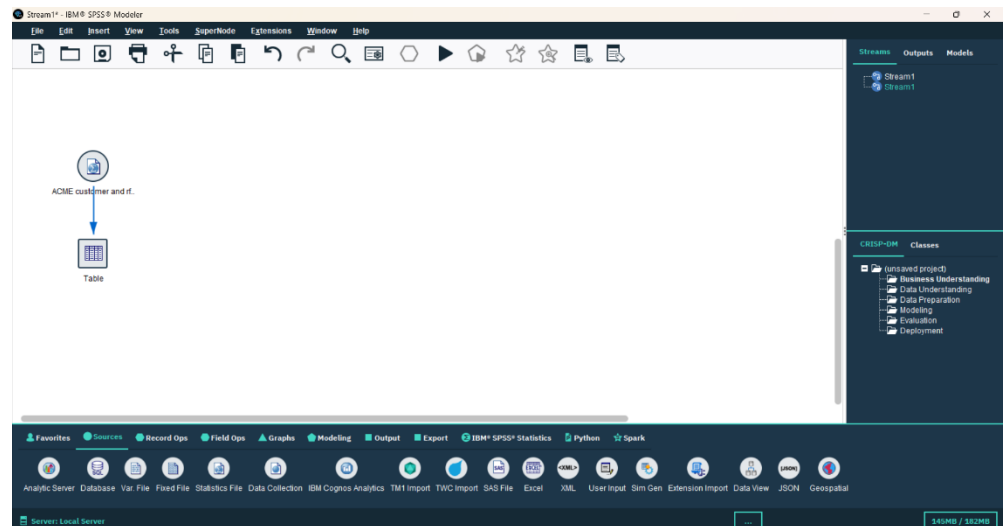
Required Tool:

The specific software required to perform the modeling and prediction is:

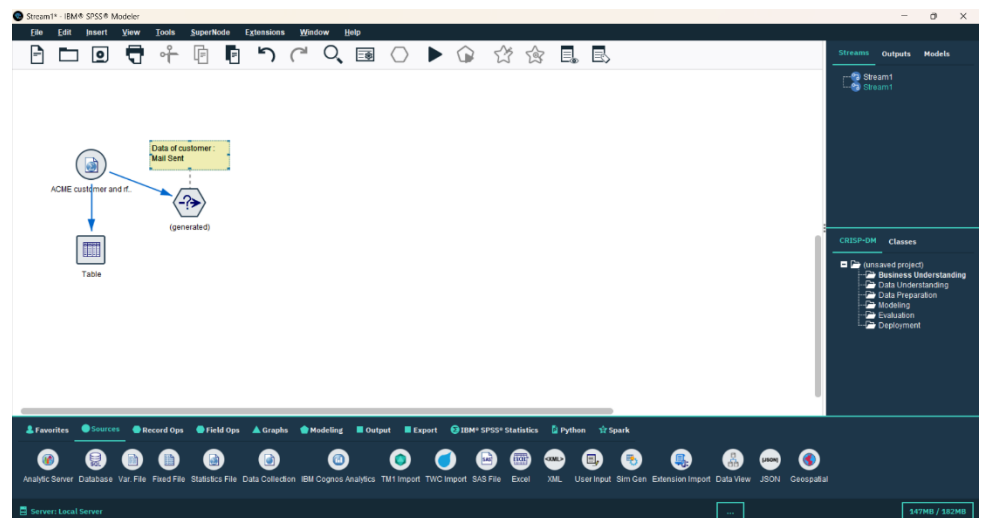
- **IBM SPSS Modeler:** This is the primary tool. It's a graphical workbench for building models without needing to write code.
- **ACME Dataset:** The required data source (in a format compatible with Modeler, like a flat file .csv,.txt, or a database connection).

