# MODUL 8
## SELF-SERVICE DATA PREPARATION IN SAS® VIYA®
## Introduction to SAS Data Studio and Plans

## THEME DESCRIPTION

Students are able to do transform data into the form required for analytics, in accordance with the SAS Data Preparation methodology address transforming the data: Prepare and Cleanse. In order to standardize your data to ensure consistency, parse data, or group data to create one best record for Analytics.

## WEEKLY LEARNING OUTCOMES (SUB-LESSONS)

CLO-3-Sub-CLO-8, Able to effectively implement the process and role of data curation in data management through the use of SAS® Viya cloud technology– C6.

Through the following learning steps:

1. SAS Data Explorer: Explore and load the CLIENT_INFO table
2. SAS Data Studio: Using Suggestions to Prepare Data
3. Add and configure the Remove duplicates suggestion
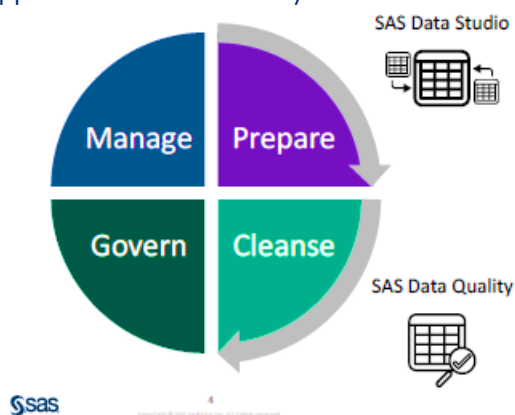4. Add the Parsing CUSTOMER_NAME suggestion
5. SAS Drive: Saved Plan overview

## PRACTICUM SUPPORT
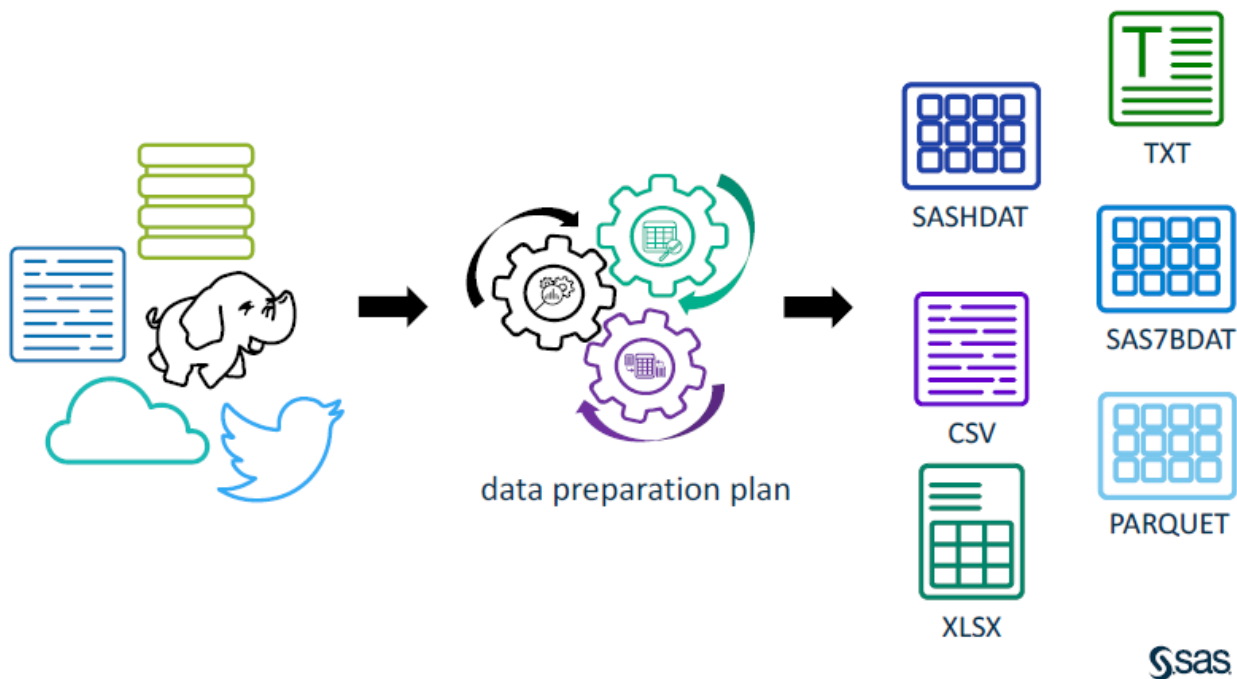
c. Windows Operating System
d. (any) Browser Application

## PRACTICUM GUIDELINES

### Exploring the SAS Data Preparation Applications

▸ Our IS-429 BDA practicum class today illustrates exploring the SAS Self-Service Data Preparation applications in the SAS Viya environment.



