



## Ontario Tech University | Seminar in SAS Machine Learning Tools

### Framing, Part I: Learning Objectives and Purpose

- Expose participants to a broad range of machine learning tools
- Show how to produce advanced analytics with code and no-code options
- Provide useful follow-up resources

### Framing, Part II: Brief Overview of SAS Academic Programs

- Our Mission
  - Helping to train the next generation of analytics talent
- How we support Academics
  - SAS Student Skill Builder
  - SAS Educator Portal
- Tools... +click here
  - Linking links
    - [https://www.sas.com/en\\_us/learn/academic-programs.html](https://www.sas.com/en_us/learn/academic-programs.html)
    - [https://www.sas.com/en\\_us/software/viya-for-learners.html](https://www.sas.com/en_us/software/viya-for-learners.html)
    - [https://www.sas.com/en\\_us/software/on-demand-for-academics.html](https://www.sas.com/en_us/software/on-demand-for-academics.html)
    - [https://www.sas.com/en\\_us/learn/academic-programs/students.html](https://www.sas.com/en_us/learn/academic-programs/students.html)
    - [https://www.sas.com/en\\_us/learn/academic-programs/educators.html](https://www.sas.com/en_us/learn/academic-programs/educators.html)

### Framing, Part III: Code versus No Code Options

- Broad software tour
  - OnDemand for Academics
    - SAS Studio
    - SAS Enterprise Miner (for a limited time only)
  - SAS Viya for Learners
    - SAS Studio
    - SAS Model Studio
    - SAS Visual Analytics
- What are the advantages of each approach?

## Framing, Part IV: The Data

- We'll be using the classic *Titanic* data set in this analysis to explore survival rates after the sinking of the "unsinkable" ship in 1912
- Variables are:

Variable Name	Details
Age	Age
Cabin	Cabin
Embarked	Port of Embarkation (C = Cherbourg; Q = Queenstown; S = Southampton)
Fare	Passenger Fare (British pound)
Name	Name
Parch	Number of Parents/Children Aboard
PassengerId	Passenger ID
Pclass	Passenger Class (1 = 1st; 2 = 2nd; 3 = 3rd)
Sex	Sex
SibSp	Number of Siblings/Spouses Aboard
Survived	Survival (0 = No; 1 = Yes)
Ticket	Ticket Number

- More details about the Titanic data set can be found [here](#).

## Framing, Part IV: Access Software

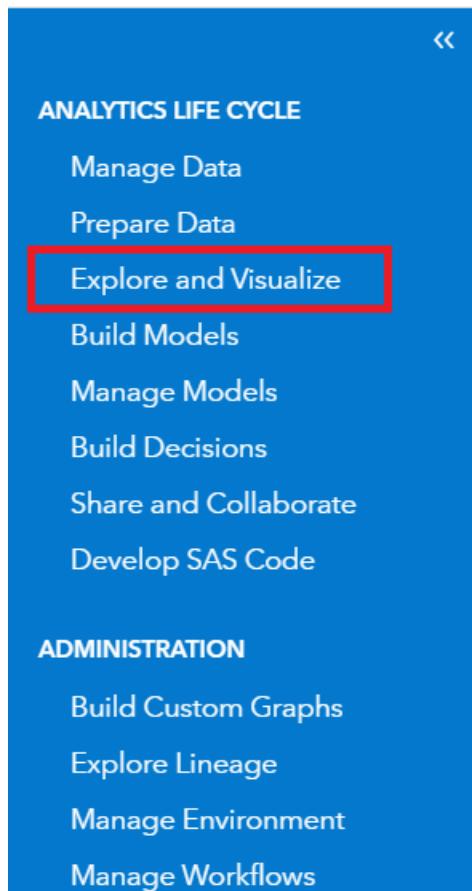
- Steps to accessing SAS Viya for Learners for the first time:
  - Navigate to the SAS Viya for Learners webpage, [https://www.sas.com/en\\_us/software/viya-for-learners.html](https://www.sas.com/en_us/software/viya-for-learners.html)
  - Click the "Access for Educators" or "Access for Students" button based on your role.
  - Log in with your SAS Profile that is linked to an academic affiliation. If you don't have a SAS Profile, click [here](#) to set one up.
  - Access the software by clicking on "Launch SAS Viya for Learners 3.5"
- For future sessions with VFL, you can access VFL by visiting <https://vle.sas.com/vfl>.
- Contact us at [academic@sas.com](mailto:academic@sas.com) if you have any questions about access.

## Part I: SAS Visual Analytics

- **Access SAS Visual Analytics**
  - How? Let the Hamburger become your friend!
    - From upper left, click



- Then select Explore and Visualize from the options available:



- **From here, Start with Data**
  - Find **Titanic** in the available data sets in SAS Viya for Learners:

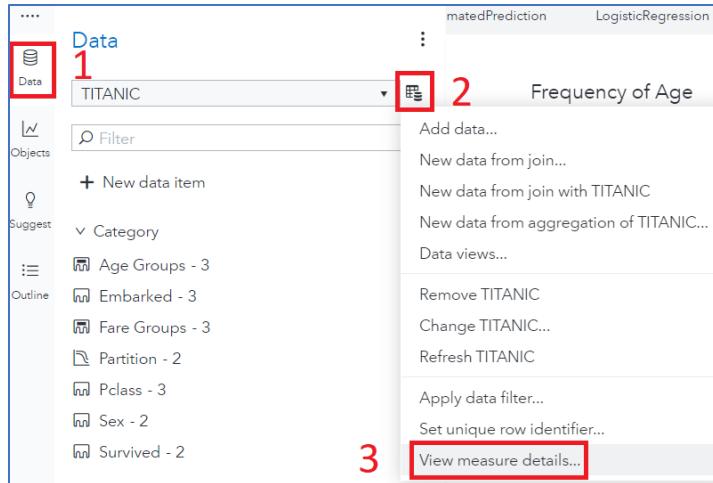
## Ontario Tech University - Seminar in SAS Machine Learning Tools - MAY2023

The screenshot shows the 'Choose Data' window in SAS. On the left, there's a sidebar with 'Available Data Sources' showing two entries: 'titan' and 'TITANIC'. The main area is titled 'TITANIC' and displays a table of 12 columns with their respective details like type, length, and format. The right side contains a summary panel with information such as 'Date profiled: 08/01/21 08:45 AM', 'Columns 12', 'Rows 1.3 K', and a 'Profile' tab. At the bottom right are 'OK' and 'Cancel' buttons.

- In the **Choose Data** window, click on *TITANIC*. Next, examine the information about the data under the **Details**, **Sample Data**, and **Profile** tabs.
- Here is sample output from the **Profile** tab:

The screenshot shows the 'Profile' tab for the TITANIC dataset. It displays a table of 12 columns with various statistics: Unique values, Null counts, Blank counts, Pattern Count, Mean, Median, Mode, Standard Deviation, Standard Error, and Minimum. For example, 'Age' has 7.49% unique values and a mean of 29.50. 'Cabin' has 63.05% unique values and a median of C23 C... . 'Embarked' has 0.23% unique values and a mode of Embark... . 'Fare' has 21.47% unique values and a mode of 8.05. 'Name' has 99.85% unique values and a median of Abbing, M.. . 'Parch' has 0.61% unique values and a mode of 0.0. 'PassengerId' has 100.00% unique values and a standard deviation of 378.02. 'Pclass' has 0.23% unique values and a lower bound of Lower. 'Sex' has 0.15% unique values and a mode of female. 'SibSp' has 0.53% unique values and a median of 0.0. 'Survived' has 0.15% unique values and a mode of Perished. 'Ticket' has 70.97% unique values and a median of CA. 2... . At the bottom right are 'OK' and 'Cancel' buttons.

- Let's further explore our data!
  - We can do this via **View Measure Details**
    - Find the tool here:



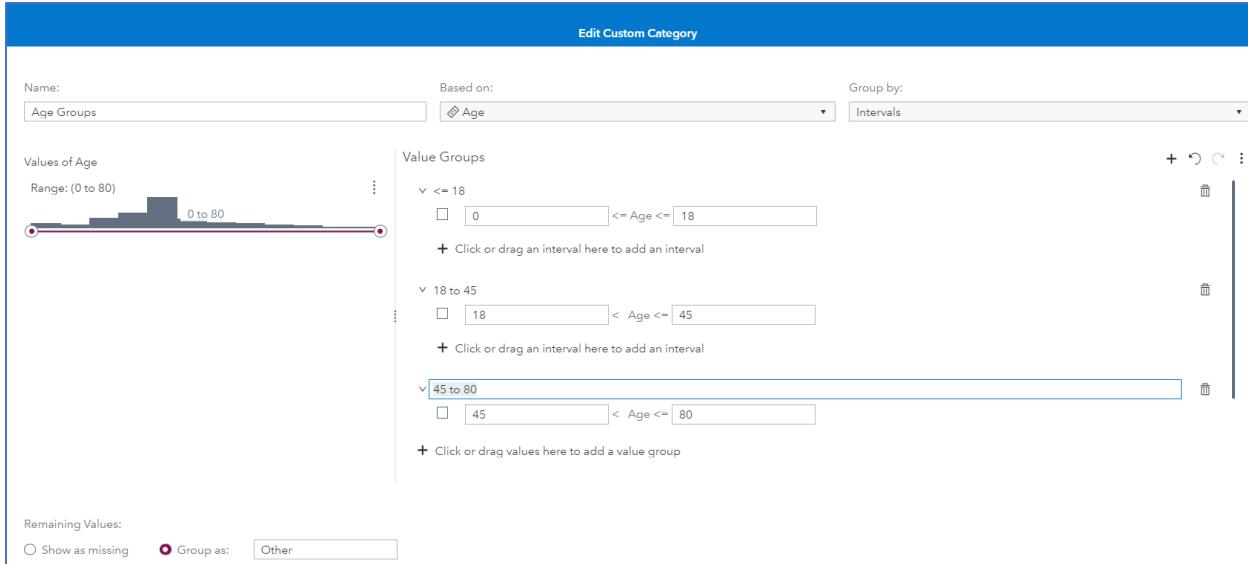
- Sample output:



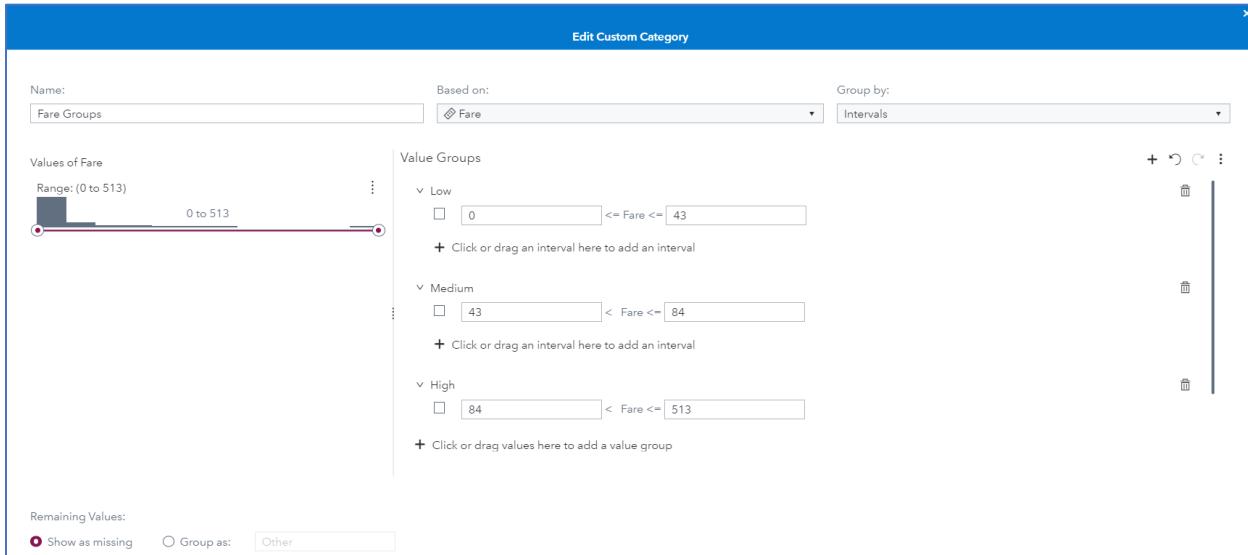
- Let's create a series of histograms in the default canvass to better understand the data
  - Label the page "DescriptiveStats" and replicate the following:



- **Helpful hints**
  - Frequency of Sex grouped by Survived
    - Grab sex and survived and then Auto Chart
    - Change **List table** to a **Bar chart** with the following settings:
      - Category = sex
      - Group = survived
      - Grouping style = Stacked
  - Frequency of Cabin
    - Change default Auto chart object to **Pie chart**
- Now, let's create New Variables → with + **New Data Item** → **Custom Category**
  - Why do this?
    - Well, the relationship between Survived and variables such as Age or Ticket Fare is likely non-linear. By creating new custom categories, we can increase the accuracy and flexibility of our models.
  - Variable 1: Age Groups
    - Group 1 | 0 to 18
      - Label = “<= 18”
    - Group 2 | 18 to 45
      - Label = “18 to 45”
    - Group 3 | 45 to 80
      - Label = “45 to 80”
  - Helpful hint before you get started:
    - Clicking on the end of the range line allows you to type in the exact value (rather than messing with the manual adjustment)



