



Interactive Modeling with Visual Statistics

Danny Modlin
Senior Analytical Training Consultant

#SASinnovate

Contents

Abstract	2
Precursory Information	3
<i>Data Dictionary</i>	3
<i>Setting Up the Report</i>	4
Logistic Regression	7
<i>Using the Logistic Regression Object</i>	7
Honest Assessment and Variable Selection	17
<i>Creating a Partition Variable</i>	17
<i>Using Fast Backwards Elimination</i>	20
<i>Switching to an Alternative Model</i>	22
Background Information	23
<i>SAS Viya Information</i>	23
<i>What is Visual Statistics?</i>	24
<i>Logistic Regression Review</i>	25
<i>Suggested SAS Courses to Learn More Information</i>	26

Abstract

In this hands-on workshop, learn to perform statistical analysis on any size data, quickly and easily, for maximum impact. SAS® Visual Statistics tasks put the power of large-scale statistical analysis at your fingertips. Learn to create generalized linear models, cluster analyses, and more!

Precursory Information

Data Dictionary

In this workshop, you use the data set **vs_bank**. This data set consists of observations taken from account holders at a large financial services firm. The accounts represent consumers of home equity lines of credit, automobile loans, and other short- to medium-term credit instruments. Appropriate data cleansing has already been applied, so we can begin with statistical modeling. The target variables relate to whether that account holder purchased a new product from the bank in the past year. The data sets contain more than 1 million rows and 24 columns. A list of variables and their labels is shown below.

Target Variable

B_TGT	New Product (Binary)
--------------	----------------------

Categorical Inputs

CAT_INPUT1	Account Activity Level
CAT_INPUT2	Customer Value Level

Interval Inputs

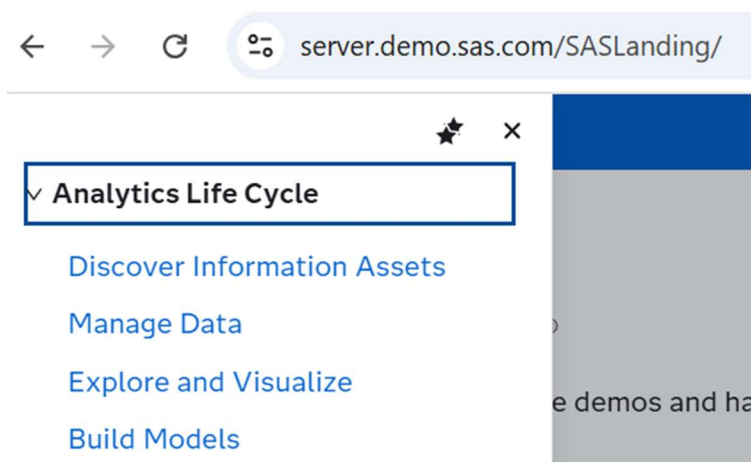
RFM1	Average Sales Past Three Years
RFM2	Average Sales Lifetime
RFM3	Avg Sales Past Three Years Dir Promo Resp
RFM4	Last Product Purchase Amount
RFM5	Count Purchased Past Three Years
RFM6	Count Purchased Lifetime
RFM7	Count Prchsd Past Three Years Dir Promo Resp
RFM8	Count Prchsd Lifetime Dir Promo Resp
RFM9	Months Since Last Purchase
RFM10	Count Total Promos Past Year
RFM11	Count Direct Promos Past Year
RFM12	Customer Tenure

Note: Other variables, not listed here, are also included in the data set. Variables with the prefix **I_** are imputed. Variables with the prefix **RI_** are imputed and replaced. Variables with the prefix **LOGI_** are imputed and log transformed. Variables with the prefix **DEMOG_** are demographic inputs.

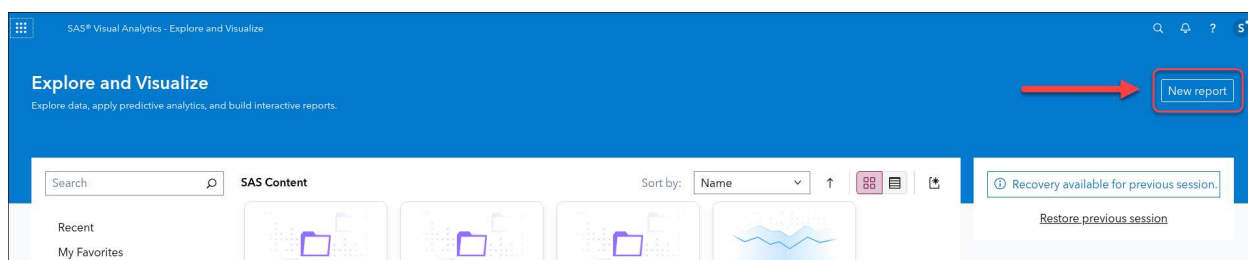
Setting Up the Report

This demonstration illustrates how to build a logistic regression model in SAS Visual Analytics. The demonstration uses the **vs_bank** data to model whether a customer contracted for at least one product in the previous campaign season. You create a binary logistic regression with both categorical and continuous explanatory variables. You then perform model validation and variable selection.

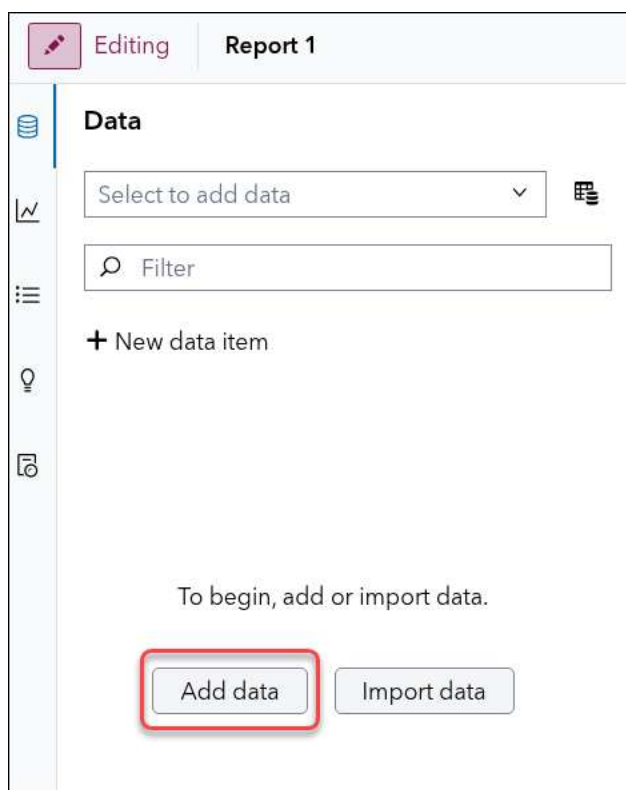
1. From the desktop, open Google Chrome.
2. From the Google Chrome toolbar, select **SAS Landing**.
3. Enter **student** in the **User ID** field.
4. Enter **Metadata0** in the **Password** field.
5. If requested to save the password, select **Save**.
6. Select **NO** when asked about assumable groups.
7. Access the applications menu in the upper left of the window and select **Explore and Visualize**.



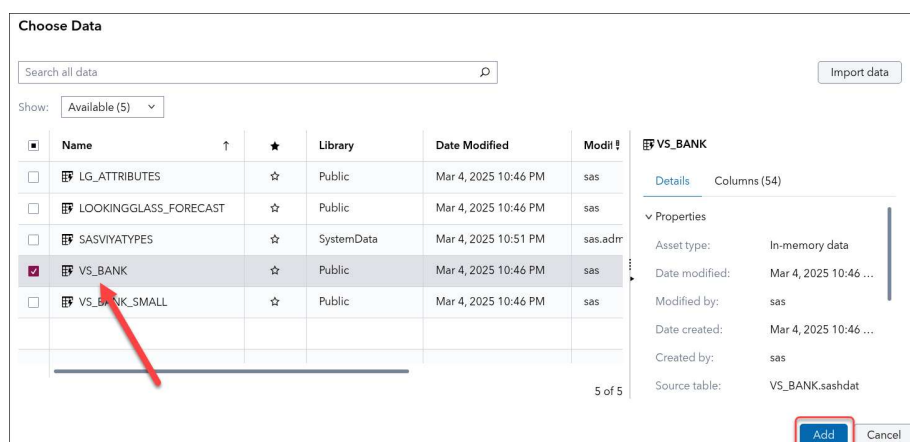
8. Load the SAS data set into CAS.
 - a. In the Explore and Visualize window, click **New report**.



