

SAS[®] Viya[®] Trial

Manage Data Guide

Data Engineering Tasks



Intro

Data and AI Life Cycle: Manage Data

A recent study by The Futurum Group showed that SAS Viya increases data and AI team productivity by 4.6x.

The analysts compared SAS Viya to alternatives in an end-to-end customer churn prediction analysis, a common use case relevant to many industries.

The first step in the data and AI life cycle is **Manage Data**. This was performed by a **Data Engineer** persona, who was tasked with evaluating and preparing the raw data into an analytical base table that could then be utilized by a Data Scientist persona.

This guide will walk you through the steps that a Data Engineer took to complete the Manage Data portion of the life cycle in SAS Viya.

Data

Data Sets

The data we use for this experiment is simulated data by using SAS code.

We'll use two data sets in CSV (comma delimited).

1. BANKING_ACCOUNT (10,088 rows & 19 columns).
2. BANKING_CUSTOMER (10,095 rows & 39 columns).

The primary key for these data sets is the "Id" column.

Our target variable (what we want to predict) is called "Churn."

The purpose for the Data Engineer is to understand, clean and prepare the data and create a modeling-ready table (aka Analytical Base Table (ABT)), which will be passed to the Data Scientist to predict which customers have the highest probability of churn so the business can take the necessary actions to try to prevent that from happening.

Data Engineer

Clean, Prepare and Govern Data

Tasks

1. Data Profiling
2. Data Sensitivity & PII Checks
3. Data Quality
4. ETL Flow

Resources

Watch before starting

- [Quick Start - Data & AI Life Cycle](#)
- [Quick Start - SAS Drive](#)
- [Quick Start - Manage Data](#)
- [Quick Start - Discover Information Assets](#)
- [Quick Start - Develop SAS Code in SAS Studio](#)
- [Quick Start - Develop Flows in SAS Studio](#)
- [Webinar - SAS Studio Flow & Steps](#)

The Data Engineer

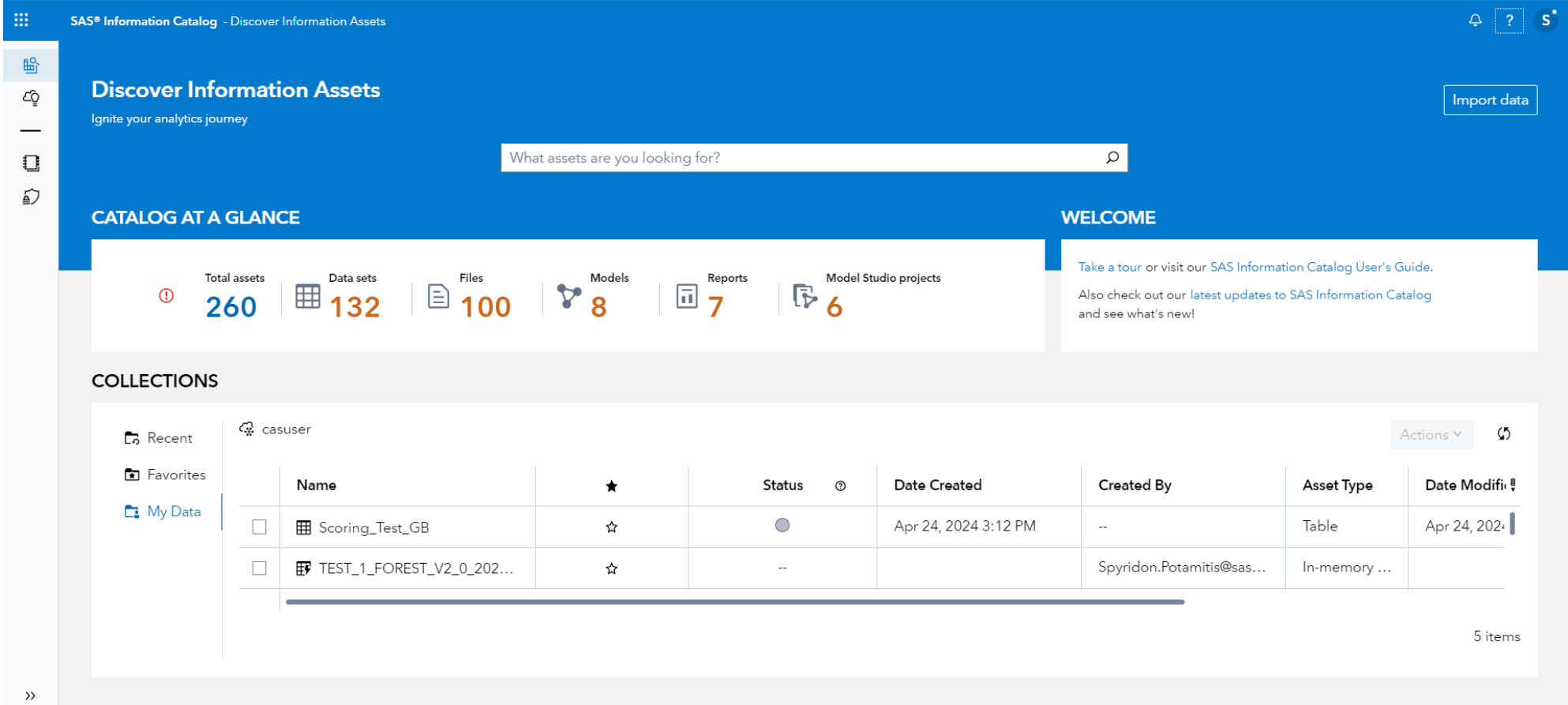
- As a Data Engineer, it is your job to upload a new data set for analysis and perform some basic data quality measures to ensure the data is ready for further use and modeling. It is also essential that the data doesn't include any personally identifiable information (PII) since this project will be accessed by many different parties who shouldn't have access to sensitive information, and you want to mitigate any risks and be regulatory compliant.
- You will work on uploading the data, profiling it to examine data quality and sensitivity and creating an ETL (Extract, Transform, Load) flow to prepare the data for modeling. The flow will create a unified data set that the Data Scientist needs to develop AI models.

Data Profiling

Information Catalog

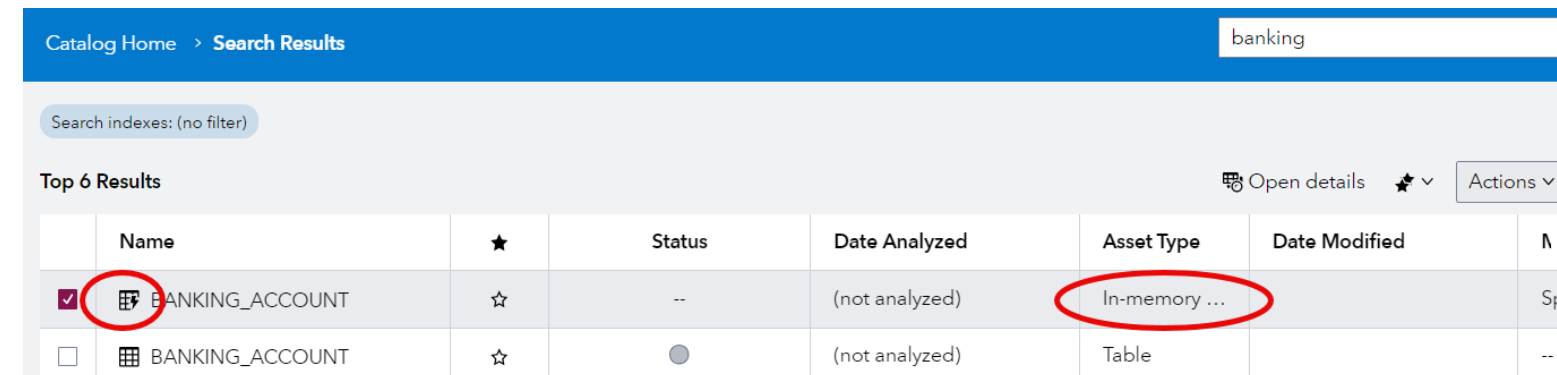
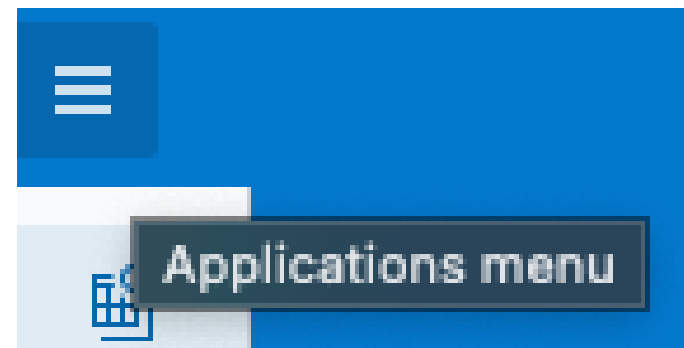
Intro

- To get more information on the data quality, we should navigate to Information Catalog and run a data profile. This will get us privacy information and descriptive statistics about the quality of our data.
- SAS Information Catalog lets you create and maintain an inventory of your information assets. Such a catalog gives you the ability to ingest, integrate and enrich metadata from the assets that are distributed across your enterprise. You can use this metadata to find and understand the relevant assets that you need to reach your business goals. An information catalog also enables data administrators to review data usage, such as when the data was created, who created the data, who modified it and when the data was modified for the last time, from a single point of access.



Data Profiling



- Select the Applications Menu (on the left of the screen) and go to “DISCOVER INFORMATION ASSETS.” From there, you can search for and select the data sets we’ve just uploaded. Filter for data sets or type “banking” in the search tab to make searching easier.
- You’ll see the names of data sets appearing twice. The ones with the thunderbolt next to their table name are already loaded into memory. We’ll only use in-memory data sets.

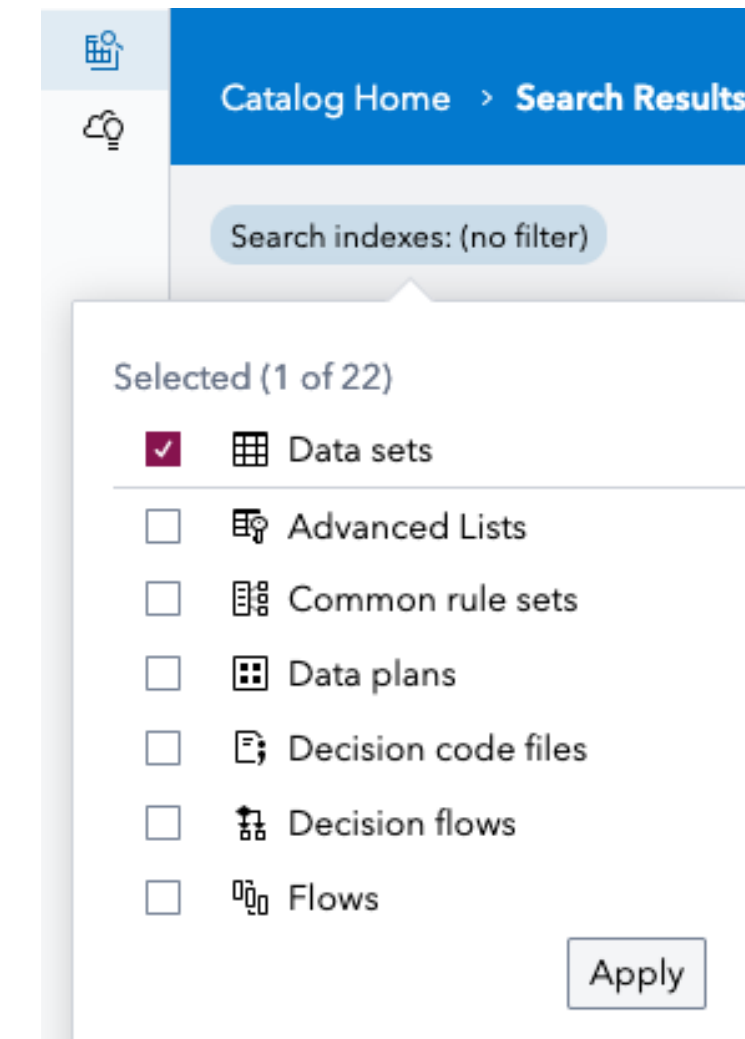


Catalog Home > Search Results banking

Search indexes: (no filter)

