## Practical #2: Data

**Objectives:**

1. Know how to preprocess data with SAS EM
2. Know how to export processed data from SAS EM to a local directory
3. Know how to use the Import node to load an external data set to SA EM

**Submission:**

- *What to submit*:  a PDF report generated by the Reporter node/tool of SAS EM, containing the information (diagram, results etc.) about the exercise done by you for this practical.
- *Deadline of the submission*: 11:59PM (Adelaide Time), Tuesday of Week 4.
- *Submission link*: "**Submission Link of Prac #2**" in **Week 3 section** on Learnonline course site.
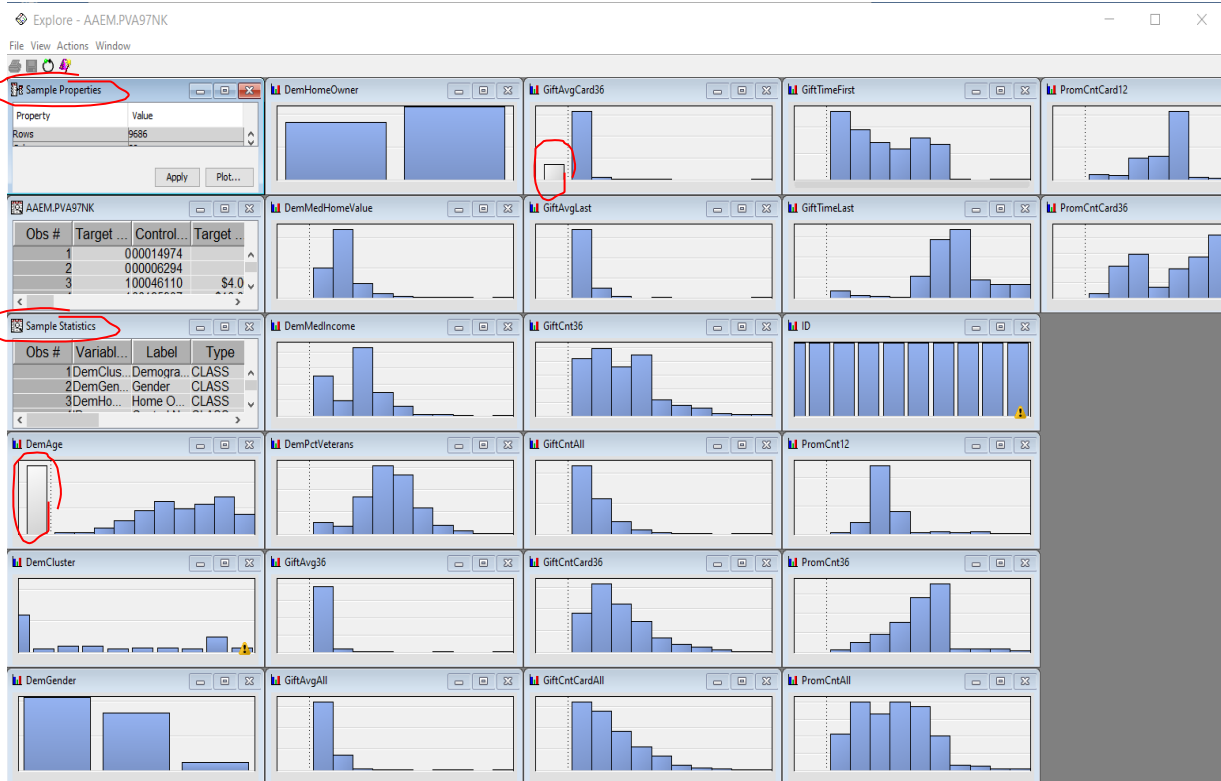- *Marks*: Prac#2 (part of the ongoing assessment of the course) is worth 2% of the total marks of the course.

**Instructions:**

### Preprocess the PVA97NK data set using SAS EM

In Practical #1, you created a project, a diagram workspace within the project, and a data source for the project using the PVA97NK data set. For this practical, you will continue to work with this data set and try a few data preprocessing functions provided by SAS EM. You will also learn how to export SAS data and import your own data file to SAS EM.

