

SPSS PRACTICAL

Name: Ganesh Agrahari

Class: BCA DS 33

Roll No.: 1230258176

Submitted To: Mr. Robin Tyagi

INDEX

SPSS Practical 6

This practical includes:

- Objective and Introduction
- Theoretical Background
- Methodology and Procedure
- Data Analysis
- Results and Interpretation
- Screenshots and Output
- Conclusion
- References

Practical: 6

Definition:

You work for a telecommunications firm and you need to cleanse and enrich a dataset in order to build models later.

Outcomes/Learning:

- Learned how to use the **Derive Node** to create new calculated fields such as billing details, segmentation, and discounts.
- Understood the use of the **Reclassify Node** to standardize and group categorical values for better analysis.
- Gained practical experience in cleansing and enriching raw datasets to prepare them for predictive modelling.

Required Tool:

IBM SPSS MODELER

Working:

In this practical, the goal is to cleanse and enrich the dataset (*telco x data.txt*) for further modelling. Fields are standardized, billing details are calculated, customers are segmented, and discounts are derived. Finally, categorical values are reclassified into consistent groups.

The main nodes used are:

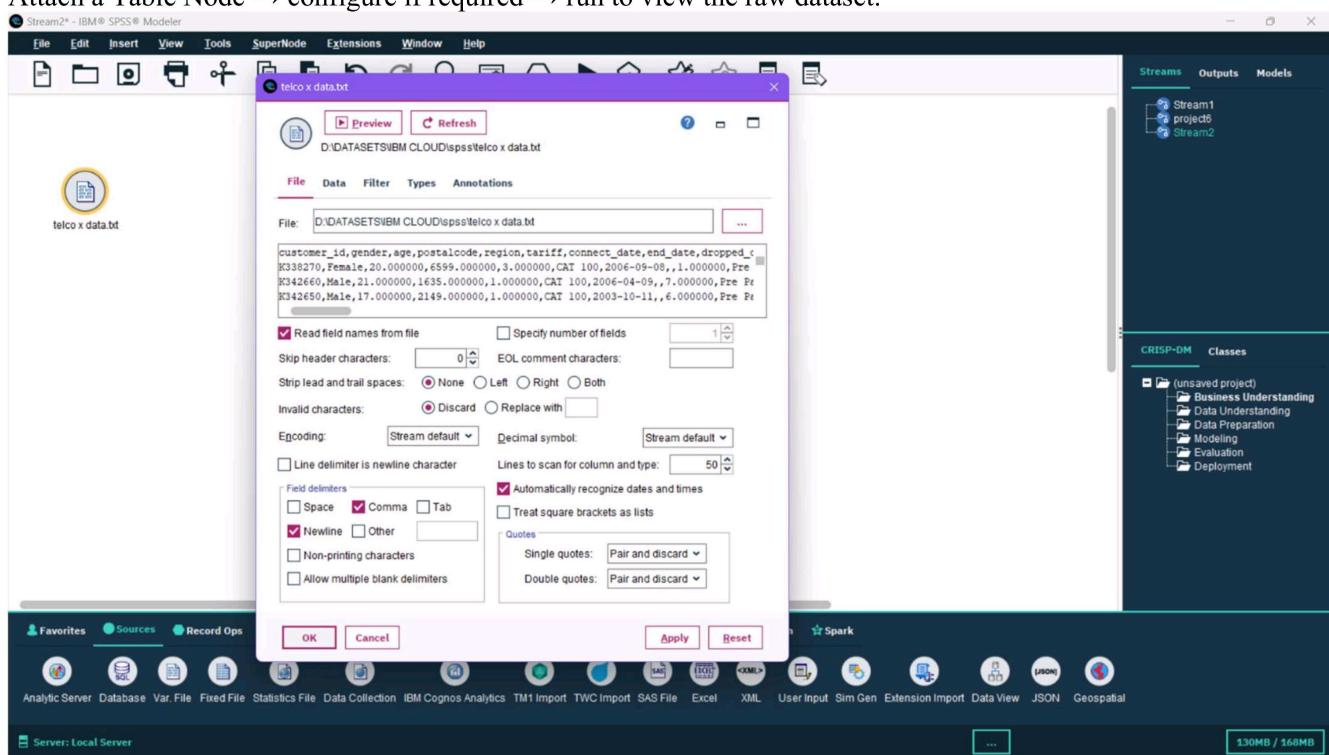
- Type Node** – to define field roles and measurement levels.
- Derive Nodes** – to create new fields such as billing, segmentation, and discounts.
- Reclassify Node** – to restructure categorical values.
- Table Nodes** – to validate results at each step.

