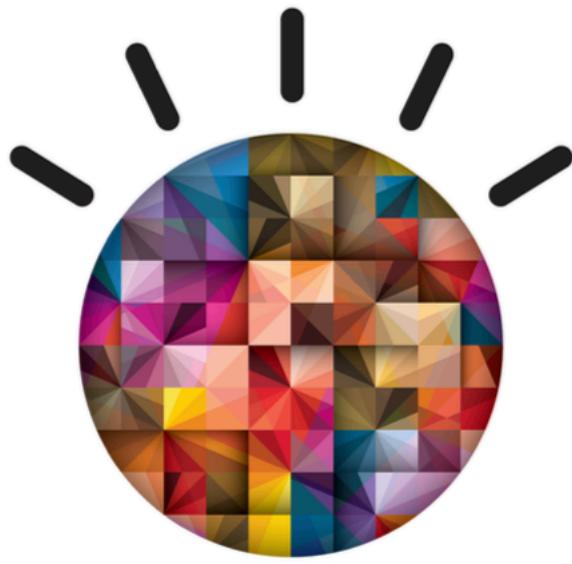


IBM
Analytics



Data Science Workshop
SPSS Modeler Guide



© Copyright IBM Corporation 2018

IBM, the IBM logo, ibm.com, and IBM SPSS are trademarks of International Business Machines Corp., registered in many jurisdictions worldwide. Other product and service names might be trademarks of IBM or other companies. A current list of IBM trademarks is available on the Web at "Copyright and trademark information" at www.ibm.com/legal/copytrade.shtml.

This document is current as of the initial date of publication and may be changed by IBM at any time. Not all offerings are available in every country in which IBM operates.

THE INFORMATION IN THIS DOCUMENT IS PROVIDED "AS IS" WITHOUT ANY WARRANTY, EXPRESS OR IMPLIED, INCLUDING WITHOUT ANY WARRANTIES OF MERCHANTABILITY, FITNESS FOR A PARTICULAR PURPOSE AND ANY WARRANTY OR CONDITION OF NON-INFRINGEMENT. IBM products are warranted according to the terms and conditions of the agreements under which they are provided.

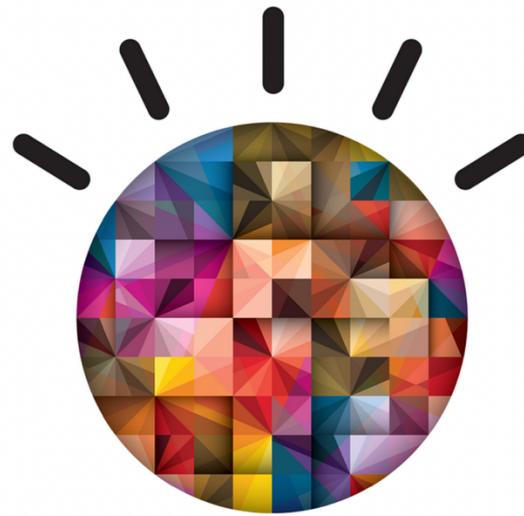
Contents

Exercise 1: Predict in 20 minutes	4
Exercise 2 – Finding Patterns and Groups.....	11

SPSS Modeler

IBM SPSS Modeler is a comprehensive data science platform, designed to bring **predictive intelligence** to everyday business problems, enabling front-line employees to make more **effective decisions** and **improve outcomes**.

This enables organizations to improve business processes and help people or systems **consistently make the right decisions** by delivering recommended actions at the point of impact. The result is a **rapid return on investment (ROI)** and the ability to **proactively and repeatedly reduce costs** and **increase productivity**.



This hands-on Data Science workshop is an instructor led session using IBM's data mining and predictive modeling software and is designed for those who are familiar with predictive analytics, as well as for beginners. Through this workshop you will experience firsthand how IBM SPSS Modeler works and how easy it is to implement predictive analytics.

Exercise 1: Predict in 20 minutes

Use Case

Goal: Approach:

Identify who has high probability to cancel the contract (Churn) to help marketing team to create a list of customer targets to receive specific marketing campaign

- Prepare data for modeling
- Define which fields to use
- Choose the modeling technique
- Automatically generate a model to identify who has cancel
- Review results

Why?

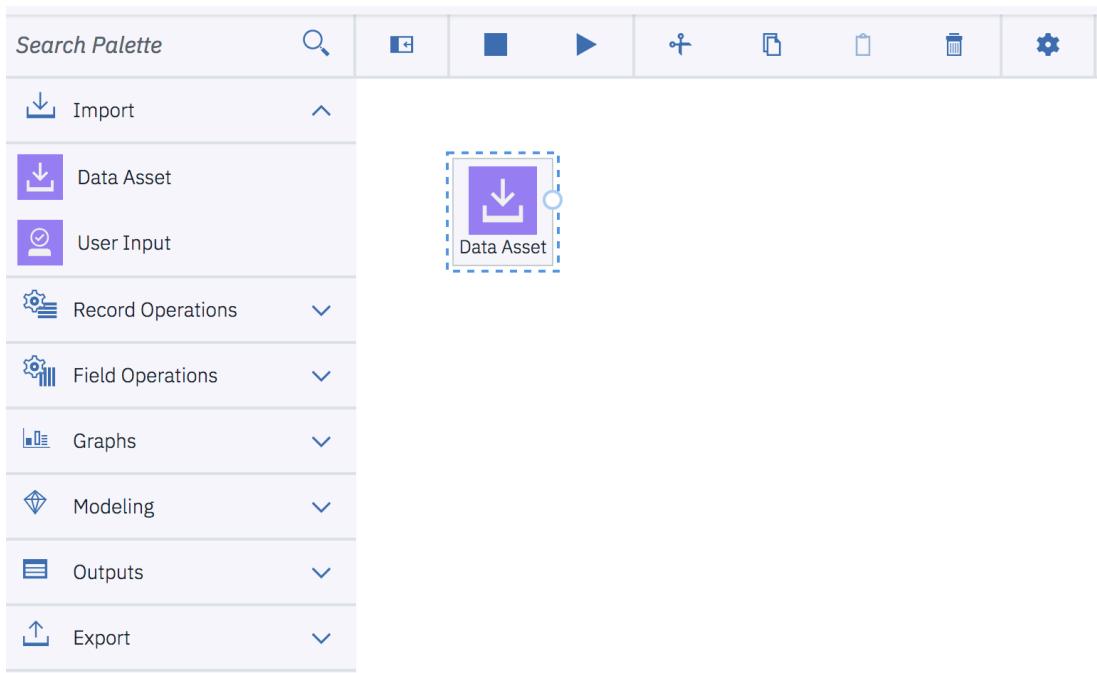
To save marketing cost and reduce the customer churn, identify those likely to churn and focus marketing efforts on those customers.

Predictive in 20 Minutes

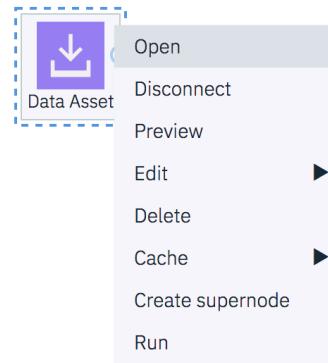
1. Open our Project and add **Customer.csv** file to **Data assets**.
2. Create a new model. Click on **+ New Modeler Flow**.
3. Type a name and description for our model and click **Create**.

The screenshot shows the IBM Watson Studio interface with the 'Modeler' tab selected. A navigation bar at the top includes links for Projects, Tools, Community, Services, Manage, Support, and Docs. Below the navigation bar, there are three tabs: 'New' (selected), 'From file', and 'From example'. A 'Name*' field is present with placeholder text 'Type name here.' and a character limit of 50. A 'Description' field is also present with placeholder text 'Type description here.' and a character limit of 500. Under 'Select flow type', 'Modeler Flow' is selected. Under 'Runtime', 'IBM SPSS Modeler' is selected. The bottom of the screen features a dark footer bar.

