# BAN210 – Final Assessment

# Using Predictive Modeling on the Breast Cancer Dataset in SAS Enterprise Miner

**April Paola Tolosa**

Prof. Uzair Ahmad
Predictive Analytics BAN210-ZBB

**Date:** April 15, 2022

# Table of Contents

## Academic Integrity Policy

# SENECA COLLEGE OF APPLIED ARTS AND TECHNOLOGY
# SENECA BUSINESS

**Academic Integrity Policy**. Seneca upholds a learning community that values academic integrity, honesty, fairness, trust, respect, responsibility and courage. These values enhance Seneca's commitment to students by delivering high-quality education and teaching excellence, while supporting a positive learning environment. The AI policy is always in effect. Note **Sections 2.3 and 2.4:**

*"…2.3 Should there be a suspected violation of this policy (e.g.…cheating, falsification, impersonation or plagiarism), the academic integrity sanctions will be applied according to the severity of the offence committed. Refer to Appendix B for the academic integrity sanctions. 2.4 Should a suspected violation of this policy be a result of, or in combination with, a suspected violation of Seneca's Student Code of Conduct and/or another non-academic-related Seneca policy, the matter will be investigated and adjudicated through the process found in the Student Code of Conduct."*

| TO BE COMPLETED BY STUDENT |
| --- |
| *SUBJECT SECTION NUMBER (e.g. QNM223 AA):*   BAN210 ZBB |
| *STUDENT NAME:*   April Tolosa |
| *STUDENT NUMBER:*   131785214 |
| *STUDENT SIGNATURE:*   April Tolosa |

I, **April Tolosa**, declare that the attached assignment is my own work in accordance with the Seneca Academic Policy. I have not copied any part of this assignment, manually or electronically, from any other source including web sites, unless specified as references. I have not distributed my work to other students.

## Introduction

In this assessment paper, the breast cancer data set will be analyzed using exploratory and modeling techniques using SAS Enterprise Miner. The goal of this paper is to show the steps and the powerful insights to be gained from the dataset using simple statistical analysis. Steps will be explained as well as the results to show ability to use the tool and to demonstrate data analysis skills.

## About the Data

The Breast Cancer data set was obtained from the University Medical Centre, Institute of Oncology, Ljubljana, Yugoslavia. It includes 201 instances of one class and 85 instances of another class. The instances are described by 10 attributes, some of which are linear, and some are nominal.

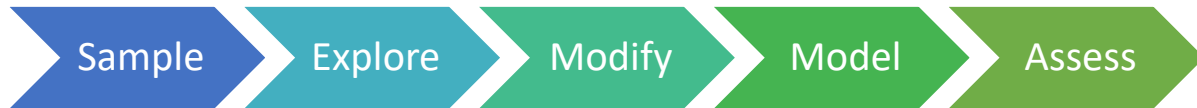| Data set Characteristics: | Multivariate | Number of Instances: | 286 | Area: | Life |
|---|---|---|---|---|---|
| Attribute Characteristics: | Categorical | Number of Attributes: | 10 | Date Donated | 1988-07-11 |
| Associated Tasks: | Classification | Missing Values? | Yes | Number of Web Hits: | 622435 |

**Attribute Information:**
1. Class: no-recurrence-events, recurrence-events
2. age: 10-19, 20-29, 30-39, 40-49, 50-59, 60-69, 70-79, 80-89, 90-99.
3. menopause: lt40, ge40, premeno.
4. tumor-size: 0-4, 5-9, 10-14, 15-19, 20-24, 25-29, 30-34, 35-39, 40-44, 45-49, 50-54, 55-59.
5. inv-nodes: 0-2, 3-5, 6-8, 9-11, 12-14, 15-17, 18-20, 21-23, 24-26, 27-29, 30-32, 33-35, 36-39.
6. node-caps: yes, no.
7. deg-malig: 1, 2, 3.
8. breast: left, right.
9. breast-quad: left-up, left-low, right-up, right-low, central.
10. irradiat: yes, no.