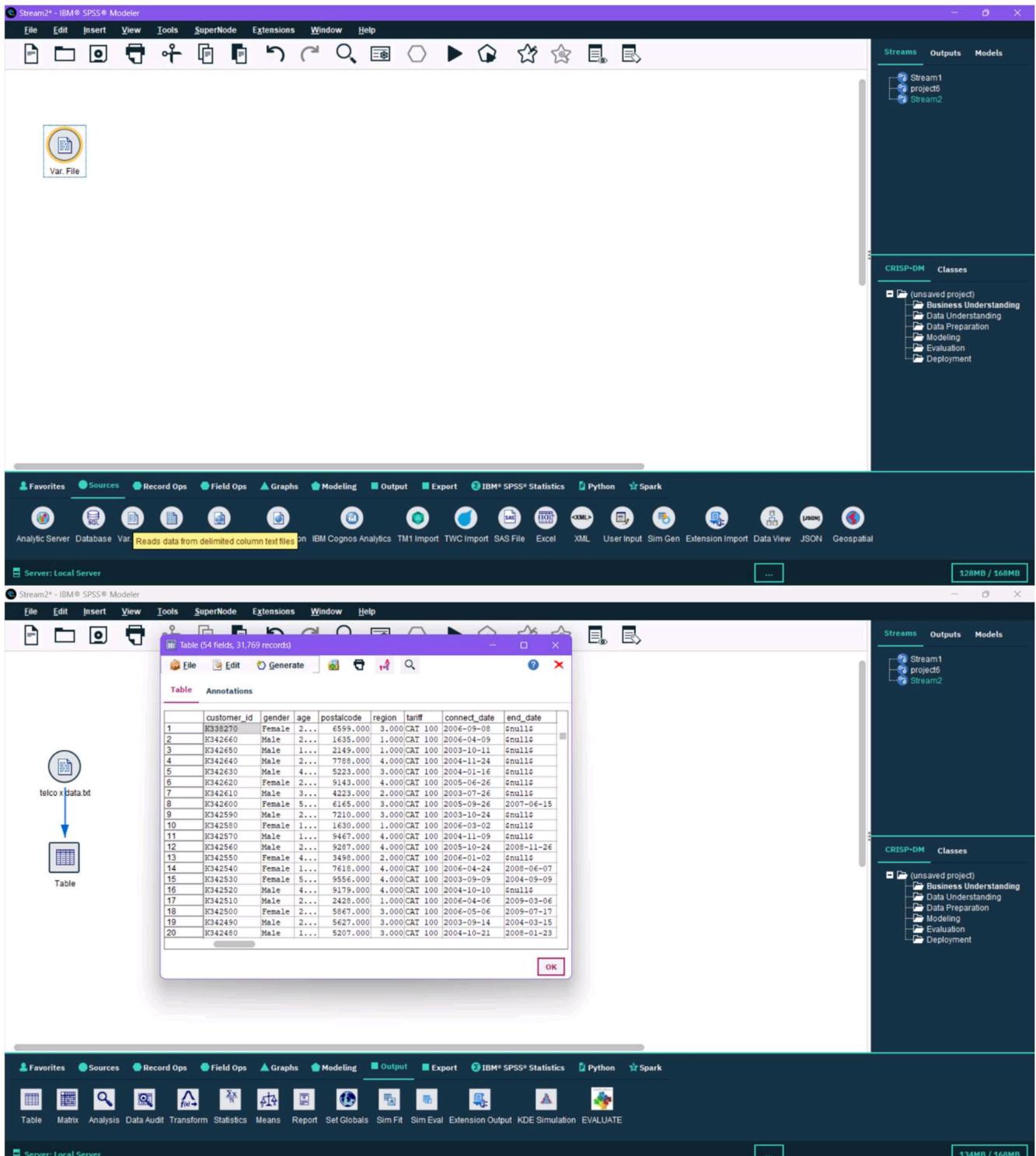
Step 1: Importing the Dataset

Open IBM SPSS Modeler → create a new stream.

Add a Var. File Node and import *telco x data.txt*.

Attach a Table Node → configure if required → run to view the raw dataset.

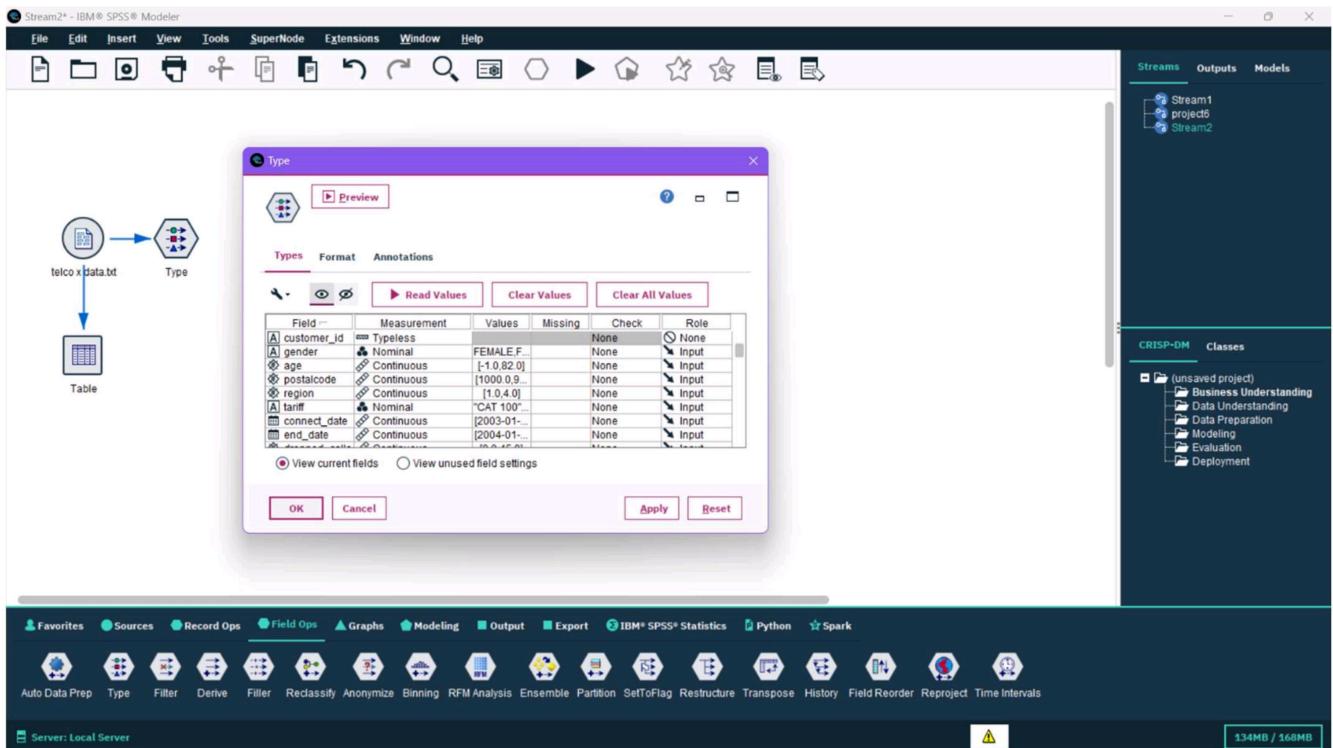




Step 2: Defining Field Properties

Connect a Type Node to the Var. File Node.