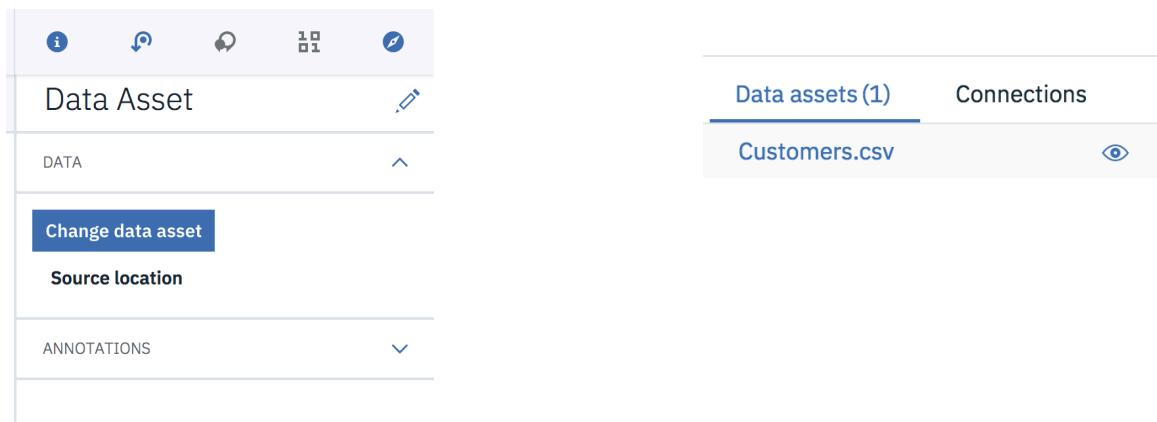
4. From Import, add the **Data Asset** node to the workflow.



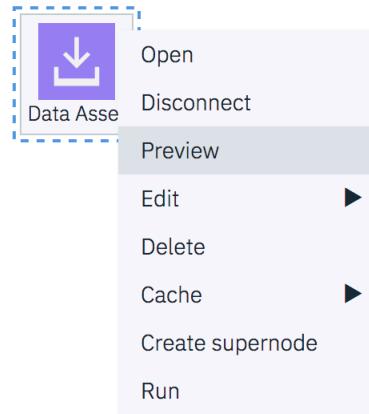
5. Access the Data Asset properties. Right click on Data Asset and select **Open**



6. Click on **Change data asset** to select the **Customer.csv** file already load in the project. Click **OK** and **SAVE**.



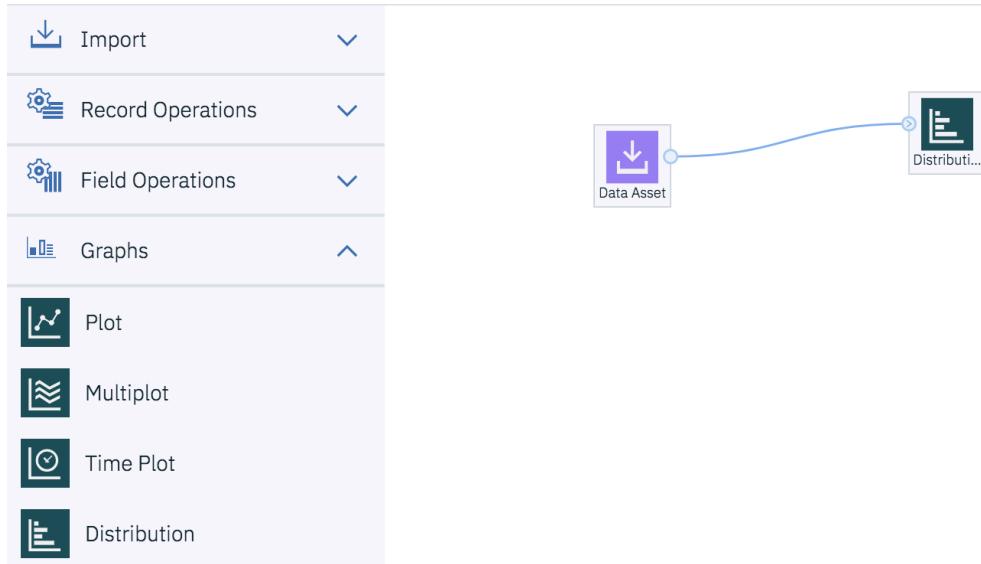
7. Right click on Data Asset and select **Preview**, to visualize the data available on file.



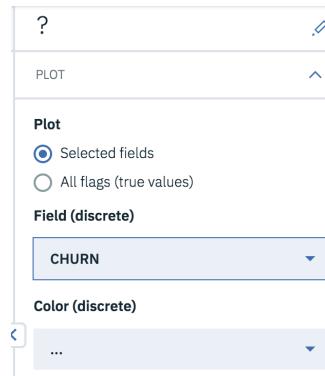
8. To go back to the workflow, click on the workflow name available on the top: My Projects / Project Name / **Workflow Name** / Preview.

	Data	Profile	Visualizations				
ID	Sex	Status	Children	Est_Income	Car_Owner	Usage	Age
2	6	M	M	2	29616	N	75.29
3	8	M	M	0	19732.8	N	47.25
4	11	M	S	2	96.33	N	59.01
5	14	F	M	2	52004.8	N	28.14
6	17	M	M	2	53010.8	N	58.87
7	18	M	M	1	75004.5	N	64.8

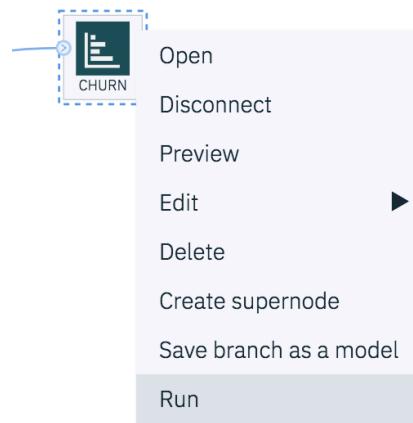
9. From Graphs options, add a Distribution node. Connect it to the data source. Click and hold the button on the first node, move the cursor to the second node and release when the cursor is on the mark of second node.



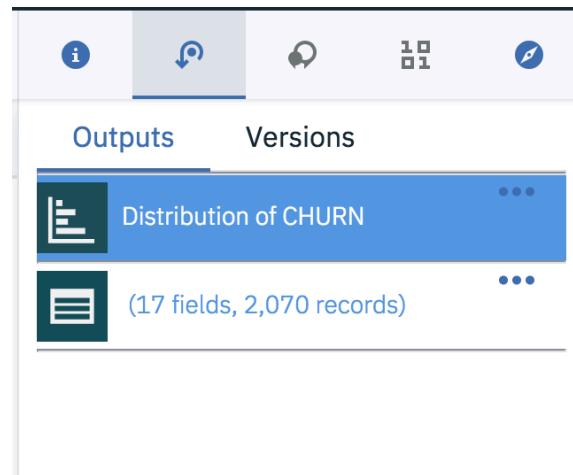
10. Double-click on the Distribution node to access the properties. At **Field (discrete)**, select **CHURN** and click save.



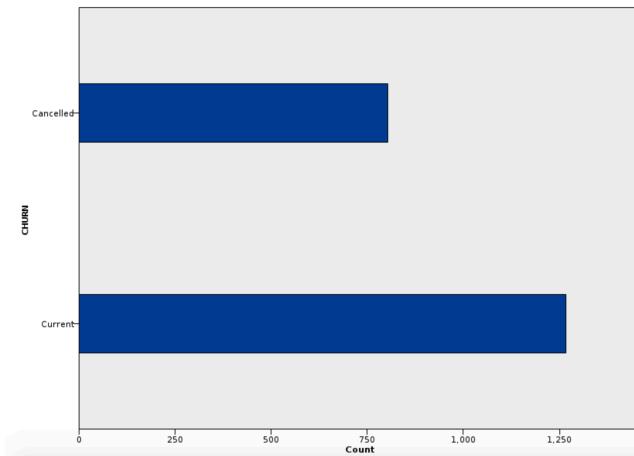
11. Note that distribution node was renamed to CHURN. Right click on it and select **Run**.