Before analysis, labels were added on the top row and the data cleaned as there were conversion errors from number to dates when the dataset was changed to a .csv file:
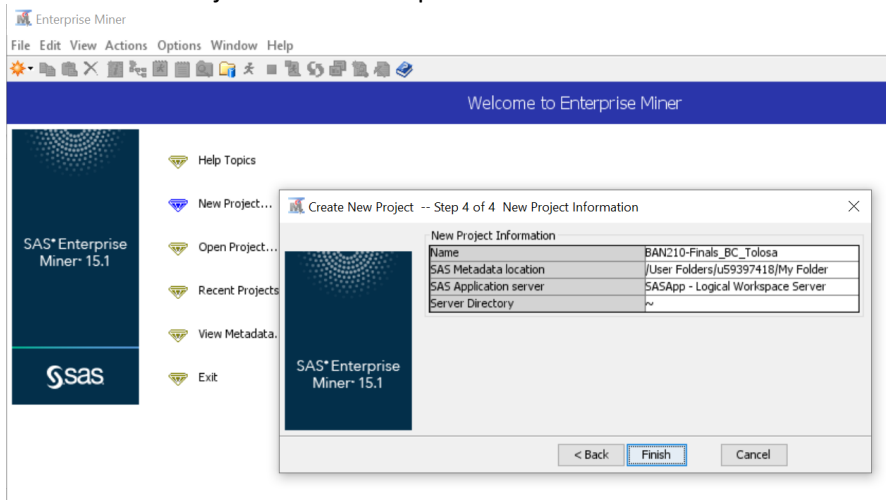
## Analysis

This section will demonstrate how the dataset is analyzed using SAS Miner guided by the data mining process learned from this course
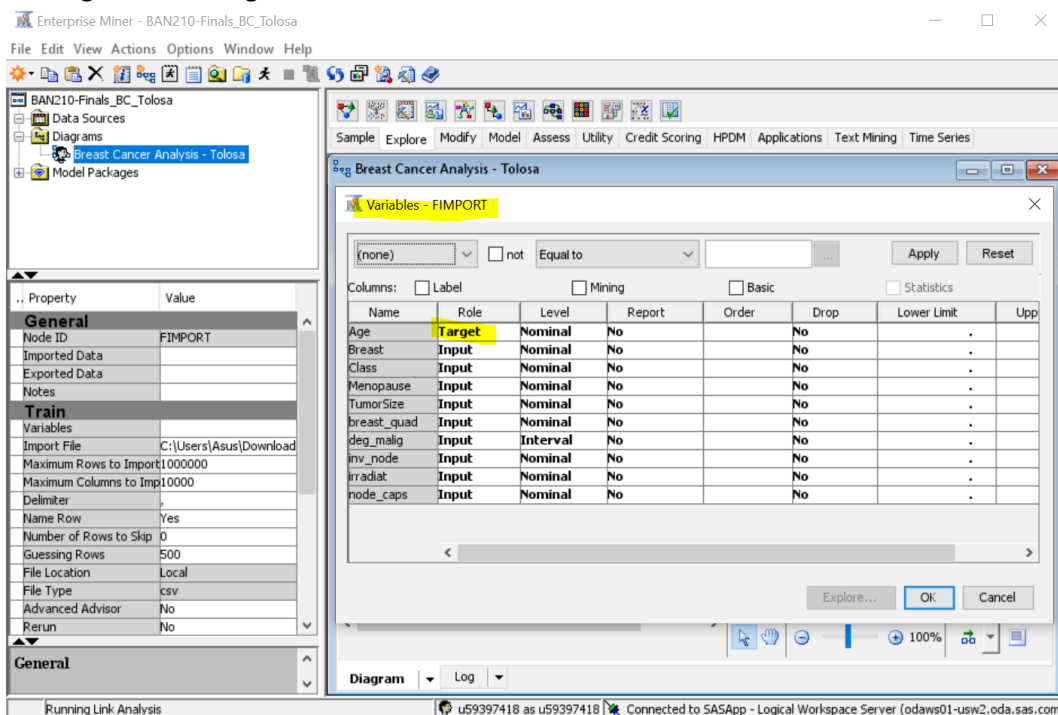
Sample → Explore → Modify → Model → Assess