- b. On the Data tab on the left side of the screen, select **Add data** to load an in-memory table to SAS Visual Analytics.

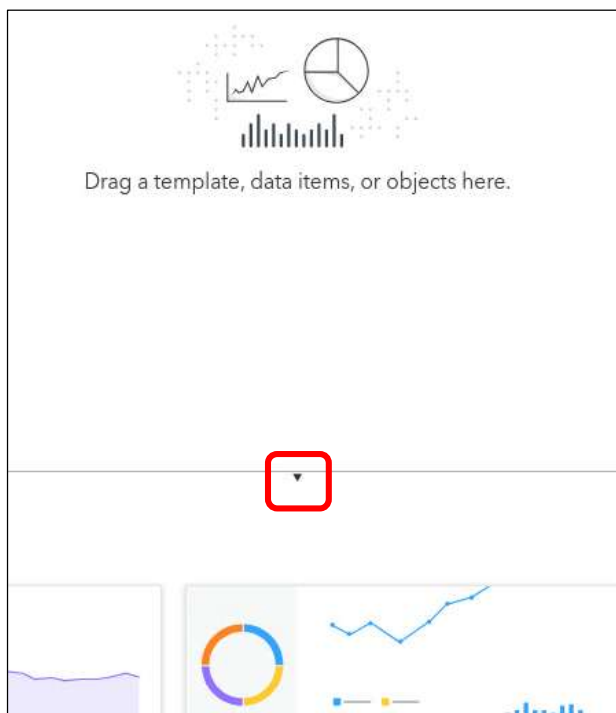



9. In the Choose Data window, select **VS_Bank > Add**.



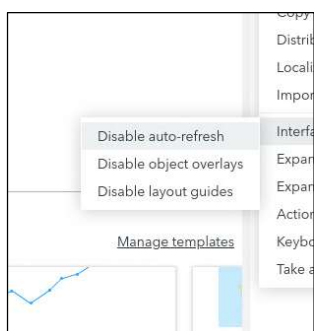
10. If prompted, close the Start your tour window.

11. Collapse the Manage Templates window if necessary.



12. From the menu bar, click  (**More**) and select **Interface options > Disable auto-refresh**.

Disabling auto-refresh enables you to set up several roles in the model before it is created. Otherwise, the model is updated anytime that a change is made.

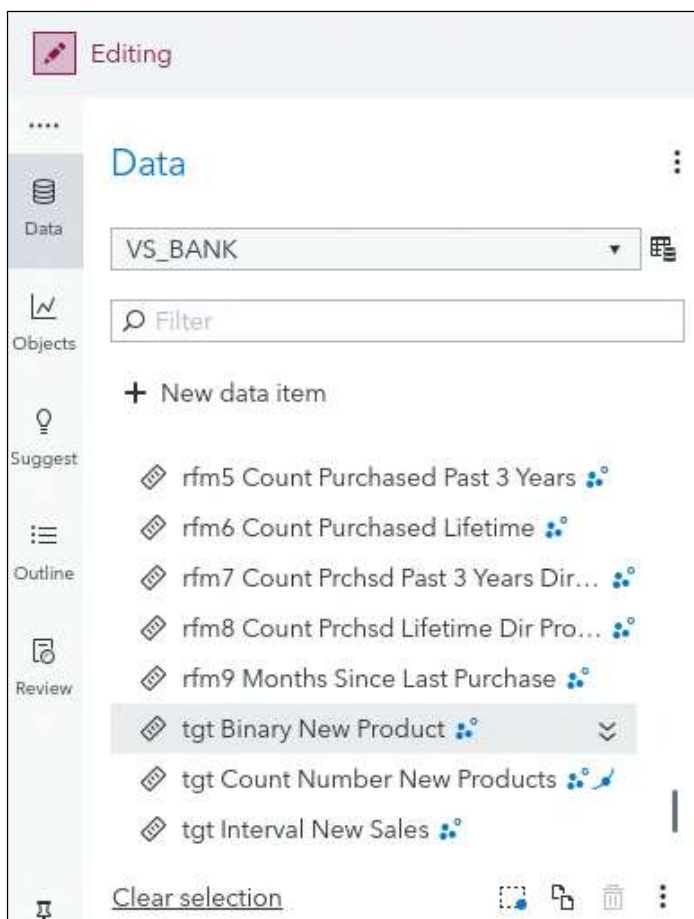


Logistic Regression

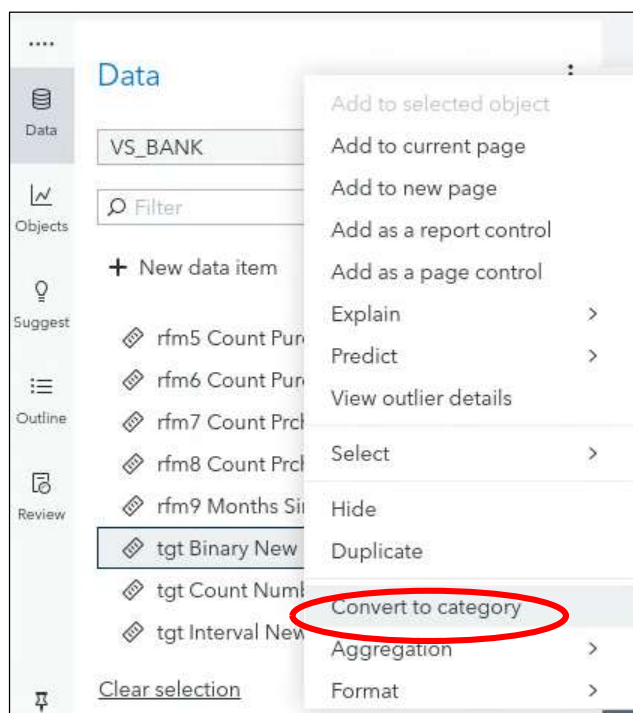
Using the Logistic Regression Object

13. On the Data tab on the left of the screen, scroll down and find the measure **tgt Binary New Product**.

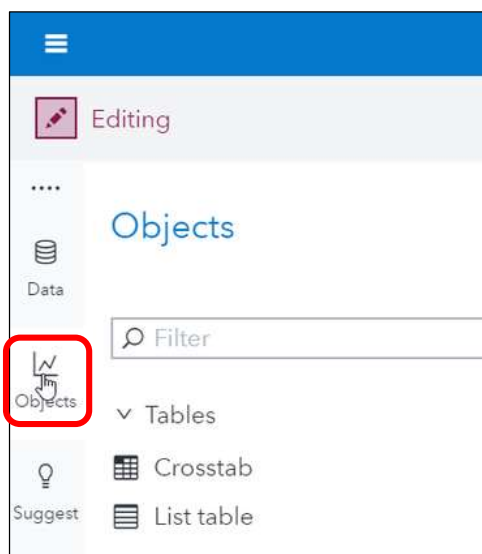
The variable **tgt Binary New Product (b_tgt)** is the primary dependent variable for categorical response modeling in this workshop. It is a binary flag that codes responders with 1 and non-responders with 0. Because it is numeric, it is treated as interval valued by default.



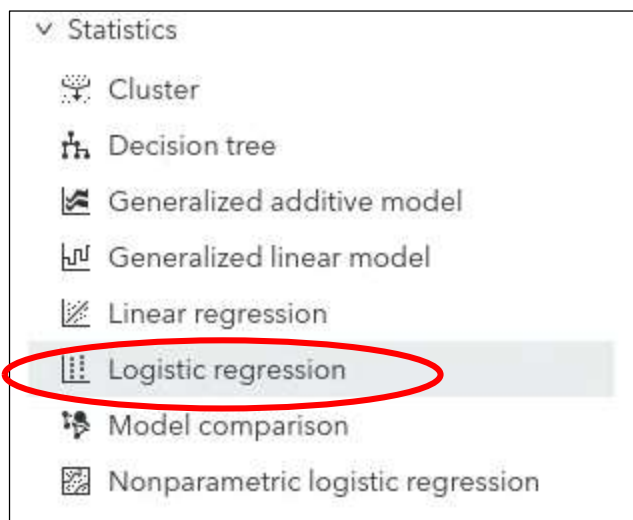
14. Right-click **tgt Binary New Product** and select **Convert to category**.