Steps:

1. Import the dataset, use the table to get all the data (12 columns and 30000 rows)



2. Filter the data for the "YES" value and generate the node, rename as training data.



- Connect that node with table to get the desired output. (10000 rows)

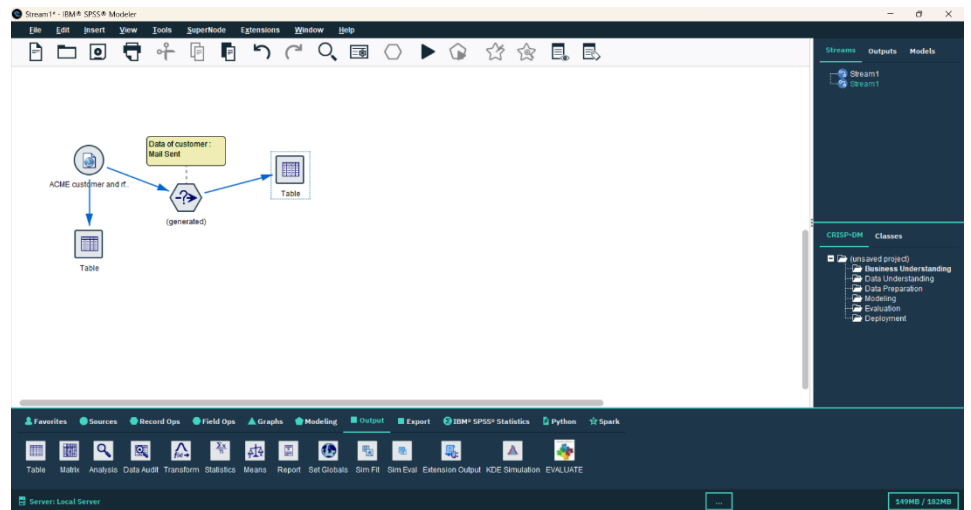
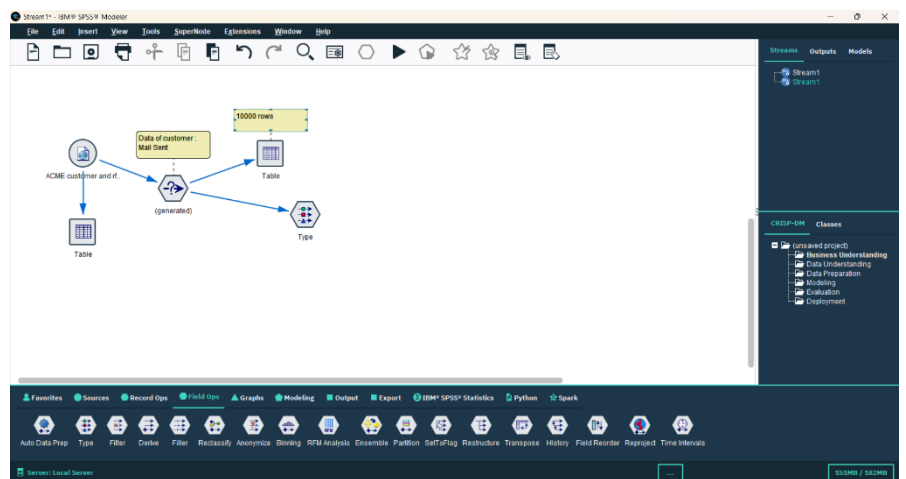