Create a New Project on SAS Enterprise Miner



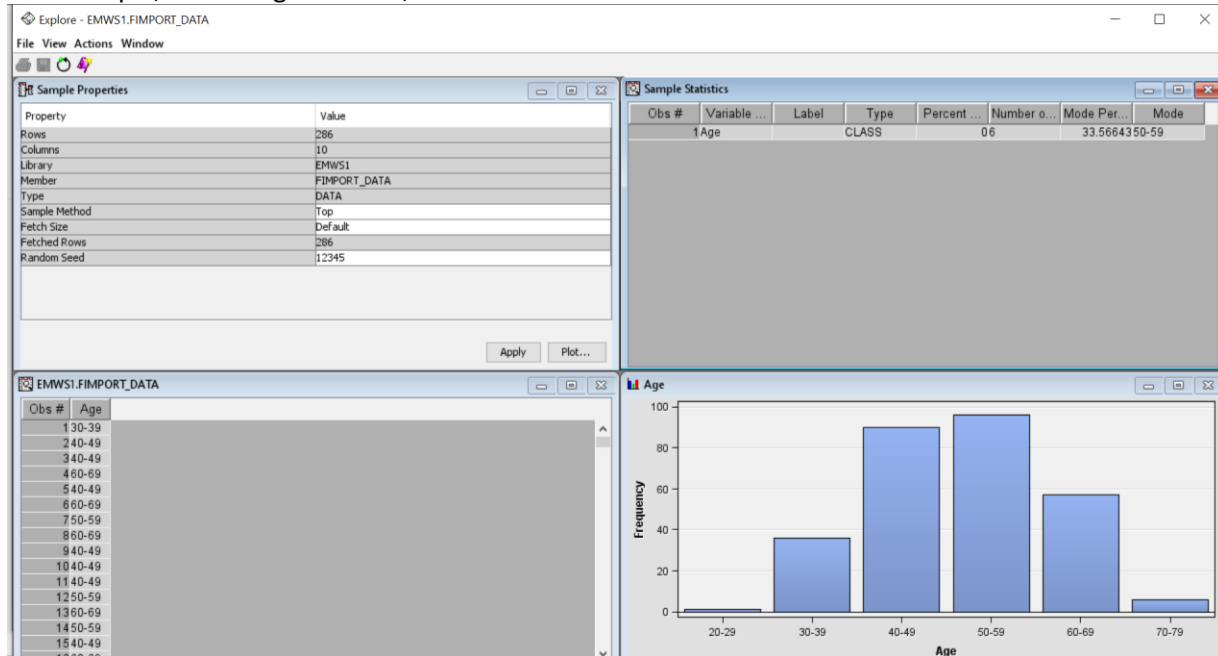Import the file using **File Import** node.
Right click on the File Import node and click Variables to open the Variables window.
Set "Age" as the Target Variable.

In this window we can click Explore to know more about the data of each variable.

For example, for the Age variable, below is a screenshot of the data:



❖ The data tells us that there are 6 classifications for age, with the age range of 50-59 being the highest frequency at 33.56%.

Exploring all variables, below are the results derived:

The File Import node is ran to check about the data. Then right click to check the Results.

❖ We can validate we have 286 observations, and 10 variables:

Results - Node: File Import  Diagram: Breast Cancer Analysis - Tolosa

File  Edit  View  Window

**Output**

```
24      The CONTENTS Procedure
25
26      Data Set Name          EMWS1.FIMPORT_DATA                    Observations           286
27      Member Type            DATA                                  Variables              10
28      Engine                 V9                                    Indexes                0
29      Created                12/04/2022 02:04:13                   Observation Length     72
30      Last Modified          12/04/2022 02:04:13                   Deleted Observations   0
31      Protection                                                   Compressed             NO
32      Data Set Type                                                Sorted                 NO
33      Label
34      Data Representation  SOLARIS_X86_64, LINUX_X86_64, ALPHA_TRU64, LINUX_IA64
35      Encoding               utf-8  Unicode (UTF-8)
36
37
38                                    Engine/Host Dependent Information
39
40      Data Set Page Size         131072
41      Number of Data Set Pages   1
42      First Data Page            1
43      Max Obs per Page           1816
44      Obs in First Data Page     286
45      Number of Data Set Repairs 0
46      Filename                   /home/u59397418/BAN210-Finals_BC_Tolosa/Workspaces/EMWS1/fimport_data.sas7bdat
47      Release Created            9.0401M6
48      Host Created               Linux
49      Inode Number               112901334
50      Access Permission          rw-r--r--
51      Owner Name                 u59397418
52      File Size                  256KB
53      File Size (bytes)          262144
```

Information about the variables is shown like the type, length, and format:

Results - Node: File Import  Diagram: Breast Cancer Analysis - Tolosa

File  Edit  View  Window

**Output**

```
56              Alphabetic List of Variables and Attributes
57
58      #      Variable        Type    Len    Format    Informat    Label
59
60      2      Age             Char    5      $5.       $5.
61      5      Breast          Char    5      $5.       $5.
62      1      Class           Char    20     $20.      $20.
63      3      Menopause       Char    7      $7.       $7.
64      4      TumorSize       Char    6      $6.       $6.
65      10     breast_quad     Char    9                           breast-quad
66      9      deg_malig       Num     8                           deg-malig
67      7      inv_node        Char    6                           inv-node
68      6      irradiat        Char    3      $3.       $3.
69      8      node_caps       Char    3                           node-caps
```
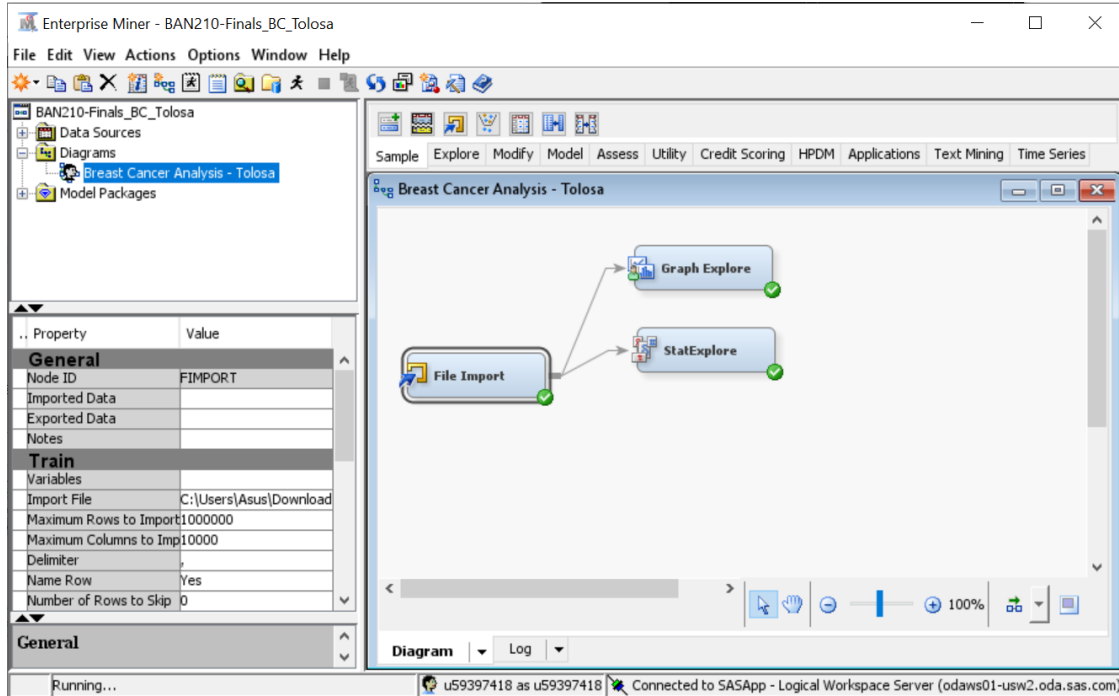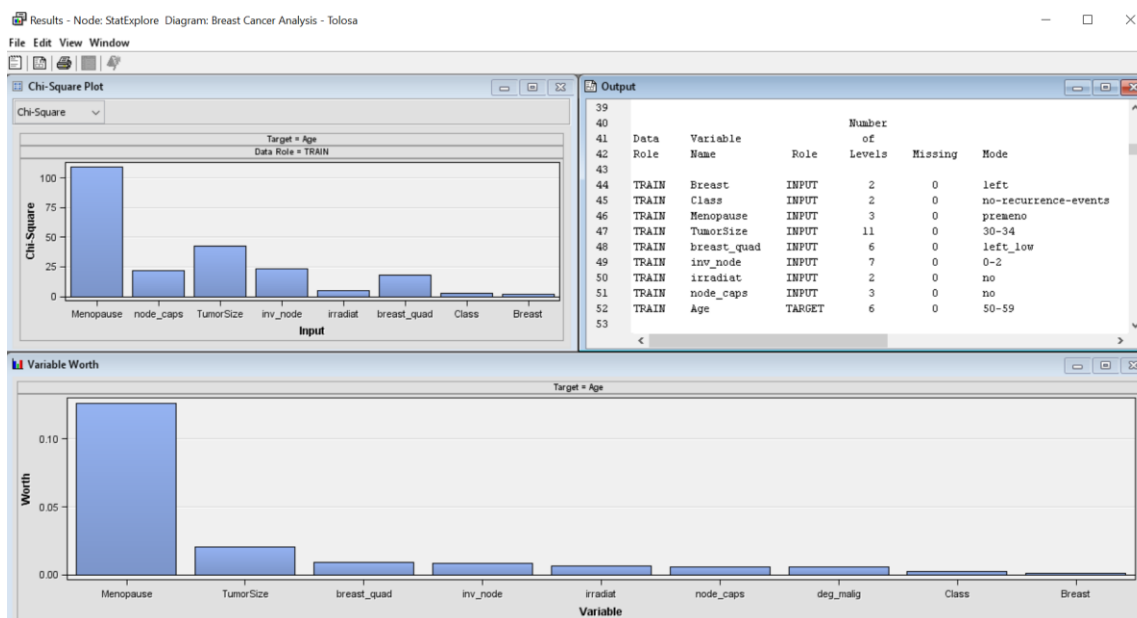
## Exploratory Data Analysis