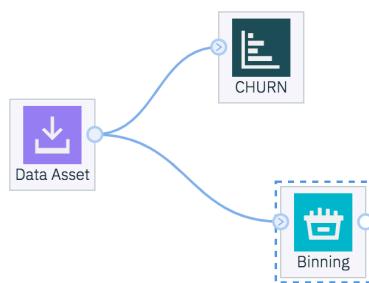
12. Double-click on the graph Distribution of CHURN to open it.



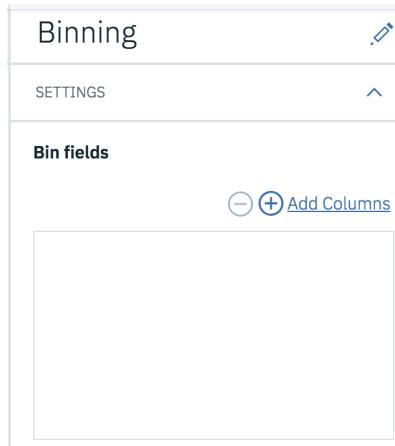
13. The resulting graph shows that of the 2070 customers in this data set, around 40% of customers have cancelled their contract. The task is to build a model to understand the relationships within the data that led to the canceled their contracts.



14. From the Field Operations, add a Binning node to the workflow and connect it to the data source. The Binning node allows you to automatically generate bins (categories) using several techniques. In this case, we will be creating categories from the continuous variable Age.
15. Click and hold the button on the Data Asset node, move the cursor to the Binning node and release when the cursor is on the mark of second node.



16. Double-click on **Binning** to open properties.



17. Click on **Add Columns** and select **Ages**. Click OK.

Select Fields for Binning

<input type="checkbox"/>	Field name ^	Data type ^
<input type="checkbox"/>	ID	integer
<input type="checkbox"/>	Children	integer
<input type="checkbox"/>	RatePlan	integer
<input type="checkbox"/>	Dropped	integer
<input type="checkbox"/>	Est_Income	double
<input type="checkbox"/>	Usage	double
<input checked="" type="checkbox"/>	Age	double
<input type="checkbox"/>	LongDistance	double
<input type="checkbox"/>	International	double
<input type="checkbox"/>	Local	double

18. Scroll down in the **Binning** properties, change the option **Bin width** to **5**. Click Save.

Bin width

5

Use the same bins for all fields

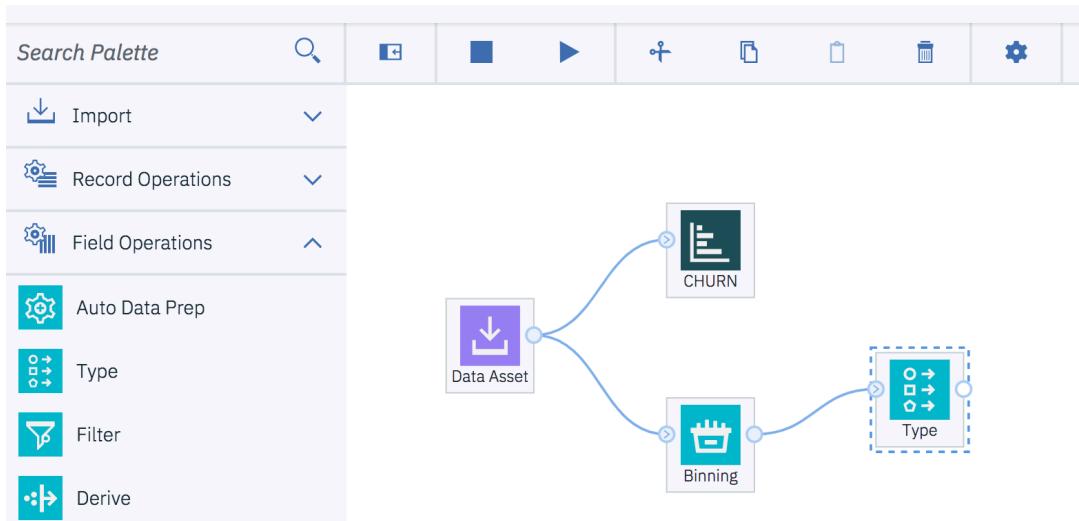
Recalculate bins

Always

If necessary

ANNOTATIONS

19. From the Field Operations, add a Type node to the workflow and connect it to the Binning node.



20. Open the Type node and click the Read Values button to scan the data as well as to display and update the range of values.

Using the drop-down box under Role, modify the following Fields:

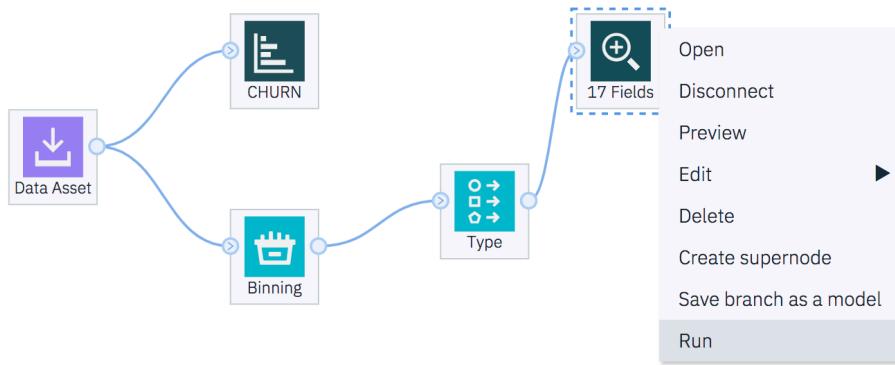
- ID = Record ID
- Age = None
- CHURN = Target

The Measurement of our Target should be set to Flag, which reflects two potential Values: Cancelled or Current. The remaining, including our new Age_Bin, will remain as Inputs in our analysis.

Save

Read Values		Clear All Values				
Field ^	Measure ^	Role ^	Value modeValues ^	Check		
Est_Inco...	Continuous	Input	Specify	96.33, 120000.0	None	
Age_BIN	Nominal	Input	Specify	1, 2, 3, 4, 5, 6, 7, 8, ...	None	
LongDist...	Continuous	Input	Specify	0.0, 59.0	None	
Children	Continuous	Input	Specify	0, 2	None	
ID	Continuous	Input	Specify	1, 3825	None	
Age	Continuous	None	Specify	12.326667, 77.0	None	
CHURN	Flag	Target	Specify	Cancelled, Current	None	
Internati...	Continuous	Input	Specify	0.0, 9.7	None	

21. From the Output, add the Data Audit node and connect it to the Type node. Right-click to select Run

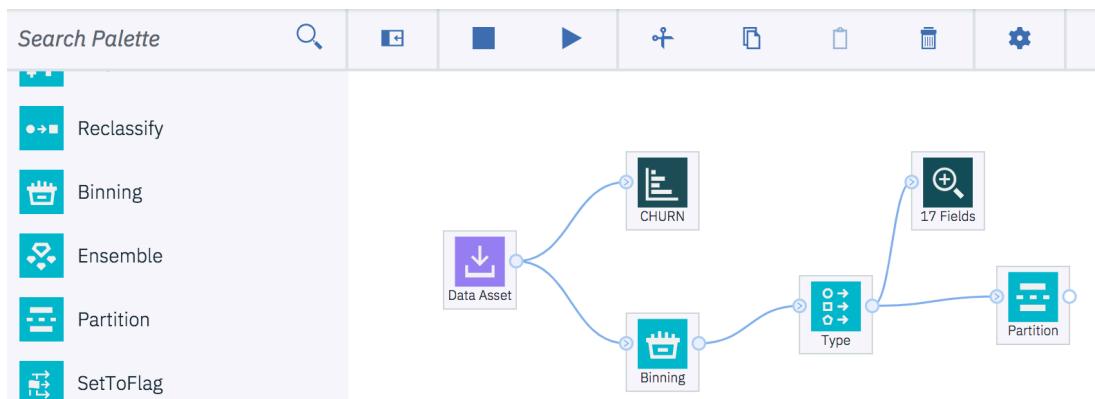


22. Double-click in the Data Audit (outputs) In the Data Audit results output, thumbnail graphs, storage icons, and summary statistics for all fields can be found. It also provides information about outliers, extremes and missing values.