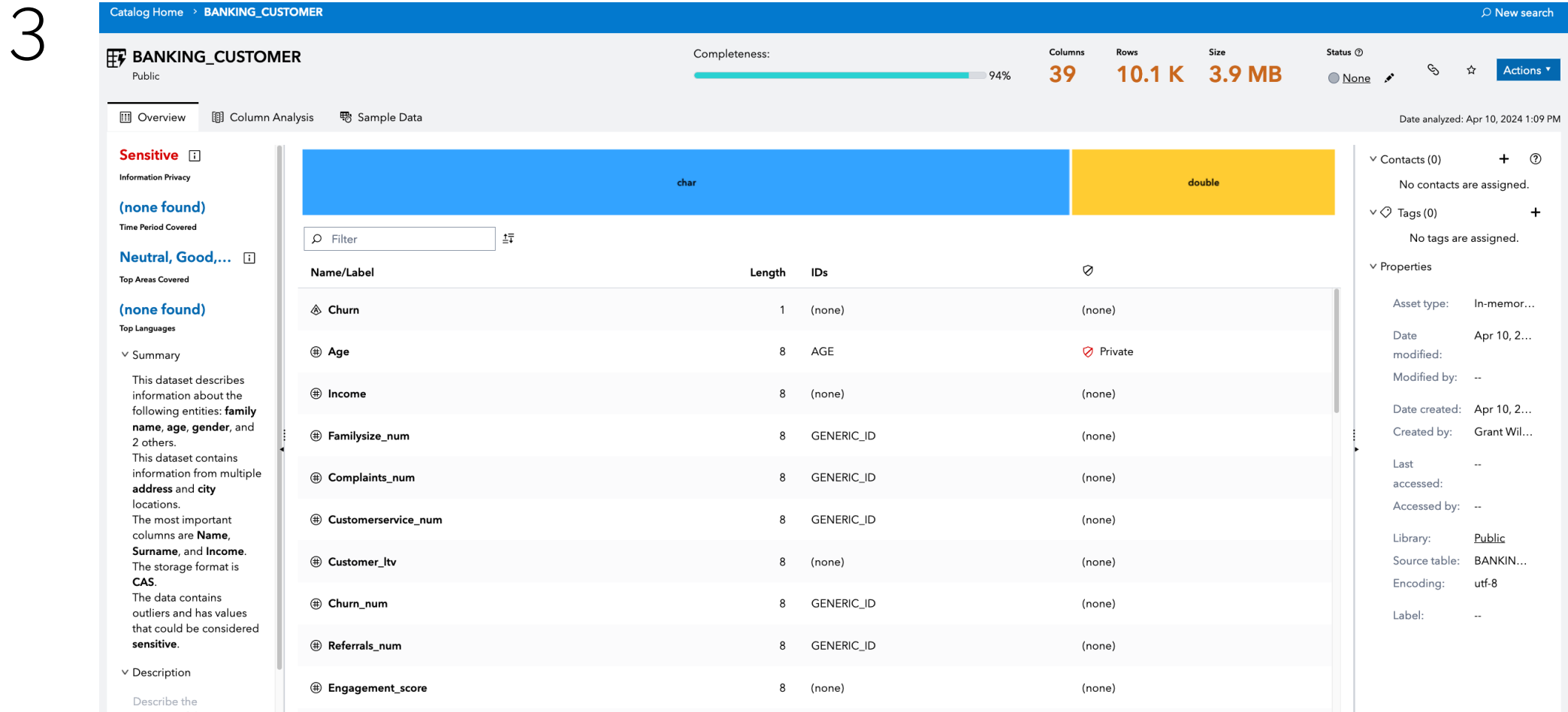
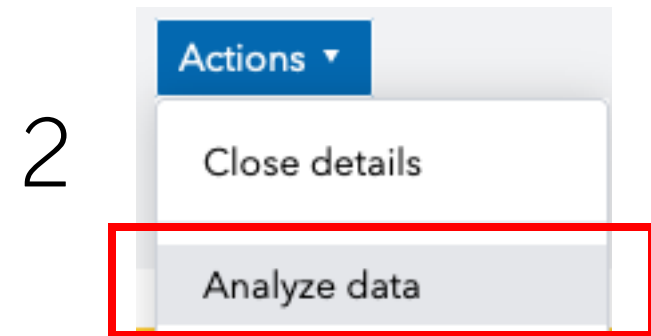
Top 6 Results Open details Actions

Name	★	Status	Date Analyzed	Asset Type	Date Modified	⌵
<input checked="" type="checkbox"/>  BANKING_ACCOUNT	☆	--	(not analyzed)	In-memory ...		Spy
<input type="checkbox"/>  BANKING_ACCOUNT	☆	●	(not analyzed)	Table		--



Data Profiling

- Let's start by selecting the BANKING_CUSTOMER data set. Navigate to the ACTIONS dropdown from the right of the screen and select ANALYZE DATA to run a data profile on the data set.
- Once complete, click on "View analysis" near the top of the screen, and you will be able to view descriptive statistics about your data set. Let's dive deeper.



Data Sensitivity & PII Checks

View Column-Level Data Sensitivity

- In the middle of the screen, we see descriptive information about our column-level data sensitivity. On the left of the screen, you should also see a summary of the findings in natural language where sensitive data and outliers are detected. Both descriptions are automatically generated by SAS Viya. The information privacy is broken into four categories (in this data set you should see all categories). Based on this information, the Data Engineer should take the necessary steps to safeguard customers' personal and sensitive data as well as investigate and treat outliers if necessary.

- None
 - There is no PII in this column.
- Candidate
 - There may be PII in this column.
 - Location, Market Condition.
- Private
 - This column contains PII.
 - Name, Surname, Age, Gender.
- Sensitive
 - There is sensitive PII in this column.
 - Marital Status.

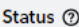

The screenshot shows the SAS Viya interface for the 'BANKING_CUSTOMER' dataset. The top navigation bar includes 'Catalog Home > Search Results > BANKING_CUSTOMER' and a 'New search' button. The main header displays 'BANKING_CUSTOMER' with a 'Public' status, a 'Completeness' bar at 94%, and statistics for 'Columns: 39', 'Rows: 10.1 K', and 'Size: 6 MB'. Below the header, there are tabs for 'Overview', 'Column Analysis', and 'Sample Data'. The 'Column Analysis' tab is active, showing a table of columns with their respective data types and sensitivity levels. The table has columns for 'Name/Label', 'Length', 'Semantic Type', 'Information Privacy', and 'Terms'. The 'Information Privacy' column is highlighted with a red box, showing 'Sensitive' for 'Gender' and 'Marital_status', and 'Private' for 'Age' and 'Education'. To the left of the table, a summary box (also highlighted with a red box) provides a natural language description of the data, mentioning entities like 'age', 'gender', 'marital status', and 'address', and noting that the data contains outliers and is considered sensitive. On the right side of the interface, there is a metadata section with details like 'Created: Apr 22, 2024 2:49 PM', 'Created by: Spyridon.Potamitis@sas.cc', and 'Accessed: Apr 22, 2024 3:39 PM'.

Name/Label	Length	Semantic Type	Information Privacy	Terms
Gender	6	Gender	Sensitive	(none)
Marital_status	9	Marital status	Sensitive	(none)
Age	8	Age	Private	(none)
Education	11	Family name	Private	(none)

Flag Data Set

- Now that we know more about the data set, we want to let our Data Engineering colleagues know that they shouldn't share this data set with other teams that shouldn't have access to customers' personal and sensitive data. To do that, we'll set a warning for this data set by clicking on the pencil icon under the "Status" tab. In there, we'll also write a comment that "personal and sensitive data was detected." This flag will be visible to everyone who has access to the data set in all SAS Viya applications. For example, if another Data Engineer wants to use SAS Viya applications such as Manage Data, Discover Information Assets or even Explore Lineage, this important info will be available to ensure proper governance of our data assets throughout the organization.

Catalog Home > Search Results > BANKING_CUSTOMER New search

BANKING_CUSTOMER Public Completeness: 94% Columns: 39 Rows: 10.1 K Size: 6 MB Status  None  Actions Date analyzed: Apr 22, 2024 3:39 PM

Overview | Column Analysis | Sample Data

(none found)

Top Languages

Summary

This dataset describes information about the following entities: **age, gender, marital status**, and 2 others. This dataset contains information from multiple **address** and **united states city** locations. The most important columns are **Surname, Name**, and **Income**. The storage format is **CAS**. The data contains outliers and has values that could be considered **sensitive**.

Name/Label	Length	Semantic Type	Information Privacy	Terms
Gender	6	Gender	Sensitive	(none)
Marital_status	9	Marital status	Sensitive	(none)
Age	8	Age	Private	(none)
Education	11	Family name	Private	(none)

Metadata:

- Created: Apr 22, 2024 2:49 PM
- Created by: Spyridon.Potamitis@sas.cc
- Last accessed: Apr 22, 2024 3:39 PM
- Last accessed by: Spyridon.Potamitis@sas.cc
- Library: Public
- Source: BANKING_CUSTOMER.sa
- Encoding: utf-8
- Label: --

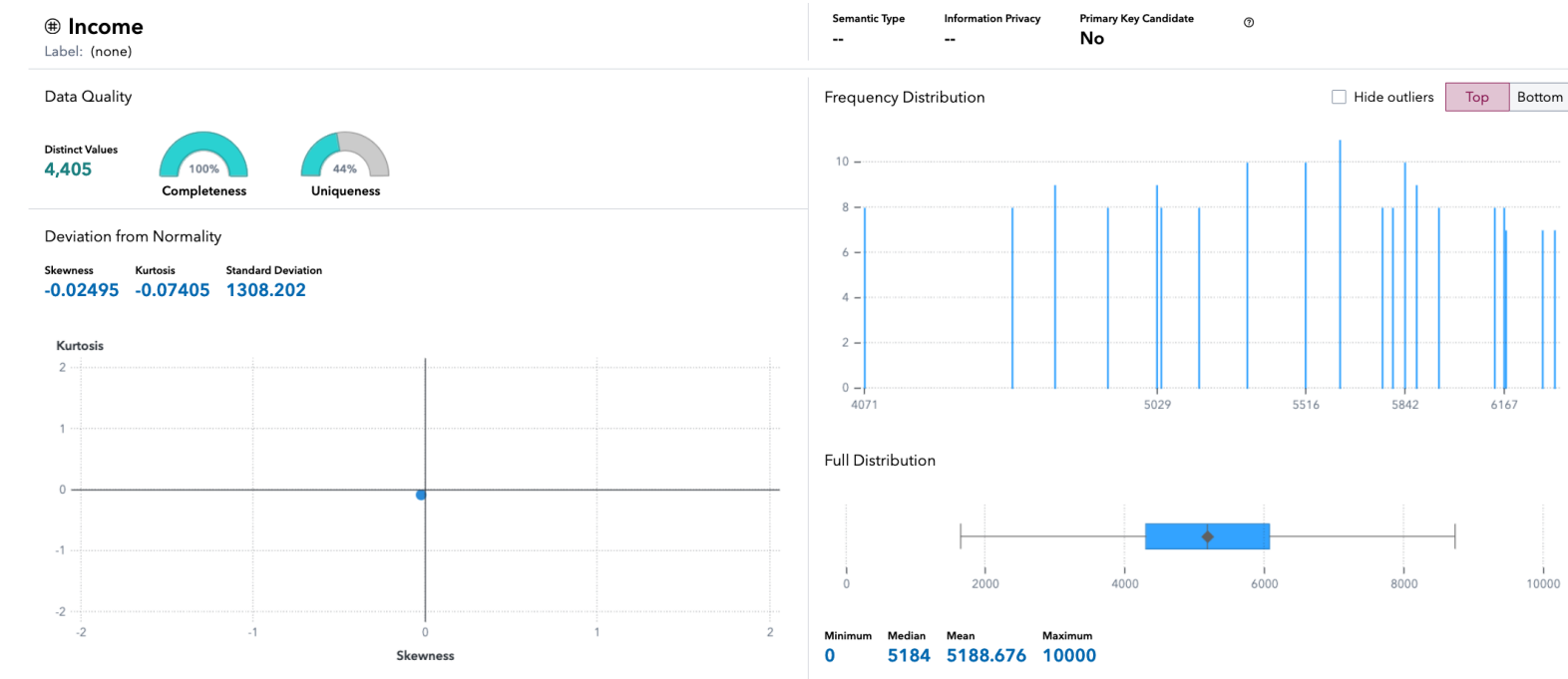


Data Quality

Get Data Quality Measures

Click a column name to view Data Quality Measures

- From the data profile, we can drill into the variables to get some metrics for data quality. Select the name of one of the numeric variables or click “Column Analysis.”
 - Skewness
 - Kurtosis
 - Completeness
 - Uniqueness
 - Distinct Values
- For numeric variables, you will see metrics. For character variables, you will see patterns.
- Very quickly we can identify that our primary key “ID” has duplicate values by examining the uniqueness indicator.