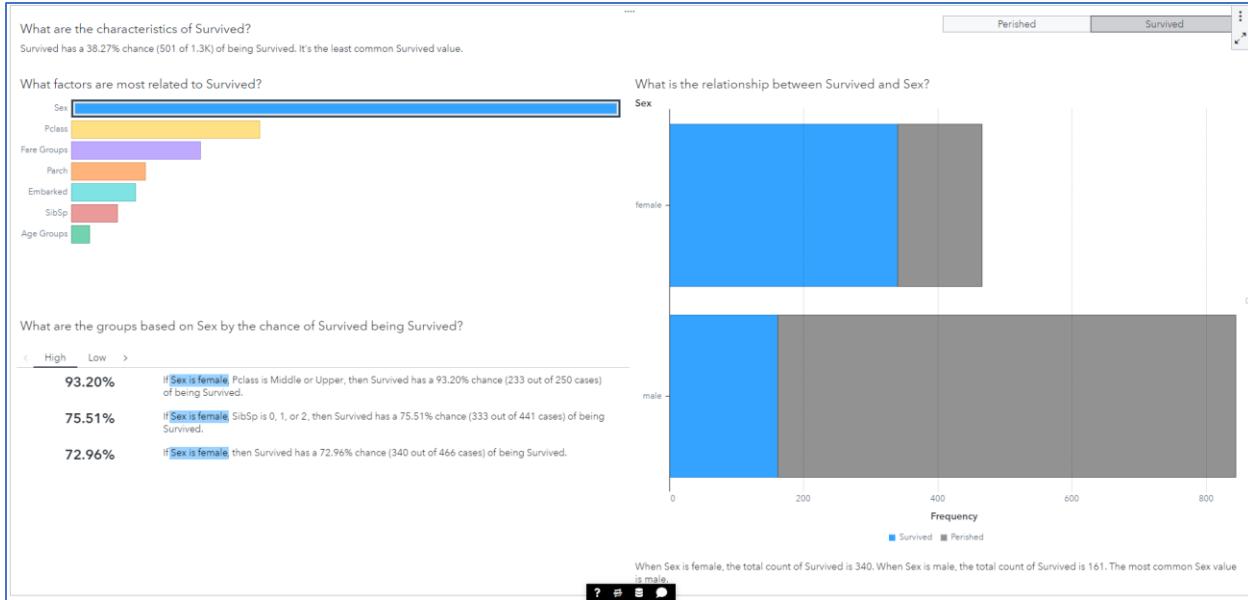
- Variable 2: Fare Groups
  - Group 1 | 0 to 43
    - Label = "Low"
  - Group 2 | 43 to 84
    - Label = "Medium"
  - Group 3 | 84 to 513
    - Label = "High"



- Finally, add **Age Groups** and **Fare Groups** to **DescriptiveStats**:
  - Sort variables as needed



- **Let's change the metadata**
  - What is metadata?
    - Simply data about the data
  - Why change it?
    - Changing the metadata helps SAS produce more useful output in our automated analyses
  - So, let's hide the following variables:
    - Cabin
    - Name
    - Ticket
    - Age
    - Fare
    - PassengerID
- **Let's see what insights we can gleam from Automated Explanation**
  - Create a new page to examine **Survived** versus **Perished**
    - Label this page "Automated Explanation"
  - Find *Automated Explanation* under *Analytics*. Drag the Object to the canvass.
    - Choose *Survived* as the *Response*.
  - Explore the output, including all the filters
    - Do you see any interesting relationships?
  - Sample output:



- **Important Question:**
  - What % of Titanic passengers survived?
    - Survival occurred 38.27% of the time
  - Why is this number important?
    - We'll use it shortly as our prediction cut-off – i.e., our probability decision threshold
- **Now on to Automated Prediction**
  - The backstory:
    - Rather than examining the bivariate relationships between Survived and individual variables – as with the Automated Explanation – let's look at all variables simultaneously via an A.I. tool called "Automated Prediction".
  - Create a new page called "Automated Prediction"
    - Find the *Automated Prediction* object under *Analytics* and add it to the canvass
      - Assign **Survived** as the **Response**
    - Try a sample prediction with the following settings:
      - Sex = male
      - Pclass = Lower
      - Embarked = EmbarkS
      - Parch = 0
      - SibSp = 0
    - What is the outcome?

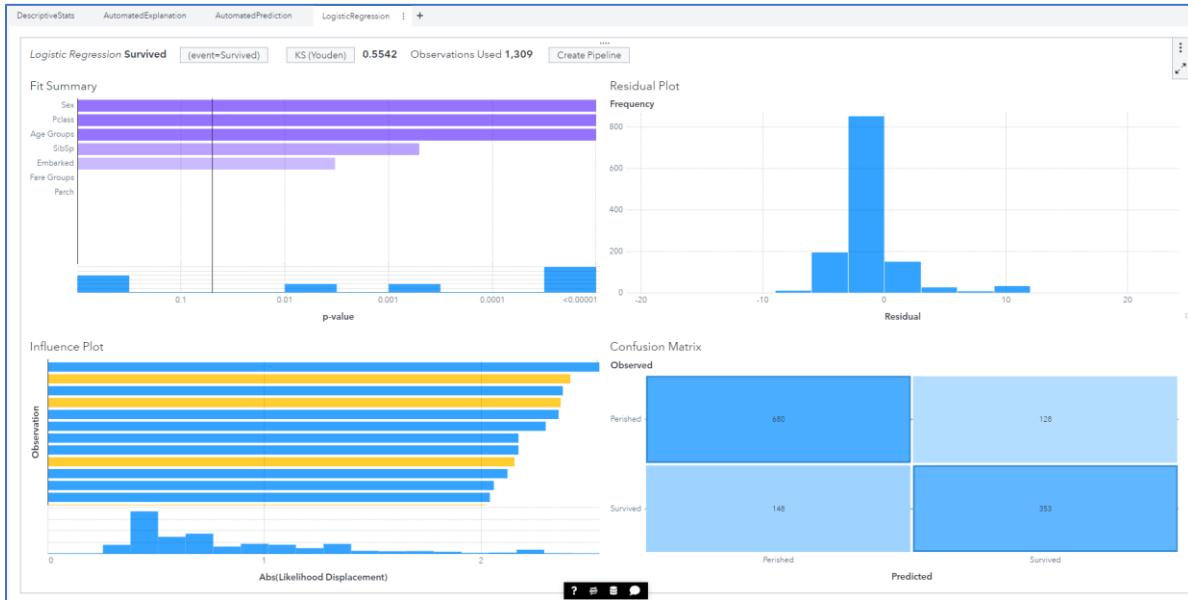
The screenshot shows a user interface for a machine learning model. At the top, there are tabs for 'DescriptiveStats', 'AutomatedExplanation', 'AutomatedPrediction', and a '+' sign. Below the tabs, a message reads: 'What values for the most important factors should be used to predict? What is the prediction for Survived?' A large blue title 'Perished' is centered. On the left, there are input fields for 'Sex' (set to 'male'), 'Pclass' (set to 'Lower'), 'Embarked' (set to 'EmbarkS'), 'SibSp' (set to '0'), and 'Parch' (set to '0'). To the right of these inputs, a note states: 'The predicted Survived, Perished, is the most common Survived value in observed cases. Most observed cases (61.73%) are Perished. The prediction is based on an automatically selected Decision Tree model.' At the bottom of the interface are several small icons.

- Use the tool to examine the Survived variable
  - What interesting trends do you notice?
- Now it is modeling time!
  - Some modeling notes to get us started:
    - For all the models, create a new page and change the label to the name of the model.
    - The models are objects found under either SAS Visual Statistics or SAS Visual Data Mining and Machine Learning
  - Start with a **Logistic Regression** model
    - Use the following Data Roles to set up your model:

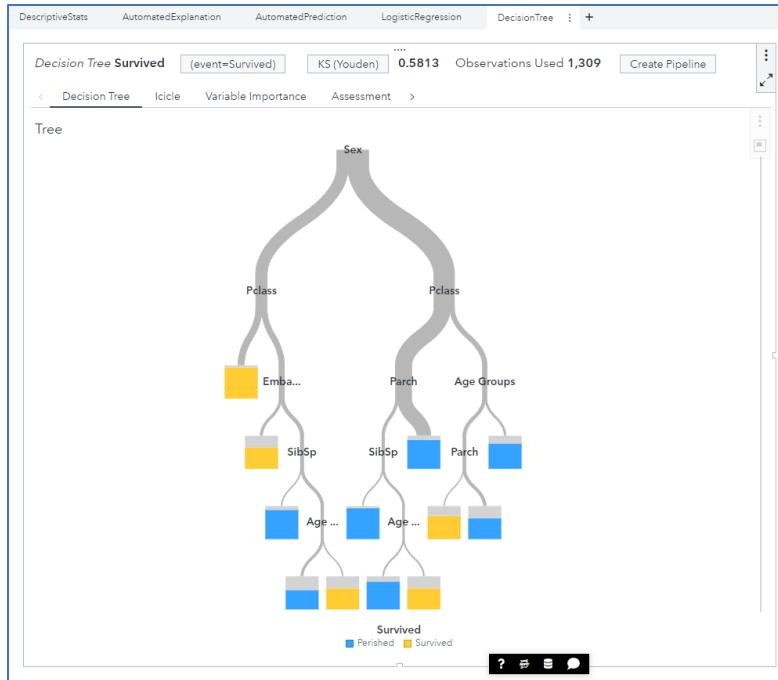
The screenshot shows a 'Data Roles' configuration window. At the top, it says 'Logistic regression - Survived 1'. Below this, there are three sections: 'Response' (with 'Survived' checked), 'Continuous effects' (with 'Parch' and 'SibSp' checked, and a '+ Add' button), and 'Classification effects' (with 'Age Groups', 'Embarked', 'Fare Groups', 'Pclass', 'Sex', and a '+ Add' button).

- Let's change a few Options
  - For variable selection method, use fast backward

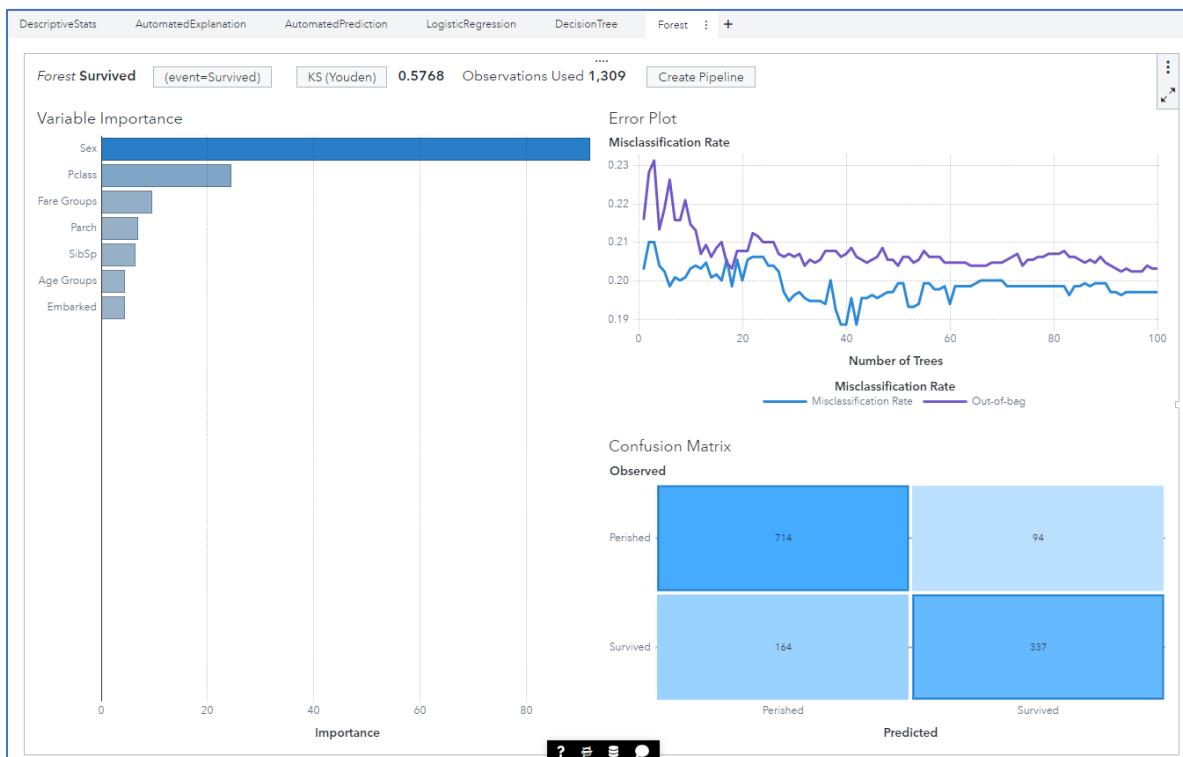
- For the residual plot → use a histogram
- Under Influence plot/variable selection plot
  - Plot to show: influence plot