Following the instruction given below to complete the practical.

1. Create a project named **Prac#2**, and in the project, create a diagram named **myPrac2** and a data source with the PVA97NK data set. Then drag the PVA97NK data source to the diagram workspace to create an Input Data node for a process flow. (Refer to the instruction in Tasks 1 to 4 of Practical #1 on how to create a project, diagram and data source.)

   **Note:** you might wish to reuse the project and data source created for Practical #1. However, it appears that SAS EM (or the SAS OnDemand setting) does not allow us to save a project with a different name, which means you will lose the project created for Practical#1 if you simply update and save it for Practical#2. It is recommended that you create a new project for Practical #2. In this way, you can also get more familiar with the process of creating a project, diagram and data source.

2. Explore the PVA97NK data set, note down the number of samples and number of variables, and check if there are missing values. (Refer to the instruction given in Step 2 of Task 5 of Practical#1 on how to explore a data set.)
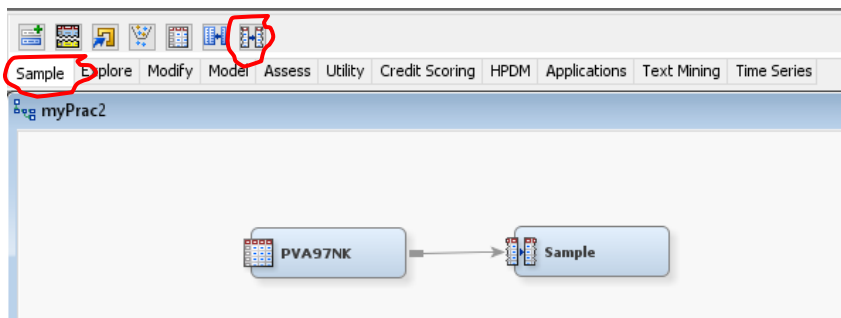
To see the number of samples and number of variables in the data set, in the **Explore** window, double click on the title bar of the **Sample Properties** window to enlarge this window. The number of samples and number of variables in the data set are shown as the **Value** of **Rows** and **Columns** respectively in the **Sample Properties** window. To restore this window to its original (small) size, double click on its title bar again.

To check if there are missing values, in the **Explore** window, double click on the title bar of the **Sample Statistics** window. In the enlarged **Sample Statistics** window, the **Percent Missing** column shows the % of missing values for each variable in the data set. Additionally, if a variable has missing values, you will see a grey bar in its histogram in the **Explore** window too.
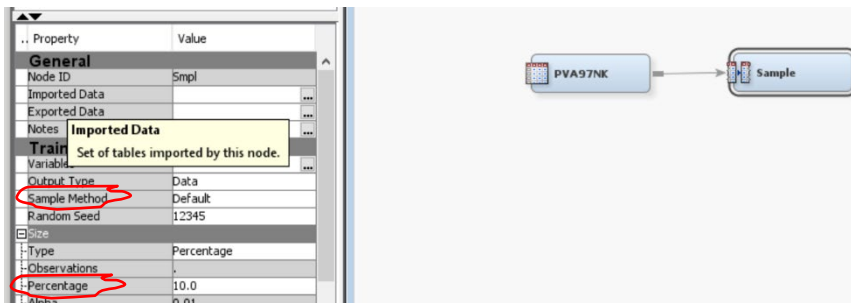
3.  Sample the PVA97NK data set using the Sample tool of SAS EM.

    To create a Sample node, click on the **Sample** tab in the tools palette. Then you can see 7 tools/nodes shown above the tabs. Drag the Sample node and drop it onto the workspace.

    Connect the Sample node with the Input Data node.



4.  Click on the Sample node created, the properties of the Sample node are shown in the **Properties** panel (on the left of the workspace). Look through the properties to know what options are available. Put mouse over a property item, you may see details of the property.