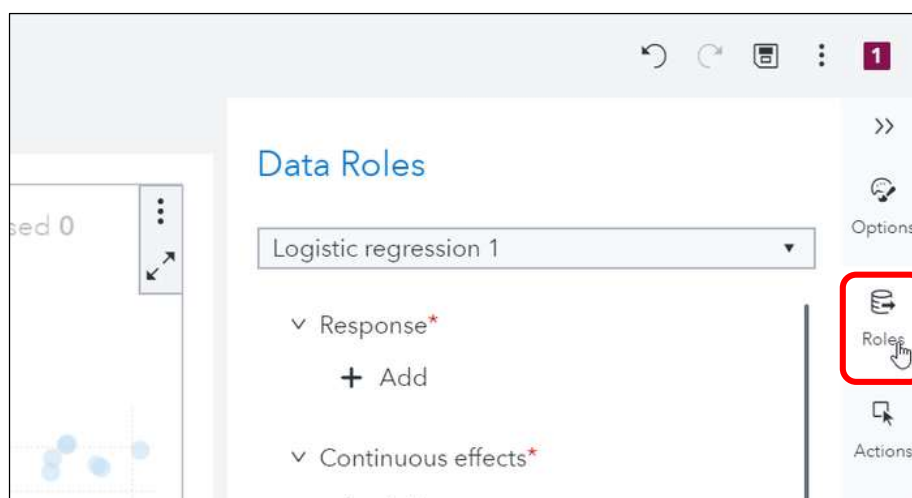
15. Click the **Objects** tab below the Data tab.



16. Scroll down and find the list of statistics.
17. Double-click **Logistic regression** or drag and drop it onto the canvas on the page.

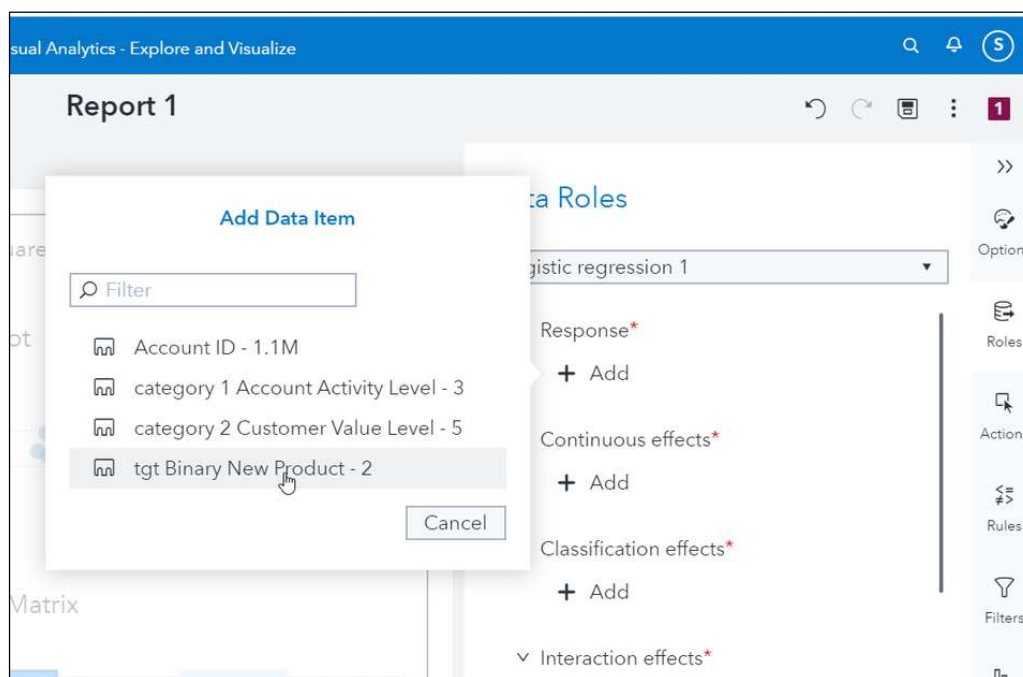


18. On the rightmost side of the screen, click the **Roles** tab.



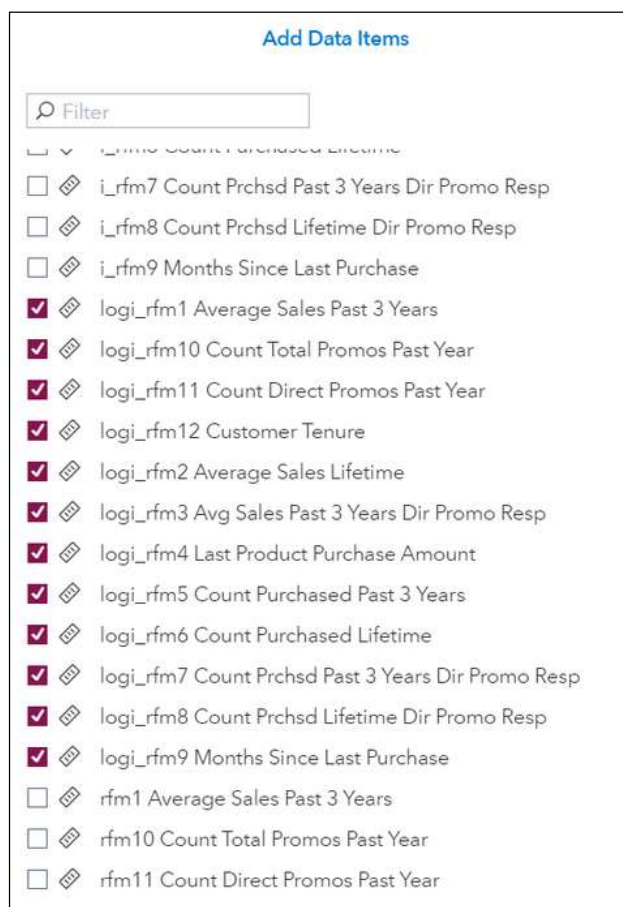
19. Next to **Response**, click **Add**.

20. Select **tgt Binary New Product** from the Add Data Item window.



21. Next to **Continuous effects**, select **Add**.

22. Select all items with a **logi_rfm** prefix.



The screenshot shows the 'Add Data Items' dialog box with a search filter. The following items are selected (checked):

- logi_rfm1 Average Sales Past 3 Years
- logi_rfm10 Count Total Promos Past Year
- logi_rfm11 Count Direct Promos Past Year
- logi_rfm12 Customer Tenure
- logi_rfm2 Average Sales Lifetime
- logi_rfm3 Avg Sales Past 3 Years Dir Promo Resp
- logi_rfm4 Last Product Purchase Amount
- logi_rfm5 Count Purchased Past 3 Years
- logi_rfm6 Count Purchased Lifetime
- logi_rfm7 Count Prchsd Past 3 Years Dir Promo Resp
- logi_rfm8 Count Prchsd Lifetime Dir Promo Resp
- logi_rfm9 Months Since Last Purchase

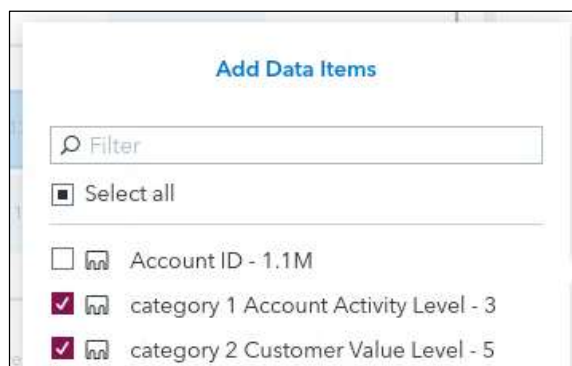
The following items are not selected (unchecked):