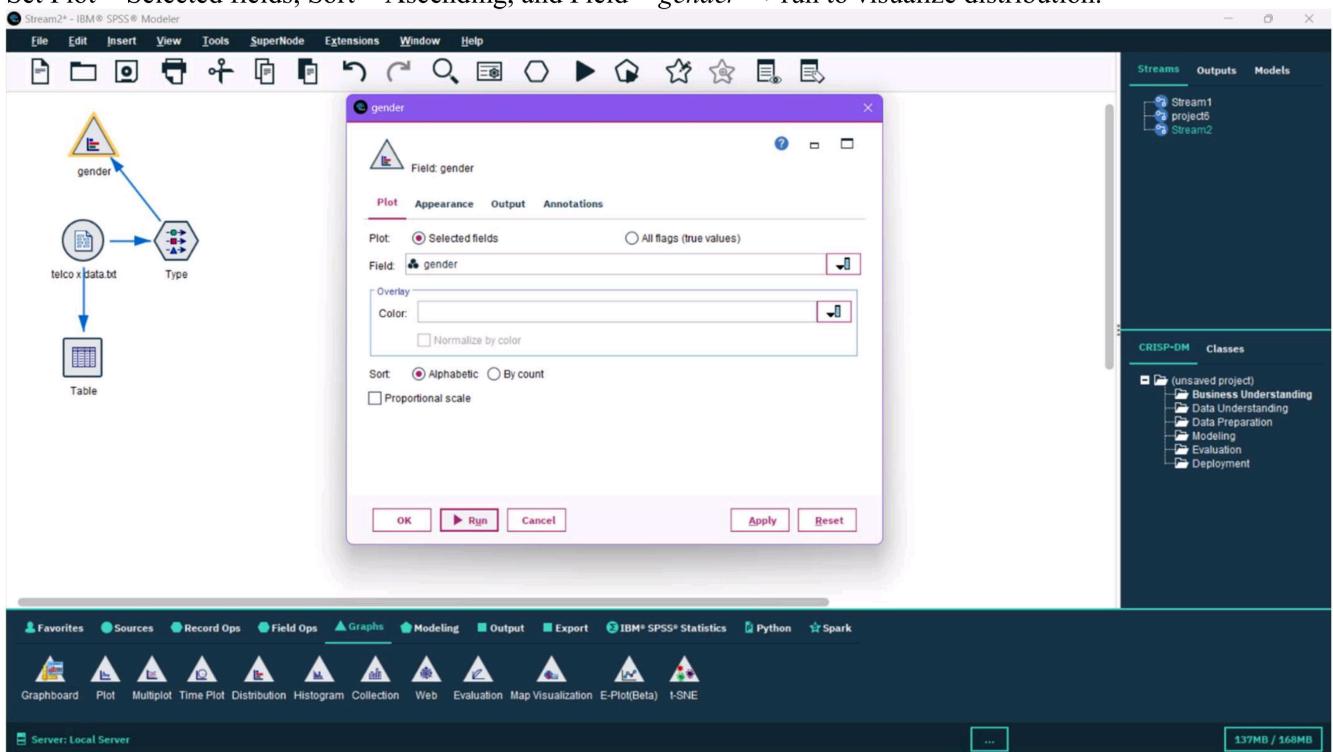
Open it → click Read Values → observe updated measurement types (Nominal, Ordinal, etc.).

