# SAS Assignment #3

Prepared By: Group #5

Joel Camacho

Mohamed Mazarul

Bradley John

The purpose of this project was to continue the knowledge enhancement of SAS Enterprise Miner. The project required us to apply decision tree nodes, partition the data and use the average square error as the model assessment statistic to answer a business problem for a supermarket with a new line of organic products. The supermarket has a customer loyalty program. As an initial buyer incentive plan, the supermarket provided coupons for the organic products to all the loyalty program participants and collected data that includes whether these customers purchased any of the organic products. Detail of these attributes is as follows.

**Business Problem**

A supermarket is offering a new line of organic products. The supermarket's management wants to determine which customers are likely to purchase these products.

**Questions & Building Model**

1. How do the purchase orders vary by gender?
2. Does age play a role in purchasing orders?

**Data description**

The ORGANICS data set contains:

- 13 variables and
- over 22,000 observations

| Variable | Role | Level | Discussion |
|----------|------|-------|------------|
| ID | ID | Nominal | Customer loyalty identification number |
| DemAffl | Input | Interval | Affluence grade on a scale from 1 to 30 |
| DemAge | Input | Interval | Age, in years |
| DemCluster | Rejected | Nominal | Type of residential neighborhood |
| DemGender | Input | Nominal | M = male, F = female, U = unknown |
| DemRegion | Input | Nominal | Geographic region |
| DemTVReg | Input | Nominal | Television region |
| PromClass | Input | Nominal | Loyalty status: tin, silver, gold, or platinum |
| PromSpend | Input | Interval | Total amount spent |
| PromTime | Input | Interval | Time as loyalty card member |
| TargetBuy | Target | Binary | Organics purchased? 1 = Yes, 0 = No |
| TargetAmt | Rejected | Interval | Number of organic products purchased |

Our first step was to import the data in EM using all the variables and set the roles as listed above then applied the required parameters to answer the business problem. The following steps were completed (see Appendices for screenshots)

1. Decision Tree nodes were created
2. A Data Partition node was added by assigning 50% of the data for training and 50% for validation

**Data Exploration**



The proportion of individuals who purchased organic products appears to be 24.3%.

**Implementations and Findings**

We created a decision tree model. We used <u>average square error</u> as the model assessment statistic. We also used <u>Subtree Assessment Plot</u> to assess the model.

- Using average square error as the assessment measure results in a tree with 29 leaves in Model 1
- Age is used for the first split
- Competing splits are Affluence Grade and Gender

For comparison purposes we created a second model that allowed an additional split at each node. In this model

- Using average square error as the assessment measure results in a tree with 32 leaves in Model 2
- Age is used for the first split
- The competing splits is Affluence Grade

**Conclusion**

Our data showed us that the solution to acquiring more customer purchases is to focus on two areas of the business. The first is to focus on women since the data shows they have a higher purchase rate. The second is to market towards middle aged customers. These customers are the top buyers and those are the ones that need to be targeted. Each of the models gave a similar final answer and could be used. We performed a model comparison (Appendix D) and the misclassification rates were nearly the same for each. Each model also performed similarly for a variety of other measures.

Appendices:

A: Overall SaS Diagram



B. Decision Tree1 Results:

B.1 – Output Only

```
File  Edit  View  Window

 Output
 46    PREDICTED      P_TargetBuy0     Predicted: TargetBuy=0
 47    RESIDUAL       R_TargetBuy0     Residual: TargetBuy=0
 48    FROM           F_TargetBuy      From: TargetBuy
 49    INTO           I_TargetBuy      Into: TargetBuy
 50
 51
 52    *------------------------------------------------------*
 53    * Score Output
 54    *------------------------------------------------------*
 55
 56
 57    *------------------------------------------------------*
 58    * Report Output
 59    *------------------------------------------------------*
 60
 61
 62
 63    Variable Importance
 64
 65                                                                              Ratio of
 66                                   Number of                                  Validation
 67    Variable                      Splitting                      Validation   to Training
 68    Name          Label           Rules        Importance        Importance   Importance
 69
 70    DemAge        Age             7            1.0000            1.0000       1.0000
 71    DemAffl       Affluence Grade 14           0.8477            0.8244       0.9726
 72    DemGender     Gender          7            0.5310            0.6371       1.1997
 73
 74
 75
 76    Tree Leaf Report
 77
 78                                     Training
 79    Node                  Training   Percent    Validation    Validation
 80    Id      Depth      Observations    1      Observations    Percent 1
 81
 82    43        5           1464        0.16       1462          0.14
 83    57        6           1271        0.05       1284          0.04
 84    54        6           1026        0.20       1043          0.23
 85    48        5            869        0.13        916          0.12
 86    56        6            659        0.01        630          0.02
 87    52        6            636        0.06        592          0.08
```
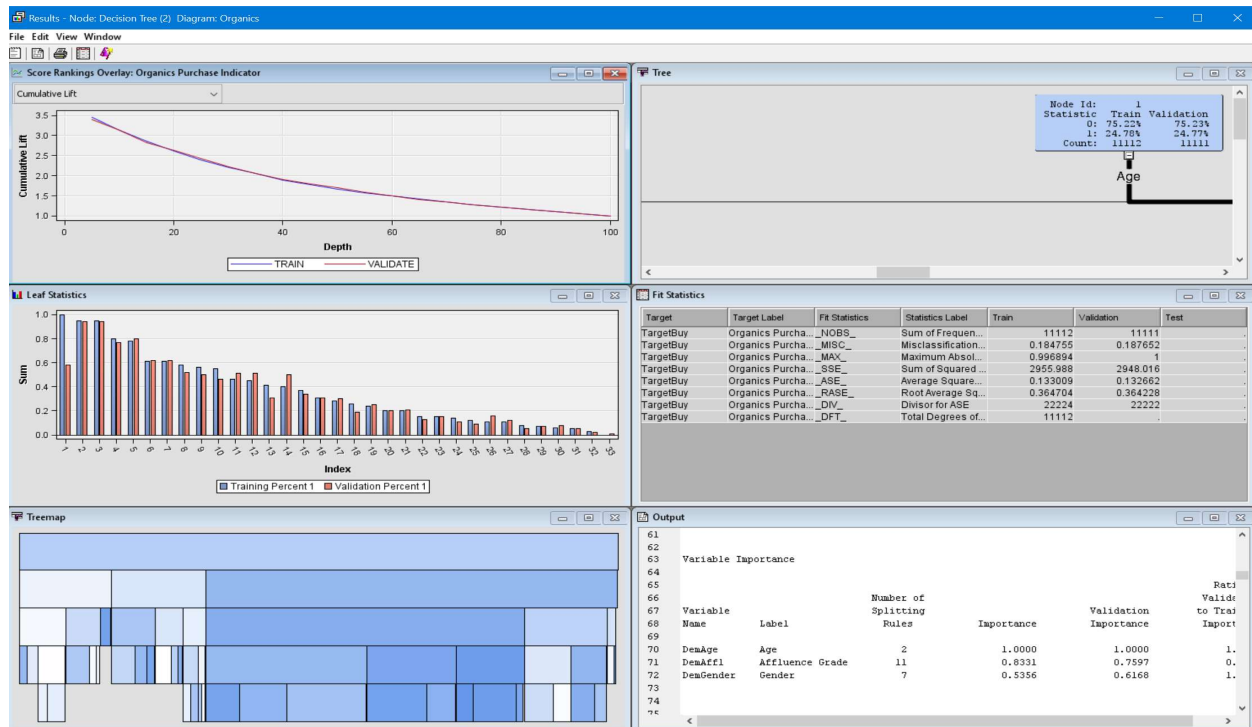
## C: Decision Tree 2