Next is the Exploratory Data Analysis (EDA) to discover patterns, spot anomalies, test hypothesis and to check assumptions with the help of summary statistics and graphical representations.

Here, the **Graph Explorer** and **StatExplorer** nodes are used.



### StatExplorer Results

The StatExplorer is a versatile tool that gives the variable distributions and statistical information of the dataset. The Chi-Square, Variable Worth, and Output is shown here.

Concentrating on the Output, more insight is gained on the statistics of how the target Age is compared with other variables.

An example below is the statistics for Age vs. Breast and Age vs. Class.

❖ We can say that for Age 40-49, 51.61% developed tumor on the left breast and 48.39% developed tumor on the right breast.

```
Output
85    Class Variable Summary Statistics by Class Target
86    (maximum 500 observations printed)
87
88    Data Role=TRAIN Variable Name=Breast
89
90                      Number
91             Target    of                    Mode            Mode2
92    Target   Level   Levels  Missing  Mode  Percentage  Mode2  Percentage
93
94    Age      20-29      1       0     right   100.0            0.00
95    Age      30-39      2       0     left    50.00   right   50.00
96    Age      40-49      2       0     left    51.61   right   48.39
97    Age      50-59      2       0     left    56.06   right   43.94
98    Age      60-69      2       0     left    51.28   right   48.72
99    Age      70-79      2       0     left    66.67   right   33.33
100   _OVERALL_           2       0     left    52.82   right   47.18
101
102
103   Data Role=TRAIN Variable Name=Class
104
105                      Number
106            Target    of                       Mode                    Mode2
107   Target   Level   Levels  Missing      Mode         Percentage      Mode2        Percentage
108
109   Age      20-29      1       0    no-recurrence-events   100.0                       0.00
110   Age      30-39      2       0    no-recurrence-events   58.33   recurrence-events   41.67
111   Age      40-49      2       0    no-recurrence-events   64.52   recurrence-events   35.48
112   Age      50-59      2       0    no-recurrence-events   72.73   recurrence-events   27.27
113   Age      60-69      2       0    no-recurrence-events   69.23   recurrence-events   30.77
114   Age      70-79      2       0    no-recurrence-events   66.67   recurrence-events   33.33
115   _OVERALL_           2       0    no-recurrence-events   67.69   recurrence-events   32.31
116
```
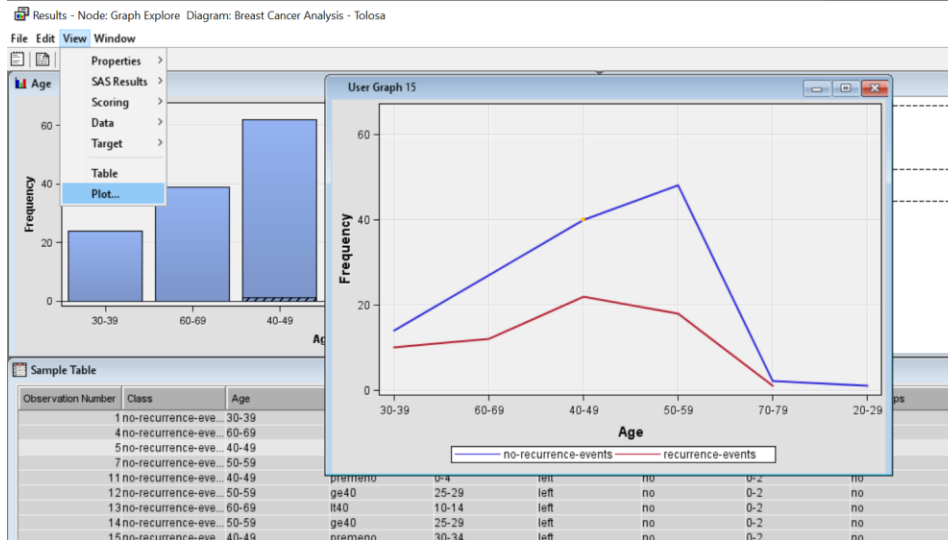
## Graph Explorer Results