- i_rfm7 Count Prchsd Past 3 Years Dir Promo Resp
- i_rfm8 Count Prchsd Lifetime Dir Promo Resp
- i_rfm9 Months Since Last Purchase
- rfm1 Average Sales Past 3 Years
- rfm10 Count Total Promos Past Year
- rfm11 Count Direct Promos Past Year

23. With the 12 data selected items, click **Apply**.

24. Next to **Classification effects**, select **Add**.

25. Click **category 1 Account Activity Level** and **category 2 Customer Value Level**.
Then click **Apply**.



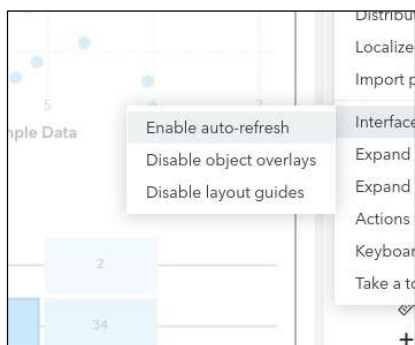
The screenshot shows the 'Add Data Items' dialog box with a search filter. The following items are selected (checked):

- category 1 Account Activity Level - 3
- category 2 Customer Value Level - 5

The following items are not selected (unchecked):

- Account ID - 1.1M

26. Create the logistic regression model by clicking **(More)** and selecting **Interface options > Enable auto-refresh**.



27. Collapse (unpin) the right pane by clicking **Roles** again on the right-side bar to maximize the display of the logistic regression model.

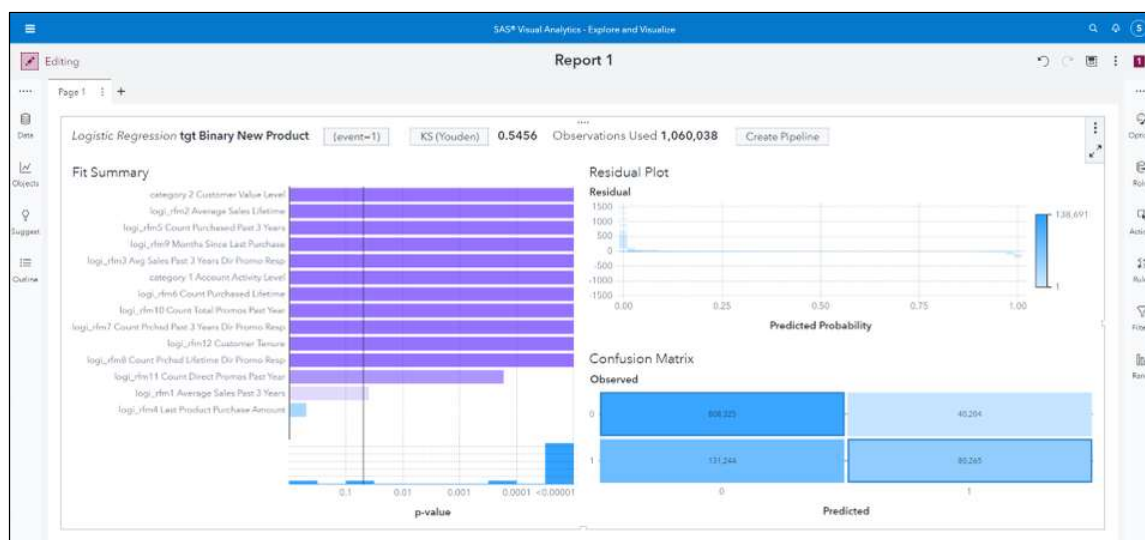
General model information appears along the top of the model, including the name of the response variable, the event of interest, the model evaluation criteria, and the number of observations used to build the model.

For our binary target variable, a value of 1 represents an account that did contract for a product during the campaign (in other words, a sale). As seen in the summary bar, the logistic regression models the event of interest (making a sale) with **event=1**.

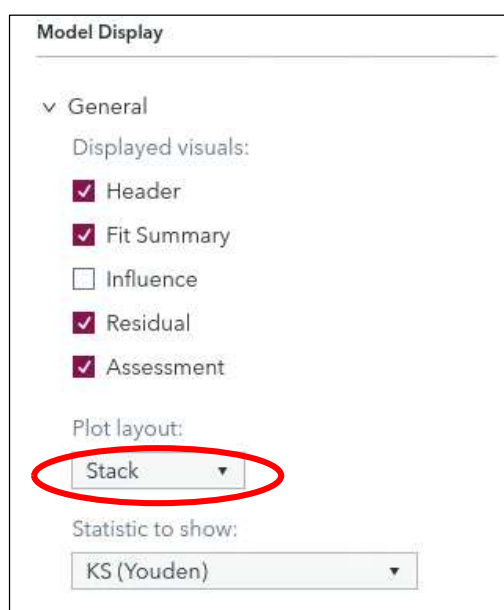
The Fit Summary pane is used to determine the most significant predictor variables that affect the response variable. The variable importance plot displays the effects on the Y axis and the p -values on the X axis. The variable importance is based on the negative log of the p -value ($-\log(p\text{-value})$). A larger $-\log(p\text{-value})$ indicates a more important variable.

The residual plot is used to assess the quality of the model and to identify outlier observations. The plot appears as either a scatter plot for smaller data or as a heat map when used with larger data.

The confusion matrix describes the performance of a classifier. It contains frequency counts of both the correct and incorrect predictions broken down by class.

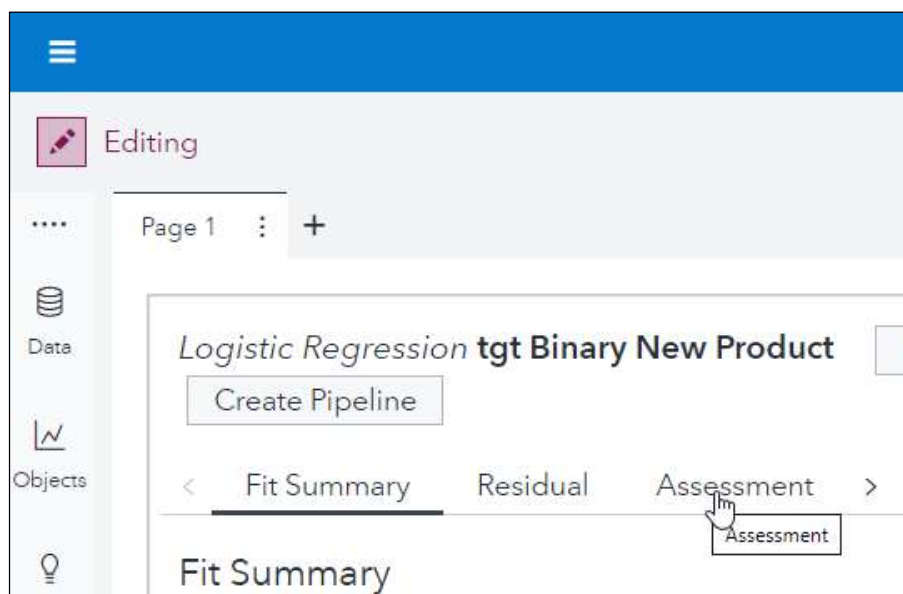


28. On the Options tab on the right side of the screen, under **Model Display**, expand **General** and change **Plot layout** to **Stack**. The model canvas appears, and **Fit Summary** is the default tab selected.