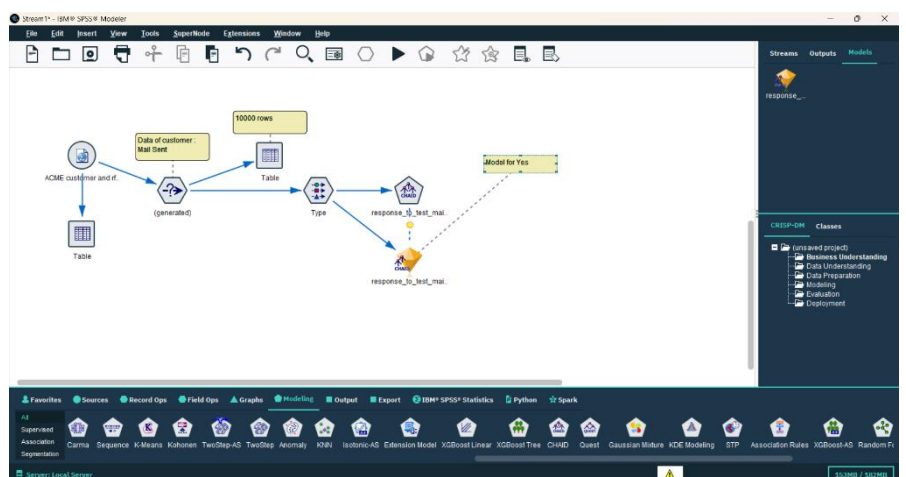
Table Annotations

	customer_id	gender	email_address	postal_code	monetary_value_01_01_2011	frequency_01_01_2011	recency_01_01_2011	has_received_test_mailing
1	723.000	male	name7502@tinet.fr	1818BO	2 medium	3 high	2 medium	yes
2	724.000	female	name25485@vnmmail.org	1132DG	1 low	3 high	1 low	yes
3	725.000	male	name15543@vnmmail.de	1803YT	3 high	1 low	1 low	yes
4	726.000	male	name28335@zigzag.be	1205WR	3 high	1 low	3 high	yes
5	727.000	female	name5354@tinet.jp	1711ON	1 low	3 high	1 low	yes
6	728.000	female	name20637@vnmmail.es	1055FG	2 medium	3 high	1 low	yes
7	729.000	female	name20636@vnmmail.es	1254MR	1 low	3 high	1 low	yes
8	730.000	female	name10414@tinet.inc	1723DG	2 medium	3 high	1 low	yes
9	731.000	male	name23372@vnmmail.inc	1713AQ	3 high	2 medium	1 low	yes
10	732.000	male	name20635@vnmmail.es	1264EC	3 high	2 medium	3 high	yes
11	733.000	female	name5356@tinet.jp	1648BT	3 high	2 medium	1 low	yes
12	734.000	female	name17582@vnmmail.de	1285XV	3 high	1 low	3 high	yes
13	735.000	female	name6388@tinet.fr	1282NB	1 low	2 medium	2 medium	yes
14	736.000	male	name10409@tinet.inc	1799IT	3 high	2 medium	1 low	yes
15	737.000	female	name13849@tinet.uk	1802UO	2 medium	3 high	1 low	yes
16	738.000	male	name25473@vnmmail.org	1971NK	1 low	3 high	1 low	yes
17	739.000	male	name13848@tinet.uk	1361RL	2 medium	3 high	1 low	yes
18	740.000	female	name23366@vnmmail.inc	1164VN	3 high	2 medium	1 low	yes
19	741.000	female	name3188@molbe.cat	1767YN	3 high	1 low	1 low	yes
20	742.000	male	name1606@lomejor.es	1681HP	1 low	3 high	1 low	yes

- Drag and drop the type node and connect with the training data node.



- Use the CHAID model and train the model with data.



- Use the table to get the result after the model is trained.

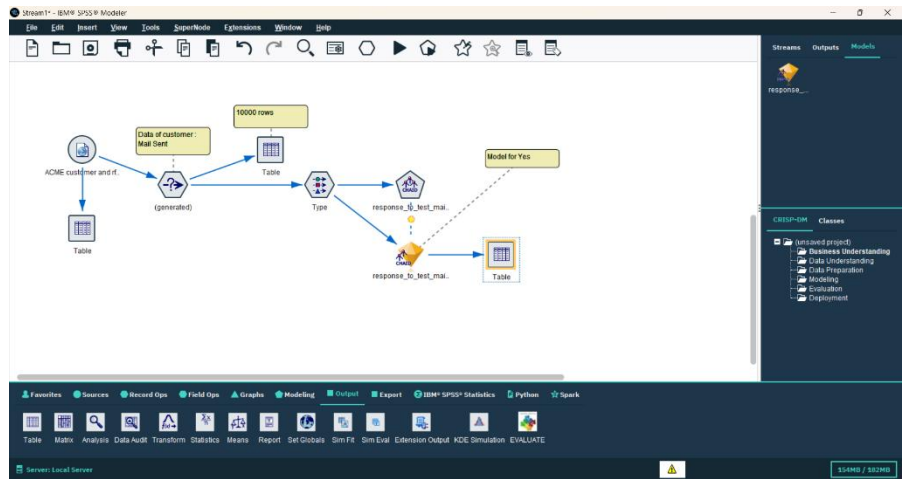
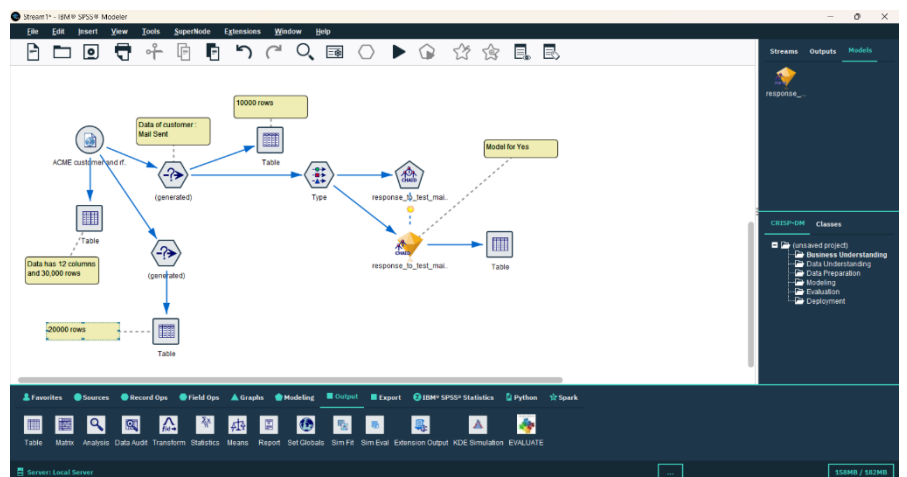