▸ The activities, and practice are about investigating the environment and SAS Data Preparation applications:
- SAS Drive
- SAS Data Explorer
- SAS Data Studio
- SAS Data Quality

▸ A **data preparation plan** is a collection of transforms that are performed on one or more tables that **prepare data for analysis and reporting**.
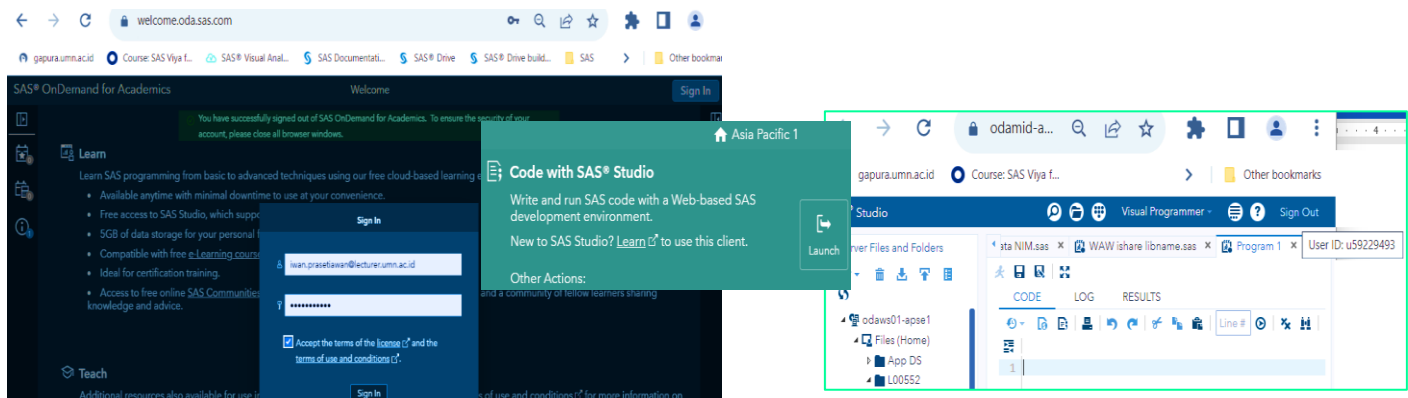
data preparation plan

> You **create plans in SAS Data Studio,** and you can save them to a SAS folder location where other users can also access them.
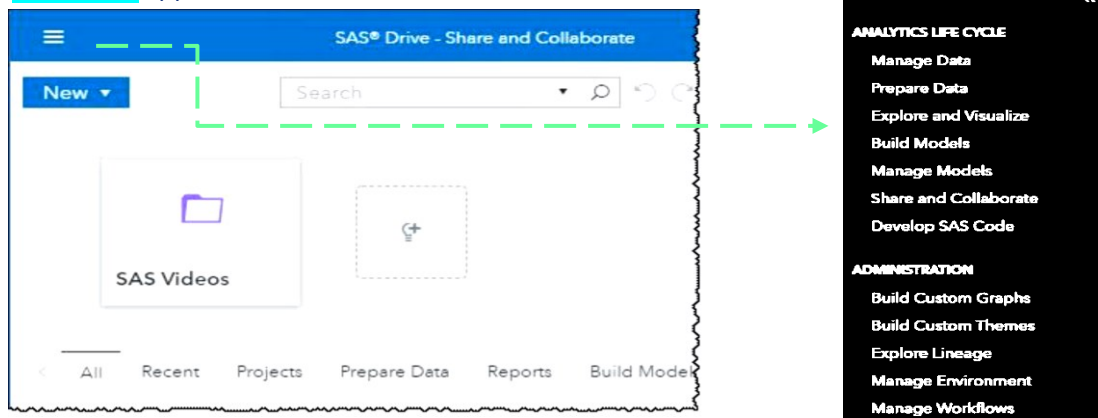
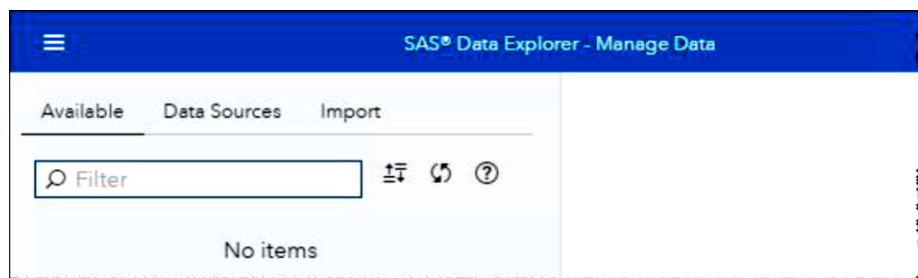## PRACTICUM STEPS