Data Audit of [17 fields]

	Field	Sample Graph	Measurement	Min	Max	Mean	Std. Dev
1	ID		Continuous	1	3825	1901.152	1094.709
2	Sex		Flag	--	--	--	--
3	Status		Nominal	--	--	--	--
4	Children		Continuous	0	2	1.148	0.843

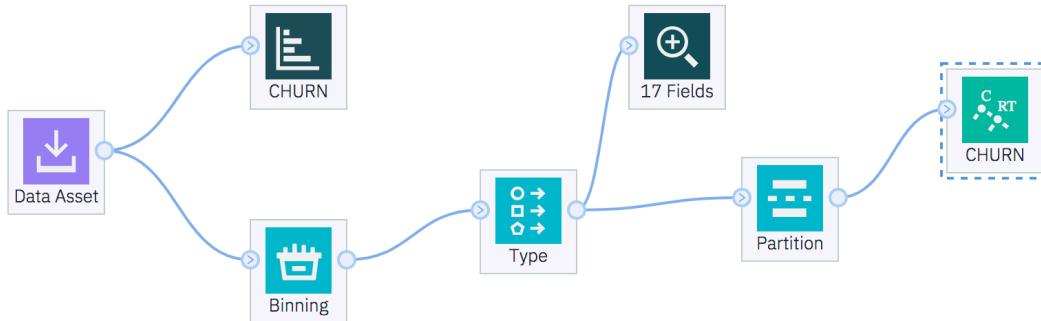
23. Now that we have explored our data, we can build a model to uncover the key drivers resulting in cancelled contract. Go back to the workflow and from the Field Operations, add a Partition node to the workflow and connect it to the Type node.



24. Double-click on the Partition node and change Training Partition to 70 and Testing Partition to 30. Click Save.

Training Partition
70
▼
Testing Partition
30
▼

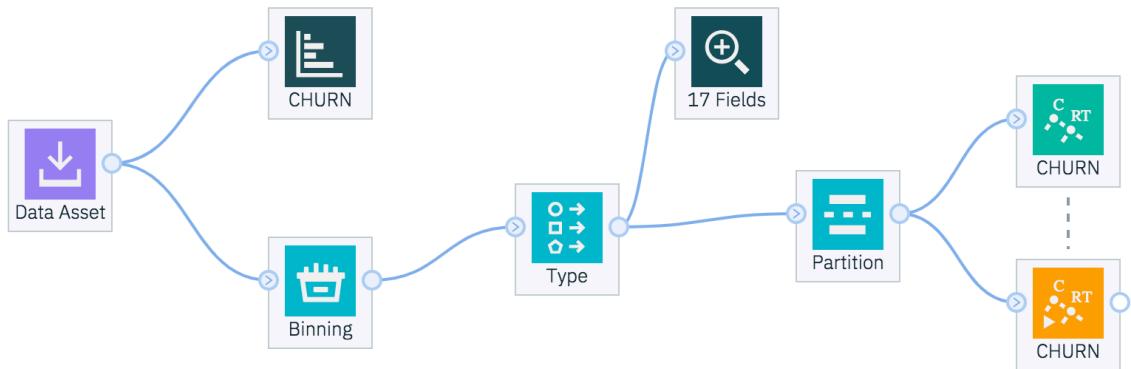
25. From the Modeling, add the C&R Tree node to the workflow and connect it to the Partition node. Note that the C&R Tree node name changes to CHURN when it is connected to the Partition node. This is because we defined CHURN as the Target Role in Step 20.



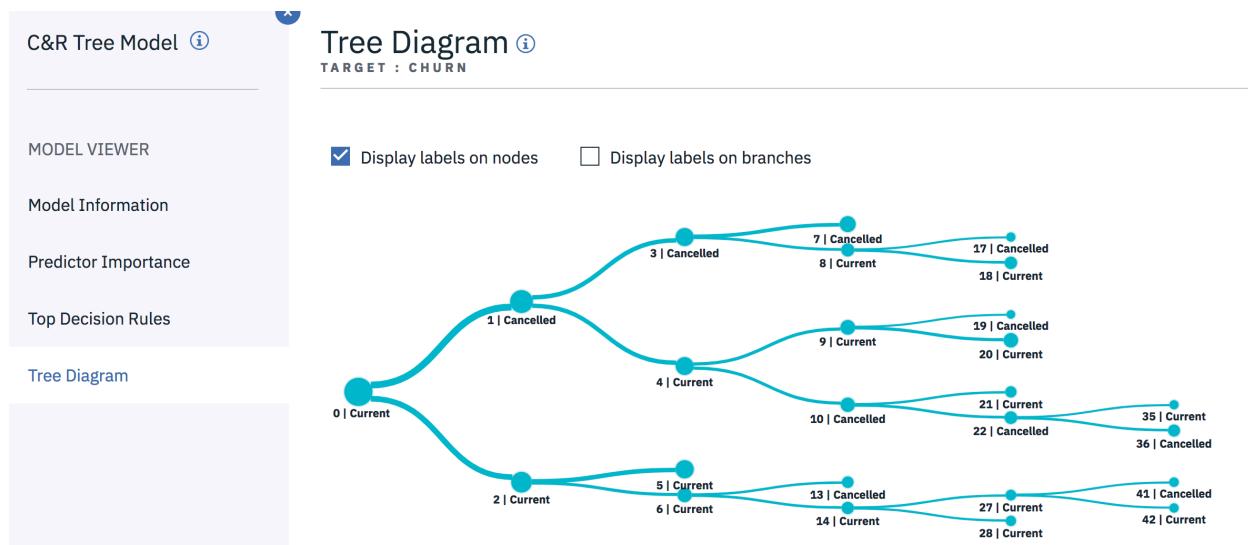
26. Open the C&R Tree node to view the settings before running the model. This example uses the defaults so closed properties and click Run.

CHURN	
	▼
FIELDS	▼
OBJECTIVE	▼
BASICS	▼
STOPPING RULES	▼
COSTS	▼
PRIORS	▼

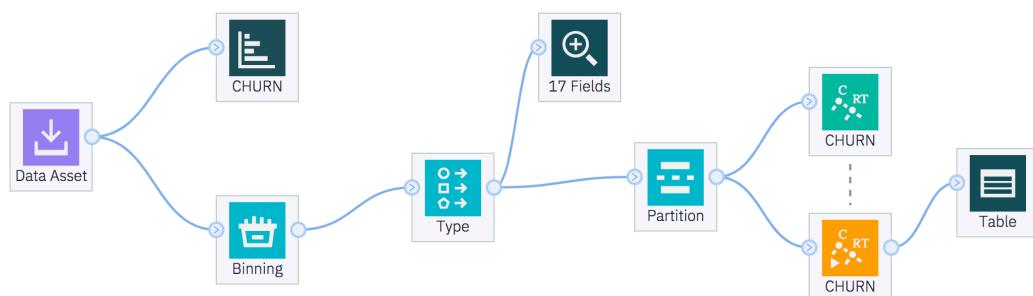
27. The C&R Tree model is generated, added to the workflow and named CHURN.



28. Right-click on the generated model and select View Model to see the outputs. The Model outputs contains Model Information, Predictor Importance, Top Decision Rules and the Tree Diagram. Click the Viewer each one.



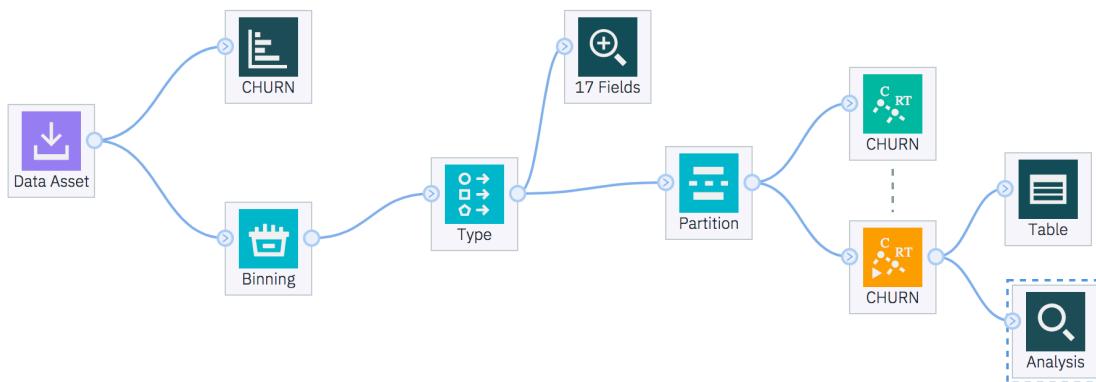
29. To view model output, from the Outputs add a Table node to the generated CHURN model and run it.



