# ANLT5050
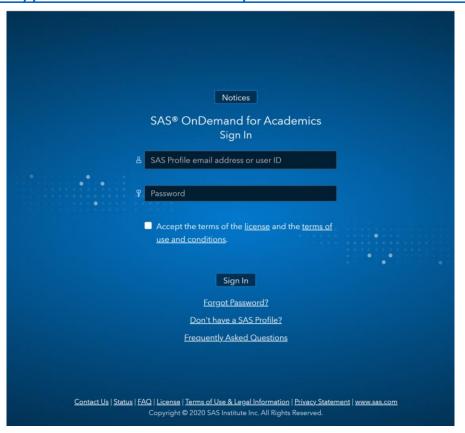
Unit 8 Assignment 2 Tutorial
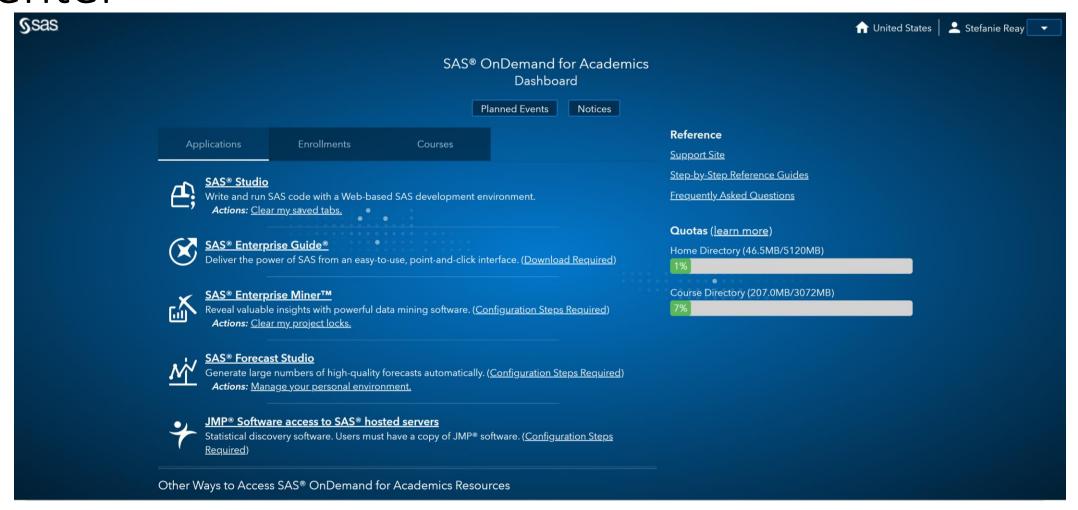
# Access the SAS OnDemand for Academics Control Center

https://odamid.oda.sas.com/SASODAControlCenter



© Stefanie G. Reay, MS, PhD, Capella University, 2021
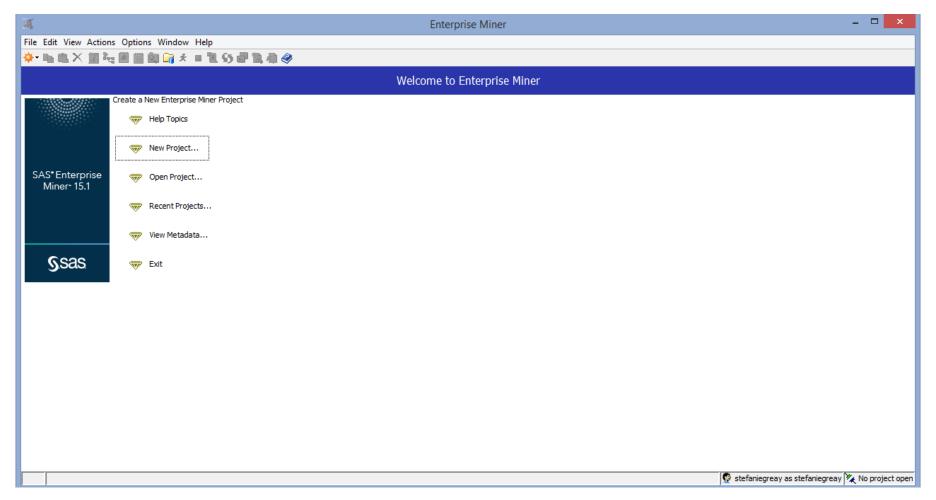
# SAS OnDemand for Academics (SODA) Control Center

# SAS Enterprise Miner Instructions

The following slides provide instructions on how to complete this task in SAS Enterprise Miner.