### Login to SAS Studio

Start, https://welcome.oda.sas.com



### Open SAS Drive Explorer

1) When you have signed in to SAS OnDemand for Academics, then open the SAS Studio application and open another tab of your browser and click SAS® Drive to go to SAS Drive page.

> *If you have problems, you can try using this url: https://v4e041.vfe.sas.com/SASDrive/ **or go to** https://vle.sas.com/course/ view.php?id=7715 then click "Launch SAS Viya for Learners3.5".*
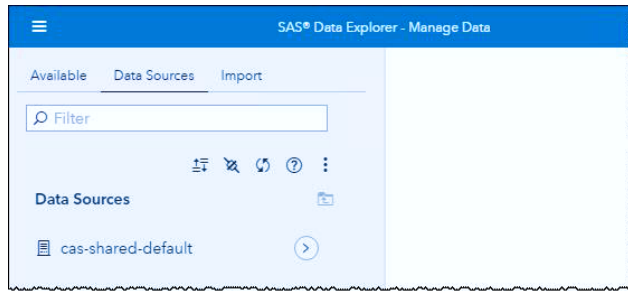
2) **SAS Drive** appears.



- ▸ With SAS Drive, you can manage all your content in one convenient location.
- ▸ You can search for an object that was created in a SAS application. You can share, tag, and preview your SAS content. Sharing takes only one click on the object.
- ▸ You decide who you want to share it with and whether they can read, edit, or share it.
- ▸ Tagging content can help you organize your work. Previewing files can help you find the right one.
- ▸ In addition to controlling your SAS content, you can drag-and-drop most content types from your local machine, and SAS Drive will upload it for you.
- ▸ For quick access to what you care most about, SAS Drive includes a collapsible area where you can place any content you choose. This paper includes screen shots and illustrative examples showing you how to manage all your SAS content in one place

3) Click ☰ (Show list of applications) to access SAS Viya applications.
4) The applications that are shown are controlled by the SAS products that are licensed and the permissions assigned to the user.

5) The list of actions in the list relate to a specific application that will be used to perform that action. For example, **selecting Manage Data opens SAS Data Explorer** to enable you to work with **data connections and data**.
6) Select **Manage Data** to **open SAS** Data **Explorer**.
   **SAS Data Explorer** appears
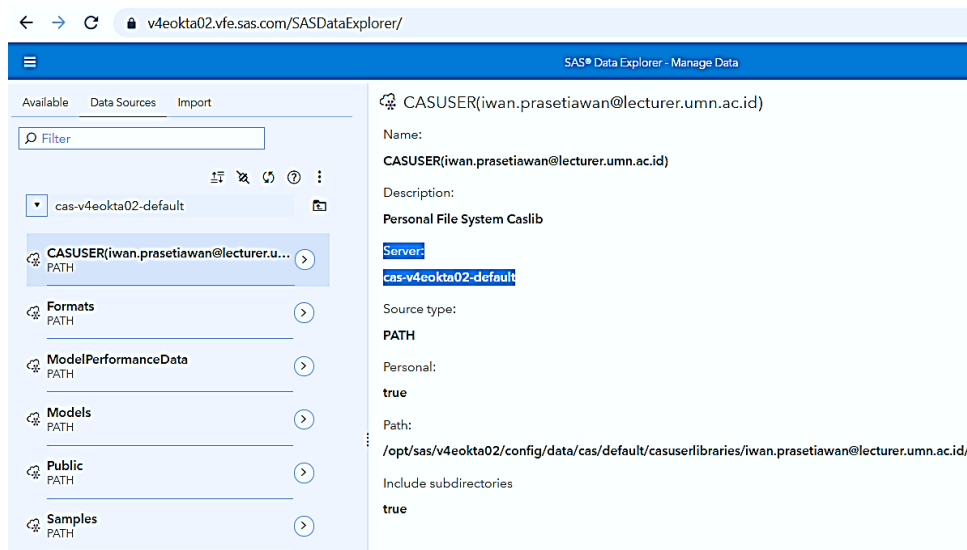