The graph explorer shows more statistics about the target variable Age.

❖ This confirms that the most frequent data is that from the 50-59 age group and the least for the 20-29 age group.
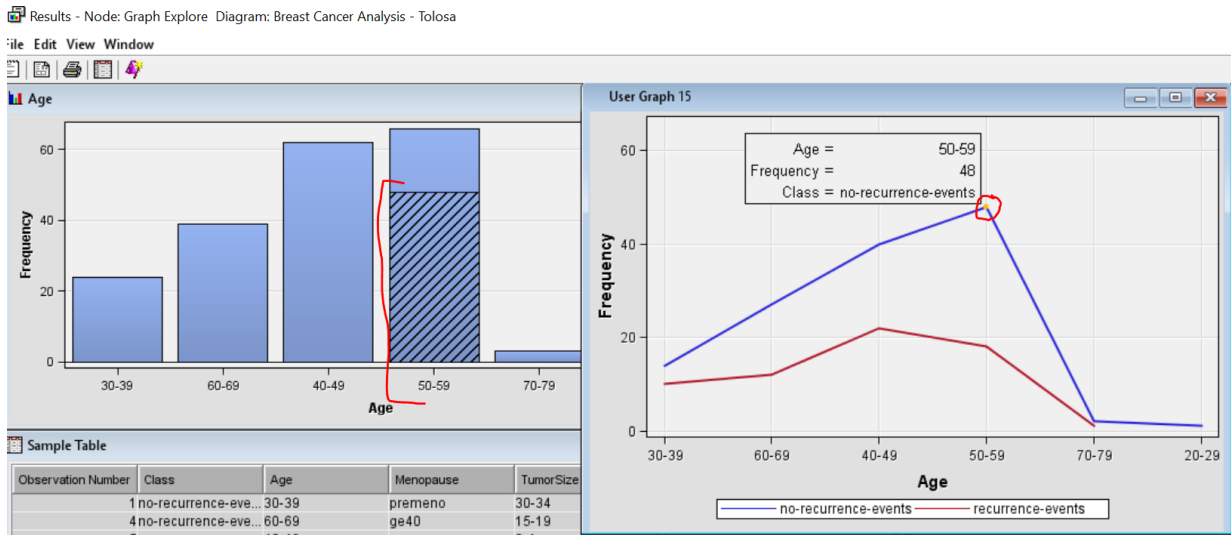
Graphs can also be made by going to View and the Plot. As an example, below is a Line graph with the Age set to Category role and the Class as the Group role