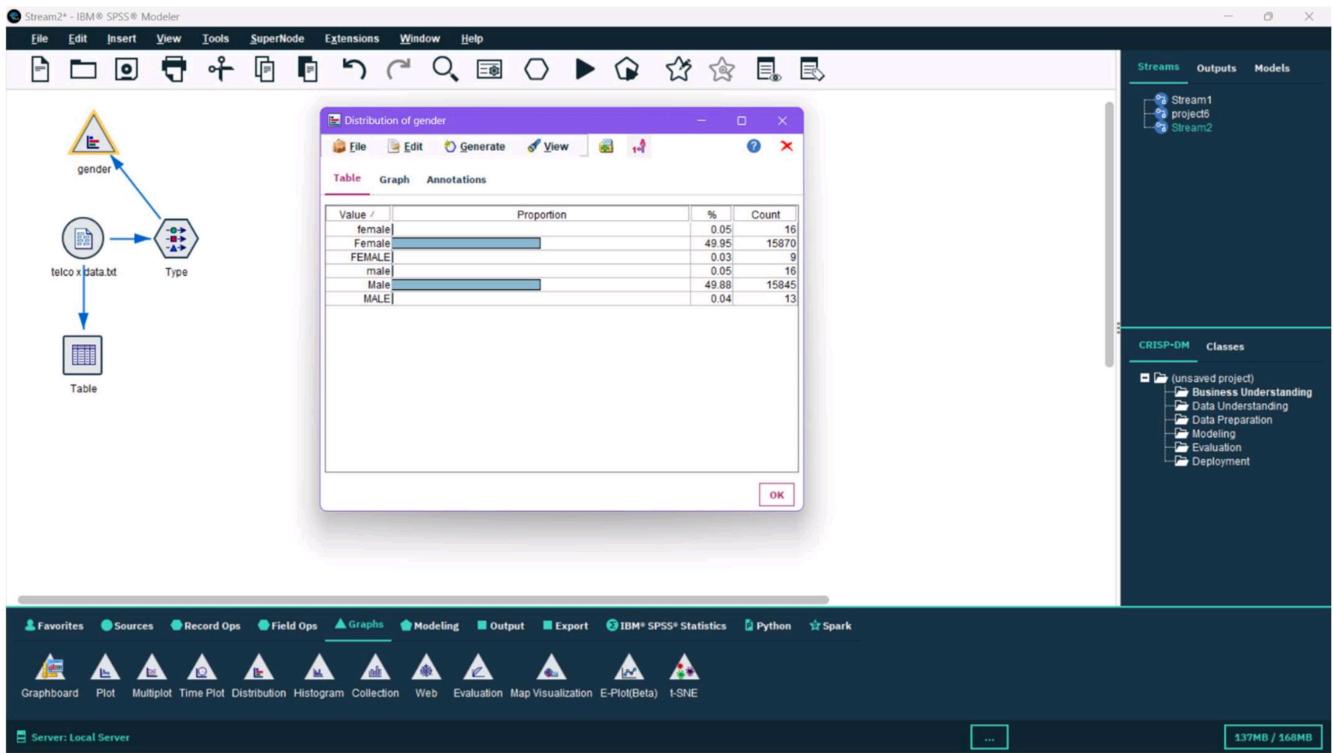


Step 3: Exploring with Distribution Graph

Attach a Distribution Node to the Type Node.

Set Plot = Selected fields, Sort = Ascending, and Field = *gender* → run to visualize distribution.



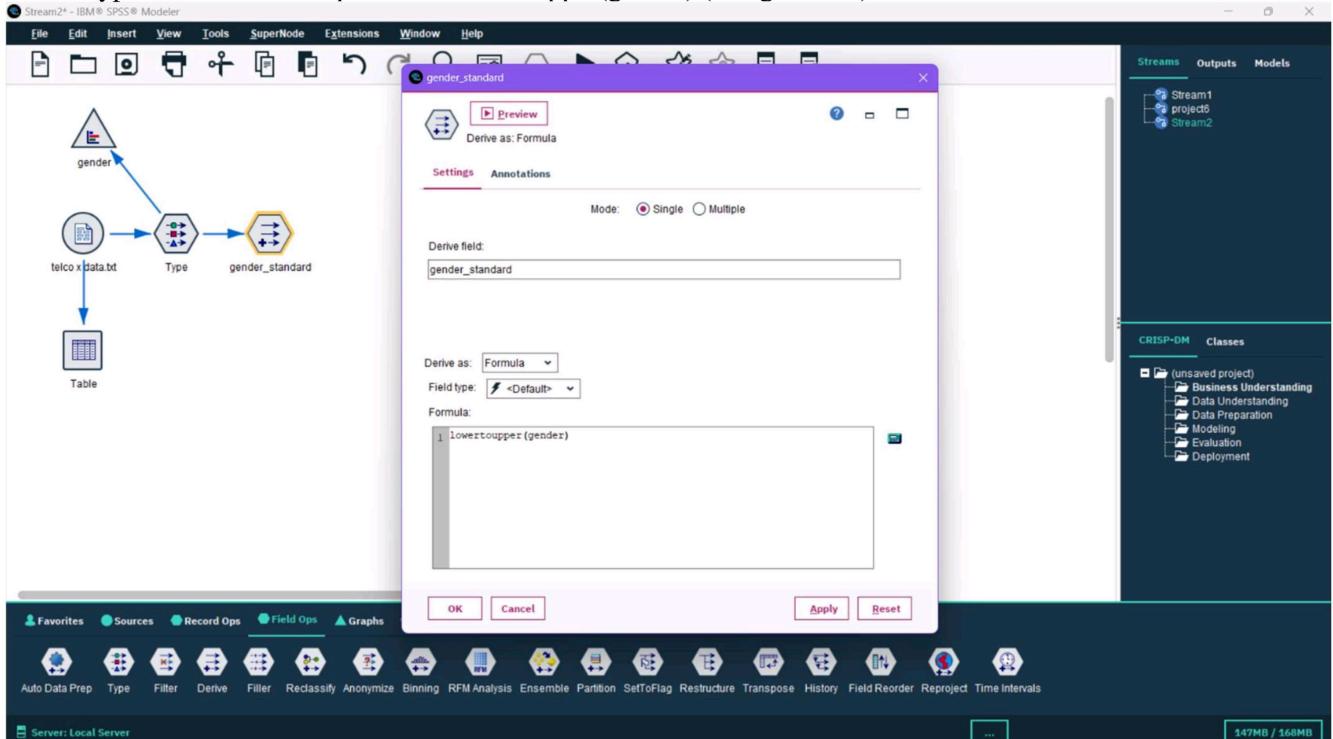


Step 4: Standardizing Gender Field

Connect a Derive Node to the Type Node.

Rename field = *gender_standard*.

Derive type = Formula → expression: lowertoupper(gender) (using CLEM).



Step 5: Visualizing Standardized Gender

Attach a Distribution Node to the *gender_standard* Derive Node.

Set field = *gender* → run to confirm corrected values.

The image consists of three vertically stacked screenshots of the Stream2+ - IBM SPSS Modeler software interface. Each screenshot shows a data flow diagram on the left and a corresponding analysis window on the right.

Screenshot 1 (Top): The data flow diagram shows a 'Table' node connected to a 'Type' node, which then connects to a 'gender_standard' node. A 'gender' node is also connected to the 'gender_standard' node. An open window titled 'gender_standard' is displayed, showing settings for a 'Plot' (Selected fields, Field: gender_standard) and a 'Table' (OK, Run, Cancel, Apply, Reset). The status bar at the bottom indicates '149MB / 168MB'.

Screenshot 2 (Middle): The data flow diagram is identical to the first. An open window titled 'Distribution of gender_standard #3' displays a table of proportions for 'FEMALE' and 'MALE'. The table is as follows:

Value	Proportion	%	Count
FEMALE	50.03	50.03	15895
MALE	49.97	49.97	15874

The status bar at the bottom indicates '149MB / 168MB'.