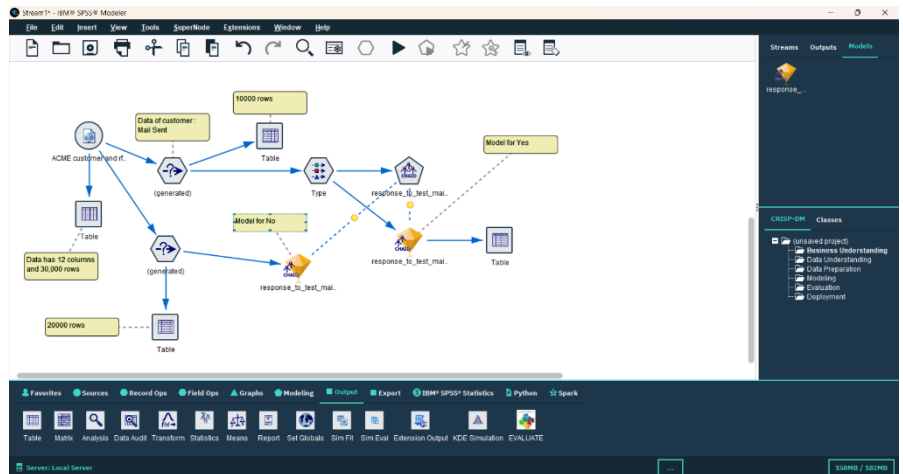
Table Annotations

	number_of_days_between_mailing_and_orderdate	ordered_within_month	\$R-response_to_test_mailing_02_01_2011	\$SRC-response_to_test_mailing_02_01_2011
1	\$null\$	nap	F	0.953
2	\$null\$	nap	F	0.993
3	\$null\$	nap	F	0.953
4	\$null\$	nap	F	0.953
5	\$null\$	nap	F	0.993
6	3.000	yes	F	0.953
7	\$null\$	nap	F	0.993
8	\$null\$	nap	F	0.953
9	\$null\$	nap	F	0.911
10	14.000	yes	T	0.625
11	\$null\$	nap	F	0.911
12	\$null\$	nap	F	0.953
13	\$null\$	nap	F	0.998
14	\$null\$	nap	F	0.911
15	\$null\$	nap	F	0.953
16	\$null\$	nap	F	0.993
17	\$null\$	nap	F	0.953
18	\$null\$	nap	F	0.911
19	\$null\$	nap	F	0.953
20	\$null\$	nap	F	0.993