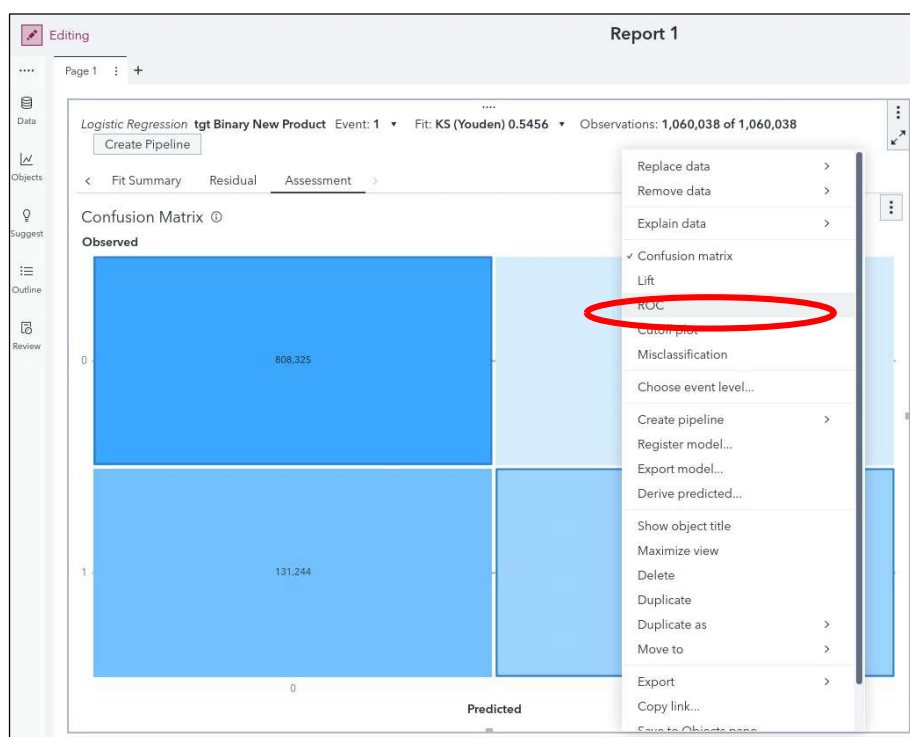
In the Fit Summary pane, one of the variables is not significant at 5%. One of the variables could be considered on the border of significance with a p -value of .0393.

29. Click the **Assessment** tab above the plot on the canvas.

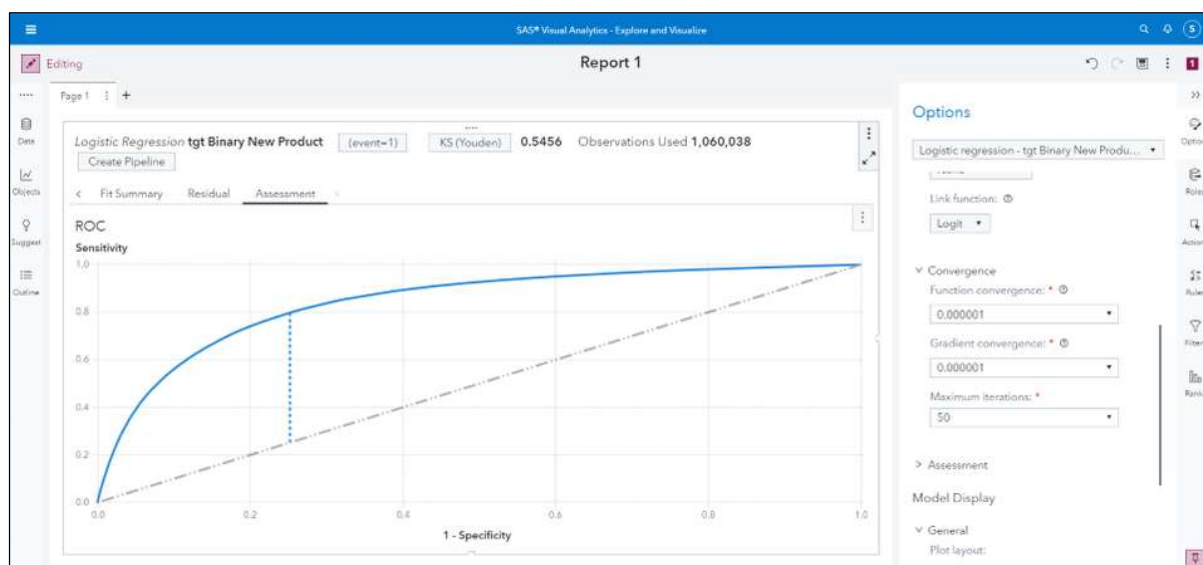


The confusion matrix shows that 80,265 customers who made a purchase were predicted to make a purchase. These are known as true *positives*. It also shows that 40,204 customers were incorrectly classified to have made a purchase (*false positives*).

30. Right-click the confusion matrix and select **ROC**.

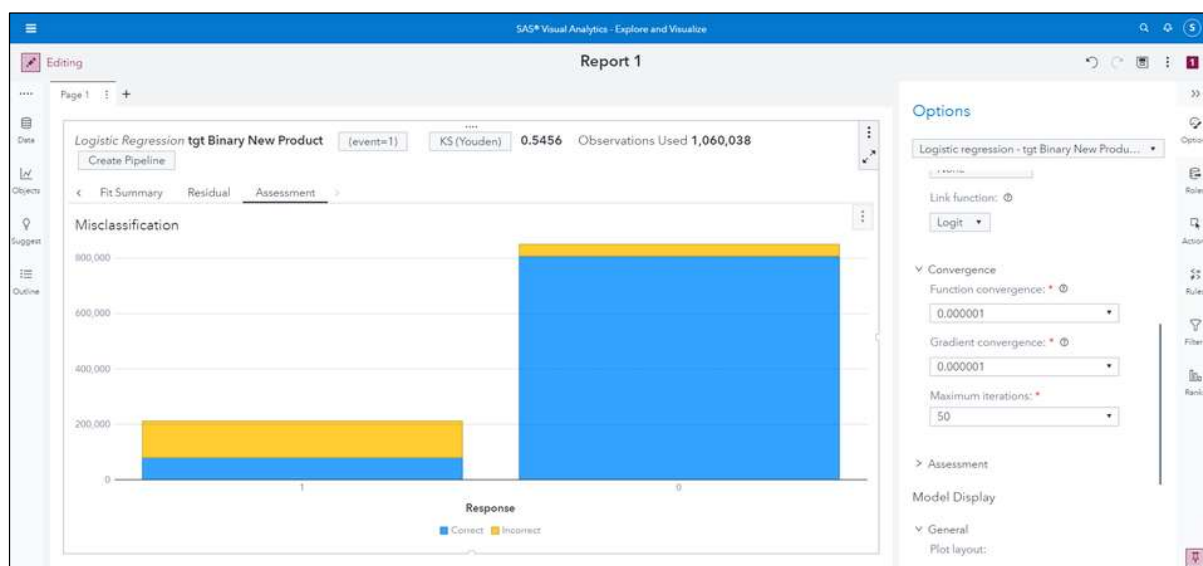


The ROC (receiver operating characteristic) chart is a graphical display that gives a measure of the predictive accuracy of a logistic model. The classification accuracy of a model is demonstrated by the area under the curve or the degree that the ROC curve pushes upward and to the left.



31. Right-click the ROC chart and select **Misclassification**.

A misclassification plot displays how many observations were correctly and incorrectly classified for each value of the response variable. This is a graphical representation of the confusion matrix.



32. On the report, click  (**Maximize**) to see the Details table.

33. Click the **Response Profile** tab to review the original distribution of the target variable.

Model Information			
Dimensions			
Response Profile			
Iteration History			
Convergence			
Fit Statistics			
Parameter Estimates			
Type III Test			
Details			
Ordered Value	Count	tgt Binary New Product	
1	211509	1	
2	848529	0	

34. Click the **Fit Statistics** tab. The Fit Statistics table displays statistics about the estimated model.

Model Information	Dimensions	Response Profile	Iteration History	Convergence	Fit Statistics	Parameter Estimates	Type III Test	Li >	□
Statistic		Value							
-2 Log Likelihood		778997.8							
AIC		779035.8							


35. Click the **Parameter Estimates** tab. The Parameter Estimates table displays the parameter estimates or coefficients of each model effect and their associated statistics.