Screenshot 3 (Bottom): The data flow diagram is identical to the first. An open window titled 'Distribution of gender_standard #3' displays a table of proportions for 'FEMALE' and 'MALE'. The table is as follows:

Value	Proportion	%	Count
FEMALE	50.03	50.03	15895
MALE	49.97	49.97	15874

The status bar at the bottom indicates '155MB / 168MB'.

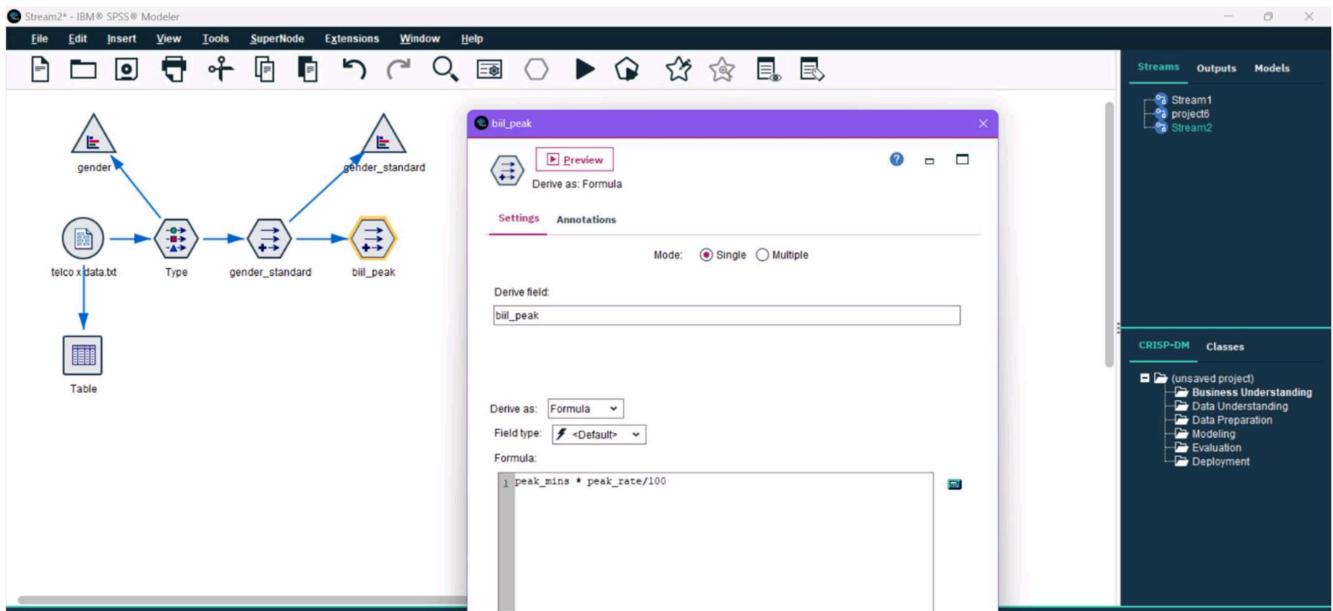
Step 6: Calculating Peak Bill

Add a new Derive Node after *gender_standard*.

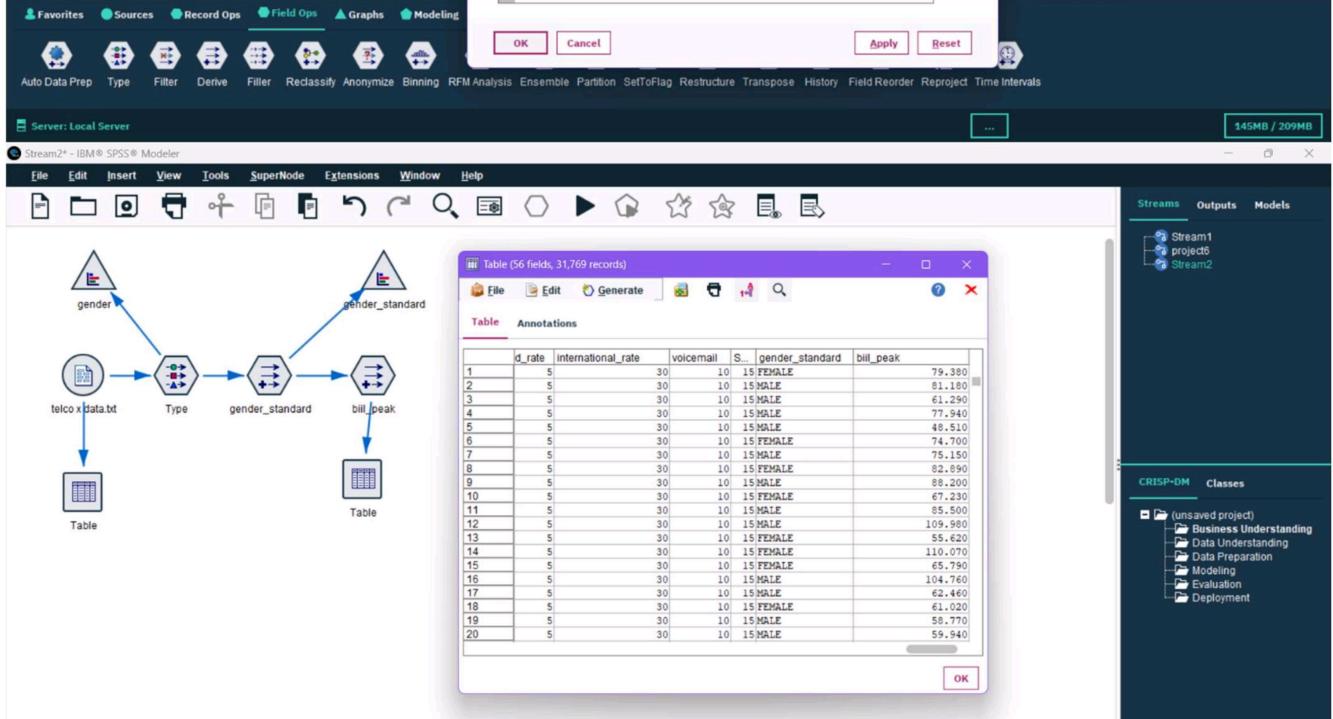
Rename field = *bill_peak*.