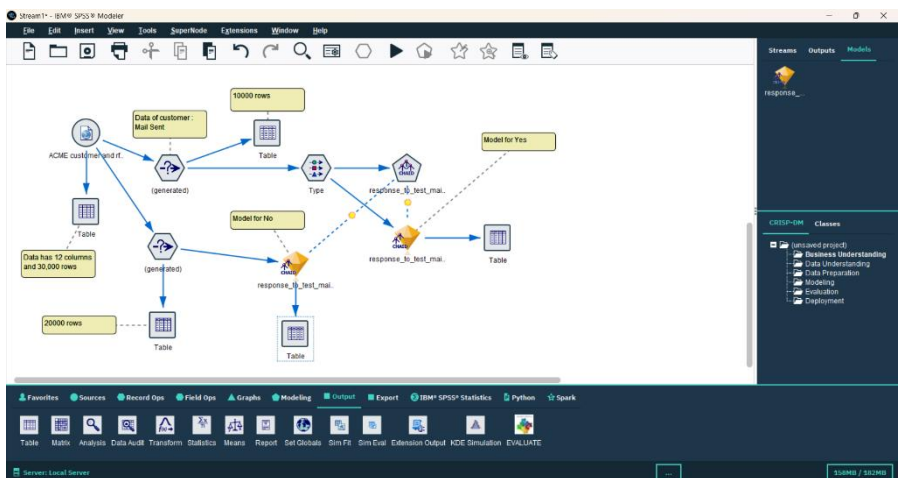
- Go back and filter the data for the "NO" value and generate the node, rename as test data for testing out model.



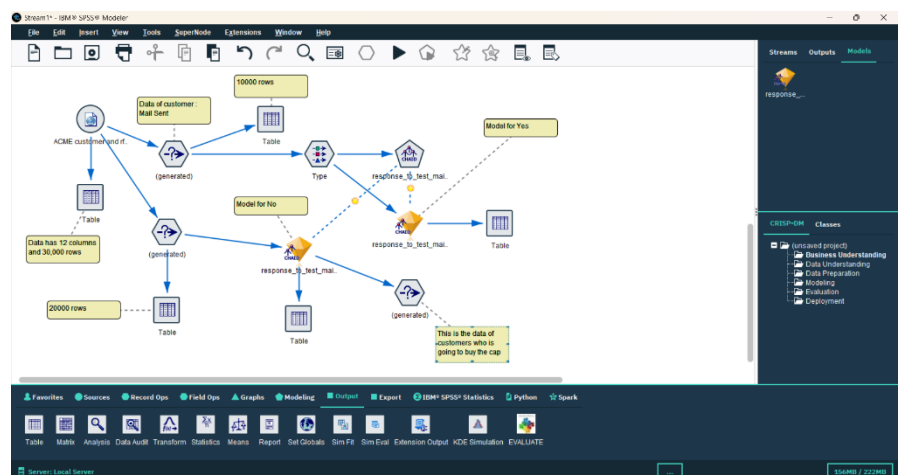
8. Make the copy of the model and connect with the test data.



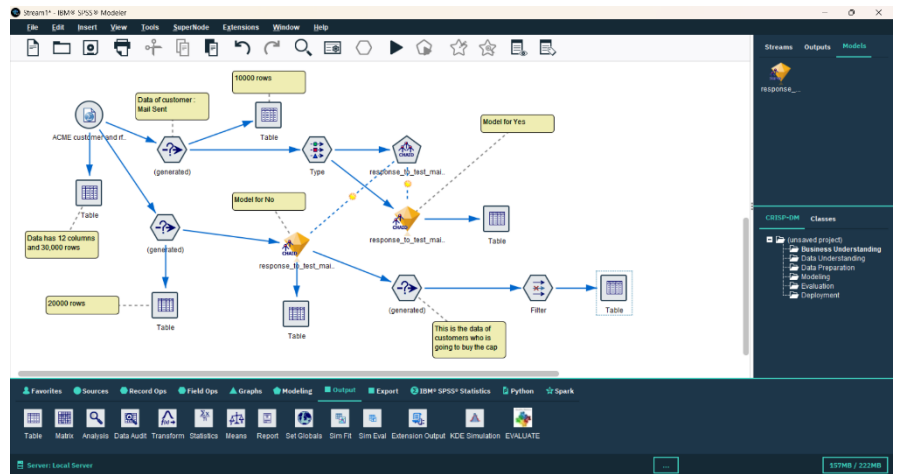
9. Use the table along with all nodes to get the result.