#	Name	Completeness	Uniqueness	Most Common Value	Least Common Value	Pattern Cour
32	Product_...	100%	0%	7	8	
33	Survey_r...	100%	0%	No	Yes	
34	Custome...	75%	0%	Good	Bad	
35	Cross_se...	100%	0%	No	Yes	
36	Market_c...	100%	0%	Neutral	Bad	
37	Name	100%	1%	Liam	Dylan	
38	Surname	100%	1%	Smith	Harrison (1 more)	
39	Id	100%	99%	9856	1	

Column Analysis

Data Quality Measures

- On the Column Analysis page, we can view Descriptive Measures, Metadata Measures and Data Quality Measures.
- First click the Column Analysis and run it. You'll be able to see the analysis by clicking the "View the analysis" button when the run is complete.
- Here, we can view Descriptive Measures, Metadata Measures and Data Quality measures to get a comprehensive view of our data.
- Opening the Data Quality Measures page gives us the following information:
 - Completeness
 - Uniqueness
 - Most and Least Common Values
 - Pattern Count
 - Semantic Type
 - Information Privacy

SAS® Information Catalog - Discover Information Assets

Catalog Home > Search Results > BANKING_CUSTOMER

A request for analysis is pending on this asset. Requested by Spyridon.Potamitis@sas.com Apr 30, 2024 12:04 PM. [Refresh view](#)

BANKING_CUSTOMER
Public

Completeness:
(Not available)

Overview **Column Analysis** Sample Data

BANKING_CUSTOMER
Public

Completeness: 94%

Columns: **39** Rows: **10.1 K** Size: **6 MB**

Overview **Column Analysis** Sample Data

Filter

Descriptive Measures Metadata Measures **Data Quality Measures**

# ↑	Name		Most Common Value	Least Common Value	Pattern Count	Semantic Type	Inform...
1	Churn	0%	0	1	--	(none)	None
2	Age	1%	47	0	--	Age	Private
3	Income	4%	5629	0	--	(none)	None
4	Family...	0%	3	6	--	Generic ID	None
5	Compl...	0%	9	20	--	Generic ID	None
6	Custo...	0%	10	0	--	Generic ID	None
7	Custo...	1%	57	0	--	(none)	None
8	Churn n...	0%	1	0	--	Generic ID	None

Column Analysis

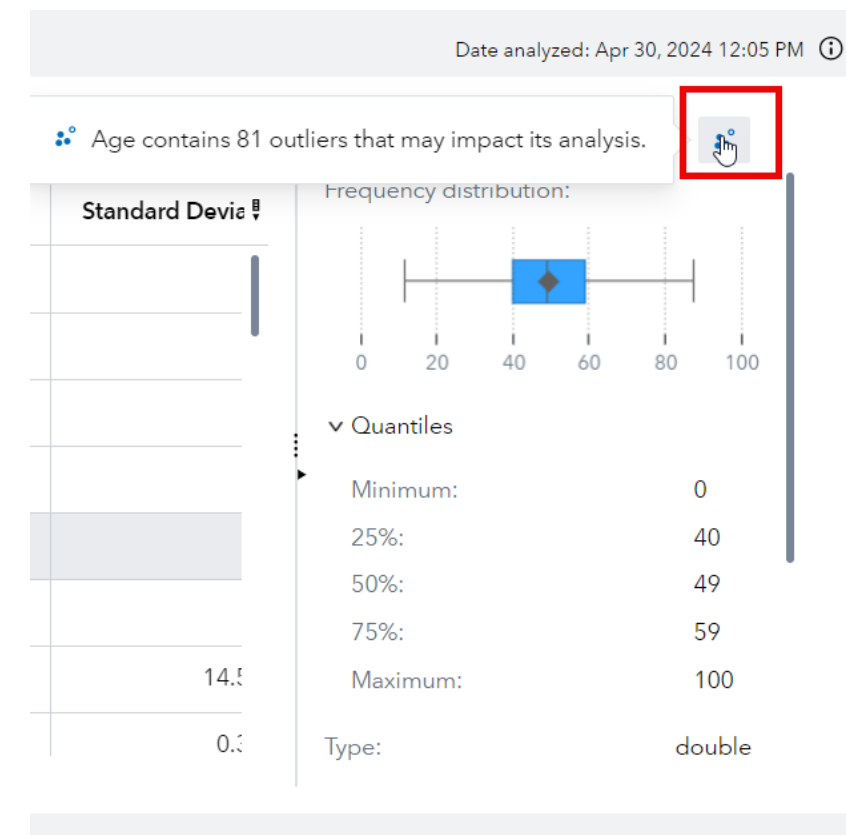
Descriptive Measures

- Let's open the "Descriptive Measures" tab for our BANKING_CUSTOMER data set.
- If we sort by our "Missing" column (by clicking on "Missing" column name), we can identify which columns have missing values. We can quickly notice a few that we'll need to assess in our data flow:
 - Loyalty Program
 - Life Event Marriage
 - Social Media Usage
 - Financial Literacy
 - Customer Sentiment
 - Digital Usage
 - Age
- When we click next to a variable name, we'll see a distribution graph, and the quantiles appear as well. Selecting the "Age" variable, for example, will give us a box and whisker plot describing the values of our selected variable and measures of central tendency. We'll also get a small icon on the top right that when we double-click on it will tell us how many outliers we have in this variable.

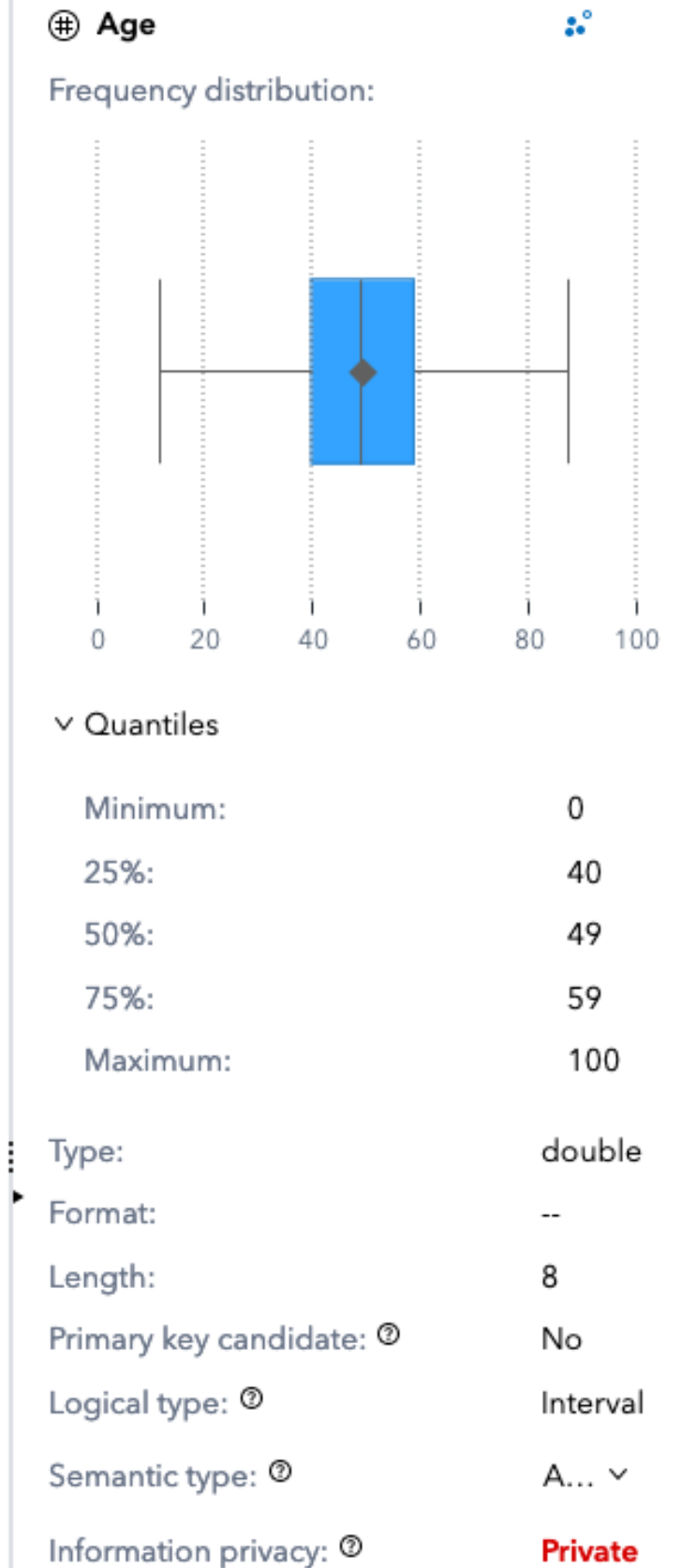
1

Name	Missing ↓
Loyalty_program	5,615
Life_event_marriage	4,972
Socialmedia_usage	3,287
Financial_literacy	2,463
Customer_sentiment	2,446
Digital_usage	1,127
Age	192

2



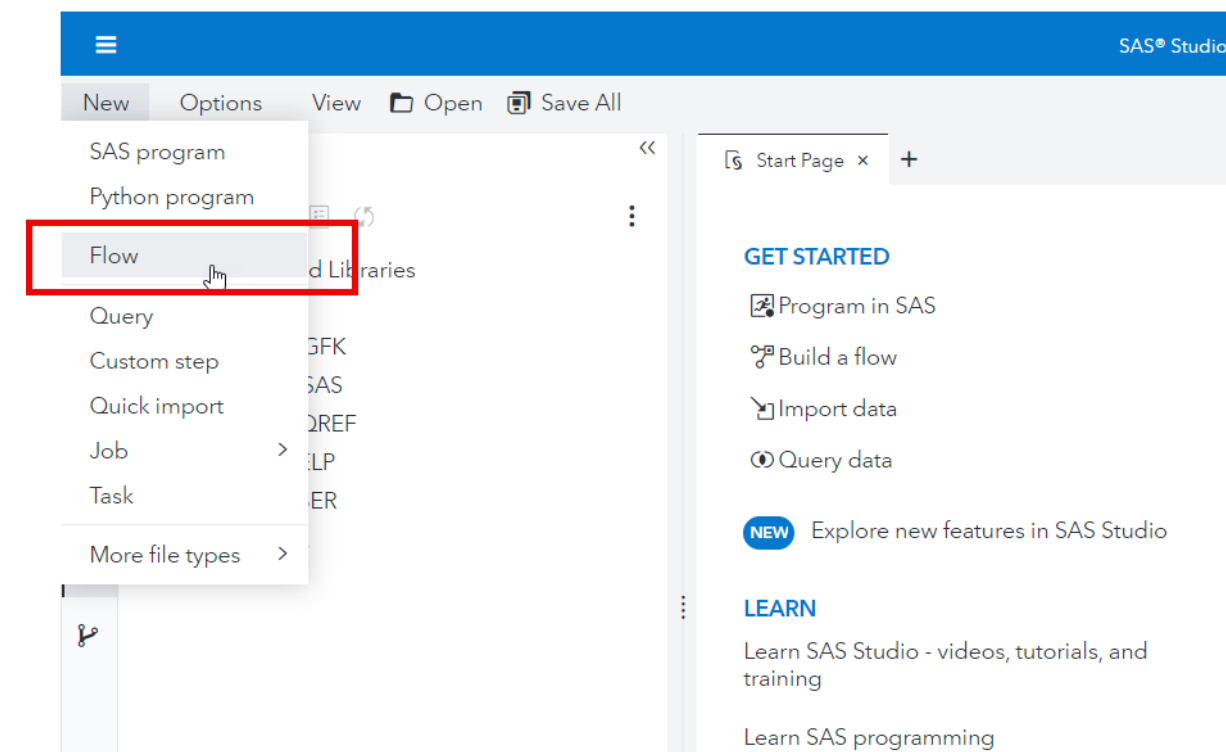
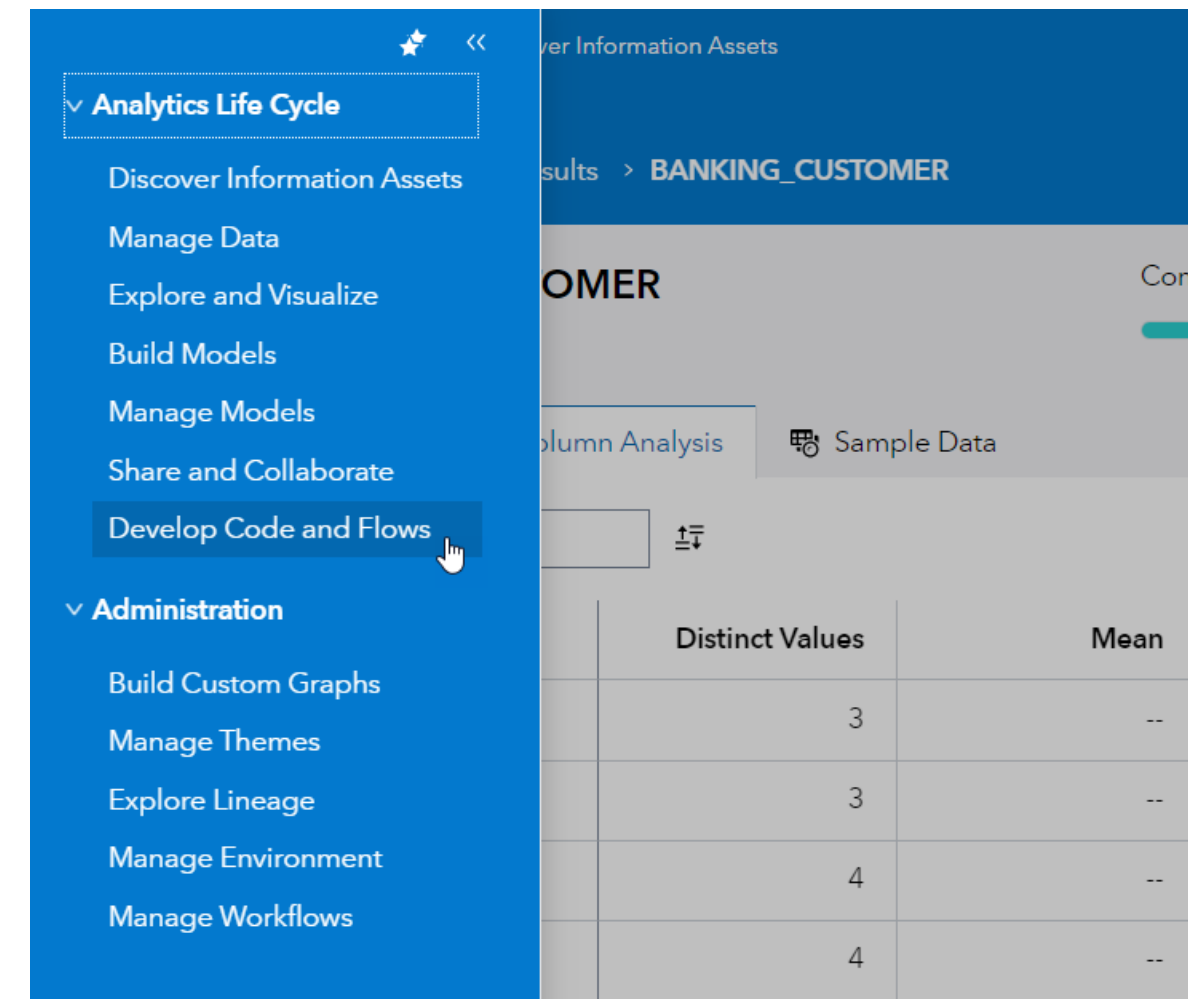
3