5. Use the random sampling method to sample the PVA97NK data set.
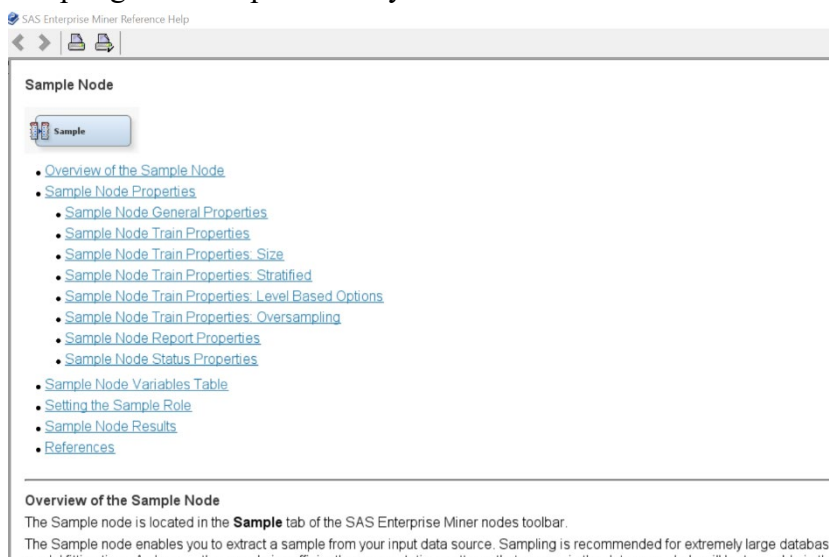
In the **Properties** panel while the Sample node is highlighted (if not, click on the node), click on the white field next to **Sample Method**, a dropdown menu will appear. Change sample method from **Default** to **Random**. Then click on the white field next to **Percentage** and change the default percentage (10.0) to 20 by typing 20 in the field.

Now right click on the Sample node, and choose **Run**, to conduct the random sampling.

In the **Run Status** window popped up, click on **Results**, the Result window will show, and you can see that the sampled data set contains 1937 observations (which is 20% of observations of the whole data set).
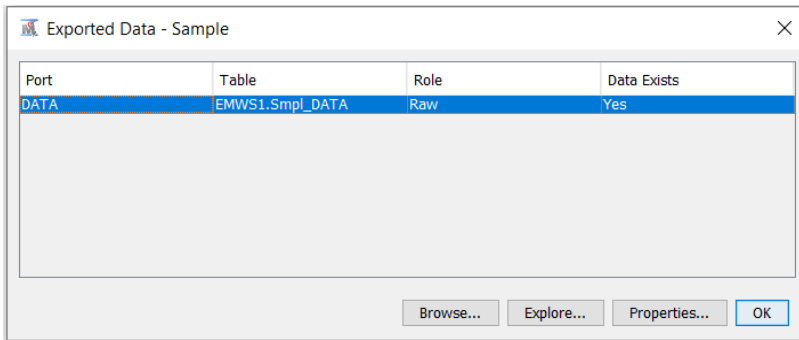