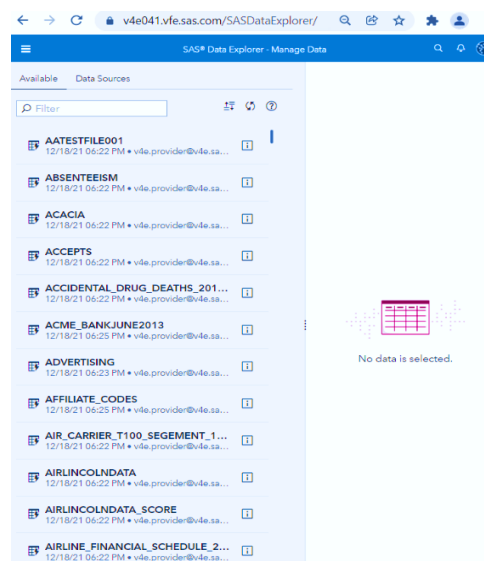1. **Explore and load the CLIENT_INFO table.**
   a. Click the Data Sources tab

b.  Only one CAS server is configured. Click (Down one level) for the **cas-shared-default server** entry to view the list of data sources that are configured on the **CAS server**.
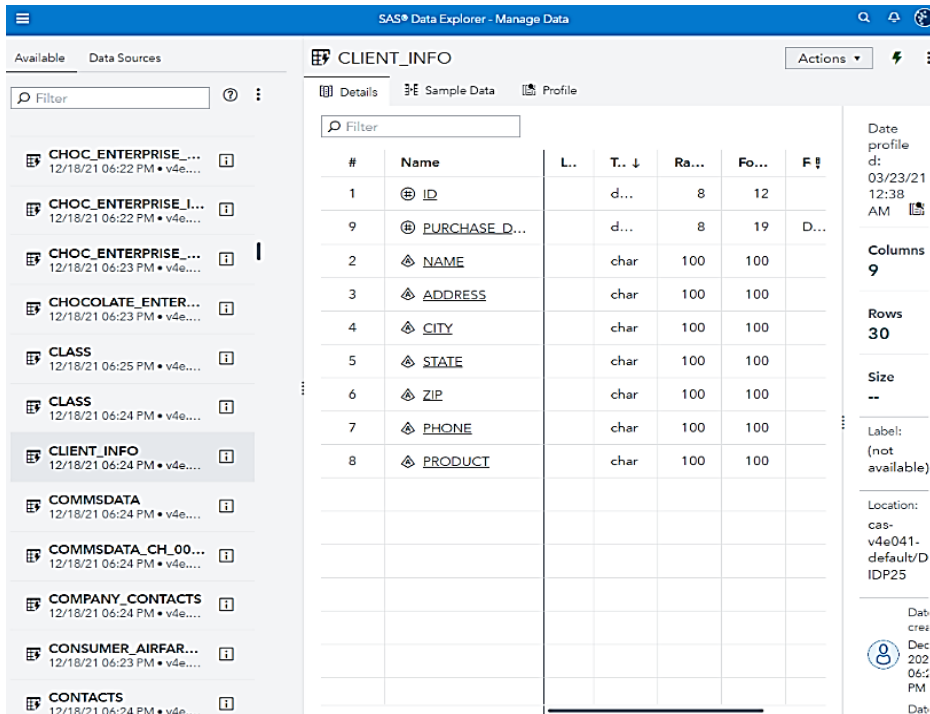**Note**: By default, users have their own CASUSER library and the Formats, Public, and Samples caslibs. Additional caslibs might be visible, depending on the SAS products that are licensed.



c.  Double-click Public to open the connection and view the data sources.  Click CLIENT_INFO.sashdat to view specifics of the table.