ETL Flow

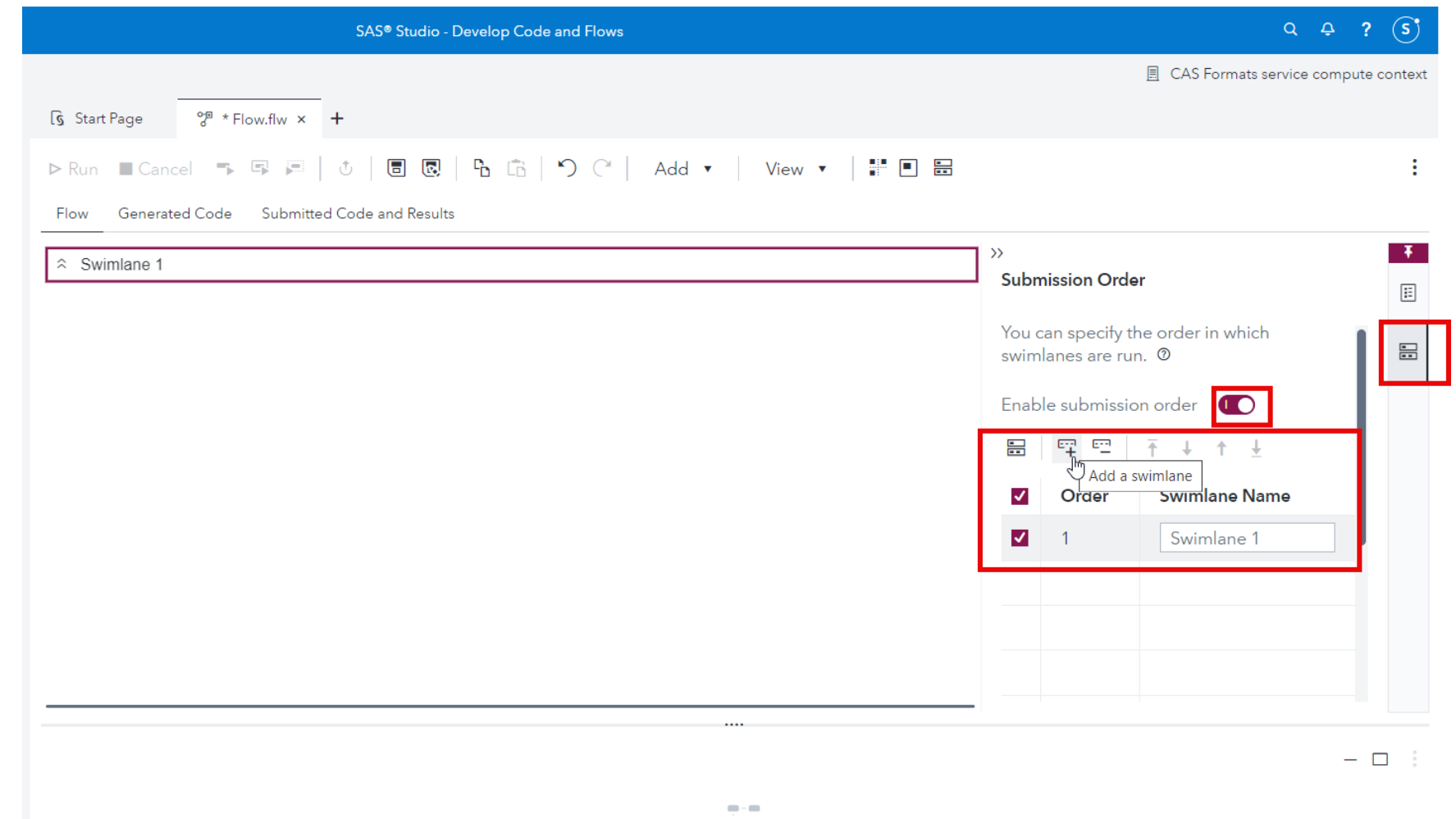
Create a Flow

- Now that we've had a chance to explore our imported data sets, let's build a custom flow to create an ETL (Extract, Transform, Load) workflow that will be easily repeatable and easy to implement in a low-code/no-code environment to automate repetitive processes.
- The end goal of our ETL process is to develop an ABT table, which is ready for the Data Scientist to use to develop models.
- Open the Applications Menu (top left of the screen) and select "DEVELOP CODE AND FLOWS."
- Open the "New" dropdown and select "Flow." This will serve as the canvas for our ETL Flow.



Create a Flow

- In examining data quality, you've noticed duplicates in the data in addition to some data sensitivity issues. We will need to clean each table and join them together.
- Start by ensuring swimlanes are enabled by opening the pane on the right and clicking the relevant button. This will enable us to run parts of our code separately and ensure we're not rerunning code needlessly. This also guarantees that we get a logical sequence in the flow with the necessary dependencies enforced.
- Click the "Add a swimlane button" as you see on the right to add a swimlane in the flow. Now you should have two swimlanes in total.



Connect to CAS

- As a first step in our flow, we need to connect to Cloud Analytic Services (CAS), which is a SAS Viya in-memory engine, and start a session.
- Let's begin by dragging a pre-built SAS Code node referred to as a "Snippet" onto our canvas. Code snippets are lines of commonly used code that you can use in an existing program or as the basis for a new program. Code snippets enable you to quickly insert SAS code into your program and customize it to meet your needs.
- Find the "Generate SAS librefs for caslibs" Snippet and drag it to Swimlane 1.
- If you click on the node you created in the flow, you will be able to see the code that is generated automatically for you.
- Run the node by right-clicking it and clicking "Run node."

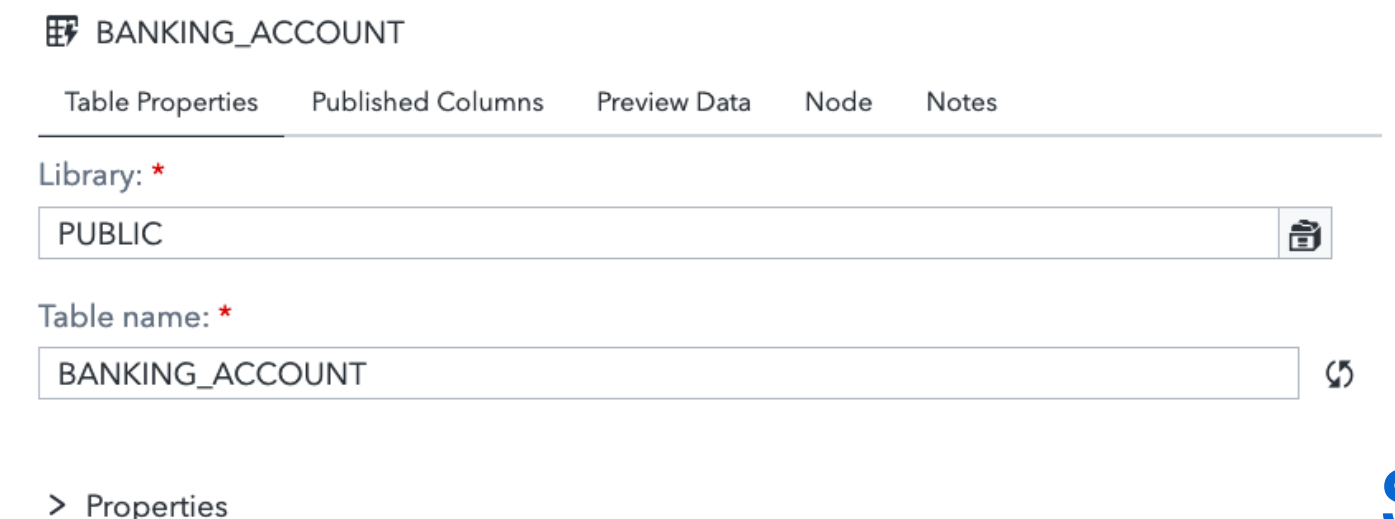
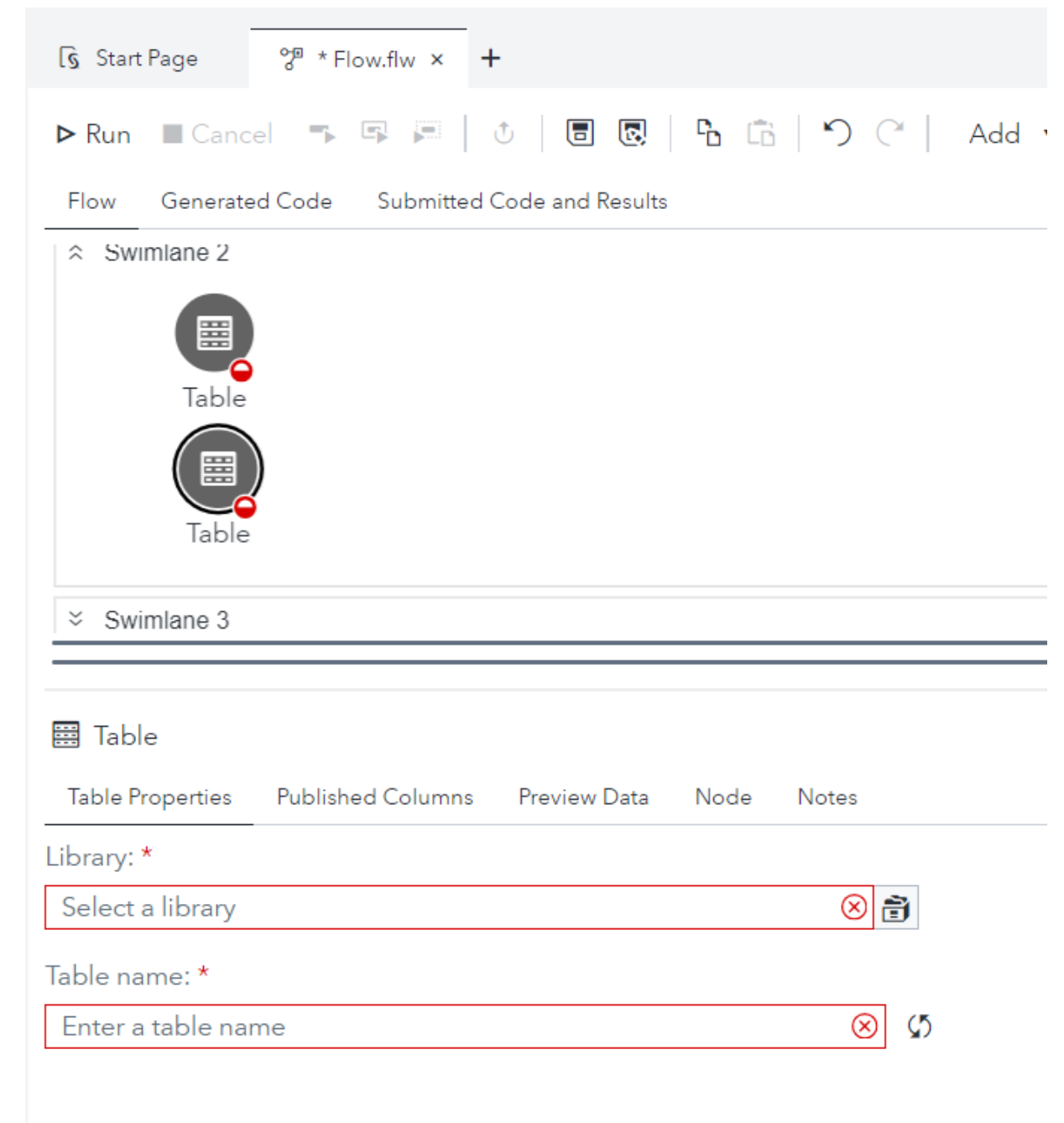
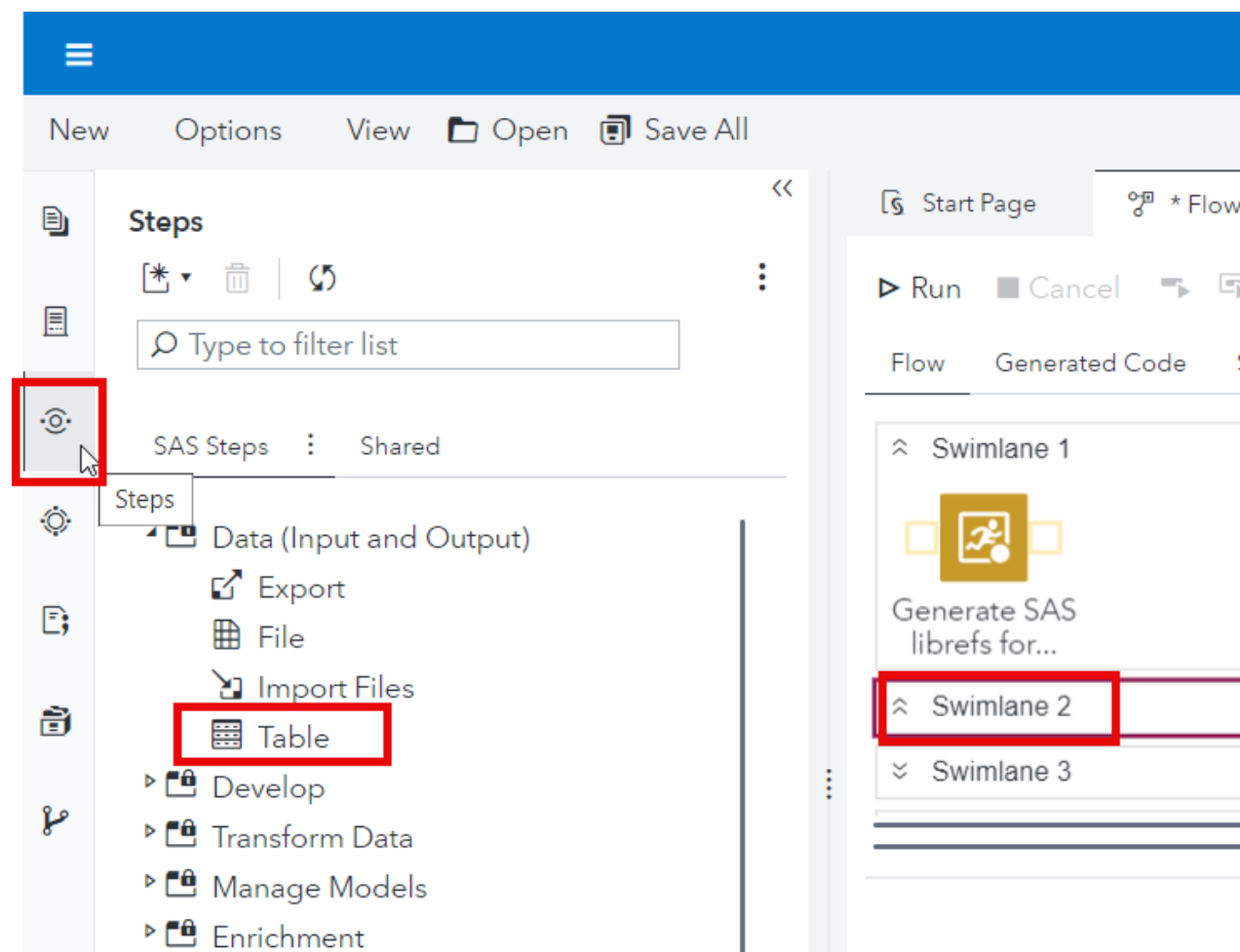
The screenshot displays the SAS Studio interface. On the left, the 'Snippets' panel is open, showing a tree view of 'SAS Snippets' and 'My Snippets'. Under 'My Snippets', the 'Cloud Analytic Services' folder is expanded, and the 'Generate SAS librefs for caslibs' snippet is highlighted with a red box. On the right, the 'Flow' canvas is visible, showing a flow with four swimlanes. The first swimlane, 'Swimlane 1', is highlighted with a red box and contains a node labeled 'Generate SAS librefs for...'. Below the flow canvas, the 'Generated Code' tab is active, showing the SAS code generated for the selected node:

```
Code Node Notes
1 /******
2 /* create a default CAS session and create SAS librefs for existing caslibs */
3 /* so that they are visible in the SAS Studio Libraries tree. */
4 /******
5
6 cas;
7 caslib_all_assign;
8
9
```