Once you have uploaded the dataset for this unit onto the SAS servers using SAS Studio, you may proceed from here using SAS Enterprise Miner.
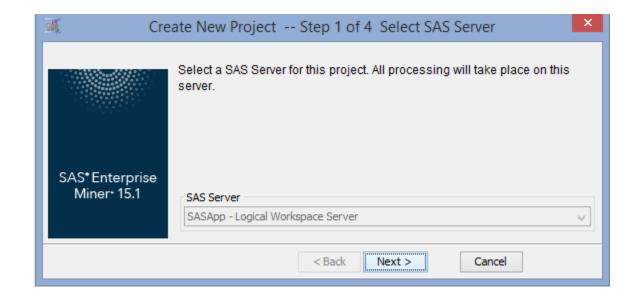
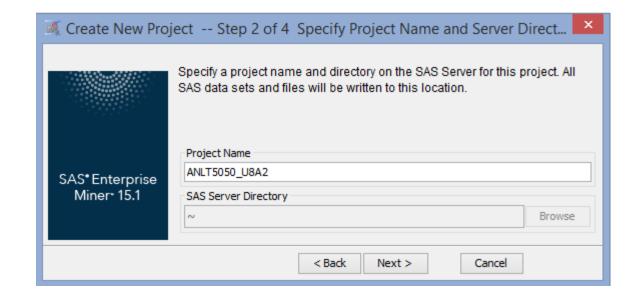Once you download and start SAS Enterprise Miner, open a new project by clicking on "New Project."
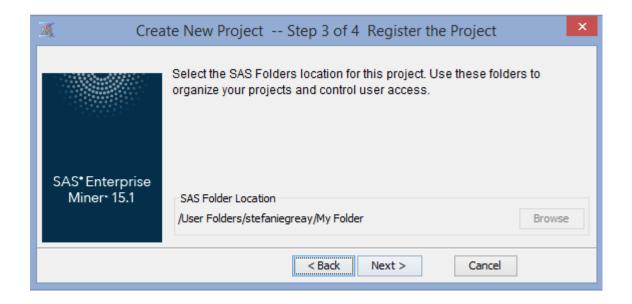
# Click "Next>" to use the default SAS Server
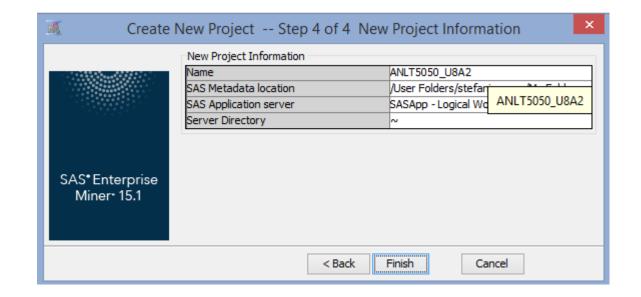
# Enter a project name and click "Next>"

# Click "Next>"

# Verify your entries and click "Finish"

# Click on the "Create Diagram" icon.

# Name your diagram and click "OK."

# Click on the project, then click on the ellipses next to "Project Start Code."

Add the library reference for where you uploaded the dataset in SAS studio, and click "Run Now."  Once it completes, click "OK."