  ▸  *If you don't find the CLIENT_INFO table in Data Sources then move to Available tab.*

d.  Click the Available tab.

▶ The Available tab enables you to see all tables that have been loaded into memory from any CAS server to which you have access. Clicked twice the CLIENT_INFO table to selected.



The Details tab shows metadata information about the table such as the column properties, table size, and number of columns.
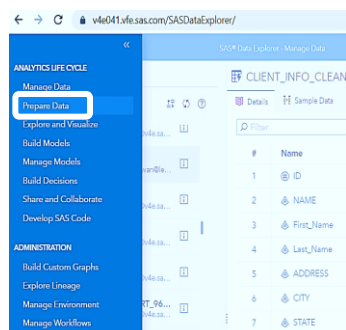
e. Click the Sample Data tab.



f. A sample data is shown displayed completely for 100 records data sample.

g. <mark>Screenshot your page of the Client_Info sample data as in the output figure A.</mark>

2. **Using Suggestions to Prepare Data.**

▶ This demonstration creates a simple plan using the Suggestions feature to quickly add and configure Data Quality transforms.

a. Clicking "**prepare data**" will take you to the **SAS Data Studio tools** for further data preparation activities.

b.  Click ▤ (Show list of applications) and select Prepare Data to open **SAS Data Studio**.

h.  Click New Plan. SAS Data Studio appears



The Choose Data window appears.

c.  Select the **BANK_CUSTOMERS** table **as the data source**. Click the Data Sources tab.
d.  Click (Down one level) for the cas-shared-default server.
e.  Click (Down one level) for the DIDP25 caslib.
f.  Right-click BANK_CUSTOMERS.sashdat and select Load to load the table into memory.
    ‣  `Ignore it, if the table is already available in CAS memory.`
    ‣  `This step shows what if you prepare a table for loading into CAS memory and will be available at available tab.`

g.  Click the Available tab. The BANK_CUSTOMERS table is displayed as an in-memory table.



h.  Double click the BANK_CUSTOMERS table. You have them now at your SAS Studio workspace.
i.  Scroll to the right and expand the CUSTOMER_NAME, ADDRESS, CITY, and STATE columns to review the data.

**BANK_CUSTOMERS**

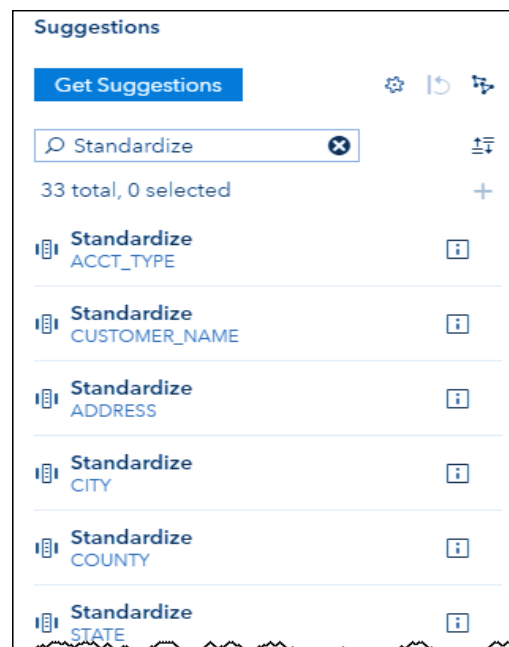| CUSTOMER_NAME | ADDRESS | CITY | COUNTY | STATE |
|---|---|---|---|---|
| Stephen seiters | 814 Blue Mound rOAD | fort WORTH | Tarrant | TX |
| Doug mAtrisciano | 12500 Ne 10th Pl | Bellevue | King | waSHINGTON |
| bridget Knightly | 66 Flint sTreet | Asheville | Buncombe | NC |
| Kennith Kirklin | 2303 21st Ave S | nashville | Davidson | TeNNeSSee |
| James Mitchell | 8253 Ronson Rd | pasadena | | CA |
| kathleen Beresnyak | 100 W 25th Avenue | San Mateo | San Mateo | CA |

**Note**: The data in the *CUSTOMER_NAME, ADDRESS, CITY*, and *STATE* columns is **mixed case** and **not in a standard format**.

j.  Click 💡(Suggestions) in the left-hand pane. Click Get Suggestions.