30. Look at the last five columns in the table. The fifth from last column is the actual outcome, whether the customer cancelled or is current; the second to last column is the prediction from the model; and the last column is the confidence in the prediction. Note that the model predicted that the customer would cancel with 78.1% confidence.

LocalBilltype	LongDistanceBilltype	CHURN	Age_BIN	Partition	\$R-CHURN	\$RC-CHURN
Budget	Intl_discount	Cancelled	3.000	1_Training	Cancelled	0.781
FreeLocal	Standard	Current	8.000	1_Training	Current	0.896
FreeLocal	Standard	Current	9.000	1_Training	Cancelled	0.781
Budget	Standard	Current	10.000	2_Testing	Current	0.810

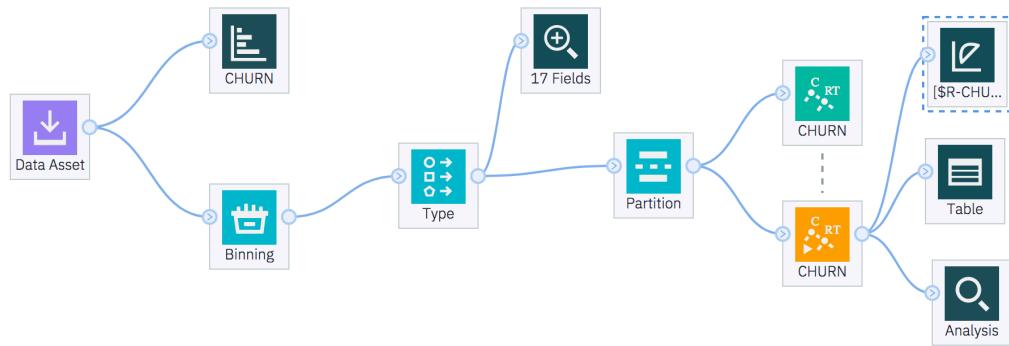
31. To see the overall accuracy of the model, from the Output, add an Analysis node, join it to the generated model and run it.



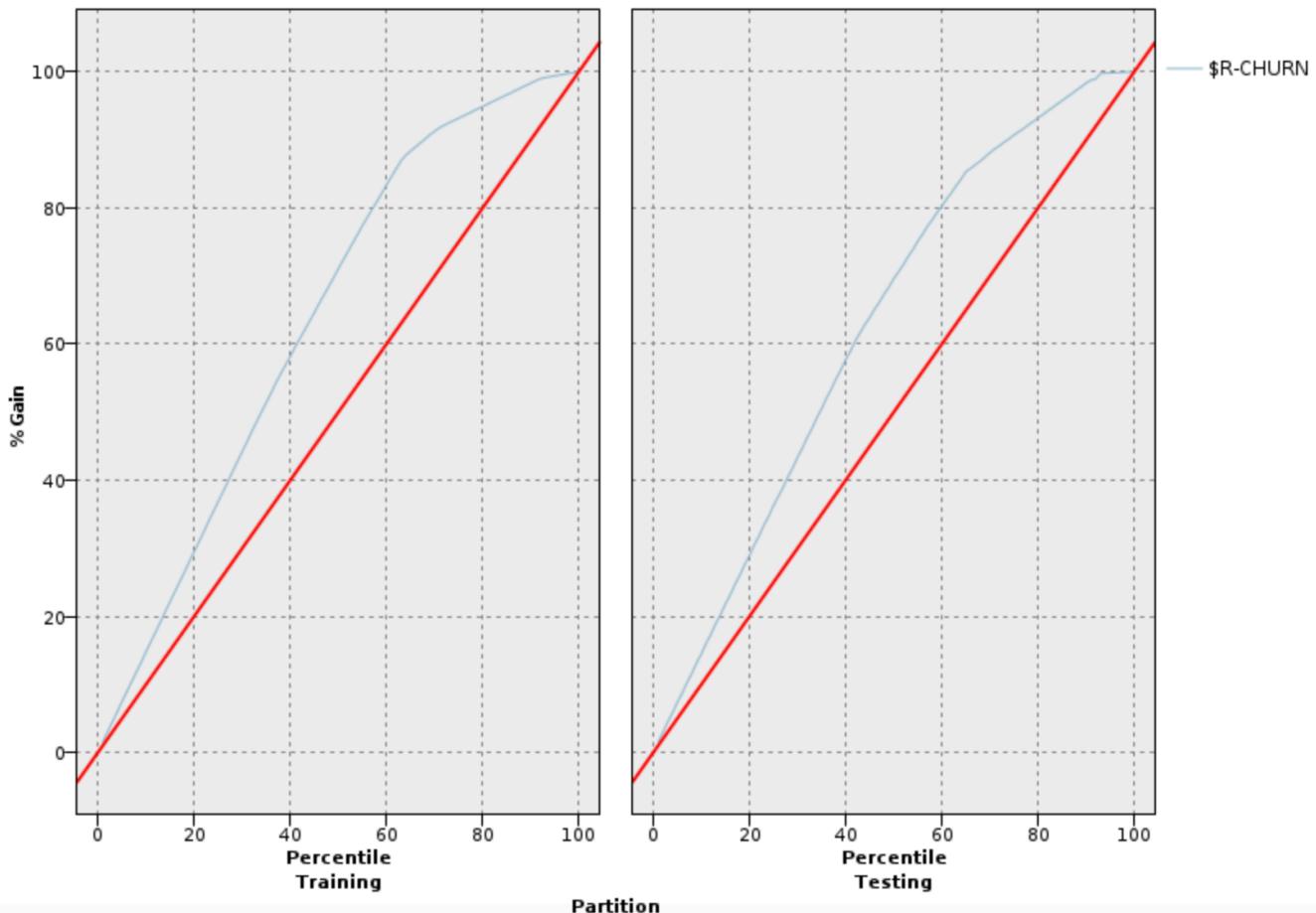
32. The resulting output indicates an overall accuracy of 77.93%. That is, the model predicted with 77.93% accuracy which customers Cancelled or Current the contract.

Results for output field CHURN				
Comparing \$R-CHURN with CHURN				
'Partition'	1_Training	2_Testing		
Correct	1,173	82.49%	505	77.93%
Wrong	249	17.51%	143	22.07%
Total	1,422		648	

33. To further evaluate the model, select an Evaluation node from the Graphs, connect it to the model node and select Run.



34. In the resulting gains chart, the red line reflects what you could expect without Predictive Analytics. The blue line, however, reflects the lift in response you could achieve utilizing Predictive Analytics. Therefore, if you were to randomly select 50% of your client base, you could expect to have captured 50% of those likely cancelled. By using Predictive Analytics, you can more effectively target those 50% of clients and capture almost 80% of those likely to cancel.



35. Finally, to see the relationships between fields, select a Matrix node from the Output and connect it to the model node. Using the drop-down menu, choose “\$R-CHURN” for Rows and “CHURN” for Columns. Select Run.

Matrix of \$R-CHURN by CHURN

\$R-CHURN	Cancelled	Current
Cancelled	580	168
Current	224	1098

Cells contain: cross-tabulation of fields (including missing values)

Chi-square = 738.419, df = 1, probability = 0

Summary

- Prepare data for modeling
- Define which fields to use
- Choose the modeling technique
- Automatically generate a model to identify who has cancelled
- Review results

Over the course of the last 20 minutes, we were able to successfully train a model by exploring SPSS Modeler’s ability to read in data, create new fields via data preparation techniques, choose and run a predictive modeling algorithm, and evaluate the results to accurately identify customer response.

Exercise 2 – Finding Patterns and Groups

Use Case

Goal: Create segments of customers

Approach:

- Merge disparate data sources, including customer data
- Define which fields to use
- Automatically generate a model to group customers
- Apply business terms to new customer groups
- Export newly created groups to a database

Why?