Deduplicate

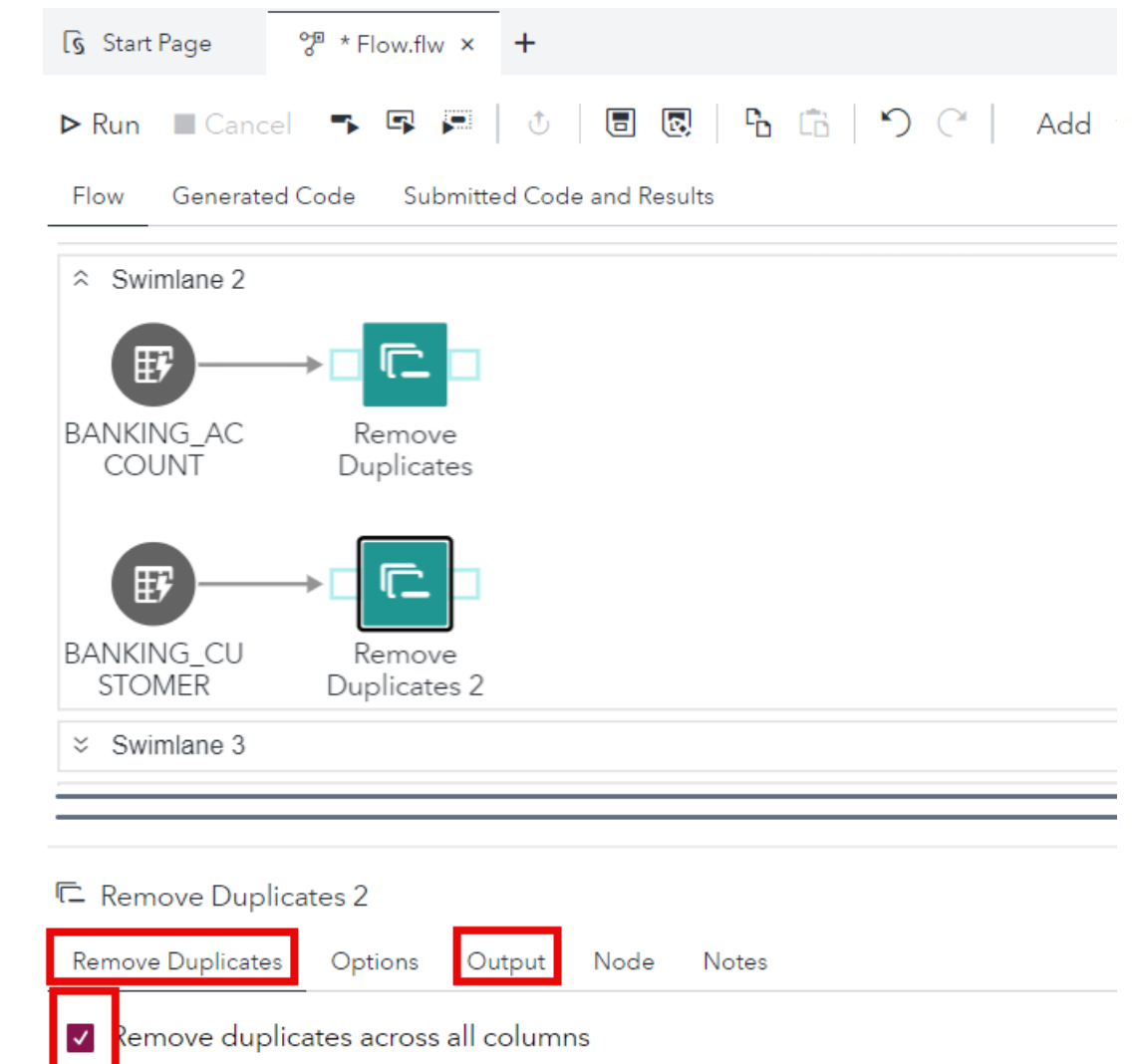
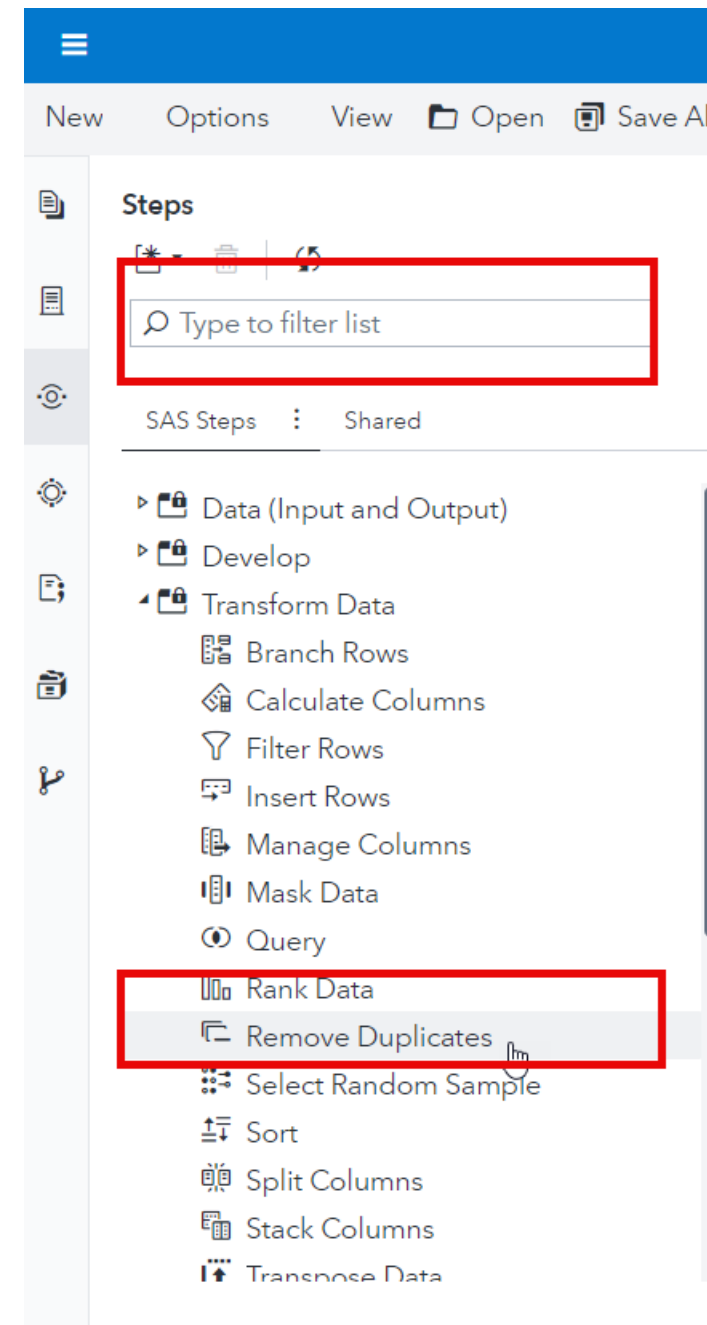
We've seen from the uniqueness analysis that input data sets contain duplicates for both the BANKING_ACCOUNT and BANKING_CUSTOMER data sets.

1. Open the "Steps" pane from the left of the screen and drag two table Steps into Swimlane 2 of the flow and select our input data sets by using the options on the bottom of the page.
2. For the first table, select "Public" as Library and "BANKING_ACCOUNT" as a table name. For the second table, select "Public" as Library and "BANKING_CUSTOMER" as the table name.



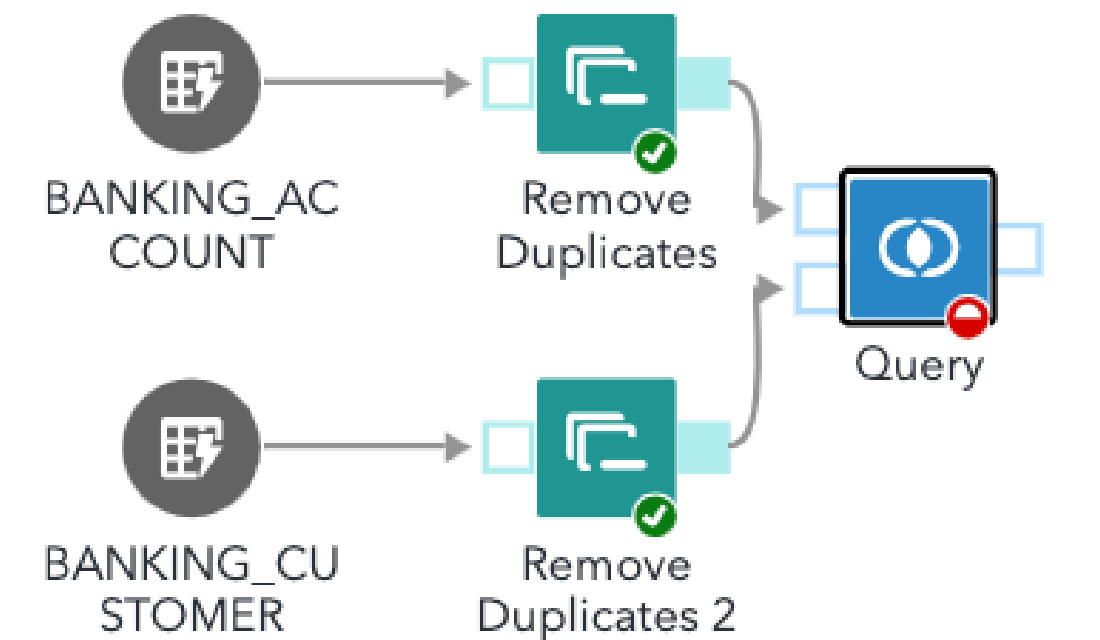
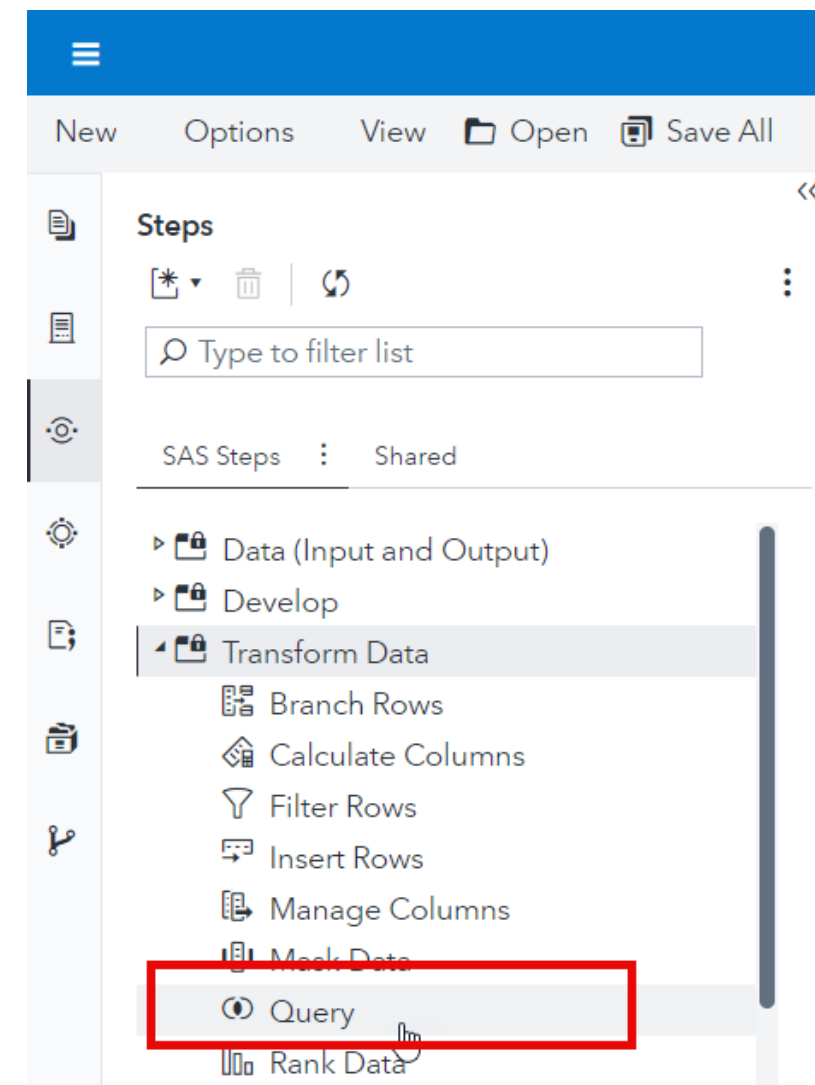
Deduplicate