The system has provided some of 32 suggestions for the selected table for you to use.

k.  **Standardize the Column**.  Enter "Standardize: in the Filter box.
    ▸ **Twelve standardizations are suggested** for this table. These columns have standardization suggestions because they were the right type of data, character, and they seemed like they contained data that corresponds to one of our standardization algorithms.
    ▸ Sometimes the suggestions make a mistake. For example, ACCT_TYPE is probably unique to this bank data and cannot be standardized by one of our default standardization algorithms. It might look enough like a phone number or address that the sort of data in the column was misclassified. However, most of the suggestions are correct. Name, Address, City, County, and State can be standardized.
    ▸ Standardization put every value in the column in a consistent format.
      For example, it applies consistent casing, abbreviations, and ordering to address values.

l.  Double-click the **Standardize CUSTOMER_NAME suggestion** to add it to the data preparation plan.



m.  The Standardize transform is added to the plan.
    ▸ Verify that the **CUSTOMER_NAME column** is selected as the <u>source column</u>.
    ▸ Verify that the selected <u>locale</u> is **English (United States)**.
    ▸ Verify that the selected <u>definition</u> is **Name**
    The suggestion gives us default settings, but we can change them.

n.  Click **Replace source** column

o.  This is the first transform in the plan. Transforms can be deleted from anywhere in the plan or moved forward or backward in the processing order.



p.  The **session table is an in-memory table**. Right now, it reflects the source table BANK_CUSTOMERS. It is not current to the plan because we have not run the plan yet. When we run the plan, the session table is updated to reflect the transforms used. Click Run.

Partial List of the Original Data     Partial List of the Standardized Data



The casing for the CUSTOMER_NAME column is now consistent. Other standardizations such as applying consistent abbreviations to name suffixes might have been applied, but we can see that the data in the first six rows are now consistently in proper case.

q. Click the X to remove the Standardize filter.
r. Click Get Suggestions in the Suggestions pane to regenerate the suggestions.



s. Enter CUSTOMER_NAME in the Filter box
   ‣ The Standardize CUSTOMER_NAME suggestion has been removed from the list of suggestions because it was already applied to the session table.
   ‣ The suggestions evaluate the session table as we update it. If we create more columns, suggestions might be generated for those columns as well.
t. <mark>Screenshot your page of Bank_Customers updated as in the output figure B</mark>.

3. **Add and configure the Remove duplicates suggestion**
   a. Double-click the **Remove duplicates suggestion**



We can choose to remove duplicates based on the values of one or more specific columns instead of using all columns. As an ID value for the table, we want **CUST_ID to be unique,** so let's remove duplicates by the CUST_ID column.

b. Clear the Remove duplicates across all columns option. Select the CUST_ID column as the column to group by.

> **2. Remove duplicates**
>
> ☐ Remove duplicates across all columns
>
> Select columns to group by: ⓘ
>
> ⊕ CUST_ID ▾ ➕

c. Click ▦ (**Properties for the Result table**) in the right pane. There are **currently 483 rows** in the table.
d. Click Run. The **new row count is 481**. Two rows with duplicate CUST_ID values were removed

> Result Table - BANK_CUSTOMERS (ses... ↻
>
> | Columns | Rows | Size |
> |---------|------|------|
> | 20 | 481 | 142.8 KB |
>
> Label:
> (not available)
>
> Location:
> cas-shared-default/DIDP

e. ==Screenshot your page of removing duplicates Cust_Id as in the output figure C==.

4. **Add the Parsing CUSTOMER_NAME suggestion**.
   a. Double-click the Parsing CUSTOMER_NAME suggestion.
      ‣ The **Parsing** transform enables us to **break up a string into smaller semantic** pieces called **tokens**.
      ‣ For example, a Name string might have smaller pieces like a first name, middle name, or last name, which all have meaning separate from the larger value.