- Better customer understanding (demographics, socio-economic, etc.)
- Tailored messages for each group/segment
- Personal and more relevant for consumers

Customer Reference

A US cable television network turns on the insights with an analytics solution that predicts success of new shows six weeks in advance.

Business challenge: This cable television network faces the challenge of managing huge volumes of information. Previously the network's research team spent a significant amount of time processing data on spreadsheets rather than analyzing it, and based decisions on a combination of experience and instinct. The company needed a large-scale analytics solution to organize this wealth of data, make sense of it, and provide answers and actionable insights.

The transformation: The solution combines television ratings data with information gathered minute by minute and viewer by viewer from a variety of channels and other sources to determine who's watching and why. Then it centralizes the data and makes it available for in-depth, predictive analytics. With insights into audience preferences gained from sophisticated statistical models, including intelligent segmentation, the network can optimize advertising revenue and viewership like never before.

IBM's implemented solution:

Accelerates analytics by extracting insights from billions of rows of audience data in seconds, instead of days.

Triples views of video-on-demand service through data-driven marketing.

Predicts success of a new show six weeks in advance of its release and adjusts marketing accordingly.

"A single day of analysis work enabled the network to design a campaign that increased the consumption of its video on demand service. Previously, that analysis would have taken weeks."

Finding Patterns and Groups

1. Create a new project and add the files Demographics.xlsx and Transactional data.csv to the Data Assets.

Data assets

0 asset selected.

NAME	TYPE
Transactional data.csv	Data Asset
Demographics.xlsx	Data Asset

2. Create a new Modeler flows on the project. Type a name and description and click Create.

Modeler

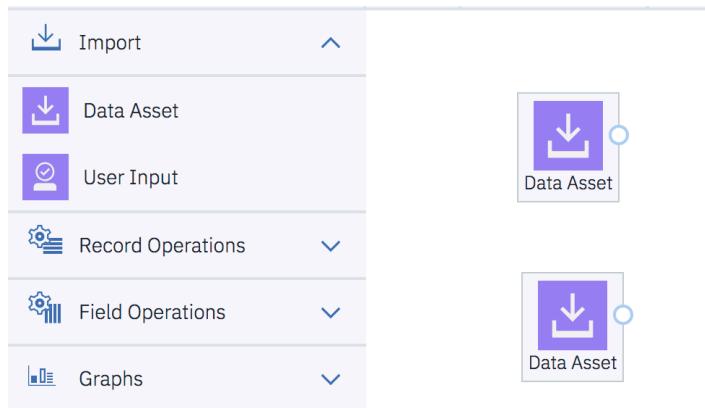
New From file From example

Name*
Customer Segmentation 29

Description
Type description here. 500

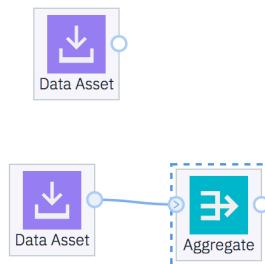
Select flow type
 Modeler Flow Neural Network Modeler BETA

3. From the Import, add two Data Assets and connect one to the Demographics data and the other to the Transactional data.



We will be joining two data files together; one file contains customer transaction data, and exists in a database, while the other file is an Excel file and contains customer demographics. Though the files are in different formats, the Merge node in IBM SPSS Modeler can very easily join the files without first requiring the analyst to translate one file and/or the other into the same format.

4. First, we need to aggregate the transactional data to Customer ID. From the Record Operations, add a Aggregate node to the workflow and connect it to the Transactional data Data Asset node.



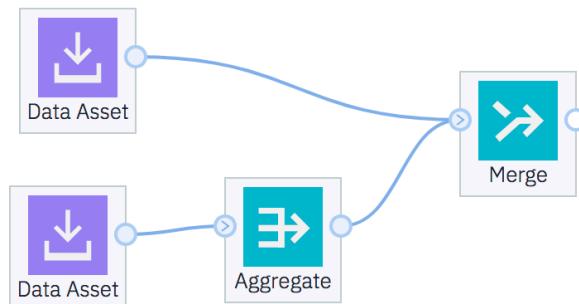
5. Double-click to the Aggregate node, in Key Fields, click Add Columns and select Customer ID.



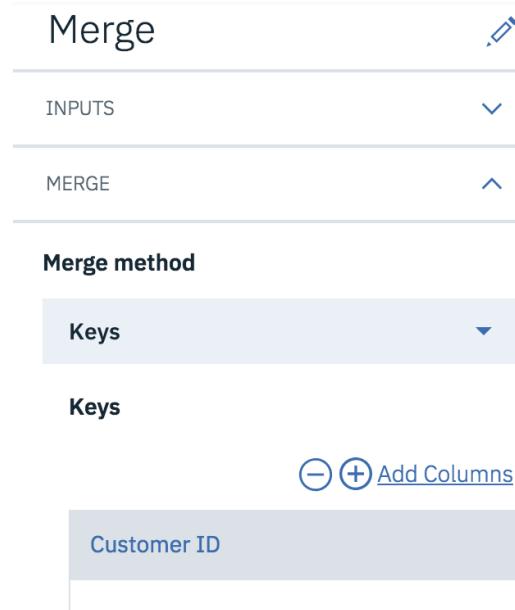
6. Click Aggregations, click Add Columns and select Total Spend and Number Transactions and click OK. Click on the edit icon and uncheck the Mean operate, click OK and Save



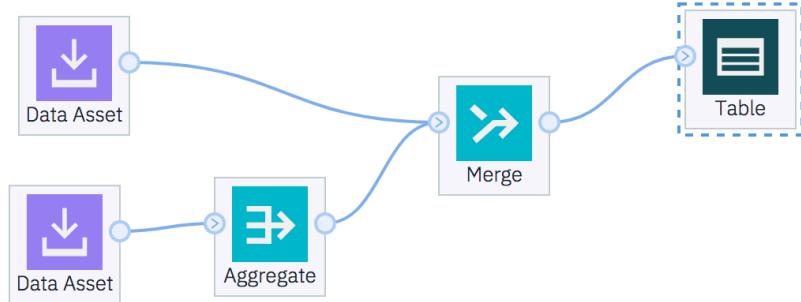
7. From the Record Operations, add a Merge node to the workflow and connect the Data Assets node (Demographics) and Aggregate node to the Merge node.



8. Open the Merge node and change the Merge method to Keys. Click on Add Columns and add Customer ID. Click Save.



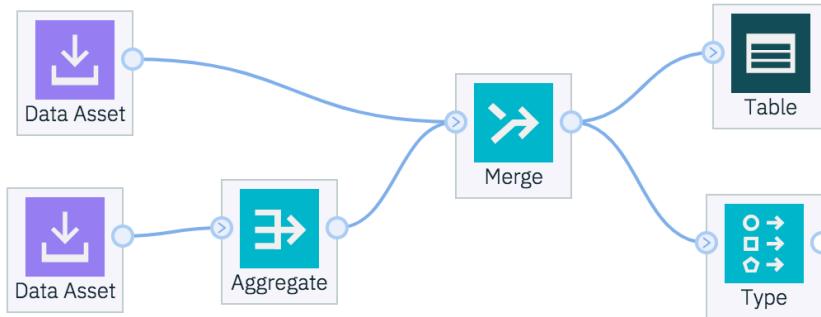
9. From the Output, add a Table node to the workflow and connect it to the Merge node.