❖ The line graph shows us that for all Age range, the tumor being no-recurring (blue line) is greater than recurring tumors.
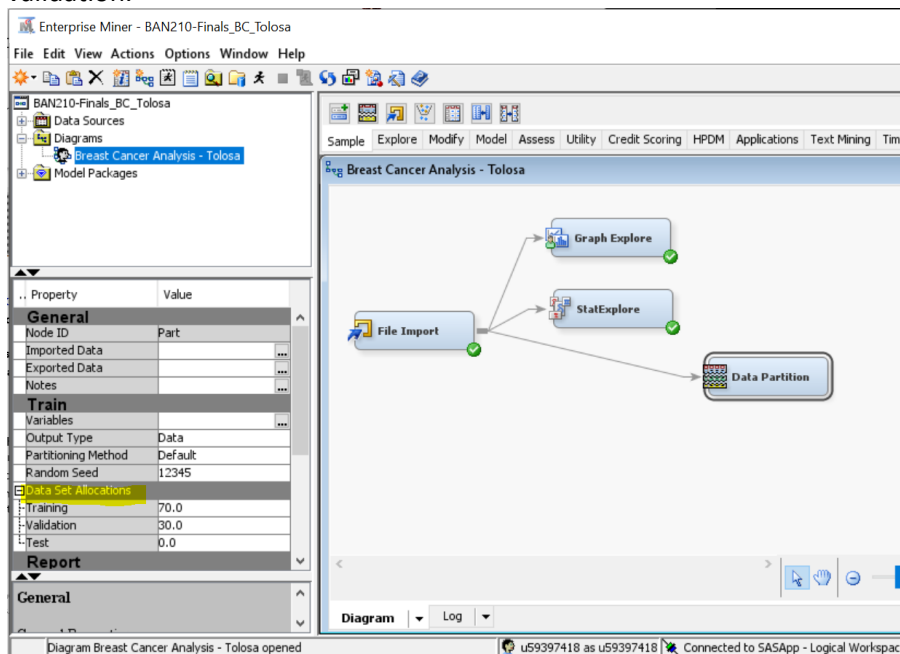


❖ There are 48 observations of no-reccurence at the age 50-59 which is the highest frequency. As this is pointed, the frequency graph on the left also shades the corresponding bar that represents this figure.

## Predictive Models

**Data Partition node** was added, and the Data Set Allocation was set to 70 for training and 30 for validation.



### K-Fold Validation

**Start Group node** is added as a way to do the **k-fold validation,** to resample the dataset to ensure generalizability of the predictive models**.** It is on **Cross Validation** mode which exports the complement of the groups specified, as opposed to the groups themselves.

**Age** is set as the target role and **Class** is set to stratification group role.

Next is the completion of the model by connecting the 3 predictive models used namely, the **Decision Tree, Regression, and Neural Network** model nodes.

These three nodes are connected to the **Model Comparison and finally the End Groups** node.



Run the End Group Node.

*Model Comparison*

First, the **Model Comparison Results** is examined— focusing on the Misclassification Rate to see how fit the model is.

❖ Based on the Fit Statistics, the **Regression** is the selected model for the no-recurring group and **Neural Network** is the selected model for the recurring group, both with the Age as target variable.

These were picked to be the best fit because they had the lowest misclassification rate and the lower, the better because it means that the forecast is closer to the actual.

The Results of the three models is shown when clicking on the End Group Results. It shows statistics of all three models and their overall results.



The Fit Statistics plot can be navigated to show the bar graph for each fit statistic label.

## Conclusion

In this report, a simple exploratory data analysis was run wherein Age vs. Breast was first analyzed and a sample conclusion was mentioned as seen from the result, that is, for Age 40-49, 51.61% developed tumor on the left breast and 48.39% developed tumor on the right breast. This is a statistical analysis which can be presented to the business about the likelihood of the location of the tumor for those age group. With insights like this, looking at the bigger picture, medical diagnosis can be shifted when the patient's age is known as well as testing and prescription that might result to savings in time and money.

Another sample statistic that can be looked at is the mode of each variable, to see their likelihood:



Analyzing the chart above, it can be said to the business that the likelihood of developing breast cancer is at the age of 50-59 (33.5%) based on the data set. The tumor is also 53% likely to be on the left breast. It is, however 70% not recurring, and 52% in the premeno stage. There is a 20.9% change that the tumor size is 30-34, 38.5% chance it is in the left lower breast quadrant, 74% chance it is in the 0-2 inv_node, and 76% not irradiat and 77% likelihood that it has no node_caps.

When it comes to predictive analysis, three models were ran, namely, the Decision Tree, Logistic Regression, and Neural Network. Age was set as Target variable and Class as the Stratification in the grouping role.

After running the model comparison and checking the Fit Statistics, the **Regression** is the selected model for the no-recurring group and **Neural Network** is the selected model for the recurring group. Knowing which model to use is important because the goal is to minimize the error or difference in the prediction to the actual, in this case, the class.

## Work Cited

Zwitter, Matjaz and & Soklic, Milan. "Breast Cancer Data Set." *UCI Machine Learning*.
        https://archive.ics.uci.edu/ml/datasets/Breast+Cancer.