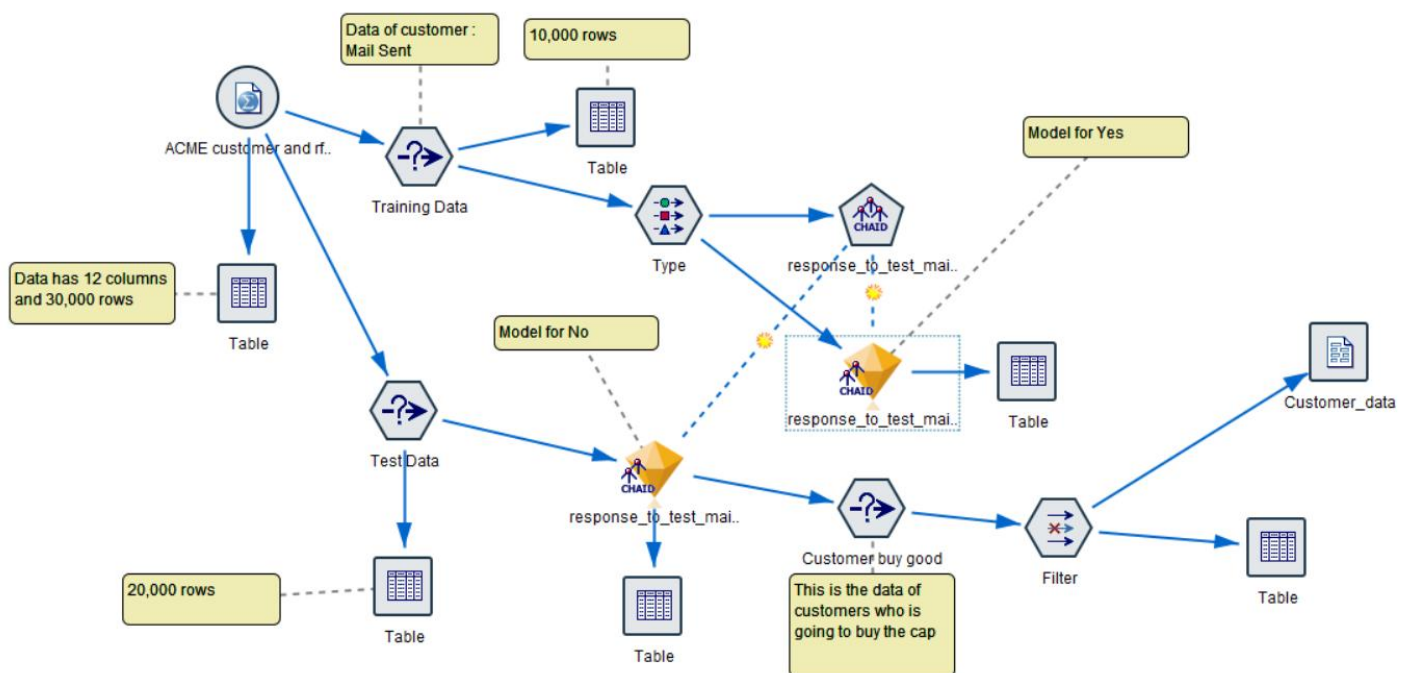
10. Again generate a new node for the "TRUE" data in the model. (Data of customers who are going to buy the cap)



11. Use the filter node to rename the specific columns and filter the data that we needed.



12. After all export the data to your local file with choosing the flat file node and specify the path where you want to save.



Q 1. How many records are in the training dataset?

Ans. There are 30000 Rows in the dataset.

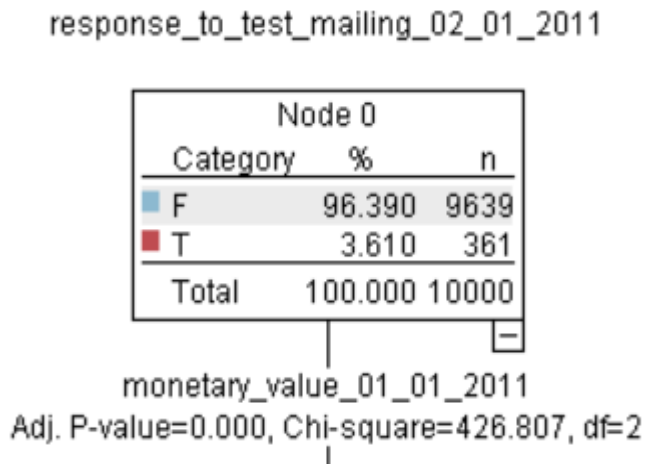
Q 2. How many fields are in the training dataset?