10. Once connected, right click on the Table node and select Run to review the data after having been merged.

Customer ID	Marital Status	Income in Thousand	Number of Children	Level of Education	Age Category_	Total Spend_Sum	Number Transactions
1.000	Married	140.000	3.000	Some High School	35 - 49	43.000	26.000
2.000	Married	225.000	4.000	High School Graduate	35 - 49	56.800	34.000
3.000	Married	225.000	3.000	High School Graduate	35 - 49	56.800	34.000
4.000	Married	210.000	4.000	High School Graduate	50 - 64	84.000	36.000
5.000	Married	150.000	3.000	Some High School	35 - 49	48.000	32.000

11. From the Field Operations, select a Type node and attach it to the Merge node.

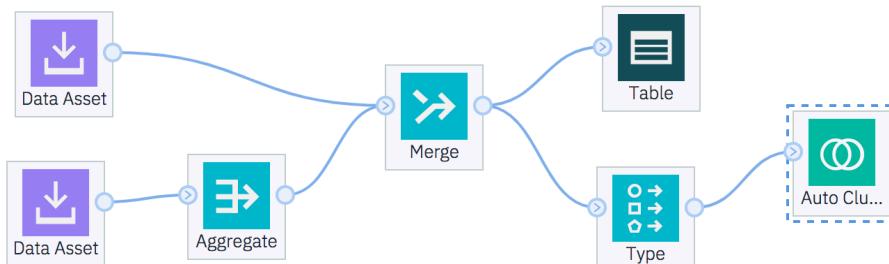


12. Double-click on the Type node to edit it. Click on the Read Values button and change the role of "Customer ID" to "Record ID". Once done, click Save.

		Read Values	Clear All Values		
Field	Measure	Role	Value mode	Values	Check
Level of E...	Nominal	Input	Specify	College Graduate, ...	None
Total Spe...	Continuous	Input	Specify	18.4, 171.0	None
Custome...	Continuous	Record	Specify	1, 157	None
Marital St...	Flag	Input	Specify	Married, Single	None
Age Cate...	Nominal	Input	Specify	19 or below, 20-34...	None
Number ...	Continuous	Input	Specify	2, 6	None

At this point we are ready to cluster the cases into segments. For this we will use the 'Auto Cluster' node. The Auto Cluster node allows you to try all of the clustering algorithms and, at your discretion, any or all of their parameters. It builds all of the models you specify and shows you the best models (3 is the default) to use with your data. So, in one step, you will have the best model(s) without having to know or guess which might work for you. This is also a nice way to see how other modeling algorithms will perform. The same holds for the other auto modeling nodes that address classification and numeric modeling techniques.

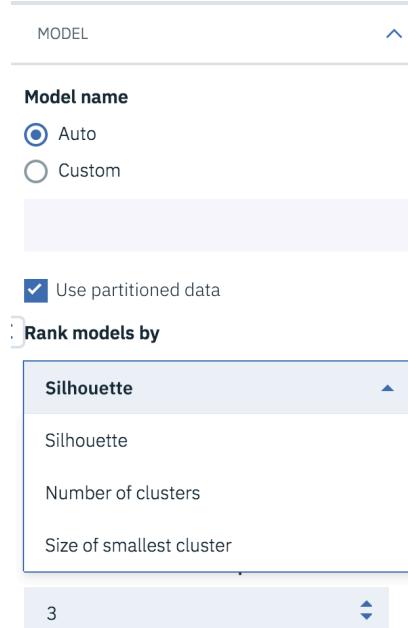
13. From the Modeling select the Auto Cluster modeling node and attach it to the Type node.



14. Double-click on the Auto Cluster node to edit it.

On the ‘Model’ options, there are various ways of ranking the quality of the models that are built. Keep the default method, ‘Silhouette’.

Also keep the default number of models to keep, 3. This means that the 3 best models, based on their silhouette measure, will be retained for our use.

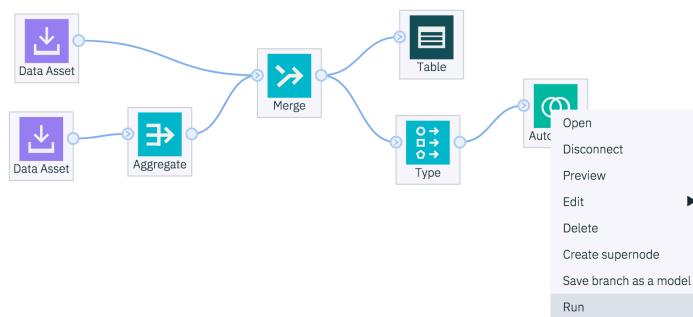


15. Scroll down on the ‘Expert’ options and click on Select models. You notice that there are 3 clustering algorithms in the list.

Select models

<input checked="" type="checkbox"/>	MODEL TYPE	SETTINGS	NUMBER OF MODELS
<input checked="" type="checkbox"/>	Kohonen	Default	1
<input checked="" type="checkbox"/>	K-means	Default	1
<input checked="" type="checkbox"/>	TwoStep	Default	1

16. Go back to the workflow and Click Run in the Auto Cluster node.



17. Right-click on the model results and select View Model to view the results of the auto-clustering analysis.

Auto Cluster - Models i

USE	ESTIMATOR	GRAPH	SILHOUETTE	BUILD TIME (MINS)	NUMBER OF CLUSTERS	SMALLEST CLUSTER (N)	SMALLEST CLUSTER (%)	LARGEST CLUSTER (N)	LARGEST CLUSTER (%)
<input checked="" type="checkbox"/>	TwoStep		0.434	< 1	4	22	0.143	47	0.305
<input type="checkbox"/>	KMeans		0.377	< 1	5	1	0.006	68	0.433
<input type="checkbox"/>	Kohonen		0.141	< 1	11	1	0.006	40	0.255

This is where you see the 3 ‘best’ models for segmenting the data. The list is in descending order by Silhouette measure like we specified. There are other important statistics about each model in the table.

The results show that the TwoStep algorithm has the best Silhouette measure followed by the K-means and Kohonen models, which were discarded.

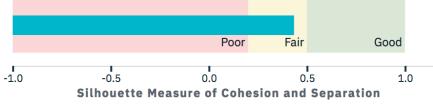
The check boxes to the left indicate that the TwoStep model will be used since it was ranked highest according to our ranking criterion. With cluster models, only one can be selected at a time; but you could choose to use any of the others by clicking the check box. For our exercise we will stay with the TwoStep model.

18. Click the first model name labeled ‘TwoStep’ to see the results of the TwoStep cluster analysis.

(←) Auto Cluster - Auto Cluster

TwoStep Clustering Model i
Cluster Quality i

EVALUATION

Cluster Quality


 Silhouette Measure of Cohesion and Separation

MODEL VIEWER

 Cluster Quality

 MODEL INFORMATION

 Predictor Importance

 Cluster Sizes

 Cluster Comparison

Cluster Quality Parameters

Overall Clustering Quality (Avg. Silhouette)	0.434
Total Within Clusters Sum of Squares	0.056
Average Within Cluster Sum of Squares	0.014

Looking at the Cluster Quality measure in the left panel, we see that the Silhouette measure (which is a measure of the clusters' internal cohesion AND how well they exclude dissimilar cases) is fair, with a value of just under 0.5. Such results are common but may also suggest that fewer and/or other variables might be needed to increase the Silhouette value.

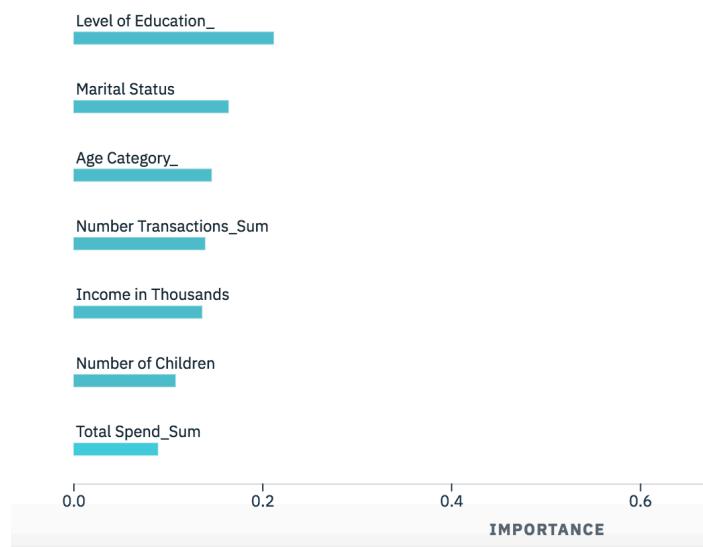
- From the menu in the right viewer select Model Information.

We specified a range of 2 to 4 clusters; and the Two Step clustering engine resolved into 4 clusters.