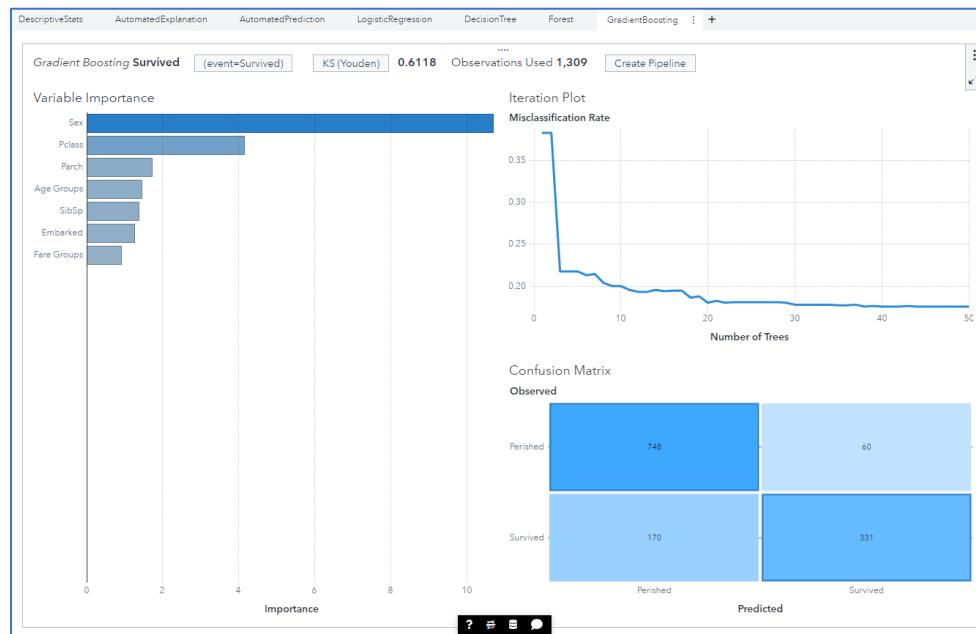
- Next up is a **Decision Tree**
  - To use your existing work, start with the Logistic Regression Object then:
    - Click alt + [ duplicate on new page as... Decision Tree ]
  - Let's explore some additional model options
    - Leaf size = 10
    - Model Display
      - Plot Layout: stack



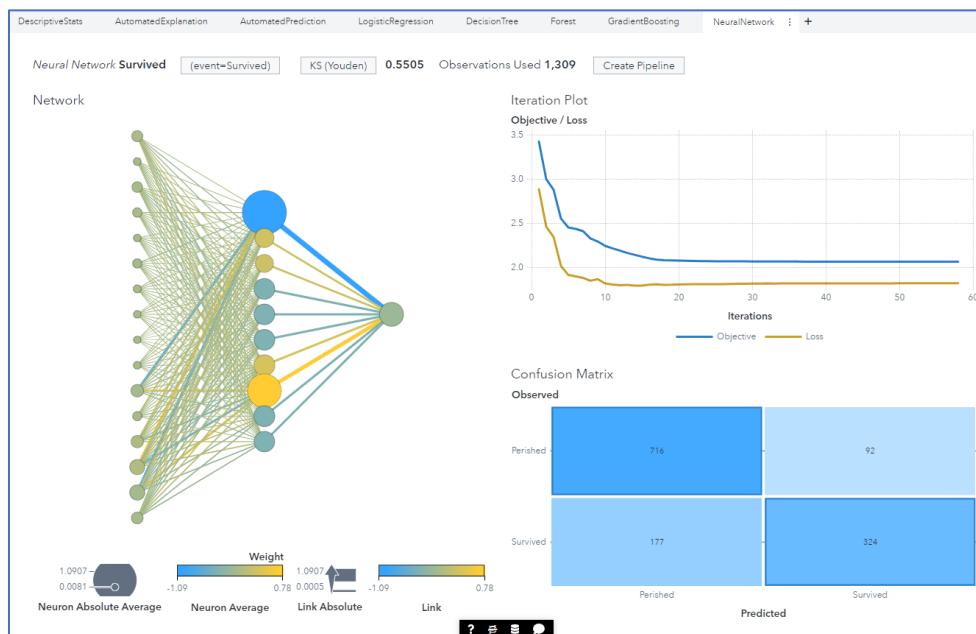
- Let's use a bunch of decision trees to create a **Forest**
  - Again, to use your previous work, start with the Decision Tree Object and then:
    - Press alt + [ duplicate on new page as... Forest]
  - Additional options to explore
    - Leaf size = 10
  - Output:



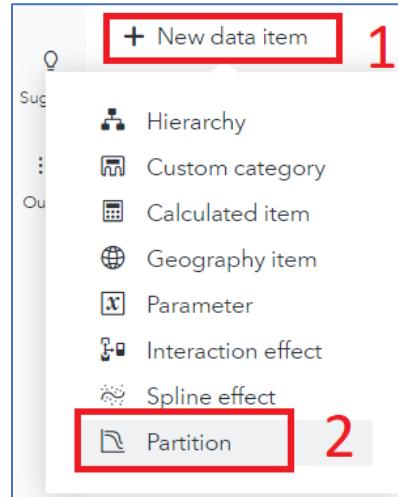
- **Gradient Boosting** is another tree-based model to explore
  - Similar approach as modeling the others
    - Forest Tree → alt + [ duplicate on new page as... Gradient Boosting ]
  - Additional Setting
    - Leaf size = 10
  - Output:



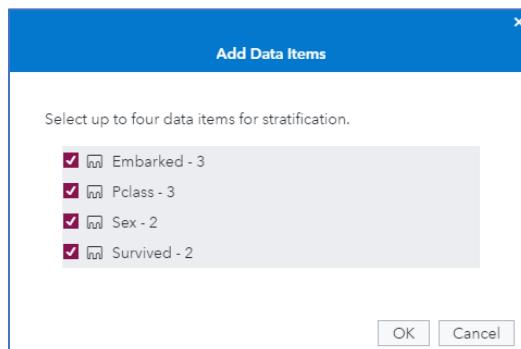
- Finally, let's end with a **Neural Network** model
  - One last time
    - Gradient boosting → alt + [ duplicate on new page as... Neural Network ]
  - Accept defaults, which produces the following:



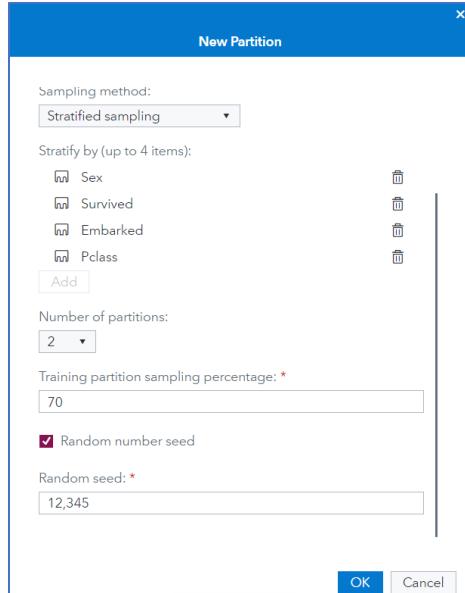
- It's time for **Model Comparison!**
  - We didn't do all that modeling for modeling's sake
    - Our goal is to find the best model, crown it champion, and do something with that model (oftentimes, predict new cases)
  - But, wait... should we have been performing honest assessment?
    - Yes! Oops!
    - Before comparing models, let's go back and fix the data!
      - +New data item ➔ partition



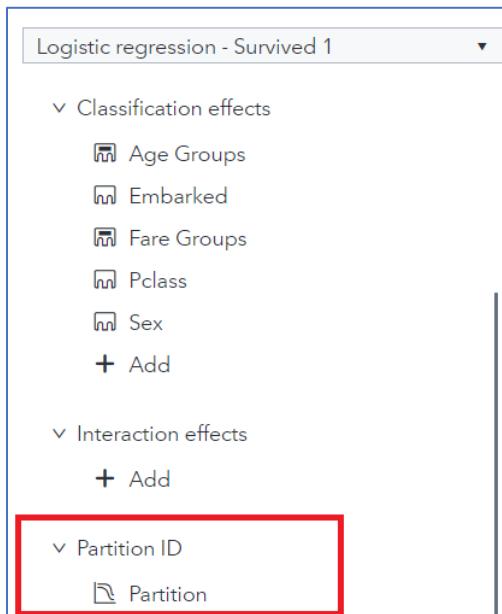
- Use a **Stratified Sampling**, adding the following variables and settings:



- Number of partitions = 2
- Training partition = 70%
- Random Seed = 12345
- Your partition will appear as follows:

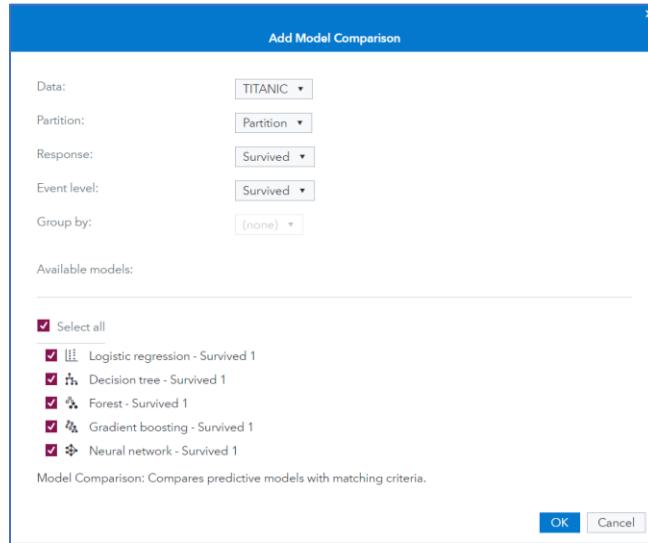


- Starting with the Logistic Regression, add the partition to the **Data Roles** in all our models and rerun!

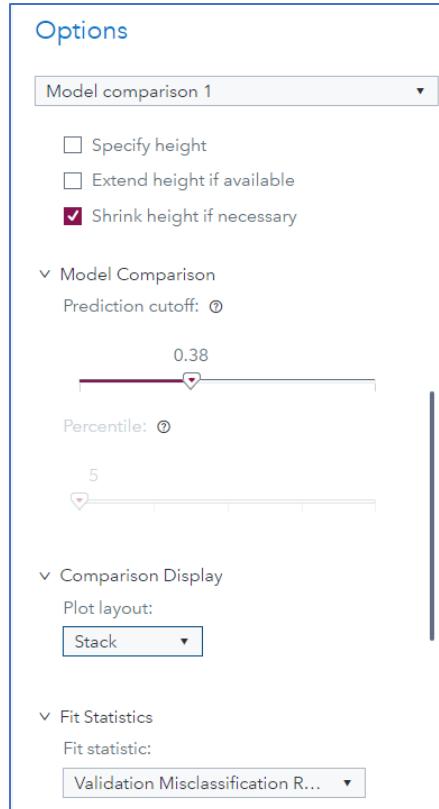


- **Which selection statistics should we choose?**
  - Discussion
    - What are you trying to optimize?
    - Which test statistic makes the most sense for your research question?
  - Settings
    - Let's choose Misclassification Rate (event) as the selection statistics
    - Use 38% as the *probability decision threshold*
- Now let's use the **Model Comparison Node** with the settings mentioned above

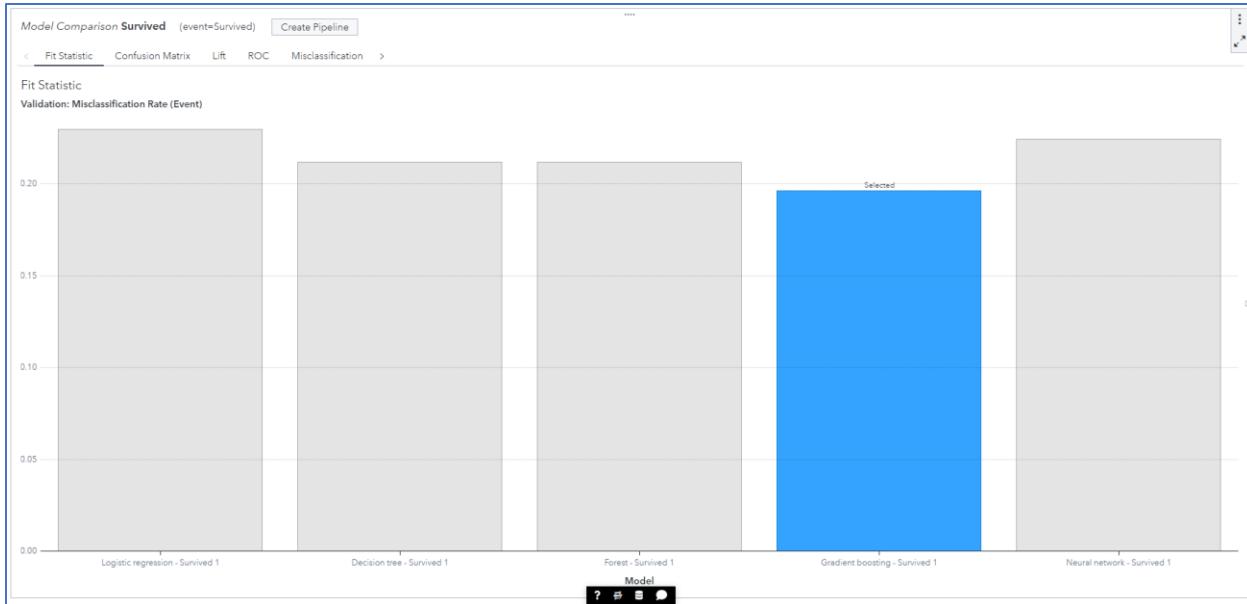
- From the options under SAS Visual Statistics, drag the *Model Comparison* object to a new page and select all available machine learning models:



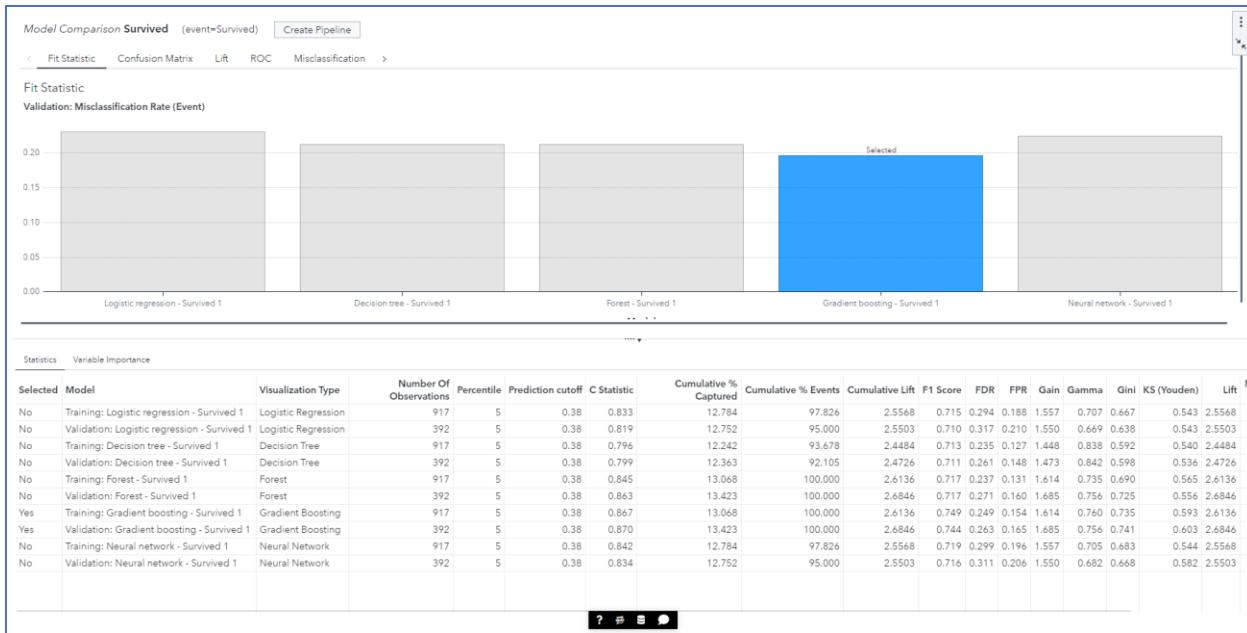
- Adjust a few of the **Options**:
  - Change the Prediction Cutoff to 0.38 (or 38%)
  - For the selection fit statistic, choose Validation Misclassification Rate
  - To make the output easier to digest, select:
    - Comparison Display → Plot layout → stack



- This yields the following:

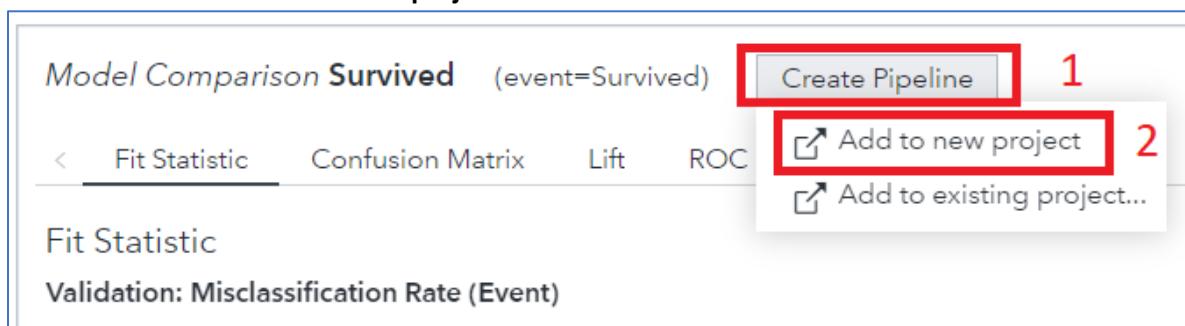


- Which model was selected by the Model Comparison object?
  - The gradient boosting has the smallest misclassification rate!
- Examine the underlying data for more details:



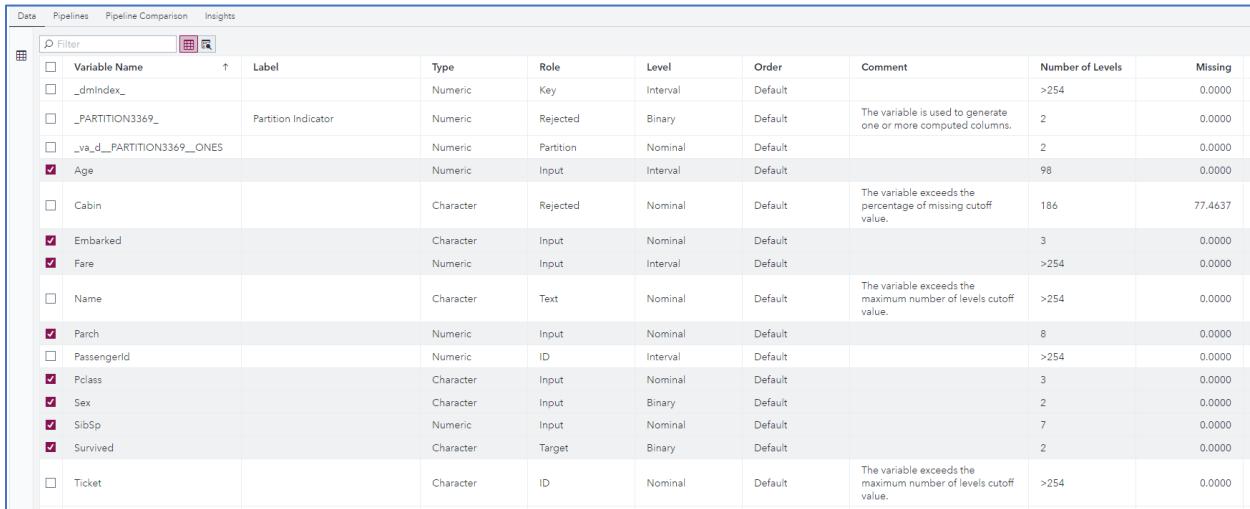
- Examine the Confusion Matrix, as well as the Lift, ROC, and Misclassification rates.
- Do you see any interesting trends?
- Promote champion model to SAS Model Studio to get us started!
  - How do we do this?
  - Well, we are just two clicks away.

- First, select **Create Pipeline**
- Then **Add to new project...**



## Part II: SAS Model Studio

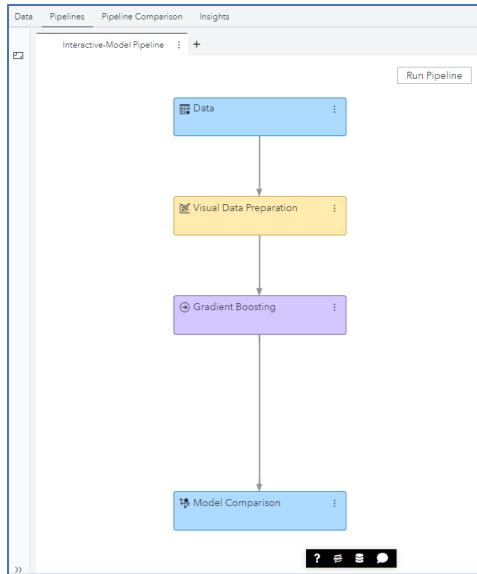
- **Welcome to SAS Model Studio!**
  - Start with a tour.
- **Metadata**
  - It's still very important.
    - So, try not to skip over it – and simply plunge into the model. Instead, check it out under the **Data** tab.
  - Why is partition rejected?
    - SAS Model Studio automatically creates training/validation/testing samples in every pipeline. So, incoming partitions are, by default, rejected.
  - Here are the main variables in our analysis:



The screenshot shows the SAS Model Studio interface with the 'Data' tab selected. The 'Variables' section displays a list of variables with their properties. Some variables have checkboxes next to them, indicating they are selected or flagged. The columns in the table include Variable Name, Label, Type, Role, Level, Order, Comment, Number of Levels, and Missing.

Variable Name	Label	Type	Role	Level	Order	Comment	Number of Levels	Missing
_dmIndex_		Numeric	Key	Interval	Default		>254	0.0000
_PARTITION3369_	Partition Indicator	Numeric	Rejected	Binary	Default	The variable is used to generate one or more computed columns.	2	0.0000
_va_d_PARTITION3369_ONES		Numeric	Partition	Nominal	Default		2	0.0000
<input checked="" type="checkbox"/> Age		Numeric	Input	Interval	Default		98	0.0000
<input type="checkbox"/> Cabin		Character	Rejected	Nominal	Default	The variable exceeds the percentage of missing cutoff value.	186	77.4637
<input checked="" type="checkbox"/> Embarked		Character	Input	Nominal	Default		3	0.0000
<input checked="" type="checkbox"/> Fare		Numeric	Input	Interval	Default		>254	0.0000
<input type="checkbox"/> Name		Character	Text	Nominal	Default	The variable exceeds the maximum number of levels cutoff value.	>254	0.0000
<input checked="" type="checkbox"/> Parch		Numeric	Input	Nominal	Default		8	0.0000
<input type="checkbox"/> PassengerId		Numeric	ID	Interval	Default		>254	0.0000
<input checked="" type="checkbox"/> Pclass		Character	Input	Nominal	Default		3	0.0000
<input checked="" type="checkbox"/> Sex		Character	Input	Binary	Default		2	0.0000
<input checked="" type="checkbox"/> SibSp		Numeric	Input	Nominal	Default		7	0.0000
<input checked="" type="checkbox"/> Survived		Character	Target	Binary	Default		2	0.0000
<input type="checkbox"/> Ticket		Character	ID	Nominal	Default	The variable exceeds the maximum number of levels cutoff value.	>254	0.0000

- One Important note:
  - The variables we created were not transferred!
- Where do they appear?
- Where did they go?
  - Wait for it...
- **Interactive-Model Pipelines**
  - The variables we created will appear here!