# Click on Actions>Add Node>Sample>Input Data

# Click the ellipses (3 dots) next to "Data Source."

# Click on "New Data Source"

# Leave it as "SAS Table" and click "Next >"

# Click on "Browse"

# Double click on the libname you just set up in the project startup code.

# Double click to select the dataset for this unit, and click "OK"

# Click "Next>"

# Verify the options and click "Next>"

# Click "Next>"

# Verify the variables and settings, adjust if necessary, and then click "Next>"

You may choose to sample the dataset here, or just keep the full dataset, then click "Next>." If you want to split into train, test, and validate, you could do this here.

# You may choose to adjust the role of the dataset, or leave it as the default, then click "Next>"

# Click "Finish" to finish the data source registration within EM.

Click "OK" to complete the process. The name of the node should then change to the name of the dataset.

# Right click on the dataset node and click "Run."

# Click on "Actions" > "Add Node" > "Model" > "Regression"

# Connect the nodes



© Stefanie G. Reay, MS, PhD, Capella University, 2021

# Right click on the dataset node and choose "edit variables."

Change the description variables to "text," the Patient_ID variable to "ID," the dates as "Time_ID" variables, the order_total_charges variable to "Target," and all others to "Input"

Second part of list (for reference).  Leave the readmit_date and readmit_discharge_date as Time ID variables.  Click "OK" once you finish editing.



| PROCEDURE_ICD_CODE | Input | Nominal | No | | No | . | . |
|---|---|---|---|---|---|---|---|
| PROCEDURE_LONG_DESC | Text | Nominal | No | | No | . | . |
| PROCEDURE_SUBCAT_CODE | Input | Nominal | No | | No | . | . |
| PROCEDURE_SUBCAT_DESC | Text | Nominal | No | | No | . | . |
| PatientAge | Input | Interval | No | | No | . | . |
| STATECODE | Input | Nominal | No | | No | . | . |
| Standard_Orders_Used | Input | Nominal | No | | No | . | . |
| ZIP | Input | Nominal | No | | No | . | . |
| admit_month | Input | Interval | No | | No | . | . |
| dx_code | Input | Interval | No | | No | . | . |
| dx_group | Input | Nominal | No | | No | . | . |
| gender | Input | Nominal | No | | No | . | . |
| i | Input | Interval | No | | No | . | . |
| icd9_target | Input | Interval | No | | No | . | . |
| op_visits6 | Input | Interval | No | | No | . | . |
| operationcount | Input | Interval | No | | No | . | . |
| order_set_used | Input | Interval | No | | No | . | . |
| order_total_charges | Target | Interval | No | | No | . | . |
| race_cd | Input | Nominal | No | | No | . | . |
| readmit_date | Time ID | Interval | No | | No | . | . |
| readmit_days | Input | Nominal | No | | No | . | . |
| readmit_discharge_date | Time ID | Interval | No | | No | . | . |
| readmit_month | Input | Interval | No | | No | . | . |
| readmit_number | Input | Interval | No | | No | . | . |

Explore...  OK  Cancel

Change the regression type by right clicking on the Regression node, and scrolling down in the properties, and selecting "linear regression" as the regression type. If you wish to include interaction effects (pairwise only), you can set the Two-Factor Interactions to "Yes". You may also opt for a variable selection method here, like Stepwise, Forward, or Backward, instead of only creating a model that contains all variables specified. Note that if you choose Two-Factor Interactions, with this many variables, you will need to edit the terms using the Term Editor to reduce them before running the node.

| .. Property | Value |
|---|---|
| **Equation** | |
| Main Effects | Yes |
| Two-Factor Interactions | No |
| Polynomial Terms | No |
| Polynomial Degree | 2 |
| User Terms | No |
| Term Editor | ... |
| **Class Targets** | |
| Regression Type | Linear Regression |
| Link Function | Logit |
| **Model Options** | |
| Suppress Intercept | No |
| Input Coding | Deviation |
| **Model Selection** | |
| Selection Model | Stepwise |
| Selection Criterion | Default |
| Use Selection Defaults | Yes |

# Right click on "Regression" node and click "Run"



© Stefanie G. Reay, MS, PhD, Capella University, 2021

# Right click on the "Regression" node and click "Results" to view the results.

# Considerations for fitting regression models

- The options for variable selection include: forward, backward, and stepwise.

