# Contents

# WHERE TO START

## 1) OPEN A NEW PROJECT:

**OPEN A NEW PROJECT → PROJECT NAME (TYPE IN THE NAME) → SAS SERVER DIRECTORY →BROWSE → FIND THE LOCATION TO SAVE YOUR PROJECT → OPEN → NEXT → FINISH**
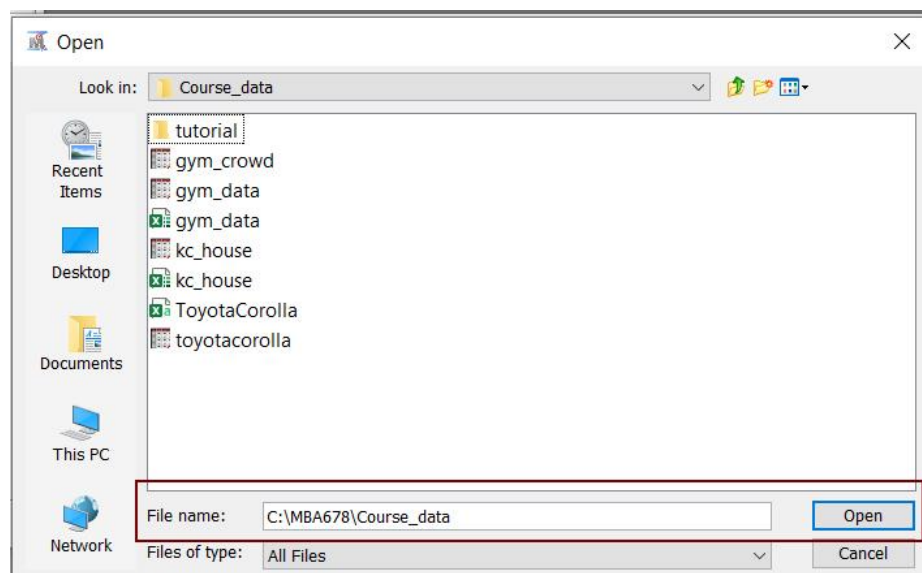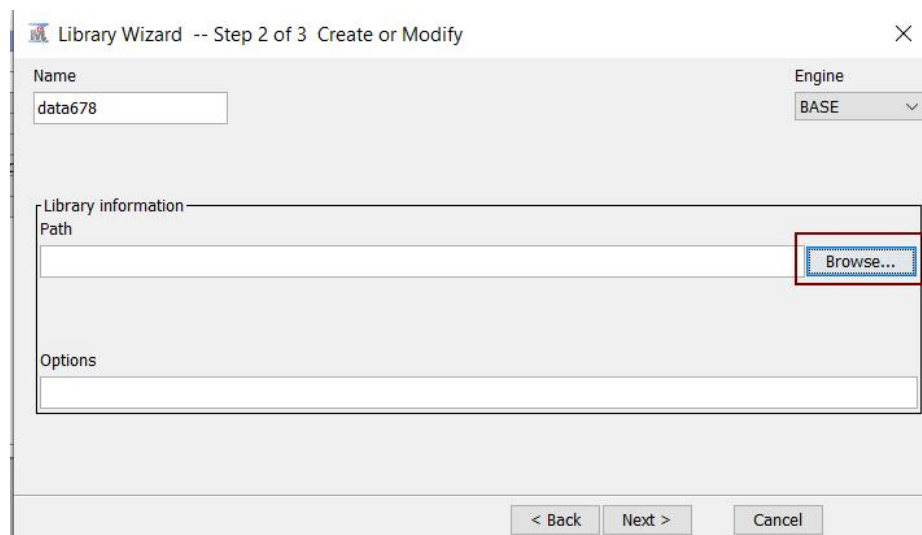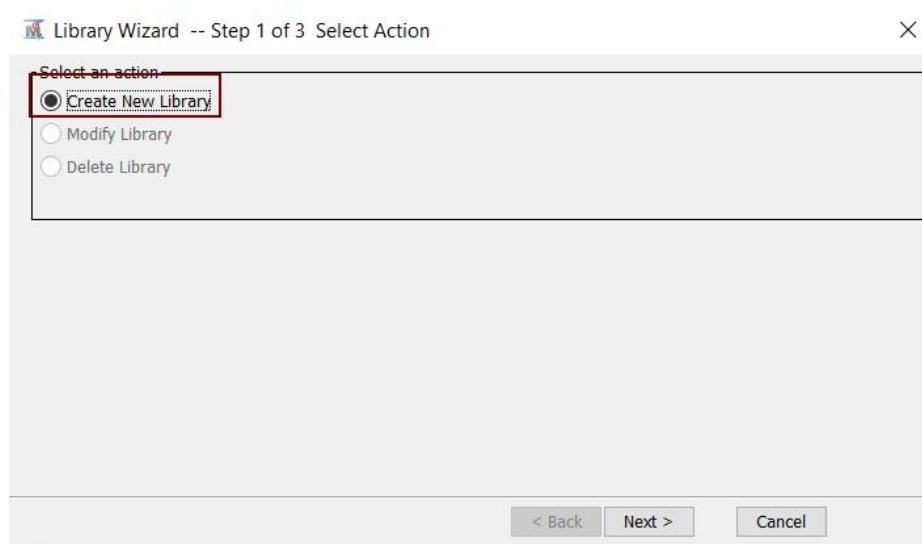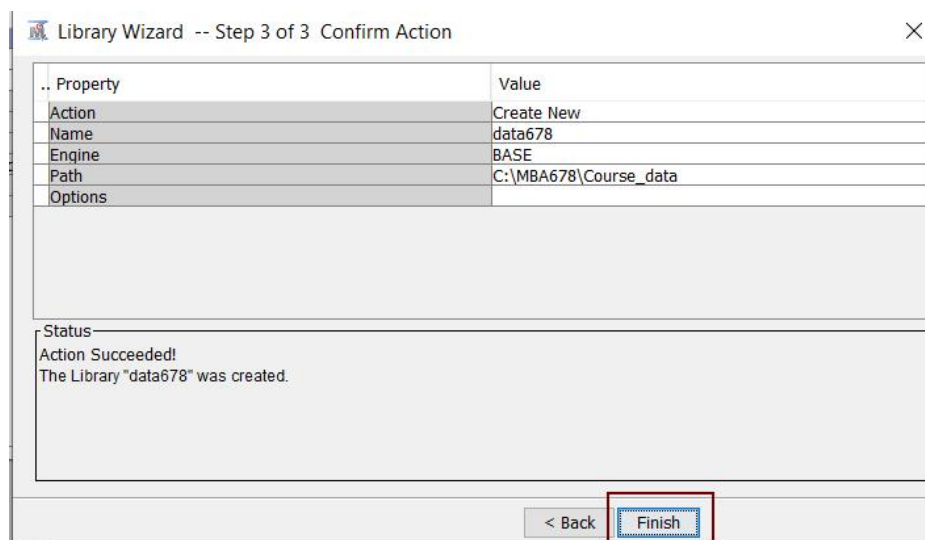
## 2) CREATE A NEW LIBRARY:

**LIBRARY → CREATE NEW LIBRARY → NEXT → NAME (TYPE IN THE LIBRARY NAME) → PATH: BROWSE (CHOOSE A LOCATION FOR YOUR NEW LIBRARY) → OPEN →
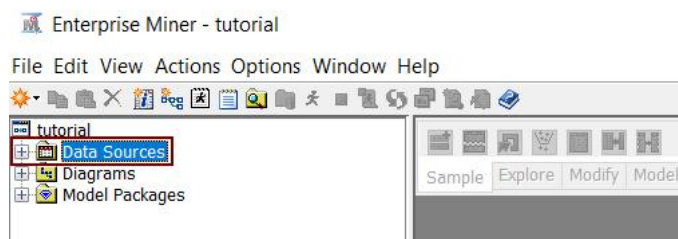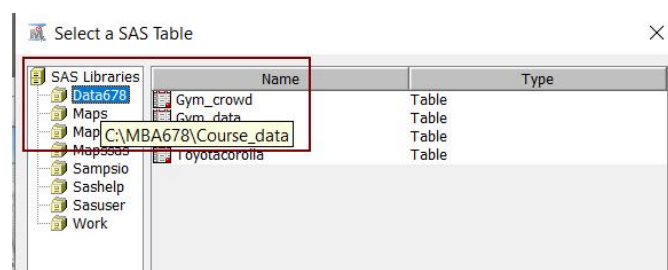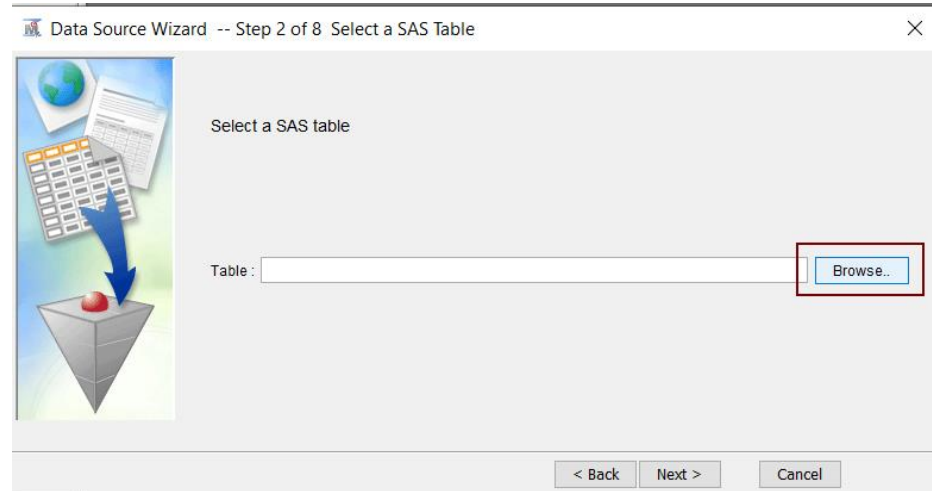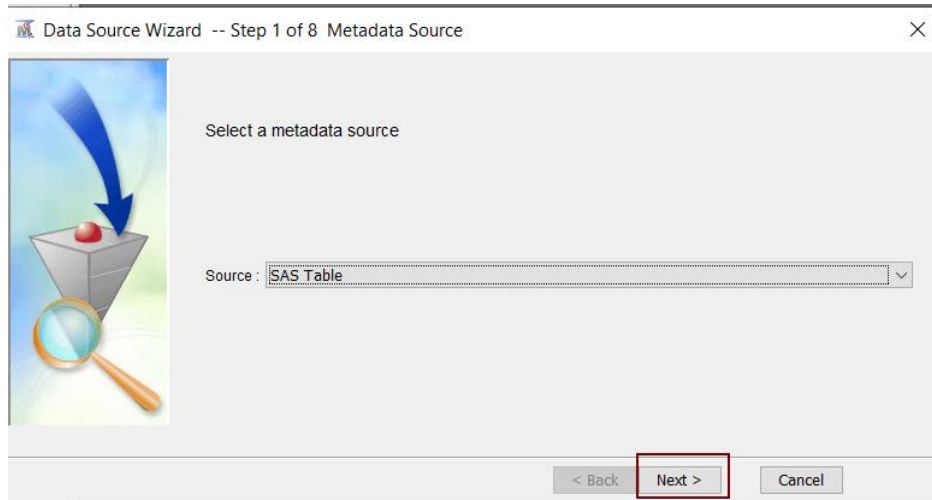NEXT → FINISH**