1. Now identify the “Remove Duplicates” step in the list under the “Transform Data” folder. You could also use the search on the top to find it quicker.
2. Drag two of these steps in the Swimlane 2.
3. Connect the table steps you see in the flow to “Remove Duplicates” steps by using your mouse. Hover over the table nodes and when you see a small hand (instead of the default arrow icon when using the mouse), hold your left click to create an arrow that will connect the two nodes.
4. Ensure you’re removing duplicates across all columns and replacing the output table with the same name (default options).
5. Right-click anywhere in the flow (Swimlane 2) and select “Run Swimlane.”



Join the Data

- Once we've deduplicated our data, we'll need to join it together.
- From the "Steps" tab, drag a Query step into the flow and connect the previous nodes ("Remove Duplicates") to the input ports of the Query node as you see on the right by using your mouse.
- Now, click on the Query step and check the options at the bottom of the page.



Join the Data

- Double-click or drag both tables (t1 and t2) onto the “Select” canvas options pane.
- Open the “Join” tab in the “Query” options and ensure that we are:
 - Joining on our ID variable. This should have been detected automatically.
 - Running a “Left Join.”
 - This won’t be the default. Select the “Venn Diagram” icon as you saw in the second graph to change to a “Left Join.”
- We conduct a “Left Join” since our customers may have multiple accounts, and we don’t want to ignore any accounts.
- Right-click on the “Query” node and select “Run node.”

Swimlane 1

Swimlane 2

BANKING_AC COUNT

Remove Duplicates

Query

Query

Options Node Notes

Columns

+ Calculated Column

Filter

t1 (Remove Duplicates)

t2 (Remove Duplicates 2)

Select Join Filter Sort Output Options

Columns Groups Filter Groups

Remove Row Convert to Aggregate

To select columns for the out

Query

Options Node Notes

Columns

+ Calculated Column

Filter

t1 (Remove Duplicates)

t2 (Remove Duplicates 2)

Select Join Filter Sort Output Options

Expression Builder Reset Join

Join 1

t1 (Remove Duplicates)

t2 (Remove Duplicates 2)

t1.Id

Inner join

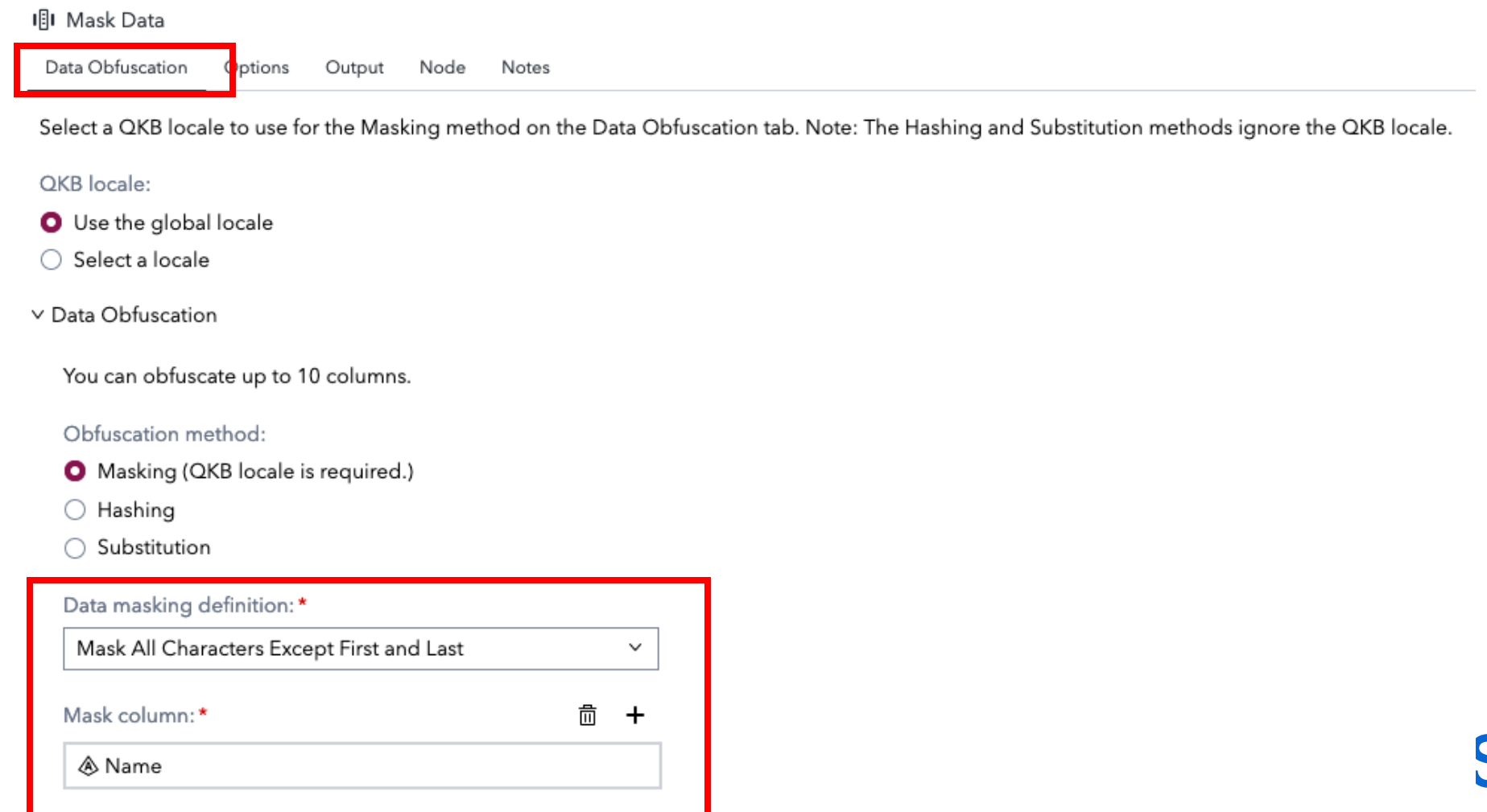
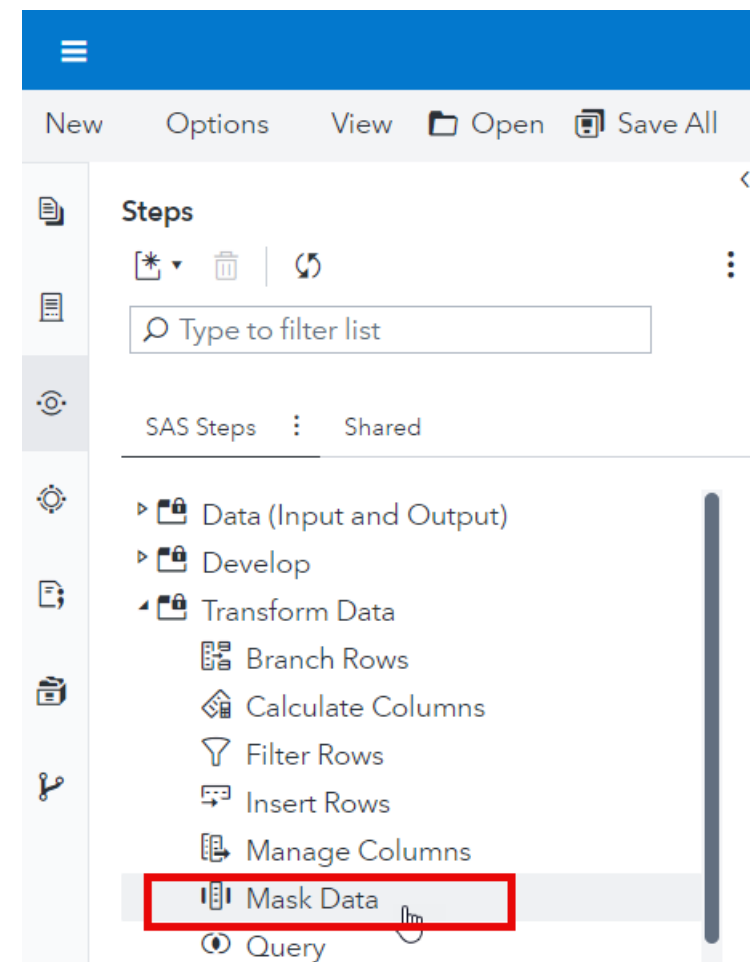
Left join

Right join

Full join

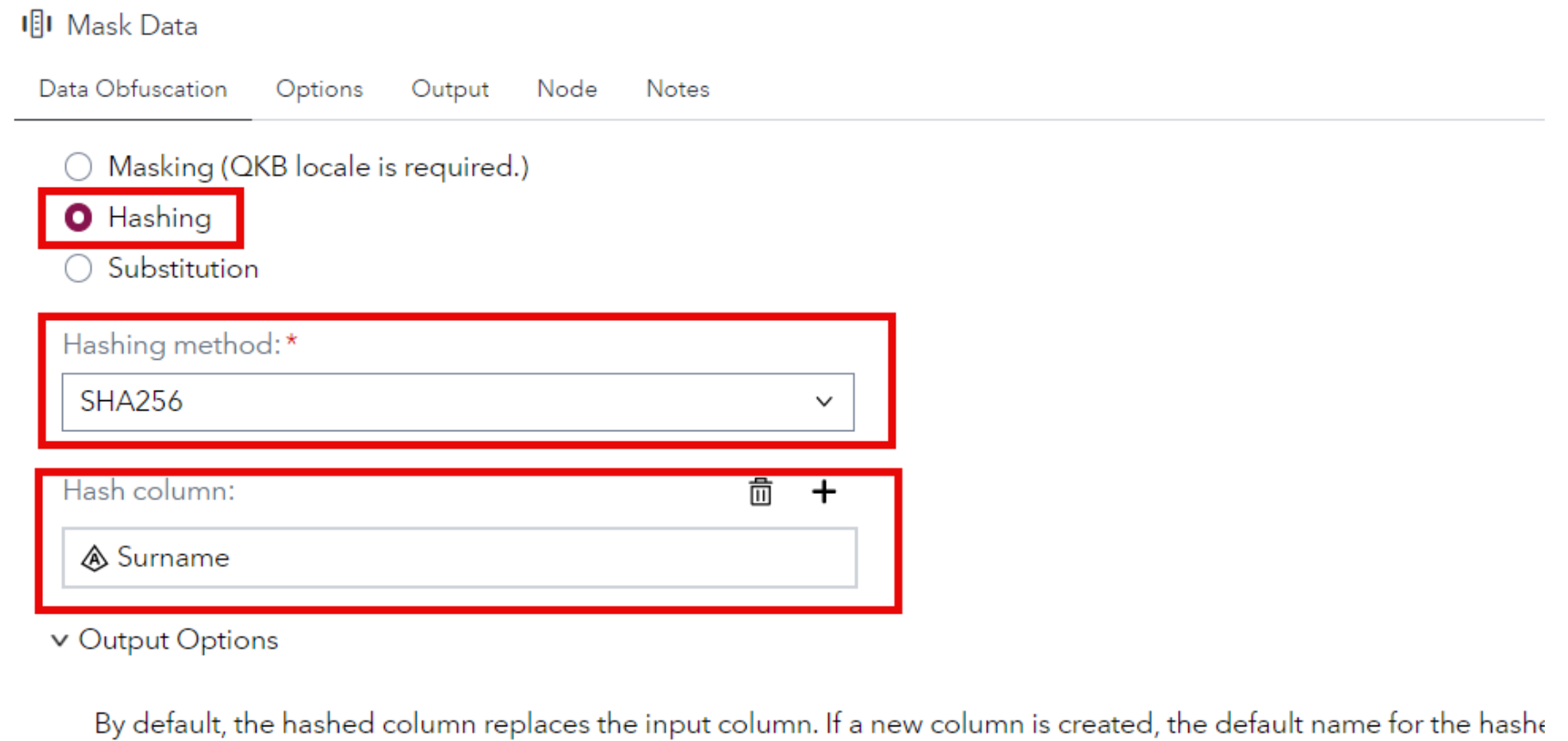
Mask Data

- In addition to joining our data, we'll need to “Mask” some sensitive and/or private information before moving the data to the Data Scientists based on our company's policies.
- Drag a “Mask Data” step into the flow and connect it to the “Query” node that you ran before.
- Recall that we have “Name” and “Surname” variables in our data that are private.
 - In the “Mask Data” options, under the “Data Obfuscation” pane in the “Mask Data” node, let's choose Masking for the “Name” column and mask all characters except the first and last and replace the column so it no longer contains the input names.