Formula = *peak_mins* * *peak_rate* / 100.

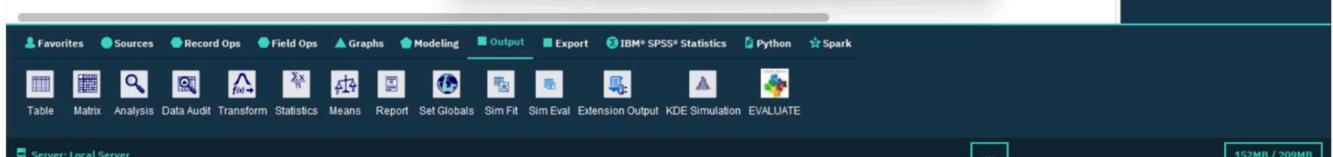
Attach a Table Node → run to view results.



The screenshot shows a data flow stream in Stream2+. The stream starts with a 'Sources' node 'telco x data.txt', followed by a 'Type' node, a 'gender_standard' node, and a 'bill_peak' node. The 'gender' field is mapped to the 'gender' node, and the 'gender_standard' field is mapped to the 'bill_peak' node. A 'Derive' dialog box is open, titled 'bill_peak', showing the formula `1_peak_mins * peak_rate/100`. The 'Derive as' dropdown is set to 'Formula' and the 'Field type' dropdown is set to '<Default>'. The 'OK' button is highlighted.



The screenshot shows the results of the 'Derive' operation in a 'Table' node. The table has 31,769 records and 56 fields. The last column, 'bill_peak', contains values ranging from 59.940 to 79.380. The 'OK' button is highlighted.



The screenshot shows the 'Output' tab in Stream2+. A 'Table' node is selected, showing the same data as the previous screenshot. The 'OK' button is highlighted.

Step 7: Calculating Off-Peak Bill

Add another Derive Node after *bill_peak*.

Rename field = *offbill_peak*.

Formula = $\text{offpeak_mins} * \text{offpeak_rate} / 100$.

Attach a Table Node → run to confirm results.

Stream2 - IBM SPSS Modeler

File Edit Insert View Tools SuperNode Extensions Window Help

Streams Outputs Models

CRISP-DM Classes

bill_offpeak

Derive as: Formula

Settings Annotations

Mode: Single Multiple

Derive field: bill_offpeak

Derive as: Formula

Field type: <Default>

Formula: l_offpeak_mins * offpeak_rate/100

OK Cancel Apply Reset

Server: Local Server

170MB / 209MB

Stream2 - IBM SPSS Modeler

File Edit Insert View Tools SuperNode Extensions Window Help

Streams Outputs Models

CRISP-DM Classes

Table (57 fields, 31,769 records)

Annotations

	onal_rate	voicemail	S.	gender_standard	bill_peak	bill_offpeak
1	30	10	15	FEMALE	79.380	15.540
2	30	10	15	MALE	81.180	4.110
3	30	10	15	MALE	61.290	4.695
4	30	10	15	MALE	77.940	8.760
5	30	10	15	MALE	45.510	6.105
6	30	10	15	FEMALE	74.700	18.765
7	30	10	15	MALE	75.150	14.535
8	30	10	15	FEMALE	82.890	11.325
9	30	10	15	MALE	88.200	1.470
10	30	10	15	FEMALE	67.230	6.120
11	30	10	15	MALE	85.500	11.550
12	30	10	15	MALE	109.380	7.515
13	30	10	15	FEMALE	55.620	4.335
14	30	10	15	FEMALE	110.070	2.070
15	30	10	15	FEMALE	65.750	6.855
16	30	10	15	MALE	104.760	0.000
17	30	10	15	MALE	62.460	3.135
18	30	10	15	FEMALE	61.020	17.295
19	30	10	15	MALE	55.770	9.630
20	30	10	15	MALE	55.940	9.615