- To see the variables, start by running the pipeline by clicking here:

**Run Pipeline**

- Examine the results of the **Visual Data Preparation** node
  - Were the variables created?
  - Yes!
  - And here is the proof:

Name	Role	New Role	Variable Level	Variable Format
Survived	TARGET	TARGET	BINARY	\$CHAR16.
Ticket	ID	REJECTED	NOMINAL	\$CHAR36.
_PARTITION3369_	REJECTED	REJECTED	BINARY	
_dmIndex_	KEY	KEY	NOMINAL	
<b>_va_c_Age_Groups</b>	INPUT	INPUT	NOMINAL	
<b>_va_c_Fare_Groups</b>	INPUT	INPUT	NOMINAL	
<b>_va_d_PARTITION33</b>	PARTITION	PARTITION	NOMINAL	
69_ONES				

CMR Variable Score Code
1 retain '_va_d_calculation2396_1';
2 retain '_va_d_calculation2396_15';
3 retain '_va_d_calculation2396_16';
4 retain '_va_d_calculation2396_17';
5 retain '_va_d_calculation2396_18';
6 retain '_va_d_calculation2396_19';
7 retain '_va_d_calculation2396_20';
8 retain '_va_d_calculation2396_21';
9 retain '_va_d_calculation2396_22';
10 retain '_va_d_calculation2396_23';
11 retain '_va_d_calculation2396_24';
12 retain '_va_d_calculation2396_25';
13 retain '_va_d_onetime2396_1';
14 retain '_va_d_onetime2396_2';
15 if (( <b>_va_c_Age_Groups</b> = 1) & ( <b>_va_c_Fare_Groups</b> = 1)) then _dmIndex_ = 1;

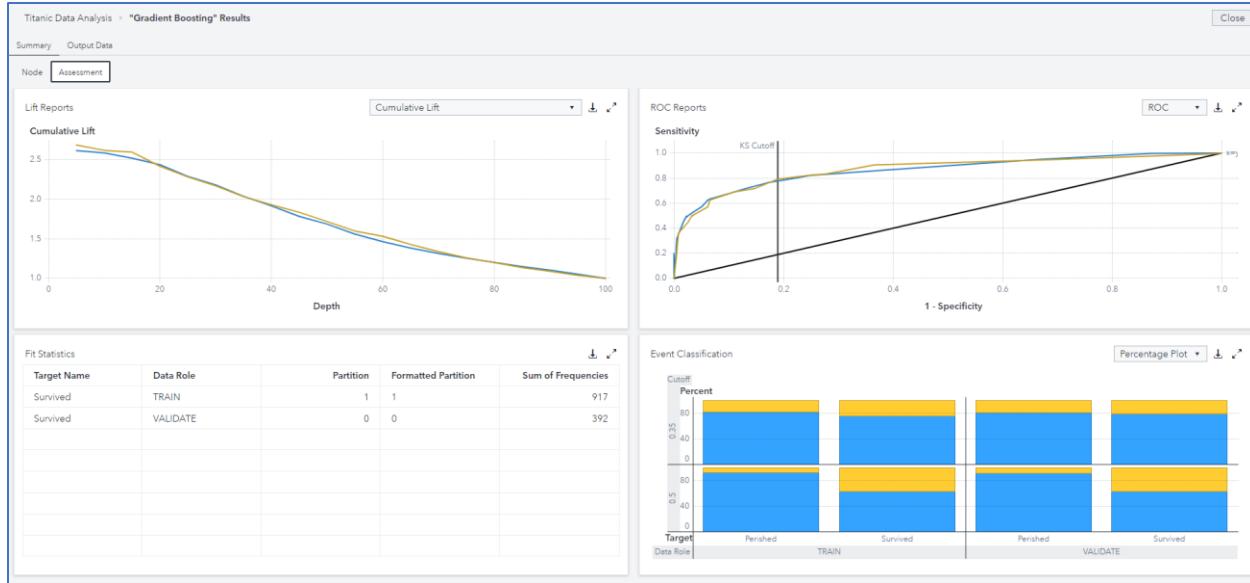
  

Properties	
Property Name	Property Value
codeLocation	mlearning
fullDatasetReconstitution	false
_omissionListing	true
dataMiningVersion	V8.5
reportingOnly	false

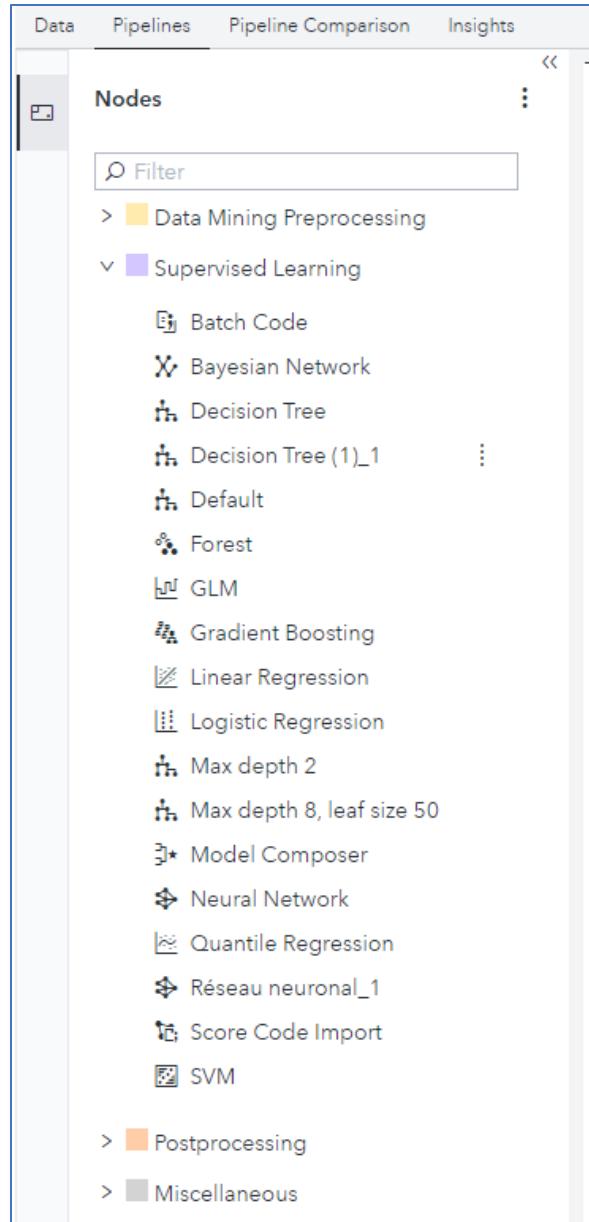
  

Output							
Age	7	double	8	4	BEST	4	0
Sex	8	double	8	1	BEST	1	0
Parch	9	double	8	1	BEST	1	0
Ticket	10	char	36	36	\$CHAR	36	0
Fare	11	double	8	8	BEST	8	0
Embarked	12	char	14	14	\$CHAR	14	0
<b>_va_d_PARTITION3369_</b>	13	double	8	12	0	0	0
<b>_va_c_Age_Groups</b>	14	char	13	13	0	0	0
<b>_va_c_Fare_Groups</b>	15	char	13	13	0	0	0

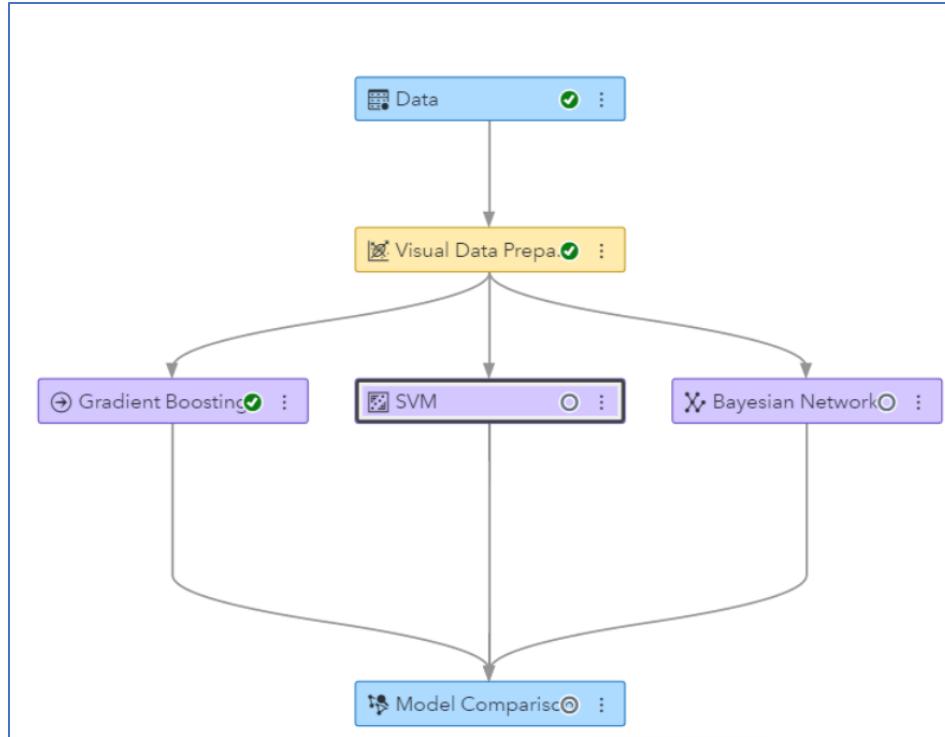
- Now, let's examine the Gradient Boosting Model



- Why stop here? Let's add more models to this pipeline!
  - General goal:
    - Can we do better than the default models of SAS Visual Analytics?
    - I think we can!
- Start by examining the types of models under the **Supervised Learning Nodes**
  - Note: your view of supervised learning models will differ slightly from mine – simply because we are (likely) on different instances of SAS Viya for Learners, and individual users can save custom nodes to the shared spaces.
  - For examples, “Max depth 2” is a user created node – whereas “Bayesian Network” is a node created by SAS. Use the SAS nodes to start.



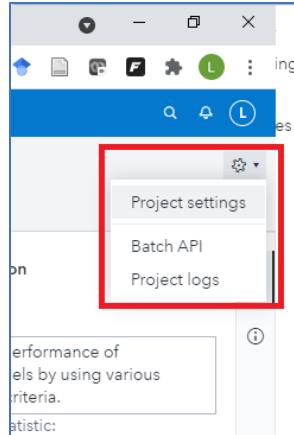
- To keep it simple – yet still expose you to new tools – let's add a SVM Model and a Bayesian network:
  - There are two general options to add nodes to an existing pipeline
    - You can drag-and-drop the node on top of Visual Data Preparation, or
    - Right-click on Visual Data Preparation and add the tool below.
  - After adding both notes, simply accept the defaults and your pipeline will appear as follows:



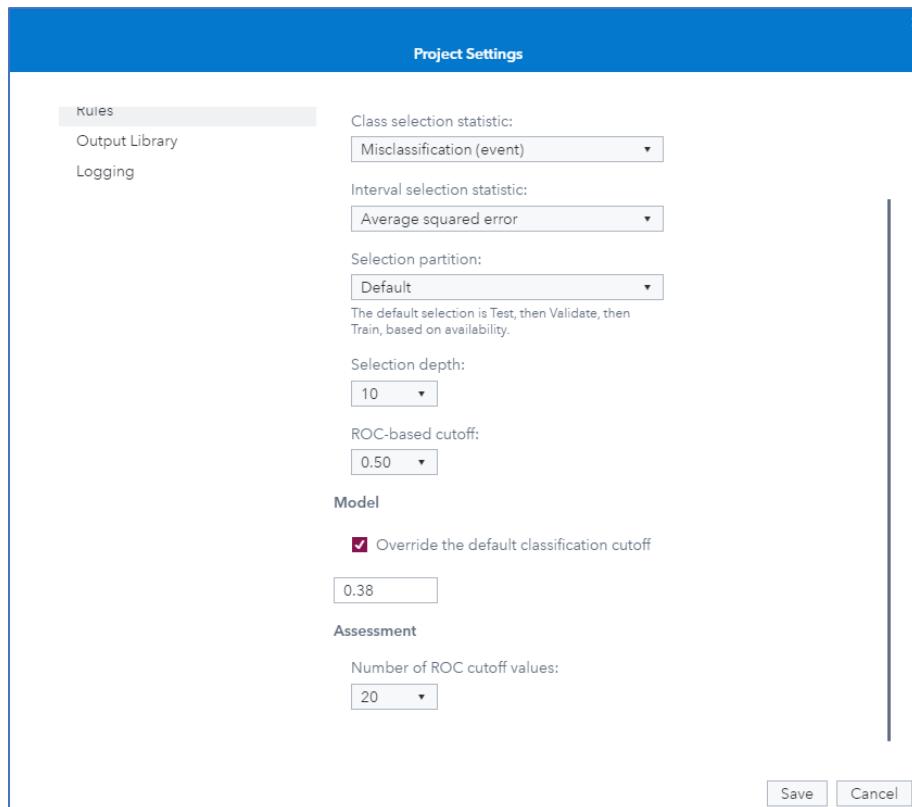
- Open the Model Comparison node and examine the models:

Model Comparison					
Champion	Name	Algorithm Name	KS (Youden)	Misclassification Rate	
☒	Gradient Boosting	Gradient Boosting	0.6087	0.1786	
	Bayesian Network	Bayesian Network	0.5323	0.2270	
	SVM	SVM	0.5241	0.2194	

- Did we beat the best model from SAS Visual Analytics?
  - No. No we did not.
  - On the plus side, it looks like the default model from SAS VA was a decent one!
- Moreover, note that we have many selection statistics to choose from.
  - We show just the KS(Youden) and Misclassification Rate in the screenshot above – but there are many more options available!
- Let's assume that we want to choose the misclassification rate.
  - We need to now adjust the global project setting to reflect this modeling choice.
  - To start, open project settings

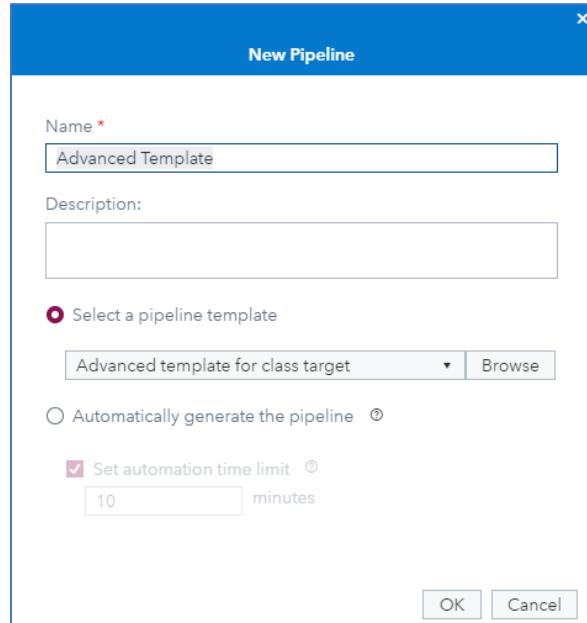


- Under Rules, select the following:
  - Class selection statistic: Misclassification (event)
  - Override the default classification cutoff = 0.38
- The final setting will appear as follows:

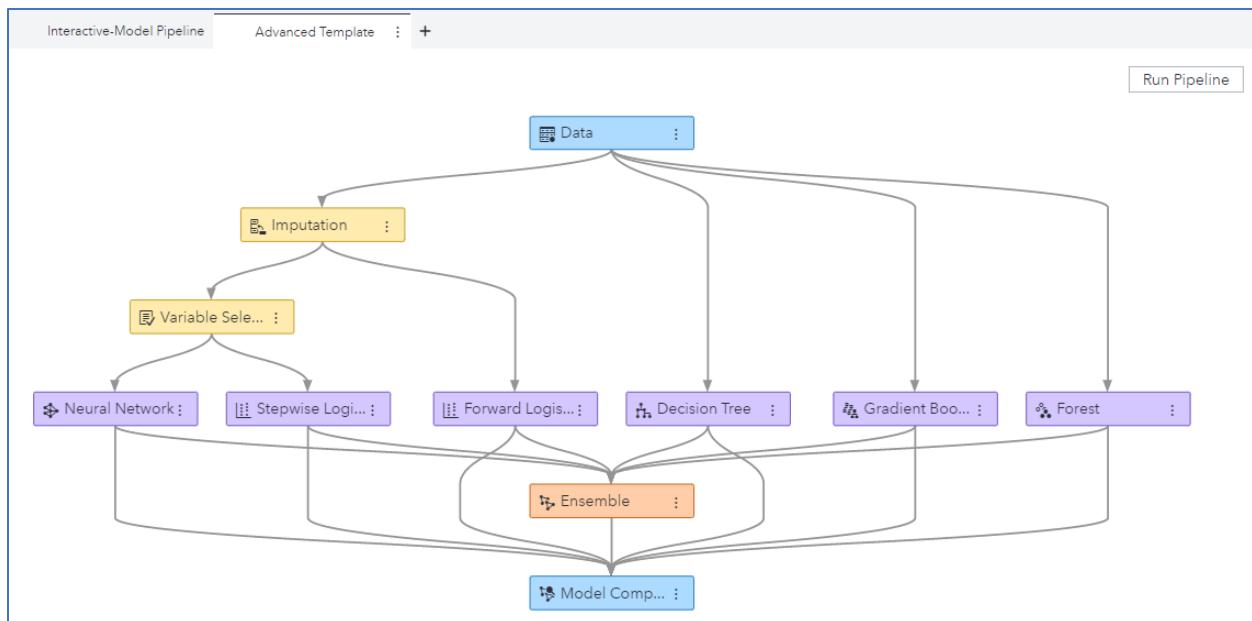


- Make sure you save before you exit this window
- Rerun your Interactive Model Pipeline (if needed)
- Now let's quickly add more pipelines – because this is the true value of SAS Model Studio!
  - We can add two pre-built pipeline templates to our analysis – and see if we can beat our SAS VA gradient boosting model
  - Pre-built templates are a great tool to speed up the modeling process

- **Template 1 – Advanced template**
  - Use the following settings:

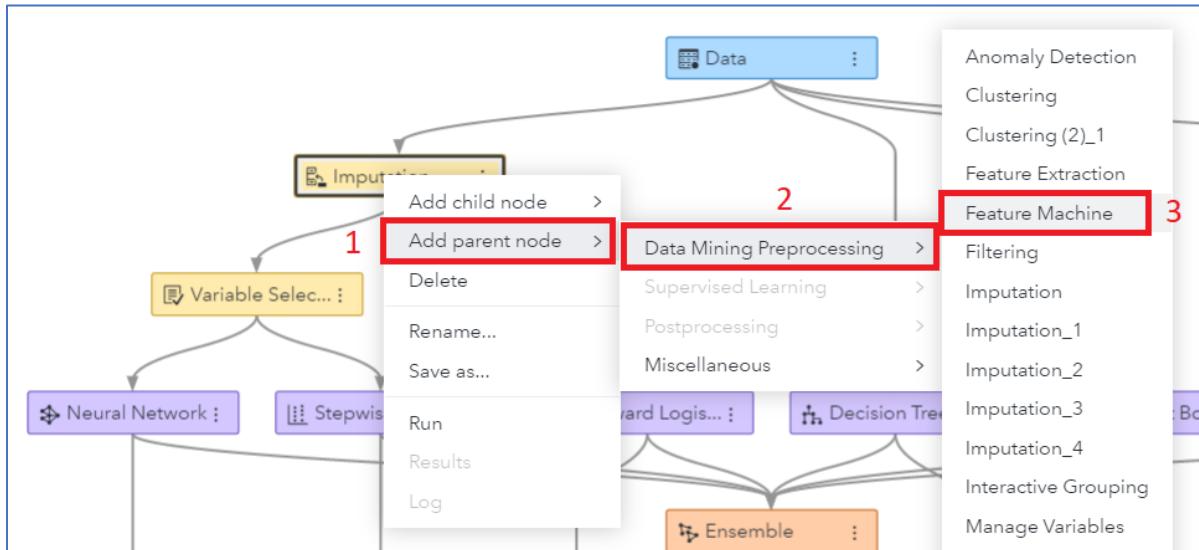


- Hint: Click on the **Browse** button to find the *Advanced template for class target*
- The new pipeline will appear as follows:

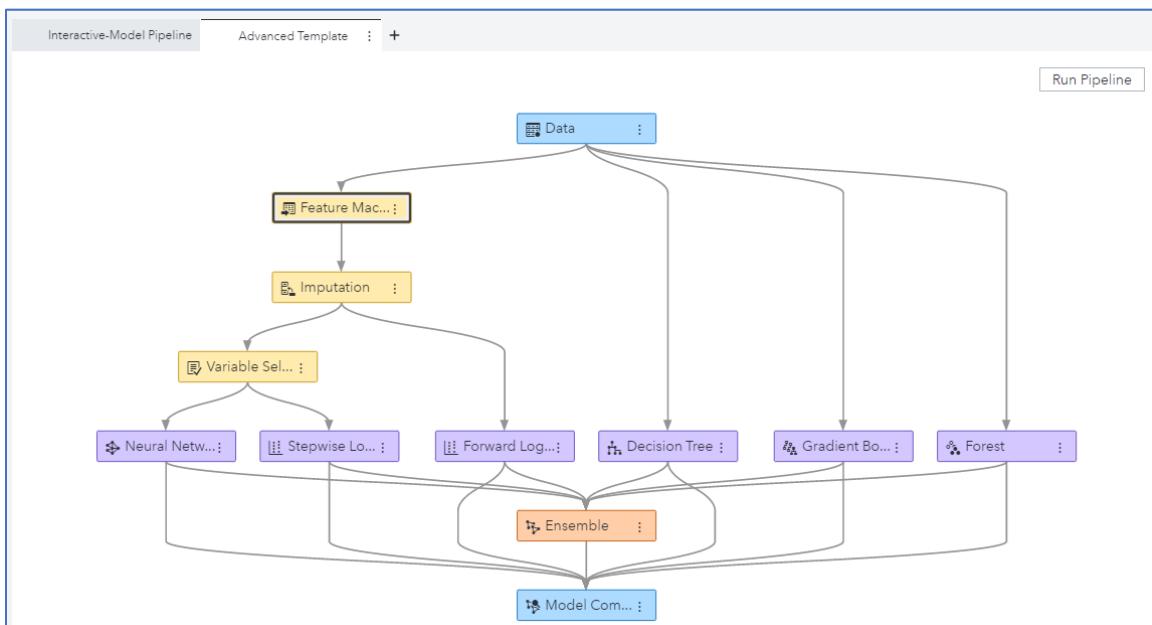


- This template is a great start...
  - But, it will not necessarily account for the fare and age variables we created.
    - Note: implicitly, tree-based models will, because split points can occur at any point on the interval. So, this is actually a significant upgrade over the variables we created.

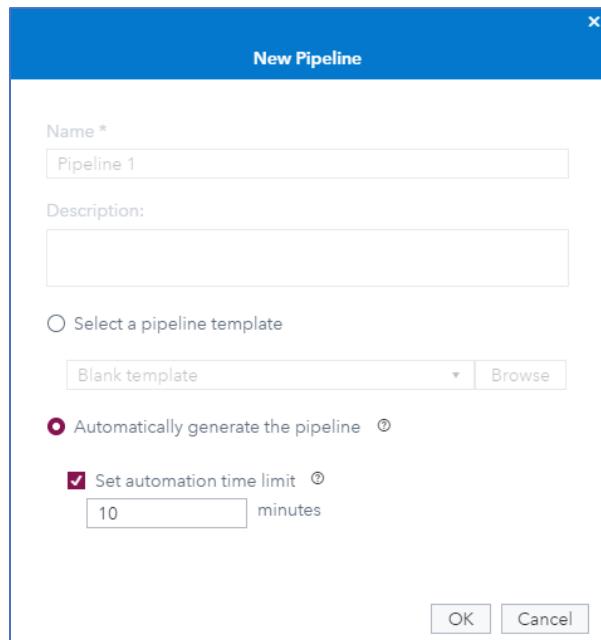
- But those are just half our models.
- For the Neural Network and the Logistic Models, we'll need to create new variables. Never fear, we can do this automatically with a **Feature Machine**!
  - Learn more about this AI tool here:  
<https://go.documentation.sas.com/doc/en/vdmmlcdc/8.5/vdmmlref/p180c9r7trn944n1ei7hilzu4cst.htm>
- You'll need to add the Feature Machine in a very precise location to have it affect our non-tree based models.
  - Start by selecting the **Imputation** node.
  - Right click, select Add parent node ➔ Data Mining Preprocessing ➔ Feature Machine:



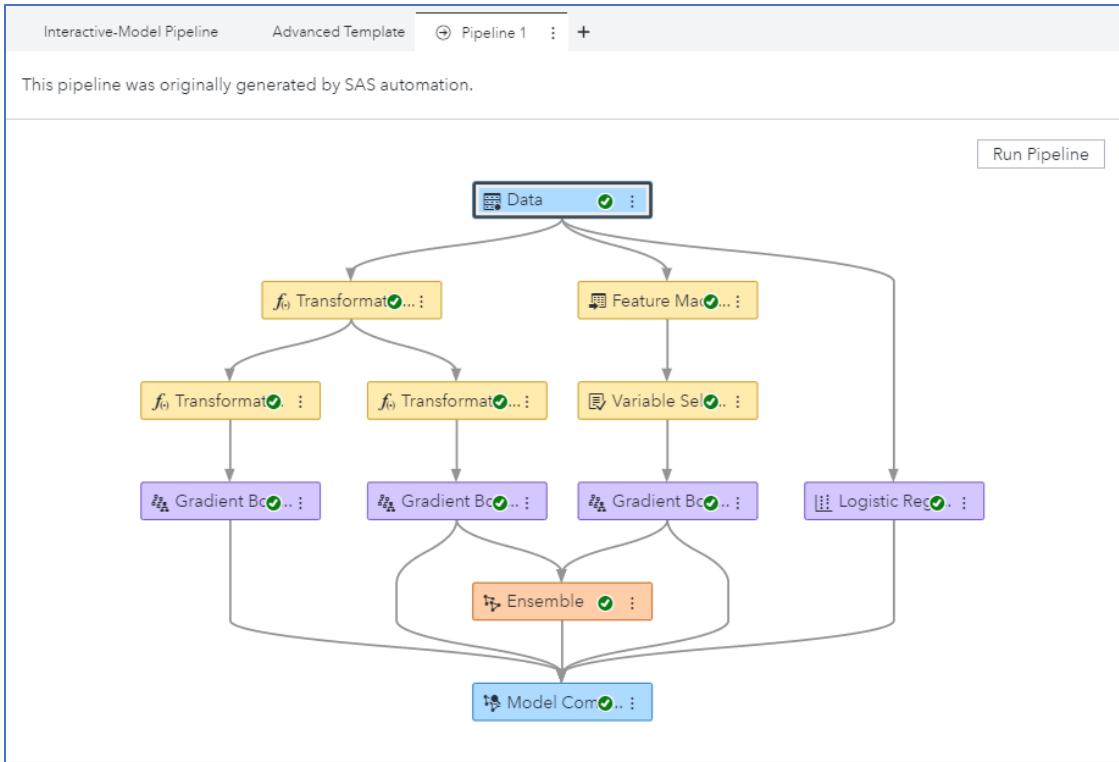
- Your new pipeline now appears as:



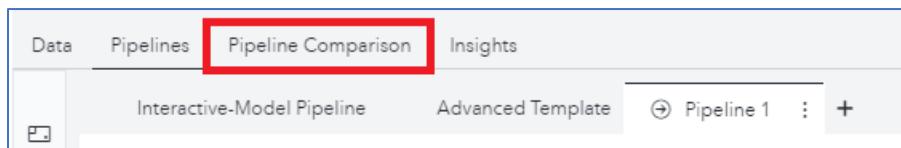
- Let it rip by clicking on *Run Pipeline*
- While that pipeline is running, create a third pipeline – this time **automatically generating the pipeline**
  - Create this **New Pipeline** by using the following settings:



- Keep the default automation time limit of 10 minutes.
- Click OK and generate the pipeline
  - Here is what mine looks like:



- Note: yours can/will look different!
- Run the pipeline once it is created
- **It's coffee break time!**
  - Take a break while the pipelines are running.
  - Don't worry, they're doing all the hard work for you.
- **Now we have pipelines to compare!**
  - While you can certainly examine the individual pipeline, most machine learning guru's only care about the predictive model performance – and want to get right to the champion models. We'll do the same.
  - Click on **Pipeline Comparison**:



- Which model has been crowned the champion? Our results:

The screenshot shows the SAS Machine Learning Tools interface with the 'Pipeline Comparison' tab selected. The table lists three models: Gradient Boosting, Gradient Boosting (1), and Forest. The champion model is highlighted with a red checkmark. Below the table are two code snippets: 'Path EP Score Code' and 'DS2 Package Code', both showing SAS code for the champion model.