Mask Data

- Scroll down to see “Additional Data Obfuscation” and create a “Hash” for the “Surname” variable so it’s not human readable but could be decoded later.
- Make sure in the “Output” pane that we are replacing the existing table.
- Don’t run the node yet, as we want to save the output table in a temporary table to check our results so far and save our flow, as well.



Mask Data

Data Obfuscation Options Output Node Notes

Masking (QKB locale is required.)

Hashing

Substitution

Hashing method: *

SHA256

Hash column:

Surname

Output Options

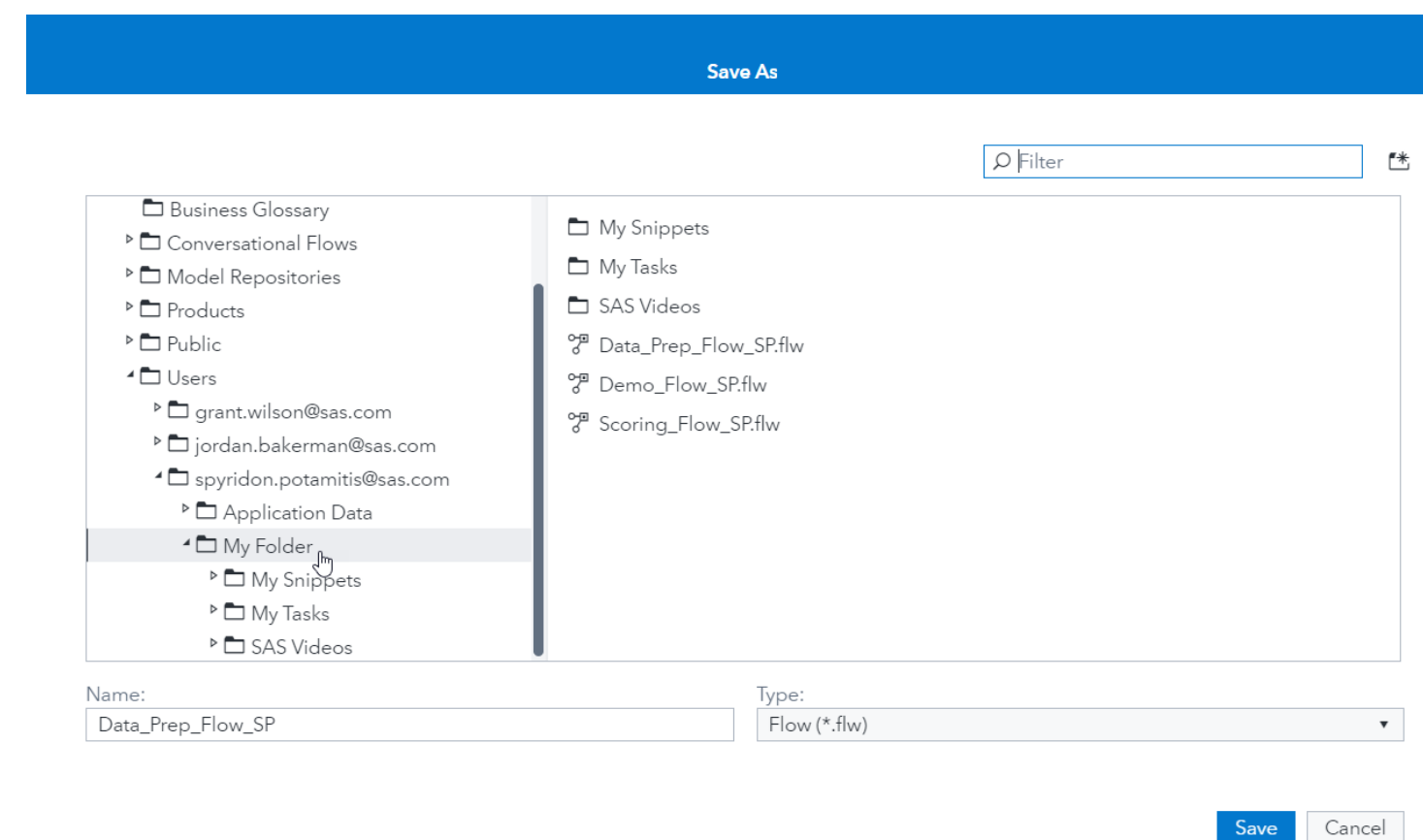
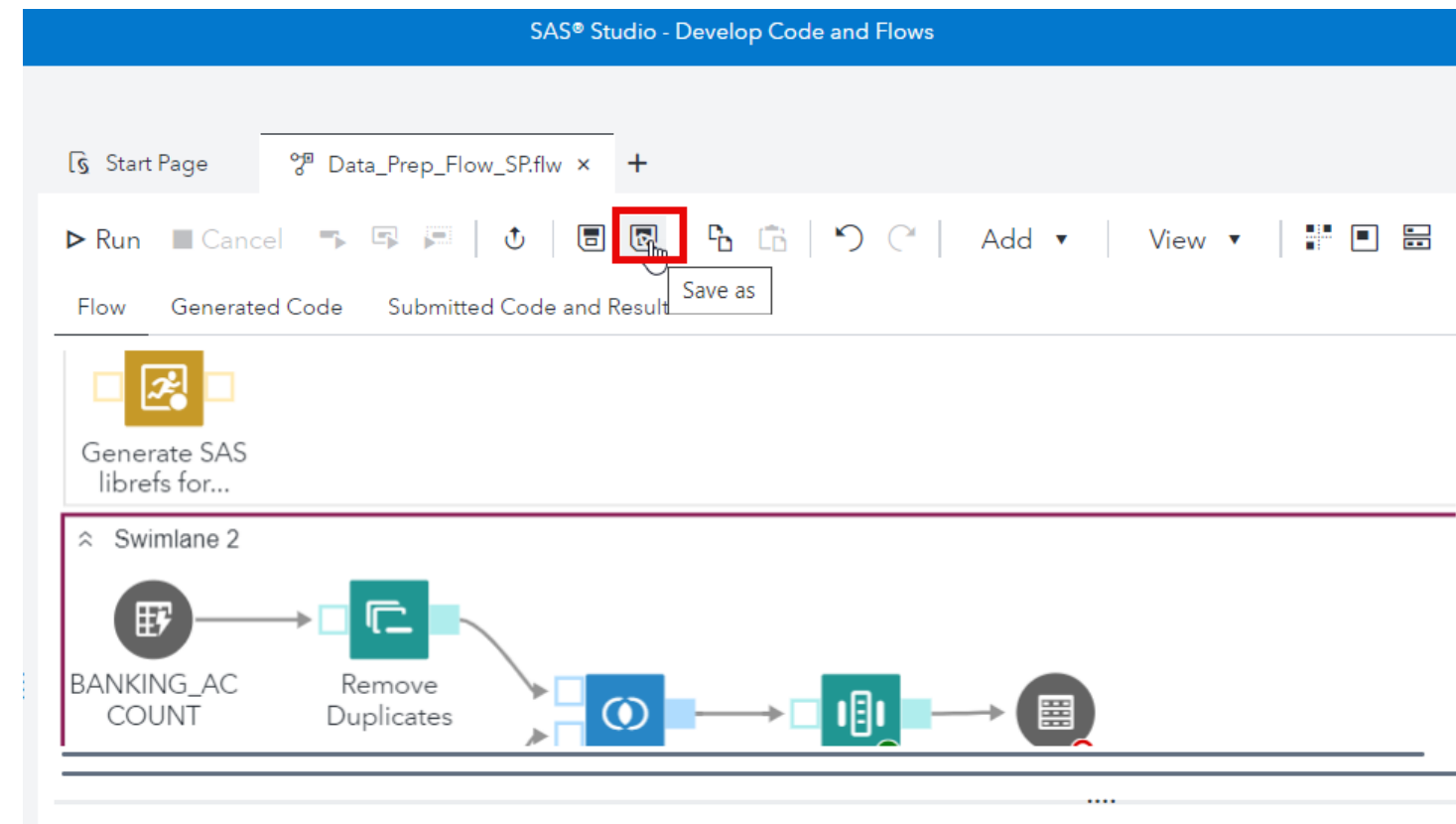
By default, the hashed column replaces the input column. If a new column is created, the default name for the hashed

Save Data and Progress So Far

Since we have performed a lot of work already, we want to save the flow we created so far and the output of our flow till now.

First save the flow by clicking anywhere on the flow and then the “Save” icon. Navigate to “SAS Content”-> “Users”-> “My Folder,” and then give an appropriate name to your flow and click save.

The next thing we want to do is to save our transformed and masked data so far in a table.



Save Data and Progress So Far

From the “Steps tab,” connect a “Table” node to the “Mask Data” node in Swimlane 2.

As this is not the final table that we want to deliver to the Data Scientist and works like a checkpoint, select in the “Table” node options as “Library” “Your User Library,” “CASUSER” and as a “Table Name,” give the name “Banking_Transformed” and click OK.

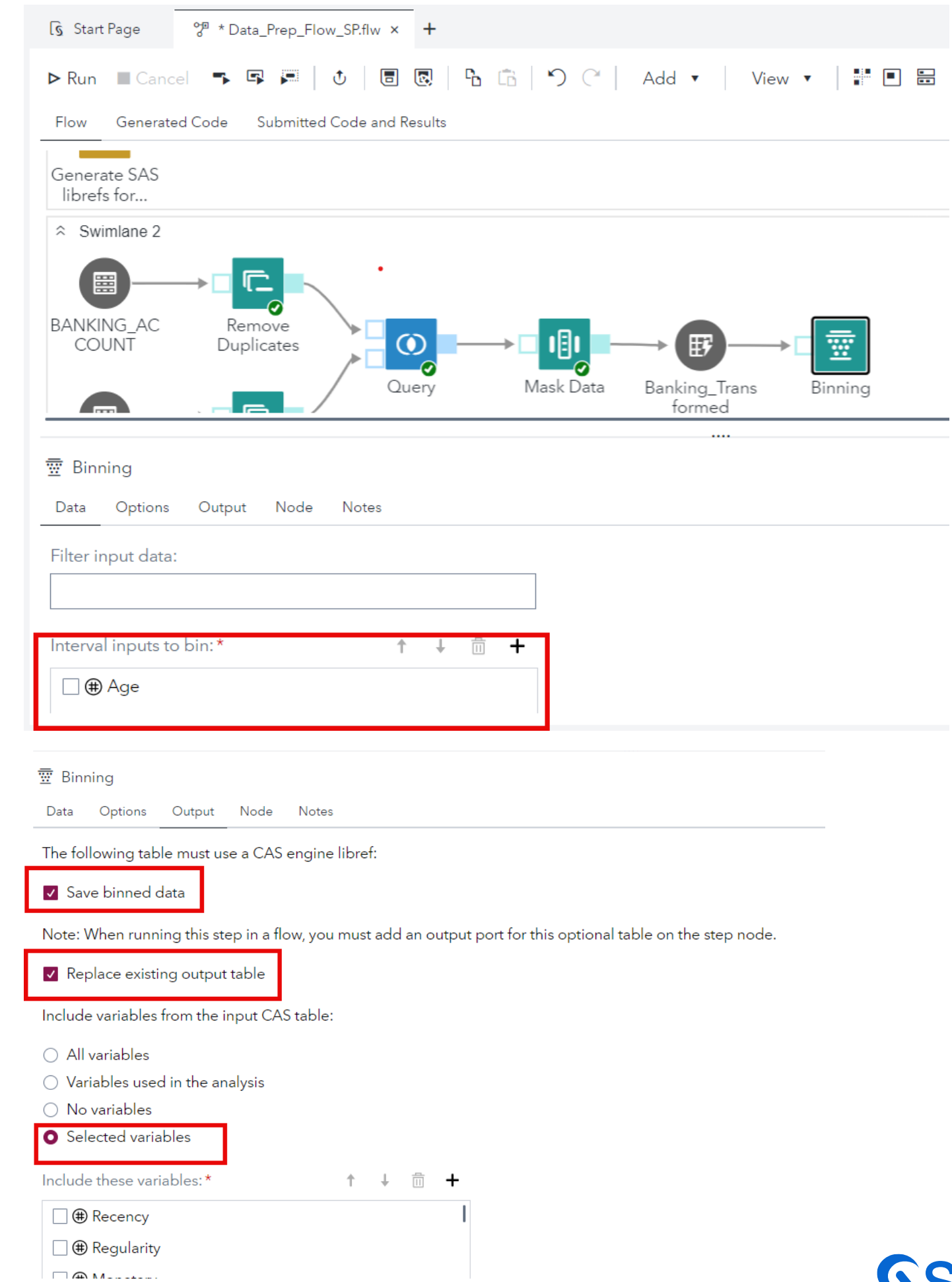
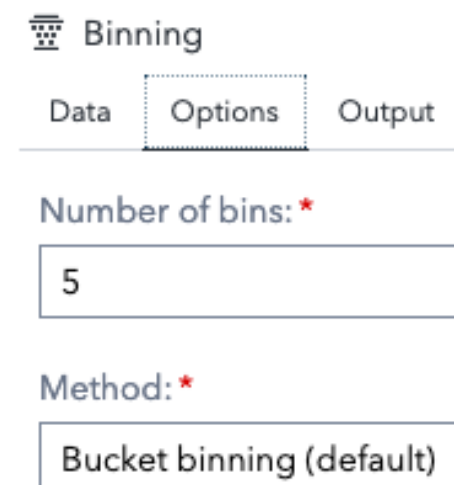
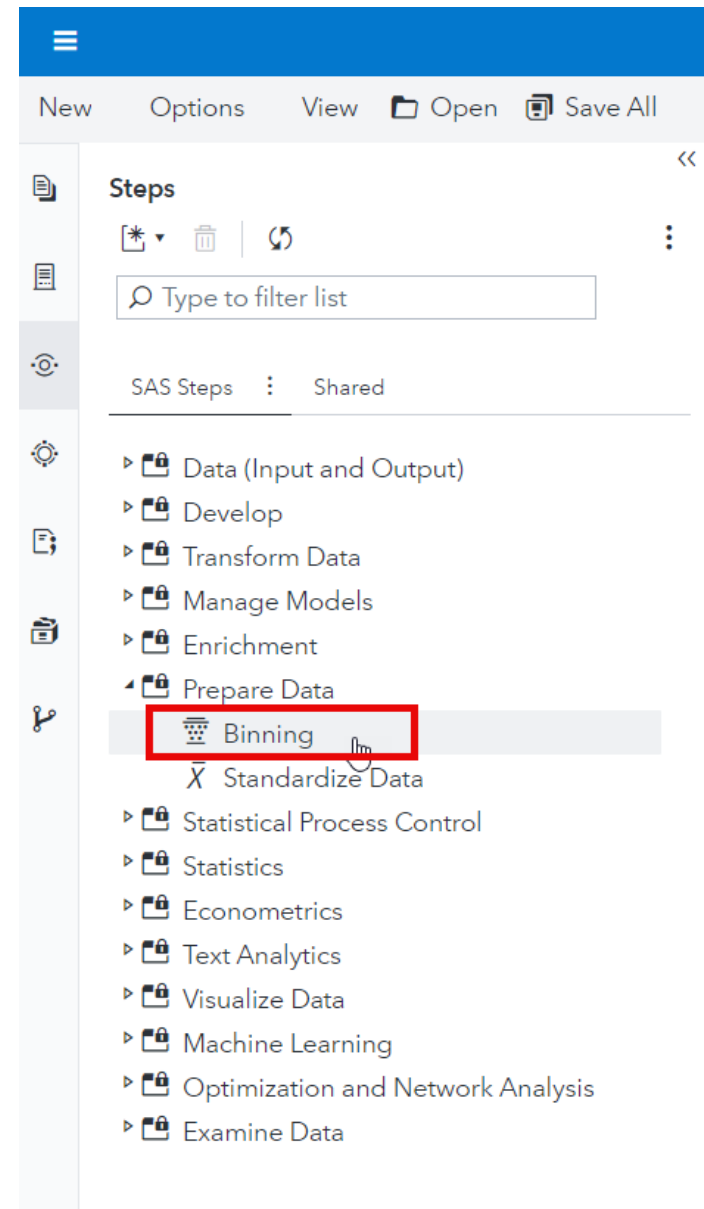
Right-click on the “Table” node you created and select “Run to Node.” After the node is run, click again on the “Table” node, and in the options below, select “Preview Data.” Check that all transformations were performed as expected.