- The process of fitting a regression model is iterative, and should not be a "set it and forget it"-type of analysis.

- SAS Enterprise Miner's regression node is a good start for analyzing several variables, but additional analysis and evaluation of assumptions should be completed outside of the regression node.

After you re-run the regression node (after making any adjustments to the settings), right click on the node and select "Results."

The Fit Statistics window shows the fit statistics for the last iteration. If you scroll down to the bottom of the Output pane, the last iteration is shown, including the variables and the tests for contributions, etc.

The Analysis of Variance table provides the F-test statistic and p-value to test for significance in the overall model's ability to predict the response variable.  The Model Fit Statistics table provides a variety of model fit statistics to evaluate the fit of the model specified.  The Type 3 Analysis of Effects provides F-test statistics and p-values for testing the contribution of each individual variable.



Analysis of Variance

| Source | DF | Sum of Squares | Mean Square | F Value | Pr > F |
|---|---|---|---|---|---|
| Model | 71 | 566486613952 | 7978684704 | 186.22 | <.0001 |
| Error | 14953 | 640656885058 | | | |
| Corrected Total | 15024 | 1.2071435E12 | | | |

Model Fit Statistics

| | | | |
|---|---|---|---|
| R-Square | 0.4693 | Adj R-Sq | 0.4668 |
| AIC | 264107.5439 | BIC | 264112.3097 |
| SBC | 264656.0018 | C(p) | -139.2217 |

Type 3 Analysis of Effects

| Effect | DF | Sum of Squares | F Value | Pr > F |
|---|---|---|---|---|
| CODE | 1 | 4236419467 | 98.88 | <.0001 |
| AT_CODE | 1 | 8485627054 | 198.06 | <.0001 |
| | 8 | 9.64134E10 | 281.29 | <.0001 |
| | 1 | 861626683 | 20.11 | <.0001 |
| | 9 | 3.91245E10 | 101.46 | <.0001 |
| TY | 3 | 5169262374 | 40.22 | <.0001 |
| HOSPITAL | 6 | 8.32923E10 | 324.01 | <.0001 |
| ICU_DAYS | 1 | 1.10245E10 | 257.31 | <.0001 |
| LENGTH_OF_STAY | 1 | 5827296625 | 136.01 | <.0001 |

The Analysis of Maximum Likelihood Estimates provides the estimate for the coefficient for a particular variable (and the intercept), as well as the standard error, t-test statistic and p-value for the t-test for significant contribution to the model. Notice that categorical variables display as individual indicators in this table (as they are treated as separate indicator/dummy variables within the model itself...this conversion is automatic within the regression node). ALL variables should be analyzed and reported on, not just the top few.