Champion	Name	Algorithm Name	Pipeline Name	Misclassification Rate (Event)	Sum of Frequencies
<input checked="" type="checkbox"/>	Gradient Boosting	Gradient Boosting	Interactive-Model Pipeline	0.179	392
<input type="checkbox"/>	Gradient Boosting (1)	Gradient Boosting	Pipeline 1	0.189	392
<input type="checkbox"/>	Gradient Boosting	Gradient Boosting	Advanced Template	0.196	392

```

Path EP Score Code
1   data sasep.out;
2     dcl package score _AS59ICLPX4ZZ1UOVWBOIBRD13();
3     dcl nchar(8) "_va_c_Age_Groups" having label n'Age Groups';
4     dcl nchar(6) "_va_c_Fare_Groups" having label n'Fare Groups';
5     dcl double "P_SurvivedSurvived" having label n'Predicted: Su
6     dcl double "P_SurvivedPerished" having label n'Predicted: Su
7     dcl nchar(32) "I_Survived" having label n'Into: Survived';
8     dcl nchar(4) "_WARN_" having label n'Warnings';
9     dcl double EM_EVENTPROBABILITY;
10    dcl nchar(16) EM_CLASSIFICATION;
11    dcl double EM_PROBABILITY;
12    varlist allvars [_all_];
13

DS2 Package Code
1   package MS_b61e43ede15945c29e0206ddf8ff2db7_12MAY2023131958346
2     dcl package score _AS59ICLPX4ZZ1UOVWBOIBRD13();
3     dcl nchar(8) "_va_c_Age_Groups" having label n'Age Groups';
4     dcl nchar(6) "_va_c_Fare_Groups" having label n'Fare Groups';
5     dcl double "P_SurvivedSurvived" having label n'Predicted: Su
6     dcl double "P_SurvivedPerished" having label n'Predicted: Su
7     dcl nchar(32) "I_Survived" having label n'Into: Survived';
8     dcl nchar(4) "_WARN_" having label n'Warnings';
9     dcl double EM_EVENTPROBABILITY;
10    dcl nchar(16) EM_CLASSIFICATION;
11    dcl double EM_PROBABILITY;
12    varlist allvars [_all_];
13

```

- Our gradient boosting model retains its crown
- But, let's compare the **diagnostic statistics** for all three models:

The screenshot shows the Pipeline Comparison table with two rows highlighted: the champion model (row 1) and the second model (row 2). Red numbers '1' and '2' are overlaid on the first and second rows respectively.

Champion	Name	Algorithm Name	Pipeline Name	Misclassification Rate (Event)	Sum of Frequencies
<input checked="" type="checkbox"/>	Gradient Boosting	Gradient Boosting	Interactive-Model Pipeline	0.179	392
<input checked="" type="checkbox"/>	Gradient Boosting (1)	Gradient Boosting	Pipeline 1	0.189	392
<input checked="" type="checkbox"/>	Forest	Forest	Advanced Template	0.196	392

- The output from this comparison:

The screenshot shows the 'Titanic IV Modeling - Compare' interface. It includes a Pipeline Comparison table, two diagnostic plots (Lift Reports and ROC Reports), and a table of Fit Statistics.

Champion	Name	Algorithm Name	Pipeline Name	Misclassification Rate (Event)	Sum of Frequencies
<input checked="" type="checkbox"/>	Gradient Boosting	Gradient Boosting	Interactive-Model Pipeline	0.179	392
<input checked="" type="checkbox"/>	Gradient Boosting (1)	Gradient Boosting	Pipeline 1	0.189	392
<input checked="" type="checkbox"/>	Forest	Forest	Advanced Template	0.196	392

**Lift Reports**

**ROC Reports**

Statistics Label	Train: Gradient Boosting ...	Validate: Gradient Boosting ...	Train: Gradient Boosting ...	Validate: Gradient Boosting ...
Area Under ROC	0.9215	0.8589	0.8736	0.8602
Average Squared Error	0.1073	0.1365	0.1298	0.1359
Divisor for ASE	917	392	917	392
Formated Partition	1	0	1	0
Gmma	0.8654	0.7447	0.7923	0.7634
Gini Coefficient	0.8431	0.7178	0.7471	0.7203
KS (Youden)	0.6857	0.5721	0.6138	0.6087
KS Cutoff	0.4000	0.4000	0.3500	0.4000

- Anything interesting in this output?
- Select **Close** when you're done comparing the models
- Finally, let's find the **Score Code** for the champion model.
  - But, wait – what does it mean to “Score”?
    - Scoring simply means that we'll apply the parameters estimated in the champion model to a new set of data.

- In other words, we'll use the old data to predict new cases – often with new data – and, typically, want to predict an event that hasn't occurred yet.

- **Step 1:** select the single model

Pipeline Comparison		
	Data	Pipelines
<input type="checkbox"/> Champion	↓	Name
<input checked="" type="checkbox"/>	★	Gradient Boosting
<input type="checkbox"/>		Gradient Boosting (1)
<input type="checkbox"/>		Forest

- **Step 2:** scroll and find the **SAS Score Code**

```

Path EP Score Code
16 _4LWVZC11FOYJ4L50QAYCH0JJP.setvars(allvars);
17 _4LWVZC11FOYJ4L50QAYCH0JJP.setkey(n'7FCF38A8D864057E3E974F3C15C |
18 end;
19
20 method _EG1QQP180KH08QW4I1HGLC0ID();
21 dcl double "_VA_D_ONETIME848_103";
22 dcl double "_VA_D_ONETIME848";
23 dcl double "_VA_D_CALCULATION848_99";
24 dcl double "_VA_D_CALCULATION848_98";
25 dcl double "_VA_D_CALCULATION848_97";
26 dcl double "_VA_D_CALCULATION848_96";
27 dcl double "_VA_D_CALCULATION848_95";
28 dcl double "_VA_D_CALCULATION848_94";
29 dcl double "_VA_D_CALCULATION848_93";
30

DS2 Package Code
31 dcl double "_VA_D_CALCULATION848_92";
32 dcl double "_VA_D_CALCULATION848_91";
33 dcl double "_VA_D_CALCULATION848_25";
34 dcl double "_VA_D_CALCULATION848_24";
35 dcl double "_VA_D_CALCULATION848_23";
36 dcl double "_VA_D_CALCULATION848_22";
37 dcl double "_VA_D_CALCULATION848_21";
38 dcl double "_VA_D_CALCULATION848_20";
39 dcl double "_VA_D_CALCULATION848_19";
40 dcl double "_VA_D_CALCULATION848_18";
41 dcl double "_VA_D_CALCULATION848_17";
42 dcl double "_VA_D_CALCULATION848_16";
43 dcl double "_VA_D_CALCULATION848_15";
44 dcl double "_VA_D_CALCULATION848_102";
45

```

- This is the SAS code that you could use to apply the entire analysis you've conducted – in SAS Model Studio – to a new data set. Two notes:
  - In SAS Viya, there is a point-and-click tool that allows direct application of the champion model on new data, and it's called **SAS Model Manager**.
    - But that's a story – and tool – for another time.
  - Secondly, I'm not sure the Titanic case is the best predictive modeling example.
    - Why?
    - Well, it seems a little inappropriate/morbid to model the survival rates of passengers on a ship that hasn't sunken – yet – using data from 1912.
- Finally, one thing we haven't spent a lot of time analyzing is **which explanatory variables matter most** in predicting *Survived* on the Titanic.
  - Why?
  - Well, again, the truth is that machine learning specialists don't really care about the RHS variables – they just care that the model predicts well, relative to all other options.
    - Therefore, we fixate on competitor models – and our goodness-of-fit statistics.
  - All that stated, this doesn't mean that others won't still be interested in the inner details of your model.
    - In other words, many in your audience will want to know that your models make sense, even if predictive accuracy is really all they should care about.

- Have no fear – there's a SAS tool for that:



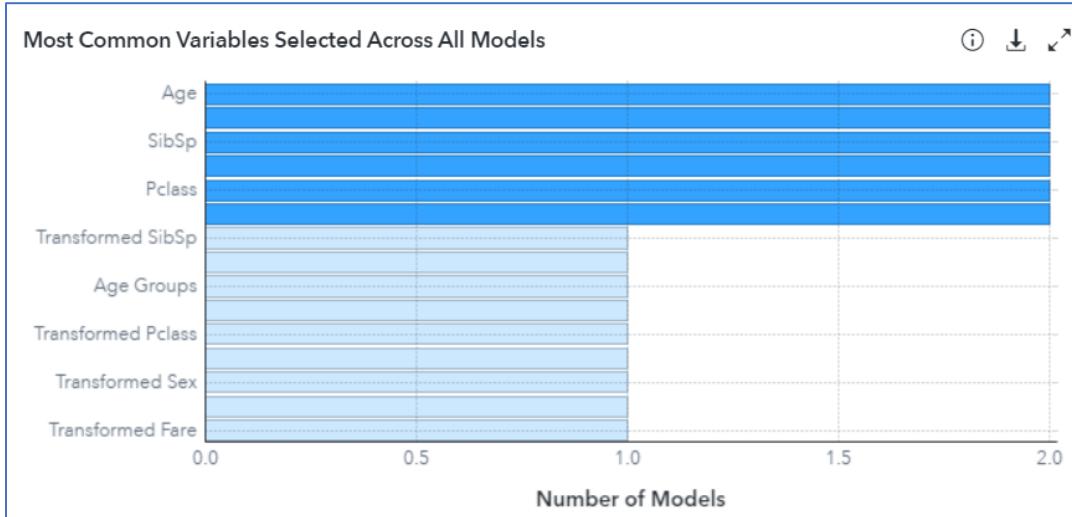
- Under the Insights tab, you'll find a whole host of information that will help with your storytelling. Some fun highlights:
  - Project Summary:

**Project Summary**

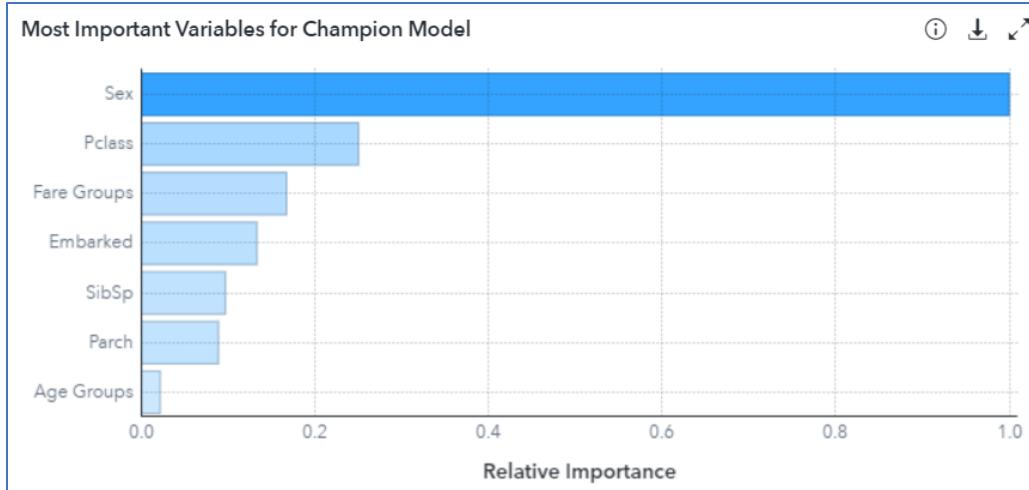
The champion model for this project is Gradient Boosting. The model was chosen based on the Misclassification Rate (Event) for the Validate partition (0.18). 82.14% of the Validate partition was correctly classified using the Gradient Boosting model. The five most important factors are Sex, Pclass, Fare Groups, Embarked, and SibSp.