6. Now open the **Help** reference of the Sample node by selecting **File → Help → Content**, then in the navigate pane of the Help page, double click on **Sample**, then click on **Sample Node**. Read the overview of the Sample node, and the content related to random sampling, and the other sampling methods provided by SAS EM.
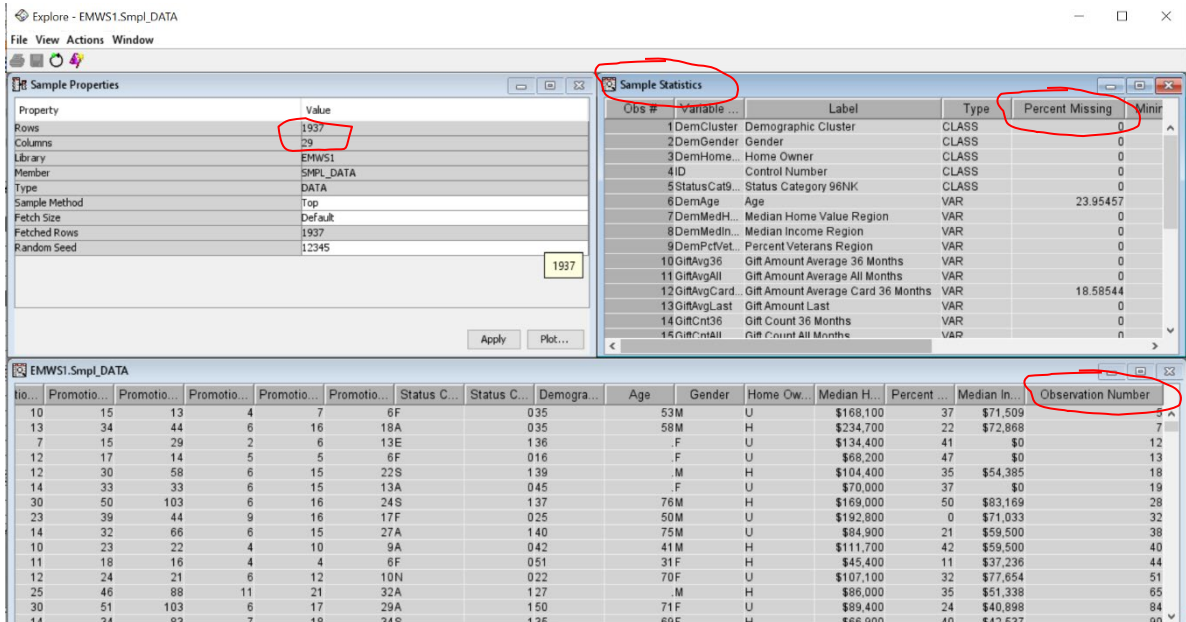


7. Explore the sampled data.

In the **Properties** panel of the Sample node (i.e. when the Sample node in the diagram is highlighted), click on the **…** (ellipse) button shown next to **Exported Data**, the **Exported Data – Sample** window will appear. In this window, click on the first row of the table (there is only 1 row in the table, in addition to the column headers), as indicated in the figure below, then click on the **Explore** button.

The **Explore** window shows:

Note that the number of samples is now shown as 1937, 20% of the samples in the original data set. Interestingly, the number of variables is 29, while the original dataset only contains 28 variables. Drag the scroll bar at the bottom, you will find that an extra column (variable) named **Observation Number** is added, which shows the record ID of each sample in the original data set. (You can expand a column to see the full column header name.)

Now have a look at the **Sample Statistics** shown in the top right window. You will see that there are missing values in the sampled data set too. Move the scroll bar on the right of the **Sample Statistics** window, you will notice that the three variables, **Age**, **Gift Amount Average Card 36 Month**s, and **Target Gift Amount** have 23.95457%, 18.58544% and 49.6128% of missing values respectively.

8. Export the sampled data set to your local computer as a **.csv** file.

   Close the **Explore** window by clicking on the **x** button at the top right corner of the **Explore** window (not the **x** in the **Sample Statistics** window). The **Exported Data – Sample** window will show again. Click on the **Browse** button (while the first row is still selected).

   A data browse window appears with the sampled data records, as indicated below. Move the scroll bars on the right and at the bottom to browse the data records. Then right click mouse anywhere on the data browse window and choose **Export to Excel**.
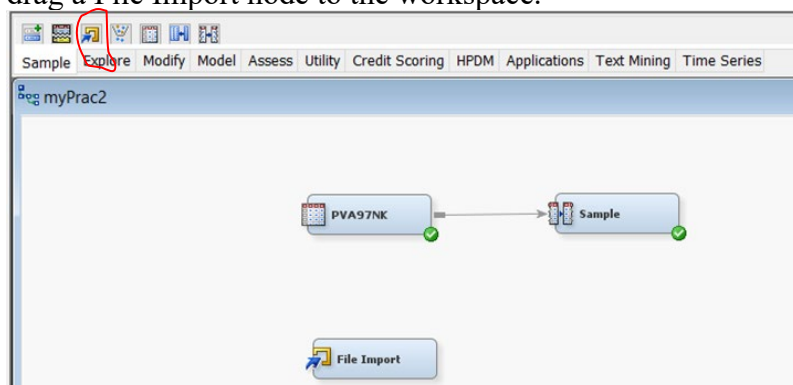
The Excel app on your computer will be opened and choose **YES** in the popup window to continue. The sampled data is now shown in the Excel app. Use the **Save As** menu of Excel to save the data set as a **.csv** file on your computer, and name this .csv file as **prac#2-sample.csv**. Remember the folder where you save this file.