Ans. There are 12 Columns in the dataset.

Q 3. How many customers were included in the test mailing?

Ans. There are 20000 customers are included in the test mailing.

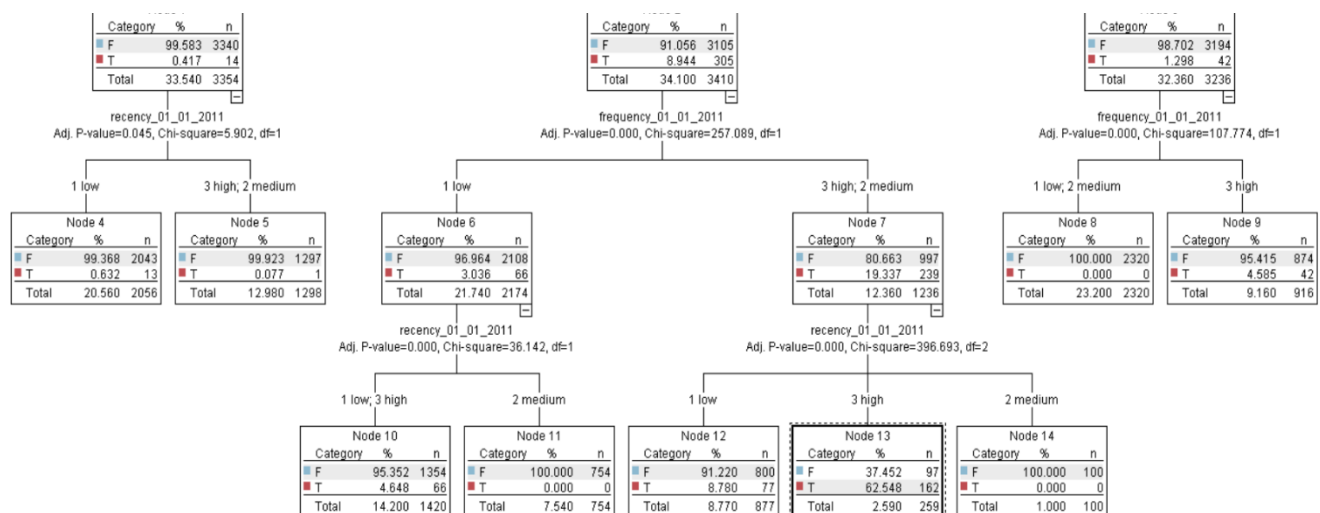
Q 4. Which field is used as the first split?



Ans. **monetary_value_01_01_2011** fields is used as the first split.

Q 5. Which group shows the highest response rate? What is the probability of responding for this group?

Ans. **recency_01_01_2011** shows the highest response rate and the probability of responding for this group is **62.548 %**



Q 6. Identify the two new fields added by the model.

Ans. \$R-response_to_test_mailing_02_01_2011 and

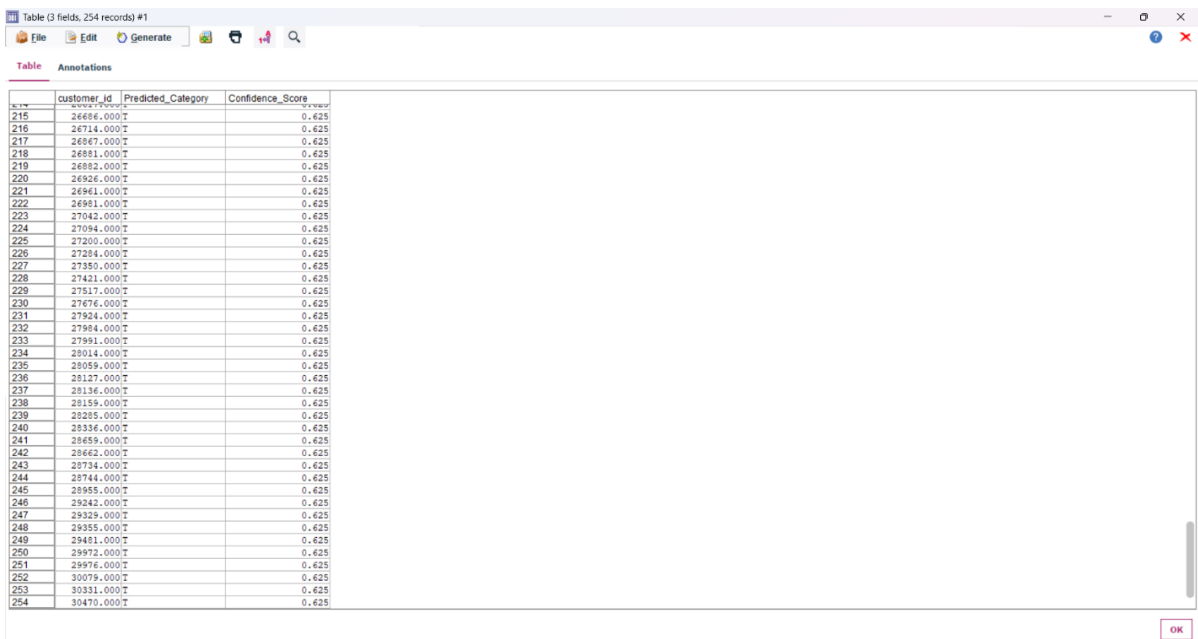
\$RC-response_to_test_mailing_02_01_2011 are two new fields added by the model.