OK

Server: Local Server

171MB / 209MB

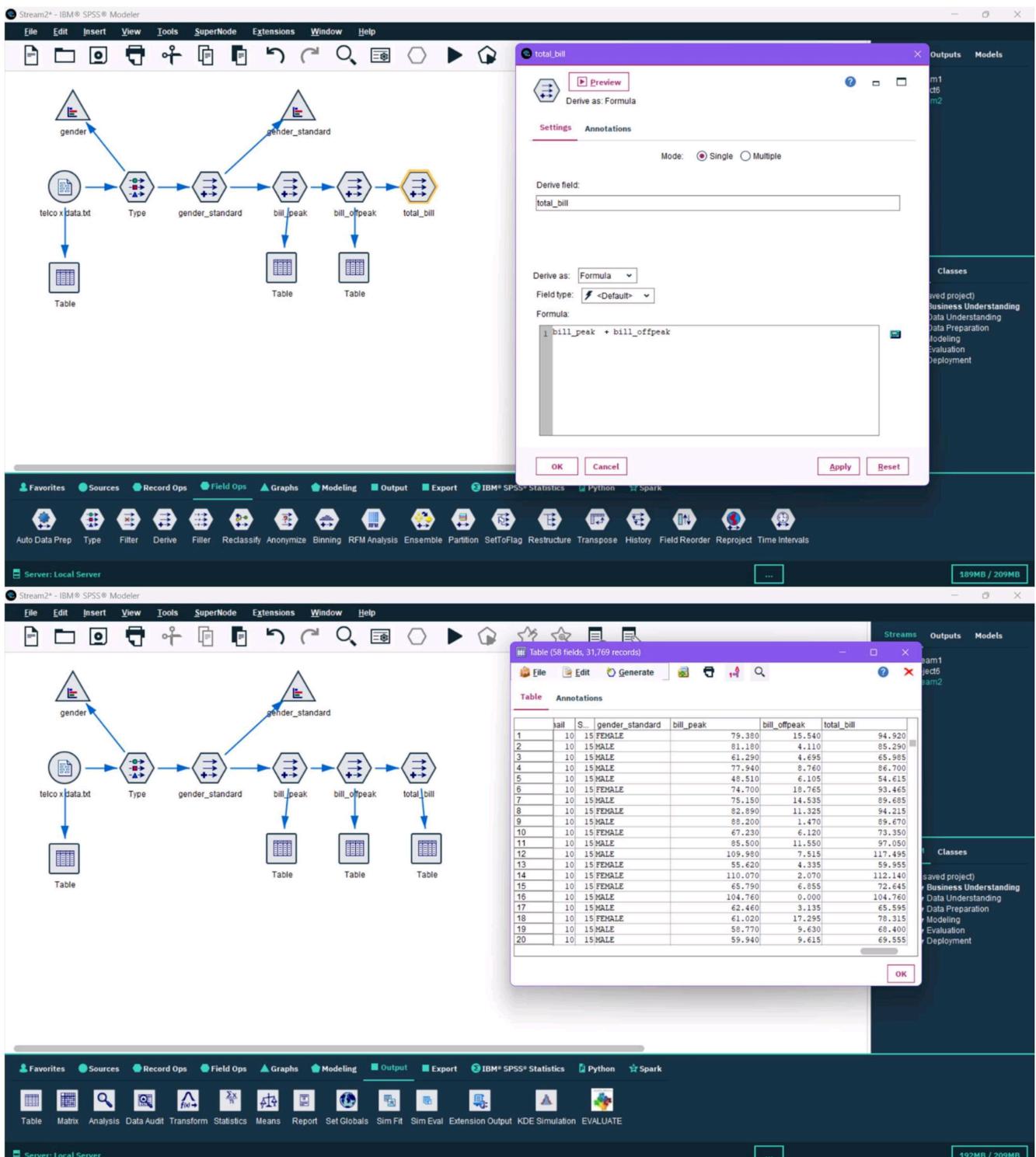
Step 8: Calculating Total Bill

Connect a Derive Node to *offbill_peak*.

Rename field = *total_bill*.

Formula = *bill_peak* + *offbill_peak*.

Attach a Table Node → run to display total bill values.