> **3. Parsing**
>
> Source column: Locale: Definition:
>
> ⚠ CUSTOMER_... ▾  English (United States) ▾  Name ▾ ⓘ ➕
>
> ⌄ Selected tokens: none
>
> Available tokens (6): Selected tokens (0):
>
> 🔍 Filter
>
> **Family Name**
> CUSTOMER_NAME_FamilyName_Parse    No items
>
> **Given Name**
> CUSTOMER_NAME_GivenName_Parse

   b. Verify that the CUSTOMER_NAME column is selected as the source column.
   c. Verify that the selected locale is English (United States).
   d. Verify that the selected definition is Name.
   e. Double-click the Family Name token to add it to the Selected tokens box.

f. Double-click the Given Name token to add it to the Selected tokens box



Notice the long column name under the token name. We will change these to simpler column names.

g. Click "Options for new columns".



In this window, you can change column properties.

h. Rename CUSTOMER_NAME_FamilyName_Parse to Last_Name.

i. Rename CUSTOMER_NAME_GivenName_Parse to First_Name.

j. Changes Length for First and Last Name to 64

k. Click OK to save the column properties.

l. Click Run to run the data preparation plan.

| ⚠ Last_Name | ⚠ First_Name |
|---|---|
| Seiters | Stephen |
| Matrisciano | Doug |
| Knightly | Bridget |
| Kirklin | Kennith |
| Mitchell | James |

There are now two new columns for First_Name and Last_Name.

**m. Screenshot your parsing page of Customer_name suggestion as in the output figure D**

5. **Change the order of the transforms**.
   a. Change the order of the transforms.



   b. If our data had a lot of duplicates, it would be more efficient to remove duplicates before doing other processing. Let's make that transform the first one to run.
   c. If necessary, click [▦] (Plan).
   d. Click the Remove duplicates transform.
   e. Click ↑(up arrow) to move it above Standardize



   Remove duplicates is now the first transform in the plan.
   f. **Save the plan**. Click [⋮] (Options) in the upper right corner to see the plan options.
   g. Click Save As to save the plan as follows:

h.  Verify that Save plan and target table is selected.
    ‣ We can open and change or run saved plans later.
    ‣ If we save the target table, we can open and use it in different applications or as a source for a new plan.
i.  Verify that sashdat is selected as the format.
    ‣ Target tables can be saved in many formats depending on the target caslib.
    ‣ For PATH caslibs, the target formats include SASHDAT, TXT, CSV, SAS7BDAT, XLSX, XLS, PARQUET, and ORC. If you do not need to save the data on disk, the Save as in-memory table only check box enables you to load the table to memory without saving it to disk.

j.  Click Save to your folder..  Click  ⋮  Options) and **Generate code for your Plan**:



    ‣ The system will automatically generate the code in SAS format and place it directly into the "download" folder in your workspace with a .txt file format.
    ‣ **Rename your plan code generated as "IS-429 Lab Week#8 Plan NIM Name.sas"**
    ‣ Select Close to close the plan.



k.  Return to SAS SAS DataExplorer
n.  <mark>**Screenhot your Saved Plan page as in the output figure E**</mark>.

➲ Finally, today's practicum is over, <mark>**collect your doc/docx and sas file into IS-429 Week#8 NIM yourName.zip**</mark> and submit immediately today to e-Learning IS-429 Practicum Week#8

# RESULTS/ OUTPUT

## A. SAS Drive Explorer: Explore and load the CLIENT_INFO table



## B. SAS Data Studio, Using Suggestions to Prepare Data: Standardize Customer_Name

## C. Add and configure the Remove duplicates suggestion



## D. Add the Parsing CUSTOMER_NAME suggestion

E. SAS Data Explorer: New Bank Customer of yours adding the DIDP25 CAS table.



**REFERENCE**

1. Anna Yarbrough. 2020. Introduction to Data Curation for SAS® Data Scientists Course Notes. SAS Institute Inc. Cary, NC, USA.
2. SAS Institute Inc. 2020. SAS® Viya® Programming: Getting Started. SAS Institute Inc. Cary, NC, USA.
3. Johnny Starling, Erin Winters, and Anna Yarbrough. 2020.Self-Service Data Preparation in SAS® Viya® Course Notes. . SAS Institute Inc. Cary, NC, USA
4. SAS® Support | Documentation
5. Other additional references are excerpts from various Online Learning/websites.