Click **OK** in the **Exported Data – Sample** window to close it.

9. Now let us create another process flow to import the saved .csv file and preprocess the data. (Note that you could connect nodes to the Sample node to preprocess the sampled data, but here we learn how to use the File Import node to load your own data, so we will create a new process flow which starts with a File Import node.)
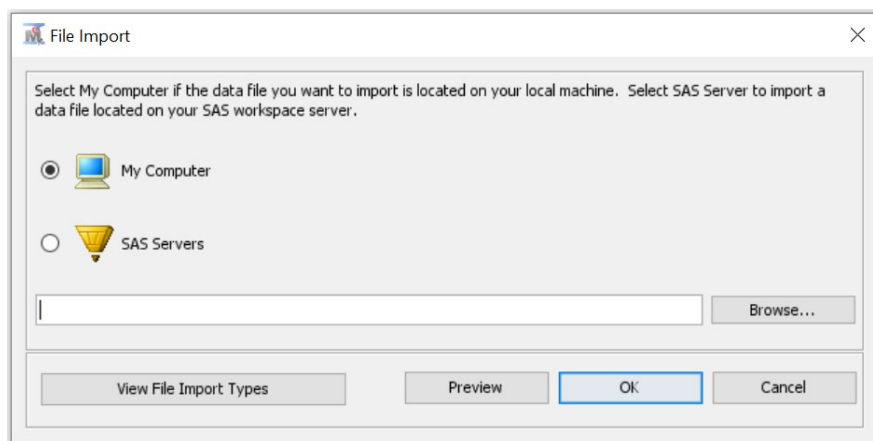
   In the same diagram workspace, make sure that the **Sample** tab in the tool palette is clicked, and drag a File Import node to the workspace.



10. While the File Important node is highlighted, in the **Properties** panel, click on the **...** button next to **Import File**.
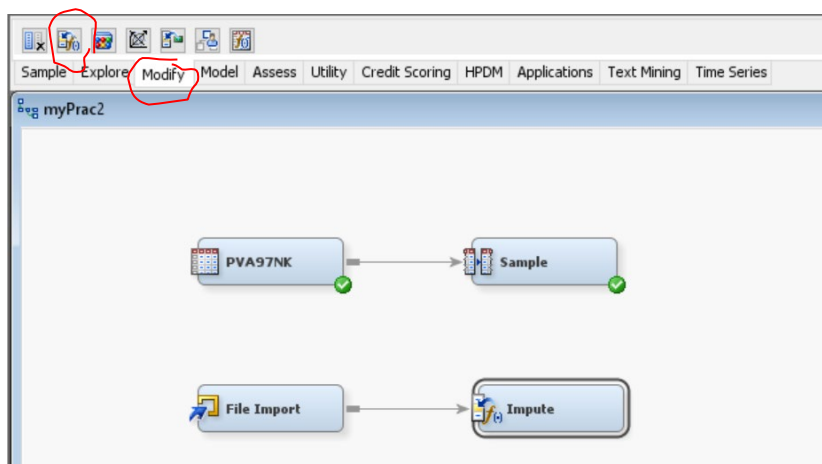


The **File Import** window shows. In the window, choose **My Computer**, and then **Browse**, and select the prac#2-sample.csv file you saved in Step 8 above, and finally click **OK** in the window to import the .csv file (the sampled data set).

**11.** Impute missing values using the Impute tool.

Recall that in **Step 7** above, we noted that in the sampled data set, three variables had missing values. In this step, we use the Impute tool to impute the missing values.