The screenshot shows the SAS Studio interface. On the left, the 'Steps' tab is active, displaying a list of steps. The 'Table' step is highlighted with a red box. The main workspace shows a workflow diagram with nodes: 'BANKING_AC COUNT', 'Remove Duplicates', 'Query', 'Mask Data', and 'Table'. The 'Table' node is highlighted with a red box. Below the diagram, the 'Table' node properties are displayed, with 'Library' set to 'CASUSER' and 'Table name' set to 'Banking_Transformed'.

The screenshot shows the 'Library and Table Name' dialog box. The 'Libraries' list on the left has 'CASUSER' selected. The 'Table' field on the right is set to 'Banking_Transformed'.

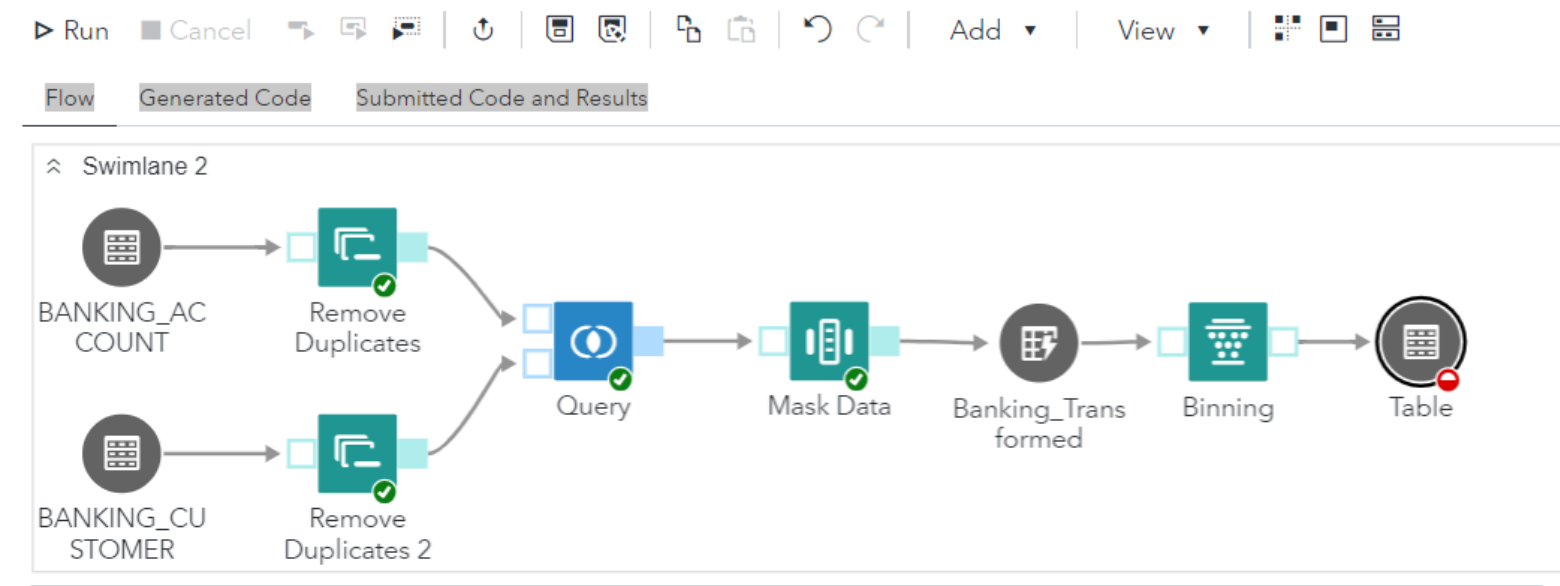
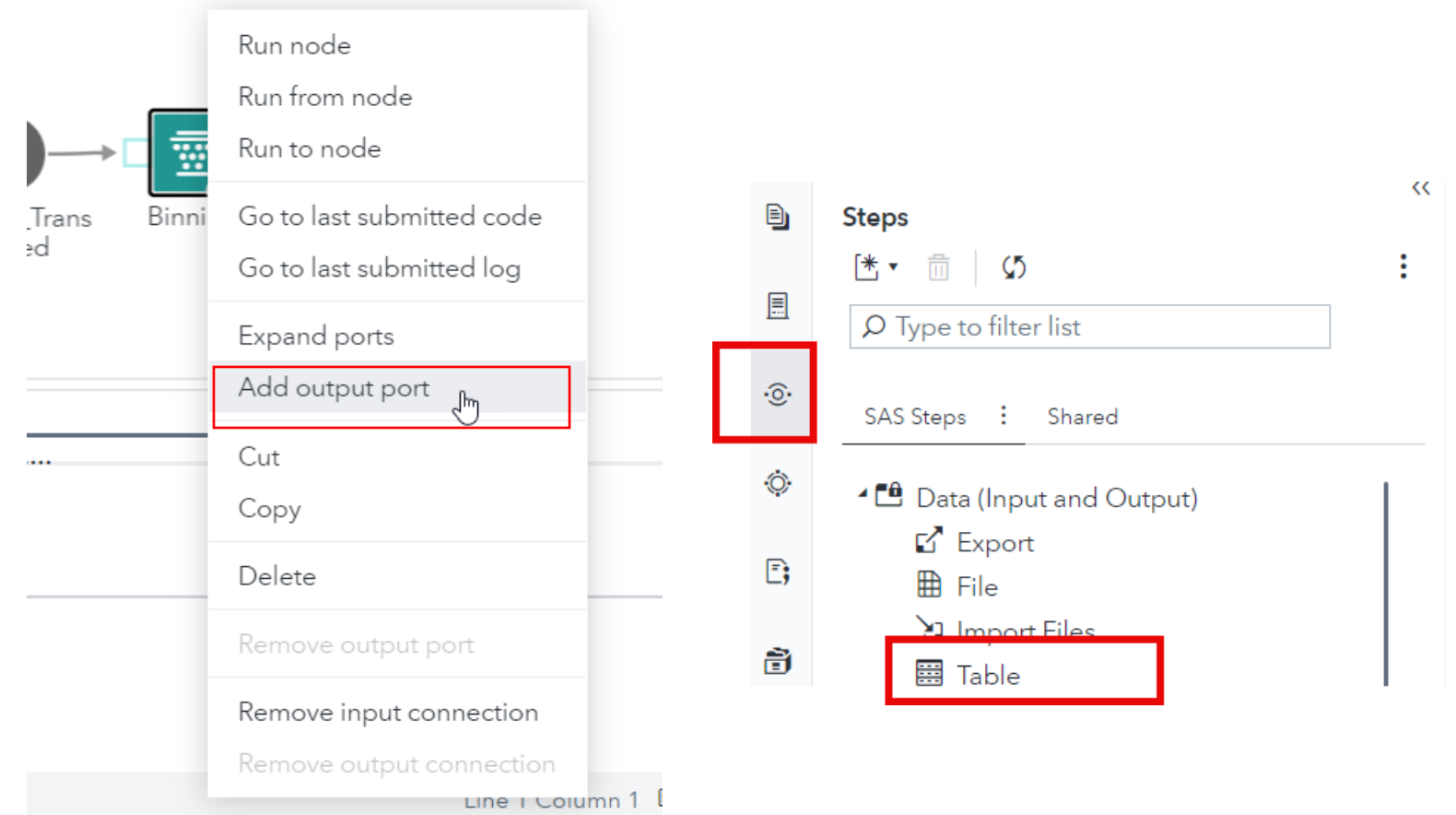
Bin Data

- Now that we've removed duplicates, joined our tables and masked private data, let's go one step further in the name of data privacy.
- From the "Steps" tab, drag a "Binning" node to our canvas and connect the "Banking_Transformed" node with it as you see on the right.
- Add "Age" as a variable to be binned. This way, we de-identify the ages of our customers, but we can still use the binned info in modeling later if we wish.
- In the options, select "5" as the number of bins.
- In the "Output" tab, select "Save binned data," "Replace existing output table" and use the "Selected variables" option to get all variables in the output data except the "Age" variable, which will be replaced with its binned version.



Bin Data

- Now we want to make sure that we save the final output to our ABT table, which is the table the Data Scientist will use.
- To do that, right-click on the “Binning” node and select “Add output port.”
- Then grab another “Table” node from the “Steps” menu.
- Connect the “Binning” node to the “Table” node as you see in the graph.
- Since this is the final table, we are going to use for modeling, use the “Table” node’s options to set the “Library” to “Public” and name the table “BANKING_ABT.” Click OK.

A screenshot of the 'Table' node properties dialog in SAS Studio. The dialog has tabs for 'Table Properties', 'Options', 'Published Columns', 'Preview Data', 'Node', and 'Notes'. The 'Table Properties' tab is active. It shows two input fields: 'Library: *' with a dropdown menu showing 'Select a library', and 'Table name: *' with a text input field showing 'Enter a table name'. Both fields have a red border and a red 'X' icon.

Bin Data

- Right-click on the “Table” node in your flow and select “Run to Node.” Now in the “Table” node you added last, “BANKING_AB T,” you can preview your data and make sure everything is OK.
- Save your flow by using the “Save” icon or use “Ctrl” and “S” to save as you would save a document.
- Your work as a Data Engineer is completed.

The screenshot shows the SAS Data Studio interface. At the top, there are tabs for 'Start Page', 'Data_Prep_Flow_SP.flw', and 'PUBLIC.BANKING_AB T'. A toolbar contains various icons, with the 'Save' icon (a floppy disk) highlighted with a red box. Below the toolbar, a flow diagram is visible, showing a sequence of nodes: 'BANKING_AC COUNT', 'Remove Duplicates', 'Query', 'Mask Data', 'Banking_Trans formed', 'Binning', and 'BANKING_AB T'. The 'BANKING_AB T' node is highlighted with a red box. Below the flow diagram, the 'BANKING_AB T' node is selected, and the 'Preview Data' tab is active, also highlighted with a red box. The preview shows a table with 10,000 rows and 60 columns. The first two rows are visible:

	⊕ Recency	⊕ Regularity	⊕ Monetary	⊕ Amount_avg	⊕ Savings_num	⊕ Creditcards_num
1	5	3	5	50	25	2
2	6	4	5	53	21	4

Thank you!