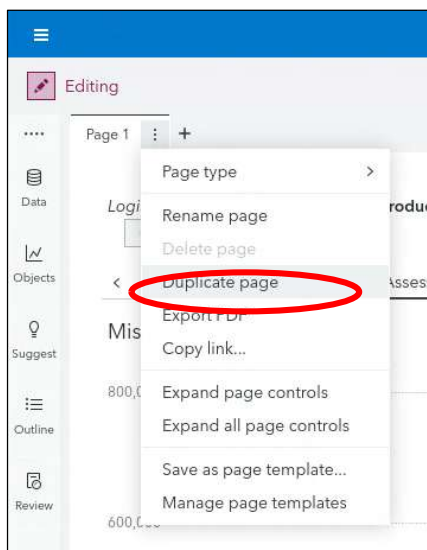
Model Information	Dimensions	Response Profile	Iteration History	Convergence	Fit Statistics	Parameter Estimates	Type III Test	Li >	□
Parameter	Estimate	Standard Error	Chi-Square	Pr > ChiSq					
Intercept	6.802414	0.069145	9678.387	<0.00001					
category 1 Account Activity Level X	0.142799	0.014336	99.22351	<0.00001					
category 1 Account Activity Level Y	0.250099	0.01745	205.425	<0.00001					
category 1 Account Activity Level Z	0	.	.	.					
category 2 Customer Value Level A	1.378718	0.008977	27750.05	<0.00001					

In the upper right corner of the report, click  (**Restore**) from the object toolbar to close the Details table and exit maximize mode.

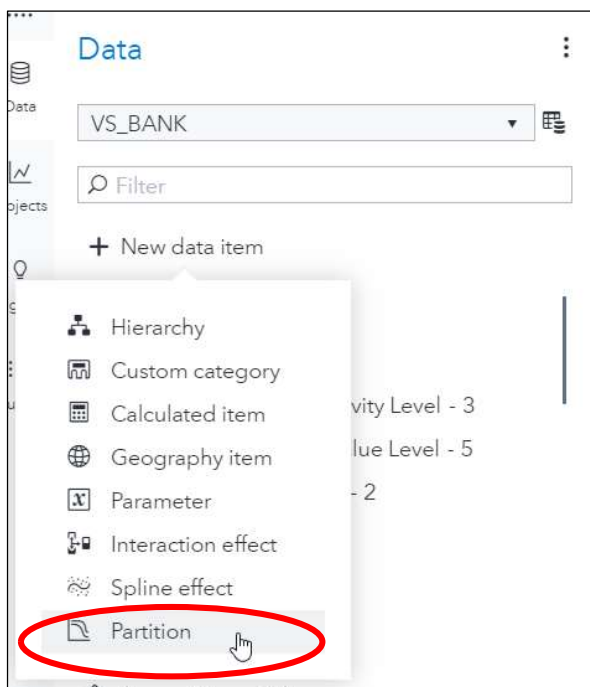
Honest Assessment and Variable Selection

Creating a Partition Variable

1. On the Page 1 tab of the Logistic Regression, click  (Page menu) and select **Duplicate page**.



2. On the Data tab, click **+ New data item** and then select **Partition** from the drop-down list.



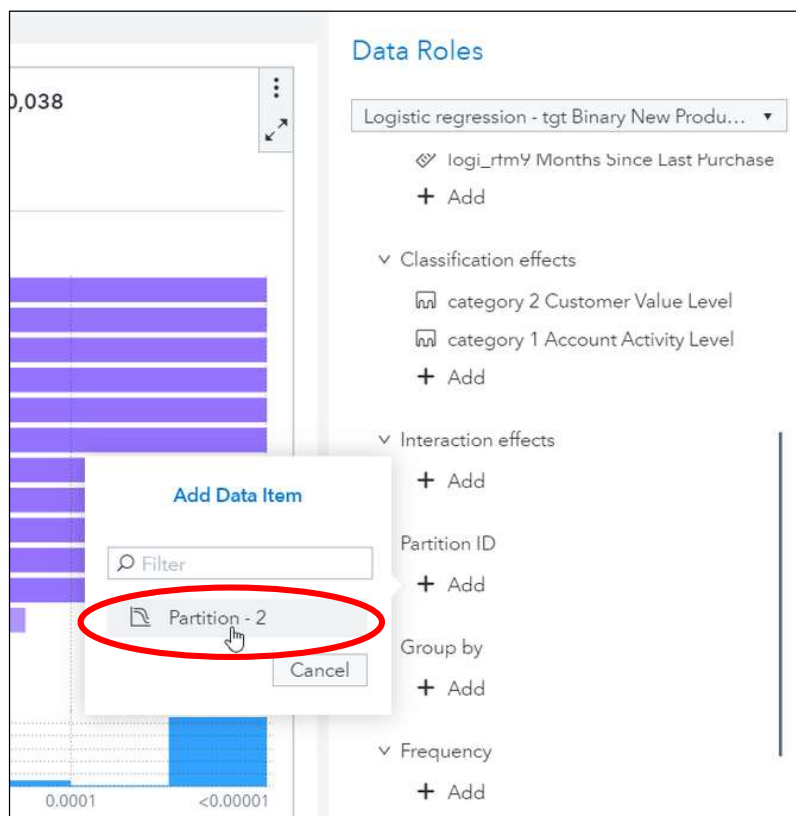
3. Create the partitions.
 - a. In the New Partition window, enter **50** in the **Training** field.
 - b. Select the box for **Random number seed** and enter **1234** in the input area.
 - c. Click **OK**.

The screenshot shows the 'New Partition' dialog box. The 'Name' field is 'Partition'. The 'Based on' section has 'Sampling' selected. The 'Sampling method' is 'Simple random sampling'. The 'Number of partitions' is '2'. The 'Training partition sampling percentage' is '50'. The 'Validation partition sampling percentage' is '50'. The 'Random number seed' checkbox is checked, and the 'Random seed' is '1234'. The 'Use unique identifiers for deterministic partitions' checkbox is unchecked. The 'OK' button is highlighted in blue.

A standard strategy for honest assessment of model performance is data splitting. A portion of the data is reserved for fitting the model, known as the *training data set*. The remaining portion, known as the *validation data set*, is held out for empirical validation. The new partition variable is added to the Category list.

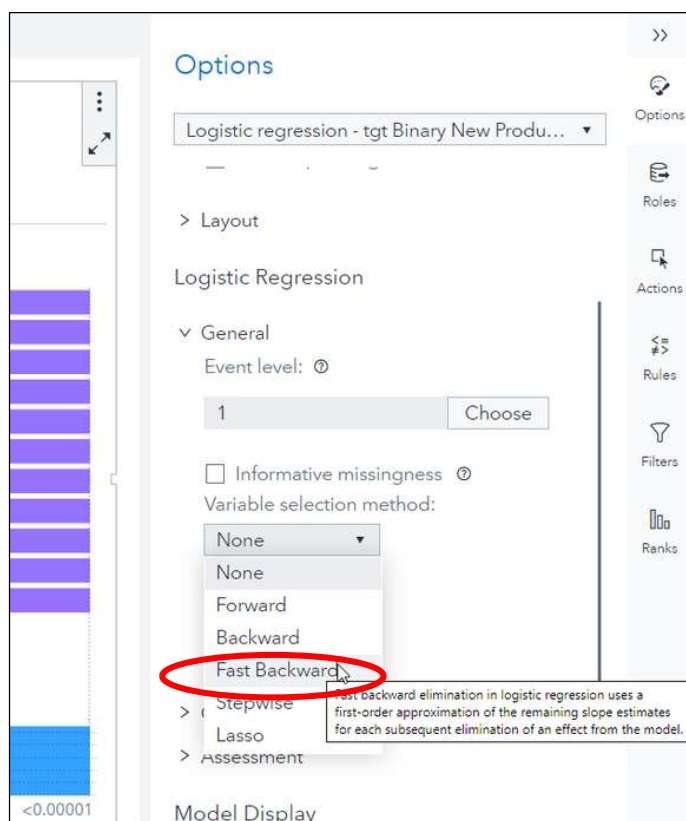
Even with a random seed specified, we might still see nondeterministic results due to the difference in data distribution and computational threads, or the walker used to sample the partition column.

4. Select the duplicated **Logistic regression** on the canvas of Page 1 (1) to make it the active object.
5. On the Roles tab on the right of the screen, scroll down and add **Partition** as the partition ID.