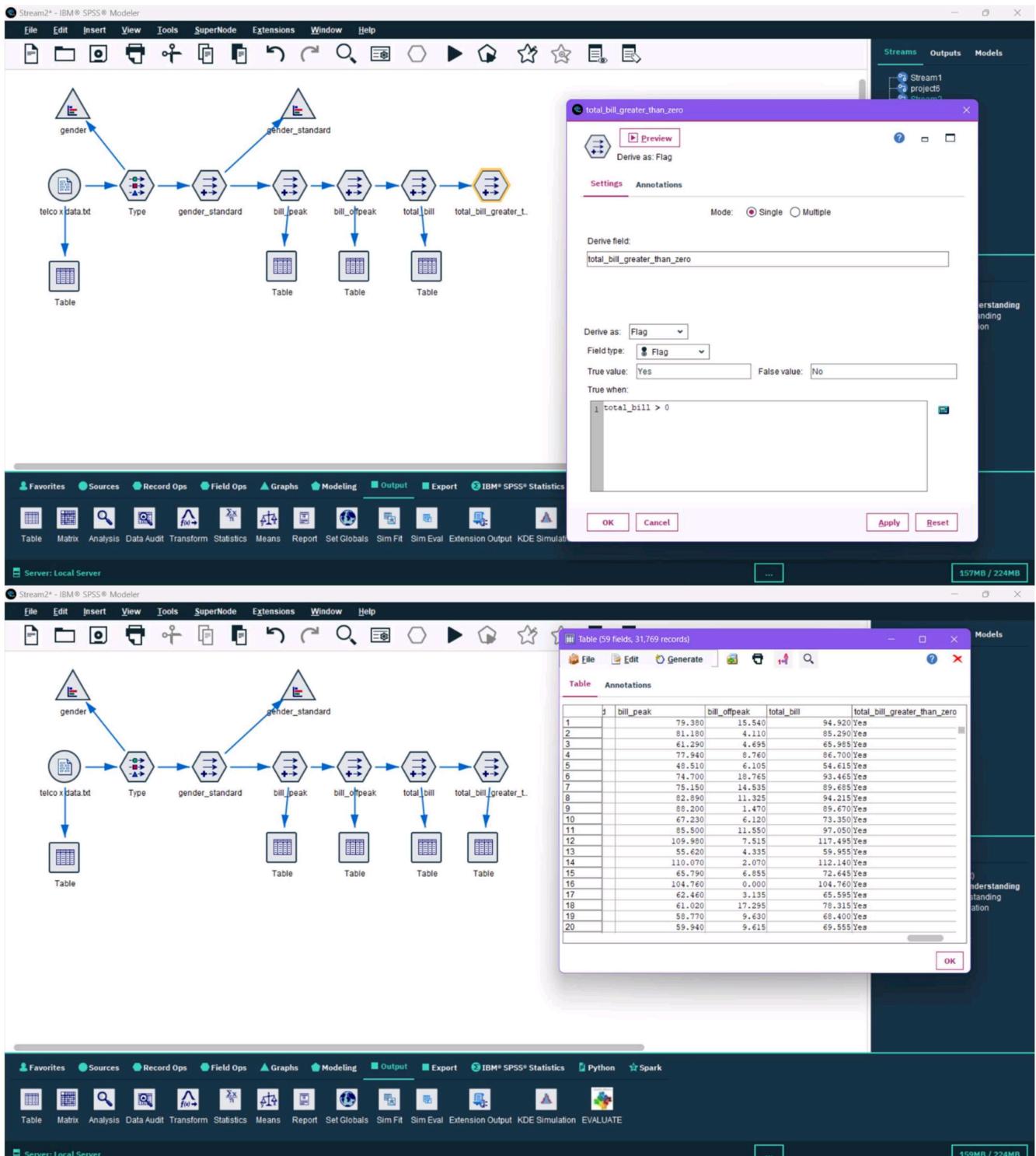
Step 9: Flagging Valid Bills

Insert another Derive Node after `total_bill`.

Rename field = `total_bill_greater_than_zero`.

Type = Flag → condition: `bill_total > 0`.

Attach a Table Node → run to validate flags.



Step 10: Customer Segmentation

Add a Derive Node after the flag.

Rename field = *segment*.

Type = Nominal and Ordinal → expression:

- *kam* if $\text{bill_total} \leq 100$
- *medium* if $100 < \text{bill_total} \leq 150$
- *jaada* if $\text{bill_total} > 150$

Attach a Table Node → run to confirm segment allocation.

Step 11: Applying Discount Rules

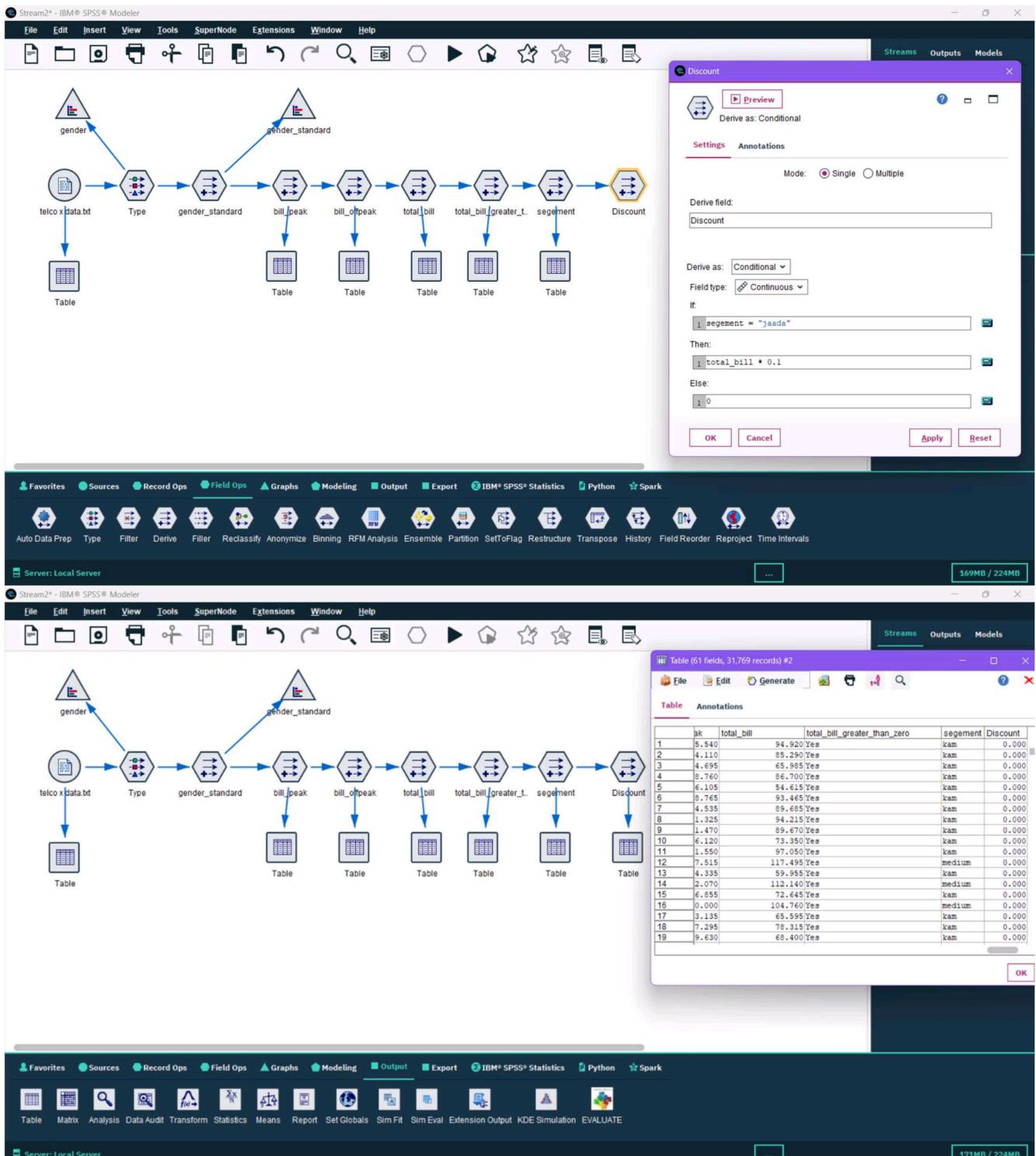
Connect another Derive Node to *segment*.

Rename field = *discount*.

Type = Conditional Continuous → expression:

If segment = "jaada" Then bill_total * 0.1 Else 0.

Attach a Table Node → run to calculate discounts.



Step 12: Reclassifying Gender Field

Finally, add a Reclassify Node after *segment*.

Reclassify field = *gender* → map values into standardized categories (M, F).

Attach a Table Node → run to generate the final cleaned dataset.