Analysis of Maximum Likelihood Estimates
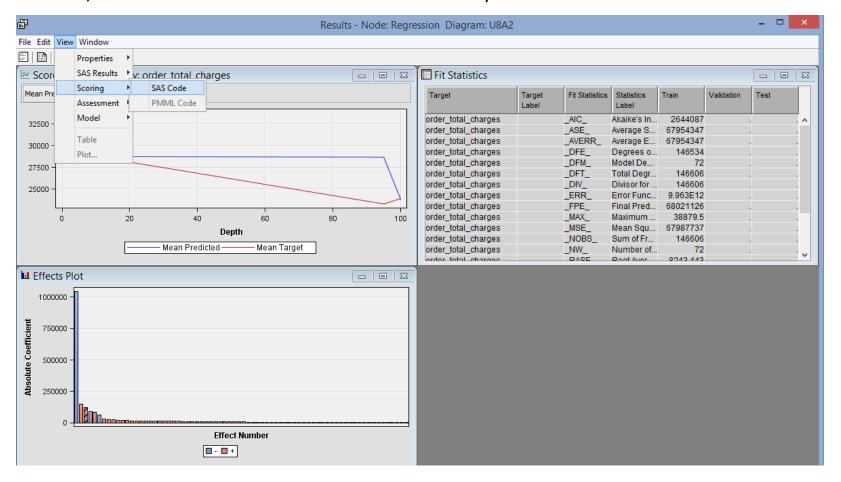
| Parameter | | DF | Estimate | Standard Error | t Value | Pr > \|t\| |
|---|---|---|---|---|---|---|
| Intercept | | 1 | -1047679 | 69504.9 | -15.07 | <.0001 |
| DIAGNOSIS_ICD_CODE | | 1 | -5153.1 | 518.2 | -9.94 | <.0001 |
| DIAGNOSIS_SUBCAT_CODE | | 1 | 7478.5 | 531.4 | 14.07 | <.0001 |
| DISCHARGED_TO | AGNST MEDICAL ADVICE AMA | 1 | -8991.0 | 491.1 | -18.31 | <.0001 |
| DISCHARGED_TO | CHG TO LTAC | 1 | -1724.2 | 1050.1 | -1.64 | 0.1006 |
| DISCHARGED_TO | HOME HEALTH AGENCY | 1 | -2790.5 | 222.3 | -12.55 | <.0001 |
| DISCHARGED_TO | HOSPICE (HOME) | 1 | 1460.9 | 423.2 | 3.45 | 0.0006 |
| DISCHARGED TO | HOSPICE MEDICAL INDAT | 0 | 0 | | | |

To get the SAS code for scoring, you can click on "View" > "Scoring">"SAS Code." This can be run against a test and validation dataset to test and then validate the model. (For example, if you split the dataset into 3 originally, for training, testing and validation, or on new data that comes in.)

# During and after the iterations, the following assumptions of regression must be validated

- There is a linear relationship between the dependent/response and independent/explanatory variables
- Multivariate normality
- No multicollinearity between the independent/explanatory variables
- No auto-correlation
- Homoscedasticity (homogeneity of variance of the residuals)

# Assumption: Linear relationship

- Linear relationships between the response/dependent and explanatory/independent variables can be checked using:

    1. A linear correlation matrix

    2. A scatterplot (or scatterplot matrix)

    3. And should also be checked for outliers using a histogram and box plot of each individual variable

# Assumption: Multivariate Normal

- Multivariate normality assumes that the residuals are normally distributed.  This can be checked using all of the following:
    1. A normal probability (or normal Q-Q) plot of the residuals
    2. A histogram of the residuals
    3. A box plot of the residuals
    4. A Komogorv-Smirnov test of the residuals

# Assumption: No multicollinearity between the independent/explanatory variables

- No multicollinearity between the independent/explanatory variables can be checked using all of the following:
  1. A linear correlation coefficient matrix (including correlations between explanatory/independent variables
  2. A scatterplot matrix showing relationships between the explanatory/independent variables
  3. Tolerance
  4. Variance Inflation Factor (VIF)
  5. Condition Index

# Assumption: No autocorrelation

- No autocorrelation can be checked using all of the following:
    1. Durbin-Watson's d test

# Assumption: Homoscedasticity (homogeneity of variance of the residuals)

- Homoscedasticity (homogeneity of variance of the residuals) can be checked using all of the following:

  1. Scatterplot of residuals vs predicted values
  2. Scatterplot of residuals vs response/dependent variable(s)
  3. Scatterplots of residuals vs each of the explanatory/independent variable(s)

  You will want to check for patterns.  Patterns of any kind suggest a lack of random spread, or, in essence, a lack of homogeneity of variance within the residuals.

# SAS Documentation Reference

The link below brings you to the SAS Documentation on the Regression Node, which has an example, including interpretation of the output.

https://documentation.sas.com/?docsetId=emref&docsetTarget=n1jqzz8cssr9m2n1ktx2iyv87q56.htm&docsetVersion=14.3&locale=en