Click on the **Modify** tab in the tools palette, then drag the Impute node and drop it onto the workspace, and connect it with the File Import node.
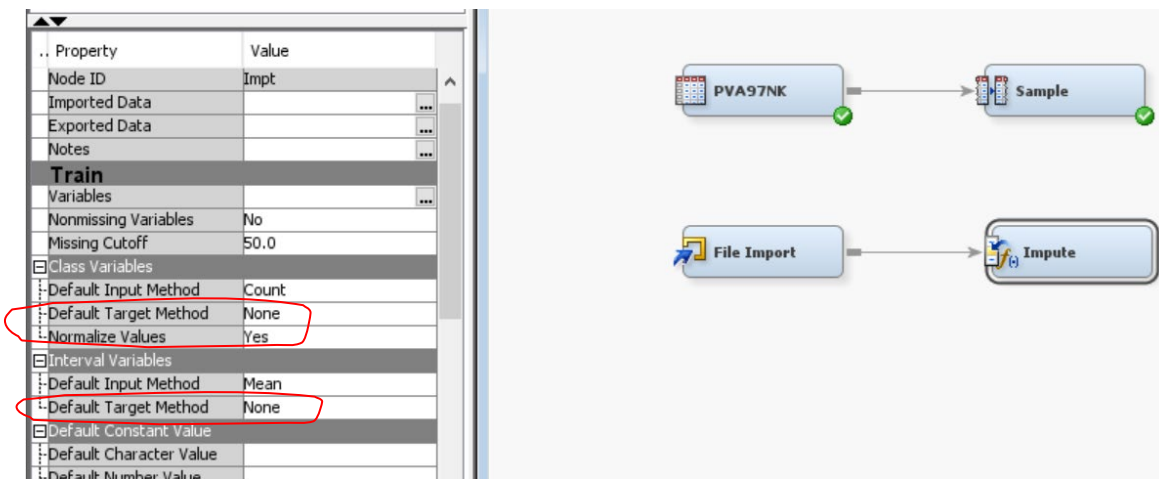


In the **Properties** panel, put your mouse over each item (or click on an item), and read the text shown in the yellow box (or in the **Property help** panel when you click on an item), to understand the meanings/purposes of these items. Also click on the white fields to see the options listed in the dropdown menu shown.
In the **Properties** panel of the Impute node, in the **Class Variables** section, change the **Default Target Method** from **None** to **Count**, and in the **Interval Variables** section, change the **Default Target Method** from **None** to **Mean.** (Note that in the PVA97NK dataset, variable **Targe Gift Amount** is set as a target variable for training a predictive model. By default, the Impute node does not imput missing values for a target variable, but only as a pratice of using the Impute node in this practical, we want to impute all missing values, so we change the default settings for target variables.)

Also in the **Class Variables** section, change **Normalized Values** from **Yes** to **No**.

Now right click the Impute node and choose **Run**. When the Run Status window pops up, select **Results** to view the Results. Note that as expected, missing values of all the three variables have been imputed using the specified method (i.e. imputing with mean).



**12.** Use the Transform Variables node to do data normalisation.

Drag a Transform node (when the **Modify** tab is selected), drop it onto the workspace, and connect it to the Impute node.
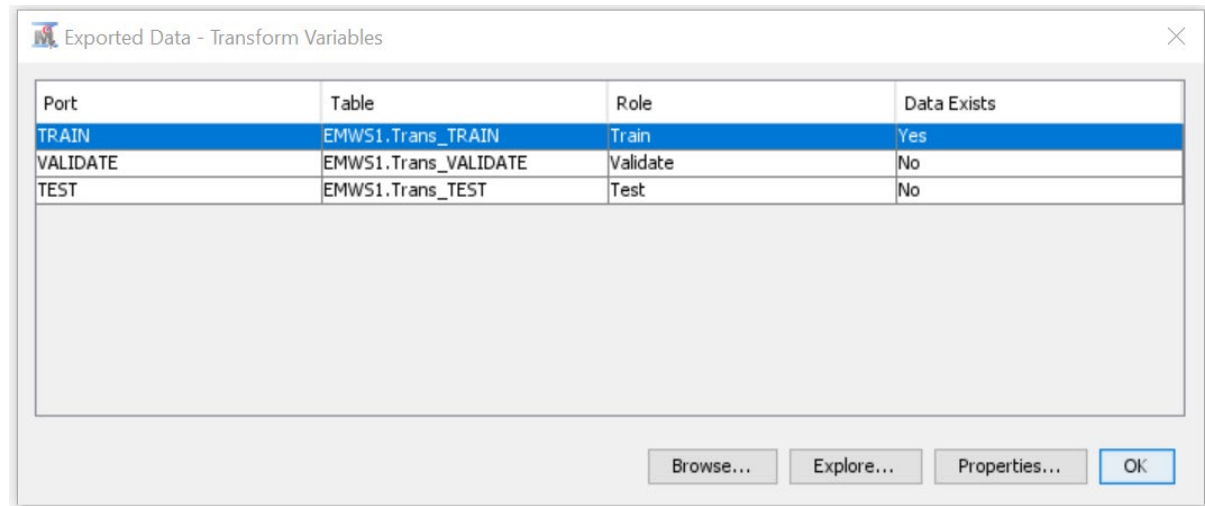