## C.1 – Model 2 Output Only

Output

```
40
41      Type             Variable        Label
42
43      TARGET           TargetBuy       Organics Purchase Indicator
44      PREDICTED        P_TargetBuy1    Predicted: TargetBuy=1
45      RESIDUAL         R_TargetBuy1    Residual: TargetBuy=1
46      PREDICTED        P_TargetBuy0    Predicted: TargetBuy=0
47      RESIDUAL         R_TargetBuy0    Residual: TargetBuy=0
48      FROM             F_TargetBuy     From: TargetBuy
49      INTO             I_TargetBuy     Into: TargetBuy
50
51
52      *------------------------------------------------------------*
53      * Score Output
54      *------------------------------------------------------------*
55
56
57      *------------------------------------------------------------*
58      * Report Output
59      *------------------------------------------------------------*
60
61
62
63      Variable Importance
64
65                                                                        Ratio of
66                                   Number of                           Validation
67      Variable                    Splitting                Validation  to Training
68      Name         Label            Rules     Importance   Importance  Importance
69
70      DemAge       Age               2          1.0000       1.0000      1.0000
71      DemAffl      Affluence Grade   11         0.8331       0.7597      0.9119
72      DemGender    Gender            7          0.5356       0.6168      1.1515
73
74
75
76      Tree Leaf Report
77
78                                   Training
79      Node               Training   Percent    Validation   Validation
80       Id     Depth    Observations    1      Observations   Percent 1
81
82       49       4         1473        0.20        1534         0.21
```

# D. Model Comparison

Model Selection based on Valid: Misclassification Rate (_VMISC_)

| Model Node | Model Description | Data Role | Target | Target Label | False Negative | True Negative | False Positive | True Positive |
|---|---|---|---|---|---|---|---|---|
| Tree | Decision Tree | TRAIN | TargetBuy | Organics Purchase Indicator | 1676 | 7978 | 381 | 1077 |
| Tree | Decision Tree | VALIDATE | TargetBuy | Organics Purchase Indicator | 1679 | 7979 | 380 | 1073 |
| Tree2 | Decision Tree (2) | TRAIN | TargetBuy | Organics Purchase Indicator | 1584 | 7890 | 469 | 1169 |
| Tree2 | Decision Tree (2) | VALIDATE | TargetBuy | Organics Purchase Indicator | 1593 | 7867 | 492 | 1159 |

```
Data Role=Valid

Statistics                                                              Tree        Tree2

Valid: Kolmogorov-Smirnov Statistic                                     0.50         0.50
Valid: Average Squared Error                                            0.13         0.13
Valid: Roc Index                                                        0.82         0.82
Valid: Bin-Based Two-Way Kolmogorov-Smirnov Probability Cutoff          0.26         0.24
Valid: Cumulative Percent Captured Response                            31.28        31.32
Valid: Percent Captured Response                                       14.15        14.28
Valid: Divisor for VASE                                             22222.00     22222.00
Valid: Gain                                                           212.54       212.94
Valid: Gini Coefficient                                                 0.65         0.65
Valid: Bin-Based Two-Way Kolmogorov-Smirnov Statistic                   0.49         0.49
Valid: Kolmogorov-Smirnov Probability Cutoff                            0.27         0.26
Valid: Cumulative Lift                                                  3.13         3.13
Valid: Lift                                                             2.83         2.85
Valid: Maximum Absolute Error                                           1.00         1.00
Valid: Misclassification Rate                                           0.19         0.19
Valid: Sum of Frequencies                                          11111.00     11111.00
Valid: Root Average Squared Error                                       0.36         0.36
Valid: Cumulative Percent Response                                     77.41        77.51
Valid: Percent Response                                                70.06        70.66
Valid: Sum of Squared Errors                                         2950.48      2948.02
```