## 3) CREATE A DATA SOURCE:

*right-click on* **DATA SOURCES → left-click on CREATE DATA SOURCE → SAS TABLE (LEAVE AT DEFAULT) → NEXT → TABLE: BROWSE → CHOOSE THE LIBRARY(DOUBLE CLICK) → CHOOSE THE FILE → OK → NEXT → STEP 3: LEAVE AT DEFAULT → NEXT → METADATA ADVISOR OPTIONS: BASIC (LEAVE AT DEFAULT) → NEXT →** <span style="color:red">**STEP 5: ASSIGN VARIABLE ROLES (!!!SET THE BINARY VARIABLE YOU ARE TRYING TO PREDICT TO TARGET AND ITS LEVEL TO BINARY) →**</span> **CONTINUE CHANGING THE VARIABLE ROLES AND LEVELS IF NEEDED → NEXT → STEP 6 (LEAVE AT DEFAULT): NEXT → STEP 7 (LEAVE AT DEFAULT): NEXT → STEP 8 (LEAVE AT DEFAULT): NEXT → FINISH**

## 4) CREATE A DIAGRAM:
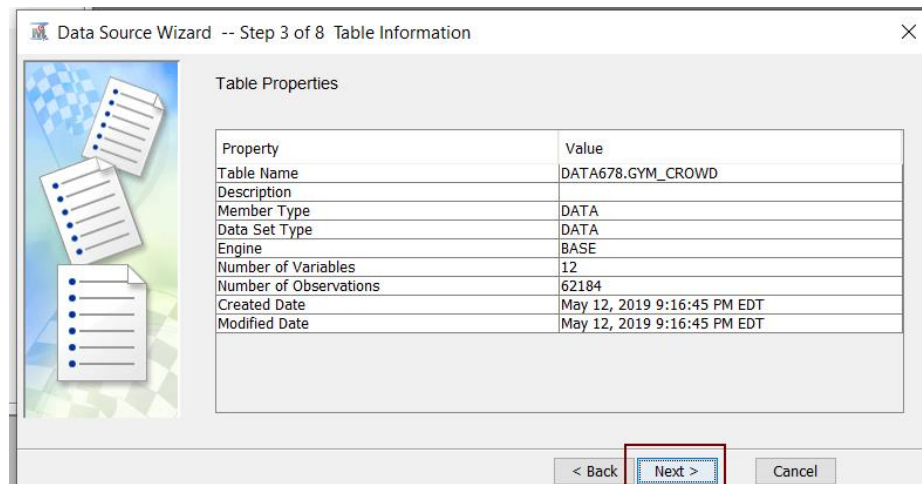**right-click on DIAGRAM → CREATE DIAGRAM → DIAGRAM NAME → OK**

# HOW TO EXPLORE DATA

## A) First way:
**Right-click on the data source and choose EXPLORE**



## B) Second way:
**OPEN THE DIAGRAM → DRAG THE DATA SOURCE INTO THE DIAGRAM → right-click on the file source you have just dragged → EDIT VARIABLES → SELECT THE VARIABLES → EXPLORE**

## PARTIONING THE DATA: CREATING A TRAIN AND VALIDATION DATASET

**SAMPLE (ON THE RIGHT TAB) → DRAG AND DROP DATA PARTITION → CONNECT THE TABLE AND THE DATA PARTITION NODES → LOOK AT THE RIGHT OF THE DIAGRAM: DATA SET ALLOCATIONS: SET TRAINING TO 70% AND VALIDATION – 30% → right-click on the PARTITION NODE → RUN → YES → RESULTS (TO VIEW THE RESULTS OF THE PARTITION)**

# LOGISTIC REGRESSION

**MODEL → REGRESSION → DRAG AND DROP THE REGRESSION TO THE DIAGRAM → CONNECT THE PARTITION TO THE REGRESSION NODE → YOU CAN EDIT VARIABLES IF YOU LIKE (ON THE LEFT; TRAIN) AND MAKE SURE THAT THE CLASS TRAGETS IS SET TO LOGISTIC REGRESSION → MODEL SELECTION: BACKWARD ETC. → Right-click on the regression node → RUN → YES → RESULTS**

**You will see 4 tables:**

- **FIT STATISTICS** will let you compare different values, i.e. MSE or some other errors for the TRAIN & VALIDATION data sets.
- **OUTPUT:** EVENT CLASSIFICATION TABLE (CONFUSION MATRIX)
- **LIFT:** CUMULATIVE AND REGULAR LIFT
- **EFFECTS PLOT:** TO EVALUATE WHICH VARIABLES HAVE THE HIGHEST IMPACT ON THE MODEL

# ASSESS A MODEL

**ASSESS → DRAG AND DROP: MODEL COMPARISON → CONNECT IT TO THE REGRESSION NODE → RUN → YES → RESULTS**

- **OUTPUT: Fit Statistics** – it will provide you with the fit statistics of each model and if you compare a few models, then it will choose the best model based on the stats; you can also find other useful statistics
- **LIFT:** for train and validation data sets and if you have a few models you will see compared lifts for different models
- **ROC CURVES**