In the **Properties** panel, put your mouse over each item (or click on an item), and read the text shown in the yellow box (or in the **Property help** panel when you click on an item), to understand the meanings/purposes of these items. Also click on the white fields to see the options listed in the dropdown menu shown.

Now in the **Properties** panel, in the **Default Methods** section, change the method for **Interval Inputs** and **Interval Targets** from **None** to **Standardize**. This will normalise an interval variable (with SAS terminology, interval variable means continuous variables) by subtracting
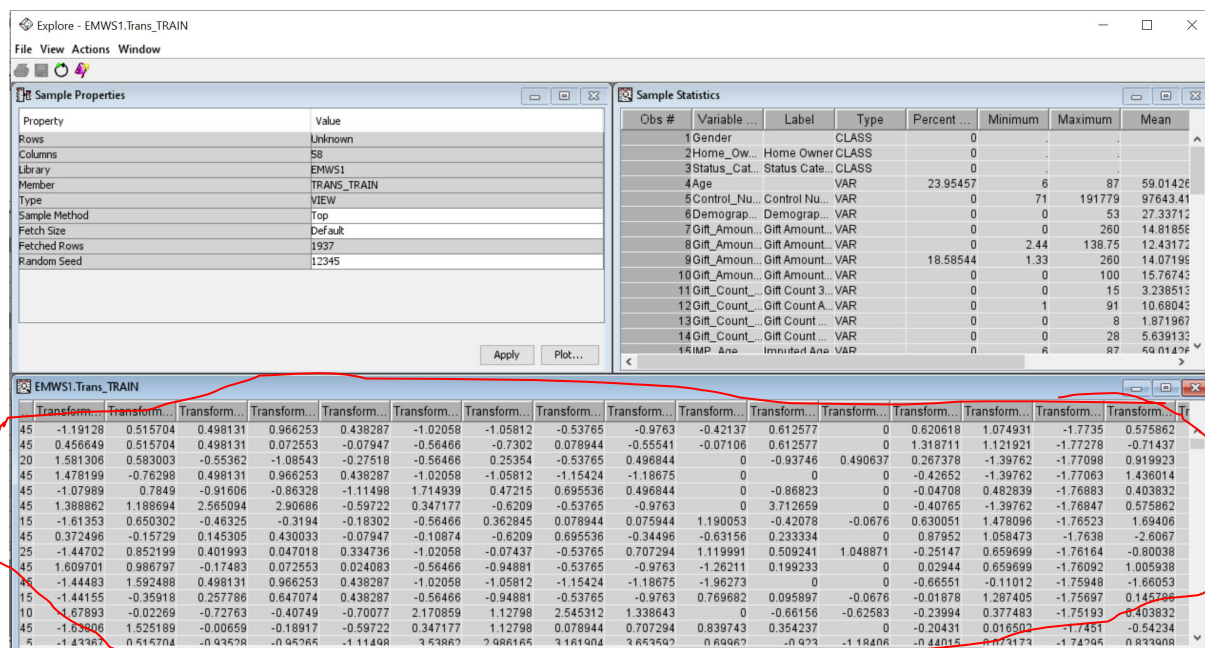
the mean and dividing by the standard deviation, which is the Z-score normalisation method we have learned in lecture. (Note that you can also specify different normalisation methods for individual variables by clicking on the **…** button next to **Variables**, then in the Variables window opened changing the method for an individual variable.)