Q 7. What do these fields represent?

Ans. The first fields represents the True/False value for the customers who are going to buy the cap or not. The second fields represents the probability of that condition is going to true.

Q 8. How many customers are predicted to respond positively (predicted = T)?

Ans. There are **254** customers are predicted to respond positively (predicted



	customer_id	Predicted_Category	Confidence_Score
214	26686.000	T	0.625
215	26714.000	T	0.625
216	26714.000	T	0.625
217	26867.000	T	0.625
218	26881.000	T	0.625
219	26882.000	T	0.625
220	26926.000	T	0.625
221	26941.000	T	0.625
222	26981.000	T	0.625
223	27042.000	T	0.625
224	27094.000	T	0.625
225	27200.000	T	0.625
226	27394.000	T	0.625
227	27350.000	T	0.625
228	27421.000	T	0.625
229	27517.000	T	0.625
230	27676.000	T	0.625
231	27924.000	T	0.625
232	27984.000	T	0.625
233	27991.000	T	0.625
234	28014.000	T	0.625
235	28059.000	T	0.625
236	28127.000	T	0.625
237	28136.000	T	0.625
238	28159.000	T	0.625
239	28205.000	T	0.625
240	28336.000	T	0.625
241	28659.000	T	0.625
242	28662.000	T	0.625
243	28734.000	T	0.625
244	28744.000	T	0.625
245	28955.000	T	0.625
246	29242.000	T	0.625
247	29329.000	T	0.625
248	29355.000	T	0.625
249	29481.000	T	0.625
250	29972.000	T	0.625
251	29976.000	T	0.625
252	30079.000	T	0.625
253	30331.000	T	0.625
254	30470.000	T	0.625

= T).

After exporting -

customers_to_contact

File Edit View H1 Edit B I G A

customer_id,Predicted_Category,Confidence_Score
5835.000,"I",0.625
5851.000,"I",0.625
6031.000,"I",0.625
6204.000,"I",0.625
6245.000,"I",0.625
6340.000,"I",0.625
6424.000,"I",0.625
6451.000,"I",0.625
6460.000,"I",0.625
6501.000,"I",0.625
6516.000,"I",0.625
6525.000,"I",0.625
6657.000,"I",0.625
6659.000,"I",0.625
6766.000,"I",0.625
6877.000,"I",0.625
6904.000,"I",0.625
7162.000,"I",0.625
7449.000,"I",0.625
7510.000,"I",0.625
7516.000,"I",0.625
7524.000,"I",0.625
7550.000,"I",0.625
7555.000,"I",0.625
7576.000,"I",0.625
7640.000,"I",0.625
7694.000,"I",0.625
7810.000,"I",0.625
7977.000,"I",0.625
8002.000,"I",0.625
8017.000,"I",0.625
8149.000,"I",0.625
8160.000,"I",0.625
8322.000,"I",0.625
8356.000,"I",0.625
8376.000,"I",0.625
8429.000,"I",0.625
8500.000,"I",0.625
8515.000,"I",0.625

Ln 16, Col 19 5,070 characters Plain text 100% Windows (CRLF) UTF-8