Using Fast Backwards Elimination

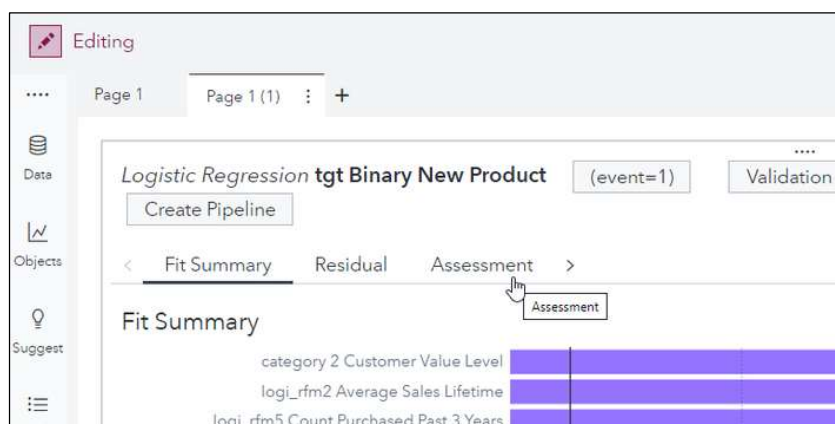
- On the Options tab, select **Fast Backward** for the variable selection method.



- Keep the significance level at **.01**.

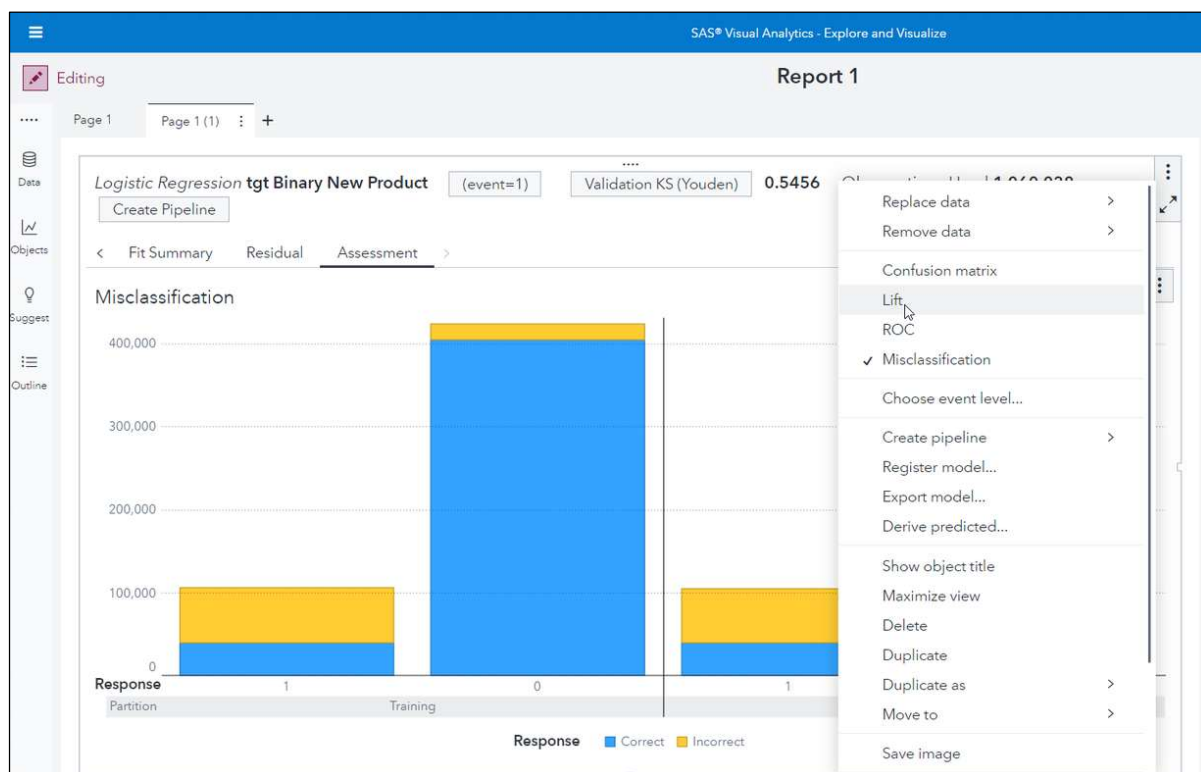
This technique is available for logistic regression models, and it uses a numeric shortcut to compute the next selection iteration quicker than the backward selection method.


- Click the **Assessment** tab on the canvas.



There are now assessment plots for both the training and validation partitions.

- Right-click the misclassification display and select **Lift**.



- On the report, click  (**Maximize**) and collapse the panes on the right to see the Details table.
- In the Details table, click the **Selection Summary** tab to verify that the variables were removed from the model during fast backward variable selection.

Note: You can select  (**More Data Tables**) to navigate through the tables.


Fit Statistics Parameter Estimates Type III Test Selection Info Selection Summary ROC Misclassification Assessment Statistics					
Control	Step	Effect Removed	Number Of Effects	Significance Level	Optimal P Value
	0		15	.	0
-	1	logi_rfm4 Last Product Purchase Amount	14	0.751888	0
	2	logi_rfm1 Average Sales Past 3 Years	13	0.620932	0


12. In the Details table, click the **Assessment Statistics** tab.

Type III Test Selection Info Selection Summary Confusion Matrix Lift ROC Cutoff Statistics Misclassification Assessment Statistics									
Partition	KS (Youden)	Misclassification Rate	Misclassification Rate (Event)	C Statistic	False Positive Rate	False Discovery Rate	F1 Score	Lift	Cun
Training	0.5451	0.1614	0.1614	0.845	0.047	0.332	0.486	3.661	
Validation	0.5463	0.1620	0.1620	0.845	0.048	0.337	0.483	3.661	

In the upper right corner of the report, click  (**Restore**) from the object toolbar to close the Details table and exit maximize mode.

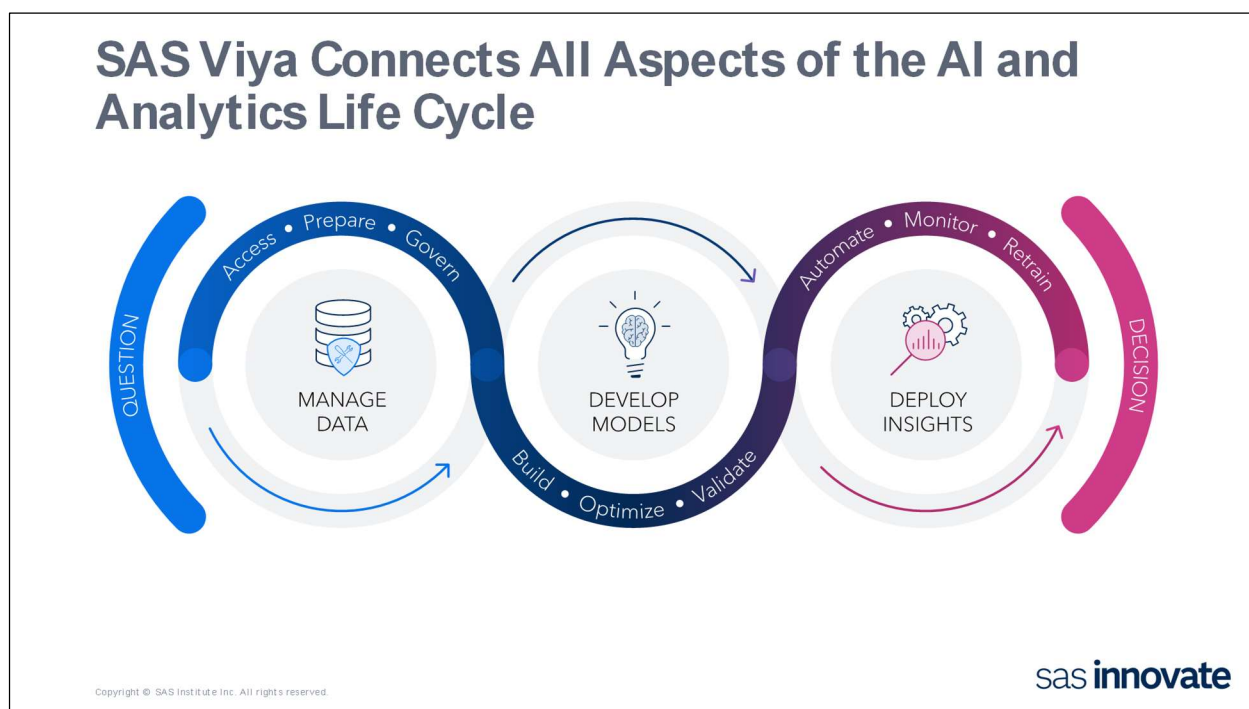
Switching to an Alternative Model

13. On the Page 1(1) tab of the Logistic Regression, click  (**Page menu**) and select **Duplicate page**.

14. From the upper right object menu of the selected Logistic Regression object on Page 1 (1), click  (**More**) and select **Change Logistic Regressions to > Decision Tree**.