Now click on the Transform Variables node and choose **Run**.

After the Run has completed, in the **Properties** panel (when the Transform Variable node is highlighted), click on the **…** button next to **Exported Data**. When the **Exported Data – Transform Variables** window shows, click on the first row, and click on the **Explore** button.
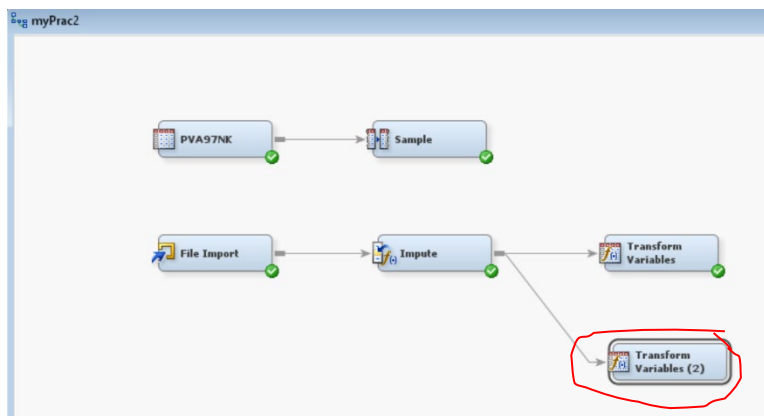


Now the **Explore** window shows. Move the scroll bar at bottom, you will see the Transformed values of the (interval) variables.



Drag a Transform Variable node, drop it onto the workspace, and connect it to the Impute node.
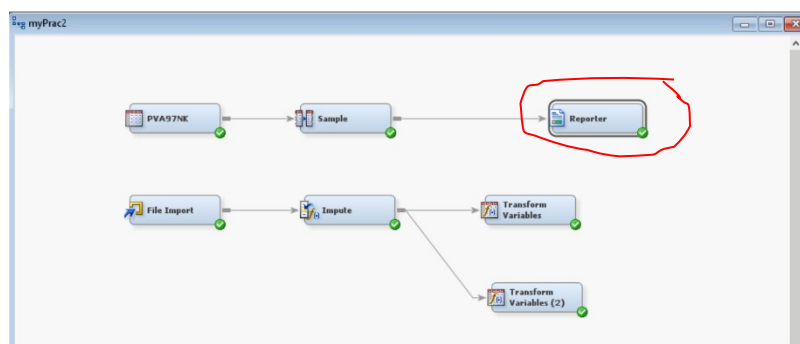
When the newly added Transform Variables node is highlighted, in the **Properties** panel, in the **Default Methods** section, change the transform methods for **Interval Inputs** and **Interval Targets** from **None** to **Range**. This normalisation method will transform a variable using the min-max method we have learned in lecture, i.e. the transformed value is equal to *(x - min) / (max - min)*, where *x* is current variable value, *min* is the minimum value for that variable, and *max* is the maximum value for that variable.

Now run this new Transform Variable node, and use the same way described above to see the transformed values, and observe the difference between the obtained transformed values from those obtained using the Z-score method.

13. Use the Reporter node to produce a report for Practical #2 submission.

Click on the **Utility** tab in the tools palette and then choose the Reporter node, drag it to the workspace, and connect it to the last node of any of the two process flows you have built.



In the **Properties** panel of the Reporter node, in the **Train** section, **change Path to All** (**Note**: this change is crucial for Practical #2 submission, otherwise the report generated will only contain the information of the process flow/path to which the Reporter node is connected, instead of the information of all the process flows.)



Right-click the Reporter node and choose Run to generate the report. After the Run has completed (for all process flows), the Run Status window will show. Click **OK**.

In the **Properties** panel of the Reporter node, click the **…** button in the **View Report** field, the report will be opened in Adobe Reader. **Save the PDF report, and submit it by following the instruction given on page 1 of the practical document.**