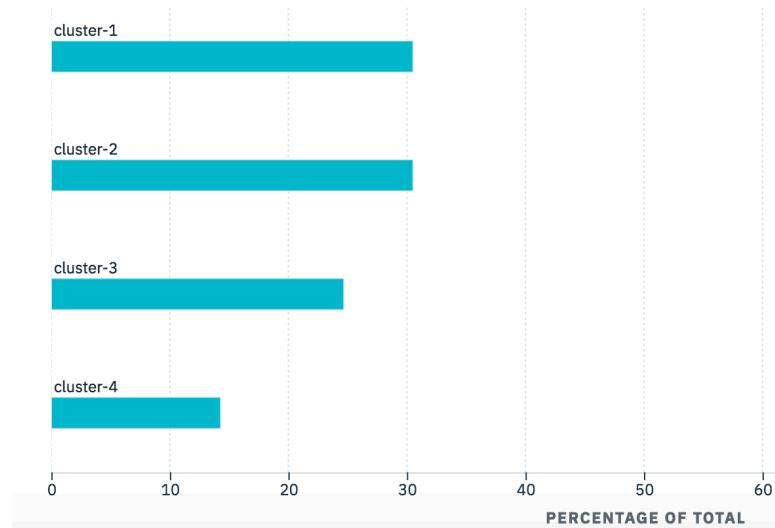
Algorithm	TwoStep	
Model Class	Distribution Based	
Number of Features	7	
Distance Measure	Log Likelihood	
Number of Clusters	4	
Number of Instances in each cluster	cluster-1	47 (30.52%)
	cluster-2	47 (30.52%)

- From the menu in the right viewer select Predictor Importance.

Now the viewer displays a graph with the variables ranked in order of importance for cluster definition. We can see that Level of Education is the most important variable, followed by Marital Status.



21. From the menu in the right viewer select Cluster Sizes.



22. From the menu in the right viewer select Cluster Comparison.



23. From the menu in the right viewer select Clusters.



24. From the menu in the right viewer select Cell Distributions (Absolute).

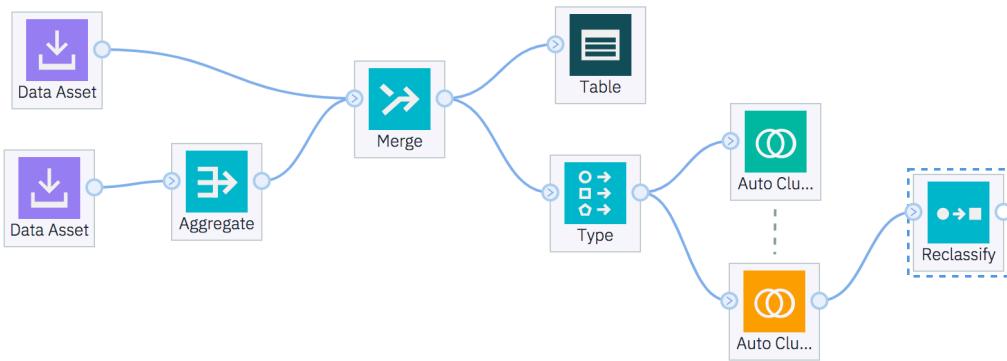
The panel of the Viewer displays the clusters in order of their size, left to right. The darkness of the shading of each variable indicates its importance in cluster definition; the lighter the shading, the less important is the variable in defining the clusters.



25. From the menu in the right viewer select Cell Distributions (Relative).



26. Go back to the workflow. From the Field Operations, drag a Reclassify node onto the workflow and connect it to the model node.



27. Double-click on the Reclassify node to edit it.

In the Reclassify Field, select the variable **\$XC-autocluster**

Click “Existing Field” just above so that it doesn’t make a new field.

Click the “**Get values**” button to populate the ‘Original value’ column for you. Enter the new values on the right, which better describe the clusters. An example is shown below, but you can enter your own labels as desired. Once completed, click Save.

Reclassify Into

New field
 Existing field

Reclassify Field

\$XC-autocluster

New Field Name

Get values Copy Clear new

> Automatically Reclassify

Values

ORIGINAL VALUE ^	NEW VALUE ^
cluster-1	Low Value Married No Kids
cluster-2	Average Married No Kids
cluster-3	Young Starters
cluster-4	High Income Families

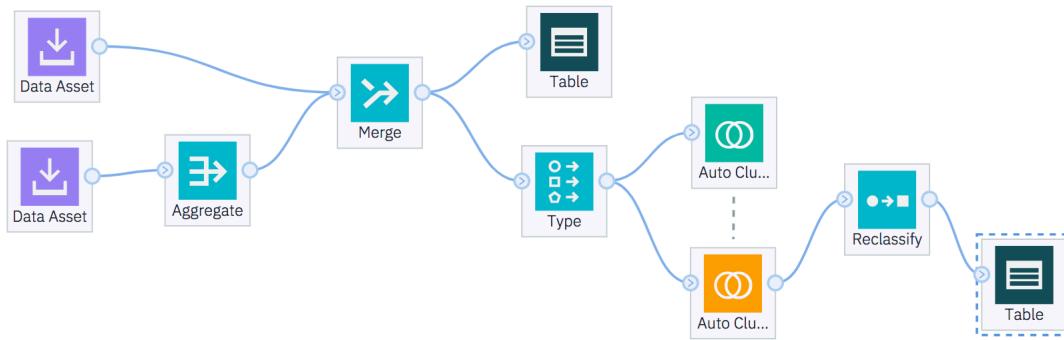
⊖ ⊕ Add Value

For Unspecified Values Use

Original value Default value

Cancel **Save**

28. From the Output, drag a Table node onto the workflow and connect it to the Reclassify node. Once connected, right-click on the Table node and Select Run.



29. The resulting table now includes a new column with the cluster assignments.

Customer ID	Marital Status	Income in Thousand	Number of Children	Level of Education	Age Category_	Total Spend_Sum	Number Transactions	\$XC-autocluster
1.000	Married	140.000	3.000	Some High School	35 - 49	43.000	26.000	Low Value Married
2.000	Married	225.000	4.000	High School Graduate	35 - 49	56.800	34.000	Average Married No Children
3.000	Married	225.000	3.000	High School Graduate	35 - 49	56.800	34.000	Average Married No Children
4.000	Married	210.000	4.000	High School Graduate	50 - 64	84.000	36.000	Average Married No Children
5.000	Married	150.000	3.000	Some High School	35 - 49	48.000	32.000	Low Value Married
6.000	Married	200.000	4.000	High School Graduate	50 - 64	68.000	36.000	Average Married No Children
7.000	Married	310.000	4.000	College Graduate	50 - 64	124.000	46.000	High Income Family
8.000	Married	170.000	3.000	High School Graduate	35 - 49	54.000	32.000	Average Married No Children
9.000	Married	193.000	3.000	High School Graduate	50 - 64	66.800	32.000	Average Married No Children

Summary

- Merge disparate data sources, including customer data from a database or CRM
- Define which fields to use
- Automatically generate a model to group customers
- Apply business terms to grouped customers
- Send new groups to file

For this exercise, we merged two data sources together in order to identify groupings within our data, we used an automated clustering technique, specifying desired parameters. The resulting clusters were reclassified into business terms and could be exported back to the file or database.



© Copyright IBM Corporation 2018

IBM Corporation
Software Group
Route 100
Somers, NY 10589

Produced in the United States of America
February 2017

IBM, the IBM logo, ibm.com, and SPSS are trademarks of International Business Machines Corp., registered in many jurisdictions worldwide. Other product and service names might be trademarks of IBM or other companies. A current list of IBM trademarks is available on the Web at "Copyright and trademark information" at: www.ibm.com/legal/copytrade.shtml.

This document is current as of the initial date of publication and may be changed by IBM at any time.

Not all offerings are available in every country in which IBM operates.

THE INFORMATION IN THIS DOCUMENT IS PROVIDED "AS IS" WITHOUT ANY WARRANTY, EXPRESS OR IMPLIED, INCLUDING WITHOUT ANY WARRANTIES OF MERCHANTABILITY, FITNESS FOR A PARTICULAR PURPOSE AND ANY WARRANTY OR CONDITION OF NONINFRINGEMENT.

IBM products are warranted according to the terms and conditions of the agreements under which they are provided.