Background Information

SAS Viya Information



SAS Viya is a cloud-enabled, in-memory analytics engine that provides quick, accurate, and reliable analytical insights. In SAS Viya, the SAS High-Performance Architecture enables the high-performance analytics engine. The CAS In-Memory Engine continues the ability to perform processing in memory and the ability to distribute processing across nodes in a cluster. The CAS In-Memory Engine adds highly efficient node-to-node communication and uses an algorithm to determine the optimal number of nodes for a given job.

SAS Cloud Analytic Services, or CAS, is a server that provides a cloud-based run-time environment for data management and analytics with SAS. By *run-time environment*, we refer to the combination of hardware and software where data management and analytics take place.

The server can run on a single machine or as a distributed server on multiple machines. The distributed server consists of one controller and one or more workers. This architecture is often referred to as a *massively parallel processing architecture*. For both modes, the server is multi-threaded for high-performance analytics.

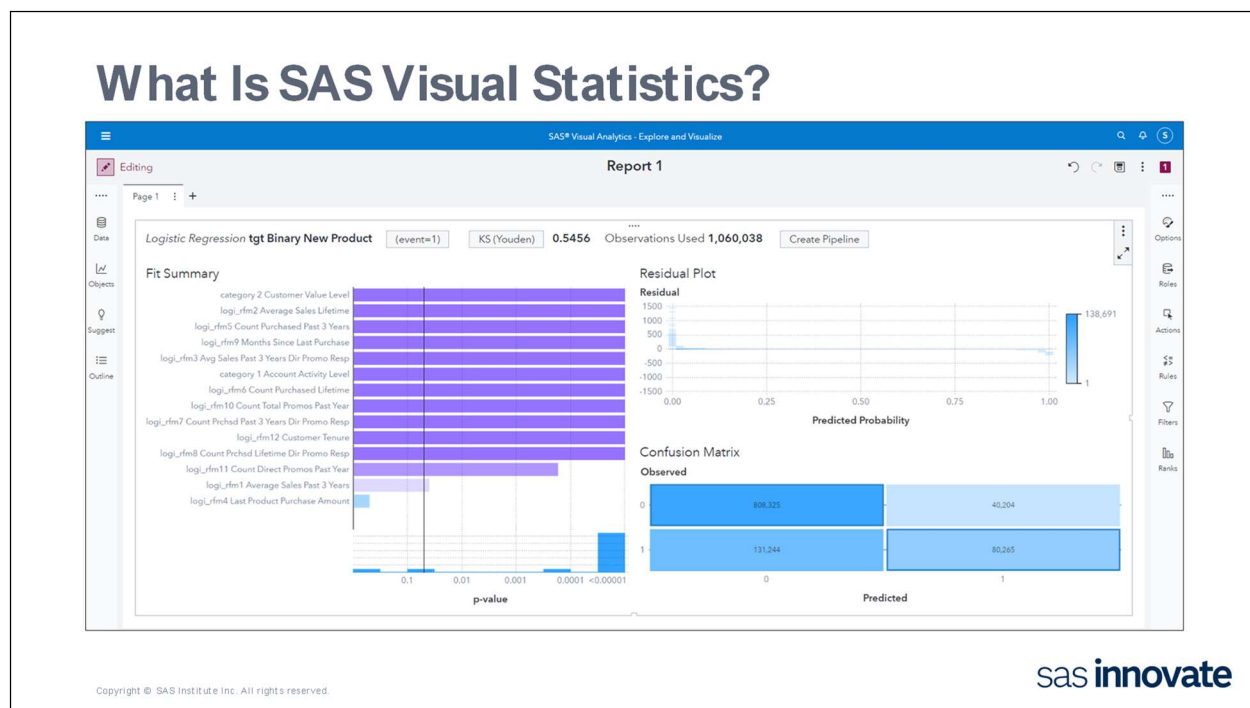
The distributed server has a communication layer that supports fault tolerance. A distributed server can continue processing requests even after losing connectivity to some nodes. The communication layer also enables you to remove or add nodes from a server while it is running.

One of the design principles of the server is to handle large problems and to work with tables that exceed the memory capacity of the environment. To address this principle, data in the server is managed in blocks. Whenever needed, the server caches the blocks on disk. It is this feature that

enables the server to manage memory efficiently, handle large data volumes, and remain responsive to requests.

You can use a variety of interfaces to interact with the CAS In-Memory Engine. These interfaces include SAS Studio, which is a browser-based interface for writing SAS code. You can also use programming interfaces for R, Python, Java, and Lua to access this CAS functionality. In addition, you can continue to submit SAS code in batch mode.

What is Visual Statistics?



SAS Visual Statistics on SAS Viya is an add-on to SAS Visual Analytics that enables you to develop and test models using a scalable in-memory engine in a common SAS Viya environment. SAS Visual Analytics enables you to explore, investigate, and visualize data sources to uncover relevant patterns. SAS Visual Statistics extends these capabilities by creating, testing, and comparing models based on the patterns discovered in analytics and visualizations. SAS Visual Statistics can export the score code, before or after performing model comparison, for use with other SAS products and to put the model into production.

SAS Visual Statistics enables you to rapidly create powerful statistical models in an easy-to-use, web-based interface. After you have created two or more competing models for your data, SAS Visual Statistics provides a model comparison tool. The model comparison tool enables you to evaluate the relative performance of two or more models against each other and to choose a champion model. A wide variety of model selection criteria is available. Regardless of whether you choose to perform a model comparison, you can export model score code for any model that you create. With exported model score code, you can easily apply your model to new data.

The following models are available in SAS Visual Statistics:

- **Linear regression** attempts to predict the value of an interval response as a linear function of one or more effect variables.
- **Logistic regression** attempts to predict the probability that a binary or ordinal response will acquire the event of interest as a function of one or more effects.
- **Nonparametric logistic regression** is an extension of the logistic regression model that allows spline terms to predict a binary response.
- **Generalized linear model** is an extension of a traditional linear model that allows the population mean to depend on a linear predictor through a nonlinear link function.
- **Generalized additive model** is an extension of the generalized linear model that allows spline terms to predict an interval response.
- **Decision tree** creates a hierarchical segmentation of the input data based on a series of rules applied to each observation.
- **Cluster** segments the input data into groups that share similar features.

Logistic Regression Review

Model Essentials: Regressions

▶ Predict new cases.	Prediction formula
▶ Select useful inputs.	Sequential selection

Copyright © SAS Institute Inc. All rights reserved.

sas innovate

Regression models enable you to characterize the relationship between a response variable and one or more predictor variables. With linear regression, the response variable is continuous. With logistic regression, the response variable is categorical. When the response variable is limited to only two categories (dichotomous), the appropriate model is binary logistic regression.

Binary Logistic Models

- Credit Scoring: Can credit score and home ownership predict loan default?

- Predictor Variables:
 - Credit Score: 300-850
 - Home Ownership: Yes/ No/ Rent



- Response Variable:
 - Loan Default: Yes/ No



Copyright © SAS Institute Inc. All rights reserved.

sas innovate

One example of a binary logistic regression model can be found in credit scoring. Can credit score and home ownership help predict the likelihood of a customer defaulting on a loan? In this scenario, one of the predictor variables is continuous (credit score) and the other happens to be categorical (home ownership) with three distinct levels. The response variable is also categorical and coded as character values.

Suggested SAS Courses to Learn More Information

[SAS Visual Analytics Learning Subscription](https://learn.sas.com/totara/program/view.php?id=31)

(<https://learn.sas.com/totara/program/view.php?id=31>)

[SAS Visual Statistics in SAS Viya: Interactive Model Building](https://learn.sas.com/course/view.php?id=445)

(<https://learn.sas.com/course/view.php?id=445>)