Results - Node: Model Comparison  Diagram: gym

File  Edit  View  Window

Output

```
46
47
48
49     Fit Statistics Table
50     Target: crowded
51
52     Data Role=Train
53
54     Statistics                                                            Reg
55
56     Train: Bin-Based Two-Way Kolmogorov-Smirnov Probability Cutoff        0.50
57     Train: Kolmogorov-Smirnov Statistic                                   0.60
58     Train: Akaike's Information Criterion                              39501.60
59     Train: Average Squared Error                                          0.15
60     Train: Roc Index                                                      0.87
61     Train: Average Error Function                                         0.45
62     Train: Cumulative Percent Captured Response                          18.64
63     Train: Percent Captured Response                                      9.22
64     Selection Criterion: Valid: Misclassification Rate                    0.20
65     Train: Degrees of Freedom for Error                               43519.00
66     Train: Model Degrees of Freedom                                       9.00
67     Train: Total Degrees of Freedom                                   43528.00
68     Train: Divisor for ASE                                            87056.00
69     Train: Error Function                                             39483.60
70     Train: Final Prediction Error                                         0.15
71     Train: Gain                                                          86.35
72     Train: Gini Coefficient                                               0.74
73     Train: Bin-Based Two-Way Kolmogorov-Smirnov Statistic                 0.60
74     Train: Kolmogorov-Smirnov Probability Cutoff                          0.45
75     Train: Cumulative Lift                                                1.86
76     Train: Lift                                                           1.84
77     Train: Maximum Absolute Error                                         0.99
78     Train: Misclassification Rate                                         0.20
79     Train: Mean Square Error                                              0.15
80     Train: Sum of Frequencies                                         43528.00
81     Train: Number of Estimate Weights                                     9.00
82     Train: Root Average Sum of Squares                                    0.38
83     Train: Cumulative Percent Response                                   87.53
84     Train: Percent Response                                              86.63
85     Train: Root Final Prediction Error                                    0.38
86     Train: Root Mean Squared Error                                        0.38
87     Train: Schwarz's Bayesian Criterion                               39579.73
88     Train: Sum of Squared Errors                                      12763.36
89     Train: Sum of Case Weights Times Freq                             87056.00
90
91
92     Data Role=Valid
93
94     Statistics                                                            Reg
95
96     Valid: Kolmogorov-Smirnov Statistic                                   0.60
97     Valid: Average Squared Error                                          0.15
98     Valid: Roc Index                                                      0.87
```

**Questions:**

### Logistic regression

1. Create a new project, assign a library, create a data source and a new diagram.
2. Open the file called gym_crowd
3. Assign variable roles and levels
4. Explore the data in two ways and answer the questions:
   a. What distribution does the number of visitors have?
   b. What is the minimum and maximum number of gym visitors?
   c. How many variables and how many observations does the dataset have?
   d. Are there any variables that have missing observations? What should you do with missing values?
   e. Which variables do you think should be excluded from the analysis since they will not contribute much to the explanation of the number of visitors to the gym?
   f. Are there any redundant variables?
   g. Which variable should you exclude from your analysis?
5. Partition the data into TRAIN (70%) and VALIDATION (30%) datasets
6. Run a full logistic regression for event = 1:
   a. Which variables are the most important in the model?
   b. What is the model's cumulative lift? Interpret it.
   c. Use a Model Comparison node to get a ROC curve. What is AUC? Interpret it.
7. Run the following logistic regressions for event=1:
   a. Stepwise
   b. Forward
   c. Backward
8. Compare all four models. Which model is the best? Which criteria have you used?
9. Using the output found after running a stepwise regression model, create a **CONFUSION MATRIX** in **EXCEL** and calculate the following:
   a. True positive rate (TPR)
   b. False negative rate (FNR)
   c. True negative rate (TNR)
   d. False positive rate (FPR)
   e. Misclassification rate
   f. Sensitivity (TRUE POSITIVE RATE)
   g. Specificity (TRUE NEGATIVE RATE)
   h. Horizontal ROC coordinates
   i. Vertical ROC coordinates
   j. Event rate
   k. Actual rate among predicted
   l. Lift
   m. Gain