<b>Project Target:</b>	Survived	<b>Project Champion:</b>	Gradient Boosting
<b>Event Percentage:</b>	38.2735 %	<b>Created By:</b>	Lincoln.Groves@sas.com
<b>Pipelines:</b>	3	<b>Modified:</b>	May 12, 2023 1:42:09 PM

- Most common variables used across all modeling options:



- Most important variables in the champion model:



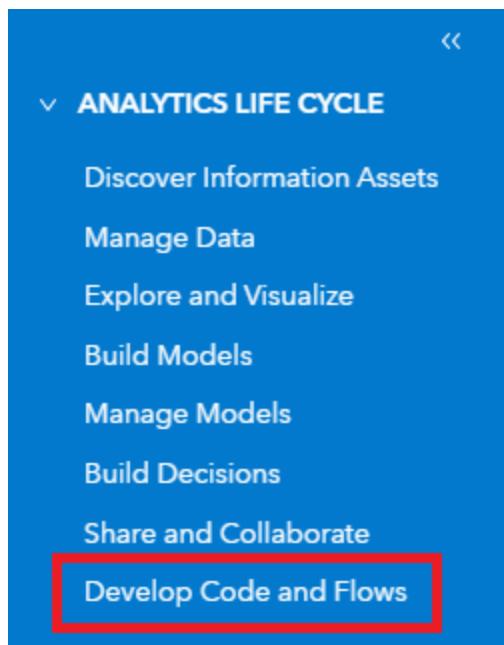
- For details on relative importance and its calculation in an example, learn more [here](#).

### Part III: Programming in SAS Studio

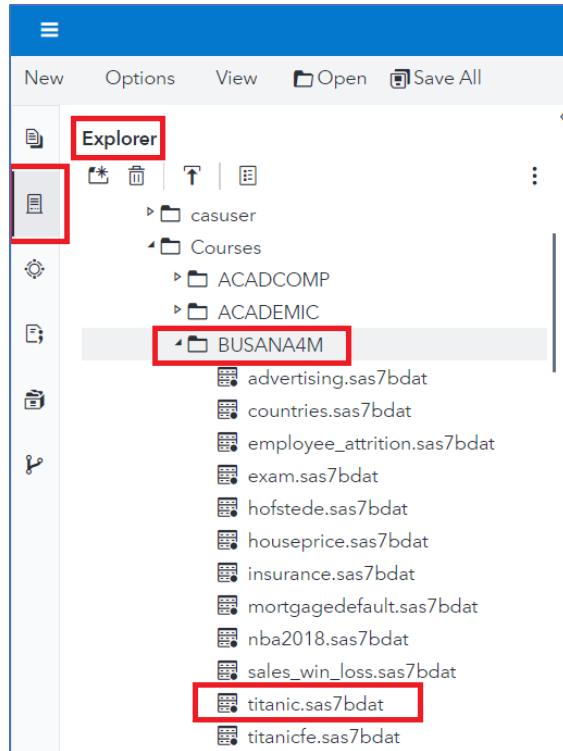
- **Overview**
  - This section is all about the SAS Studio Tasks
  - We will explore the following data lifecycle tools:
    - Prepare Data
      - Examine Data
    - SAS Viya Prepare and Explore Data
      - Transform Data
    - SAS Viya Machine Learning
      - Supervised Learning
- **How to find SAS Studio?**
  - Return to the Hamburger in the upper left-hand side!



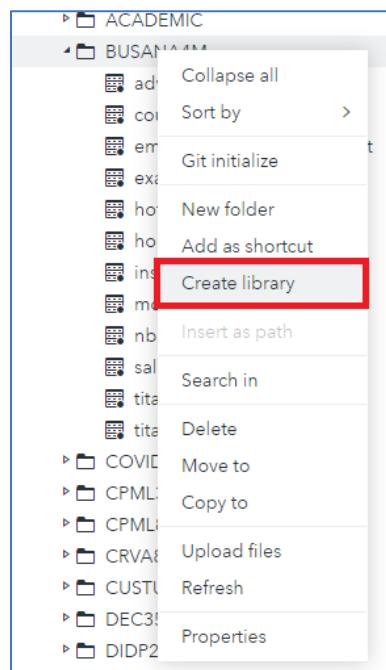
- Select *Develop Code and Flows*



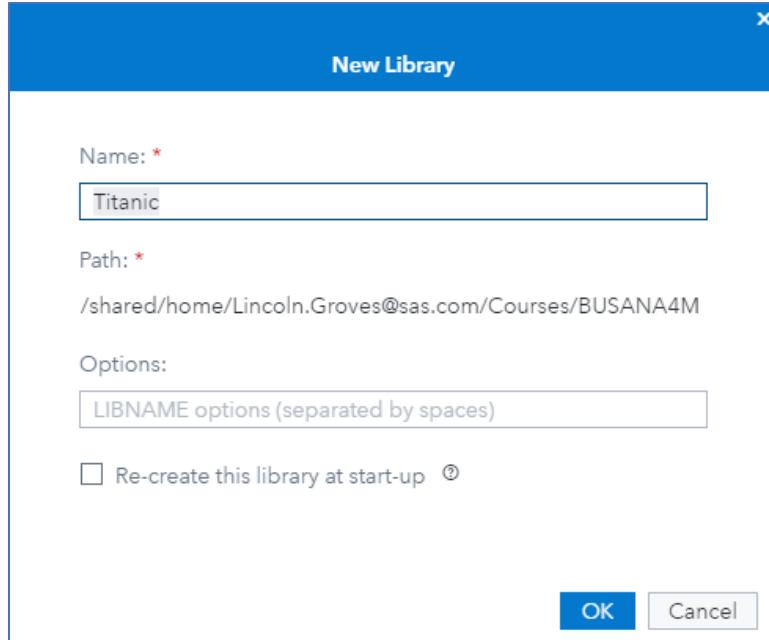
- **We need to first prepare the data**
  - **Step 1** – find the Titanic Data in SAS Studio



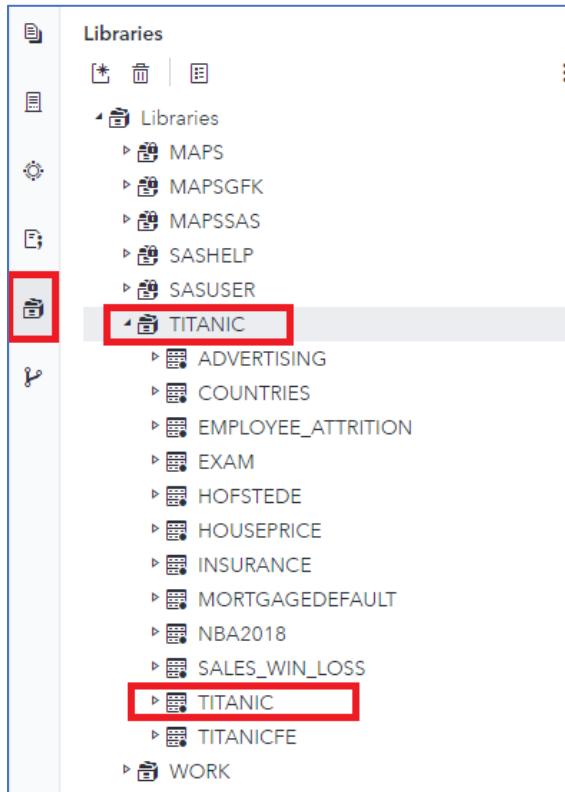
- Step 2 – create a SAS Library
  - Wait, what is a SAS Library?
    - Put simply, it tells the software where your data sets are stored
  - To create a SAS library, right-click on the BUSANA4M folder and select “Create library”:



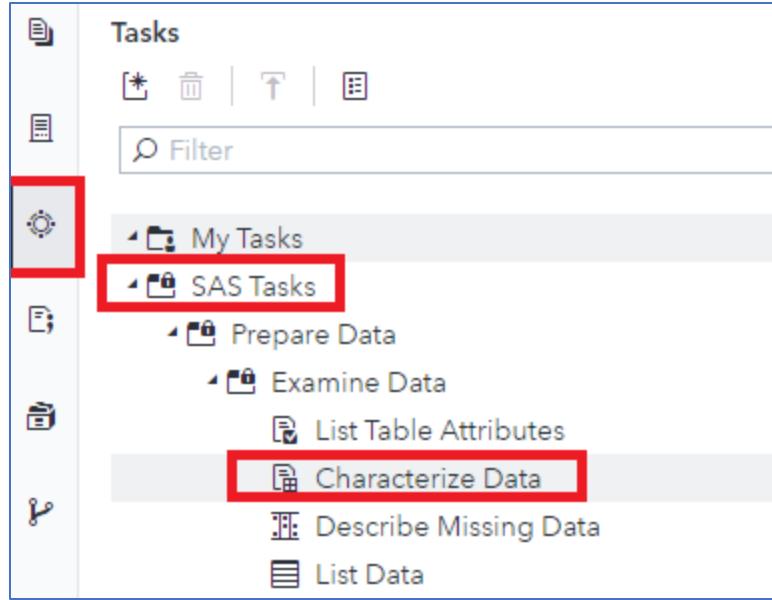
- Use the following general settings in your new library:



- Ensure that your library has been initialized by checking that your library can be accessed:



- Better understand your data with the Characterize Data Tasks
  - Find it here:



- In the Characterize Data Task, do the following:
  - Select **titanic.titanic** as the Data
  - Select all variables for Automatic Characterization
- See the settings and SAS Code below:

The screenshot shows the SAS Studio Task Editor for the 'Characterize Data' task. On the left, the 'DATA' tab is selected, showing 'TITANIC.TITANIC' as the source. The 'AUTOMATIC CHARACTERIZATION' tab is active, with all variables ('PassengerId', 'Survived', 'Pclass', 'Name', 'Sex') selected for automatic characterization. On the right, the 'Code' tab displays the generated SAS code:

```

1 /* 
2 * Task code generated by SAS® Studio 5.2
3 *
4 * Generated on '11/15/21, 12:58 PM'
5 * Generated by 'Lincoln.Groves@sas.com'
6 * Generated on server 'pdcesx23113'
7 * Generated on SAS platform 'Linux LIN X64 3.10.0-862.9.1.e17.x86_64'
8 * Generated on SAS version 'V.03.05M0P111119'
9 * Generated on browser 'Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML, like Gecko)'
10 * Generated on web client 'https://v4e055.vfe.sas.com/SASStudioV/main?locale=en_US&launchedFromAppSwitcher'
11 */
12 /*
13 ods noproctitle;
14
15 /** Analyze categorical variables ***/
16 title "Frequencies for Categorical Variables";
17
18 proc freq data=TITANIC.TITANIC;
19   tables Survived Pclass Name Sex Ticket Cabin Embarked / plots=(freqplot);
20 run;
21
22 /** Analyze numeric variables ***/
23 title "Descriptive Statistics for Numeric Variables";
24
25 proc means data=TITANIC.TITANIC n nmiss min mean median max std;
26   var PassengerId Age SibSp Parch Fare;
27 run;
28
29 title;
30
31 proc univariate data=TITANIC.TITANIC noprint;
32   histogram PassengerId Age SibSp Parch Fare;
33 run;
34

```

- Run the code by clicking and examine the output as follows:

## Ontario Tech University - Seminar in SAS Machine Learning Tools - MAY2023

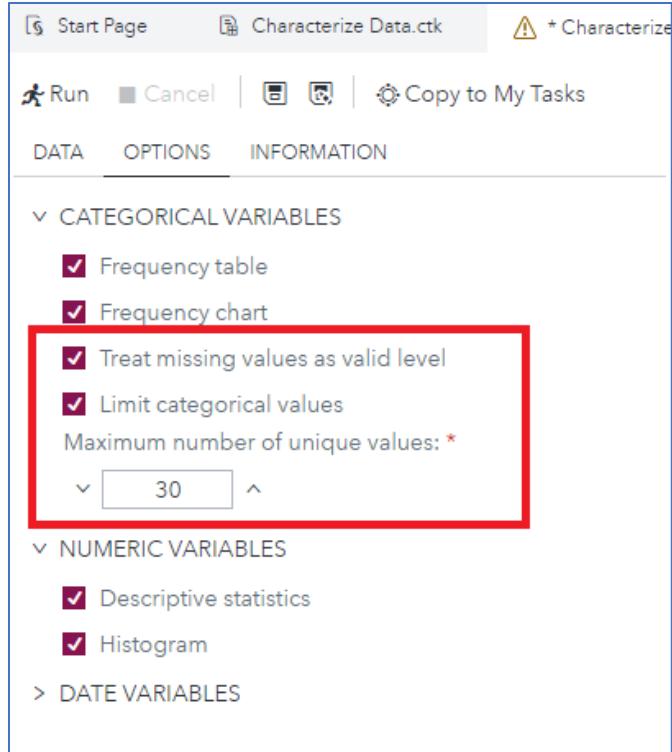
The screenshot shows the SAS Studio interface with the 'Characterize Data' task open. The left sidebar lists the DATA, AUTOMATIC CHARACTERIZATION, and CUSTOM CHARACTERIZATION sections. In the AUTOMATIC CHARACTERIZATION section, 'PassengerId' is marked as a key variable. The right pane displays the results of the FREQ procedure, including a bar chart titled 'Distribution of Survived' and a frequency table for Pclass.

Survived	Frequency	Percent	Cumulative Frequency	Cumulative Percent
Perished	808	61.73	808	61.73
Survived	501	38.27	1309	100.00

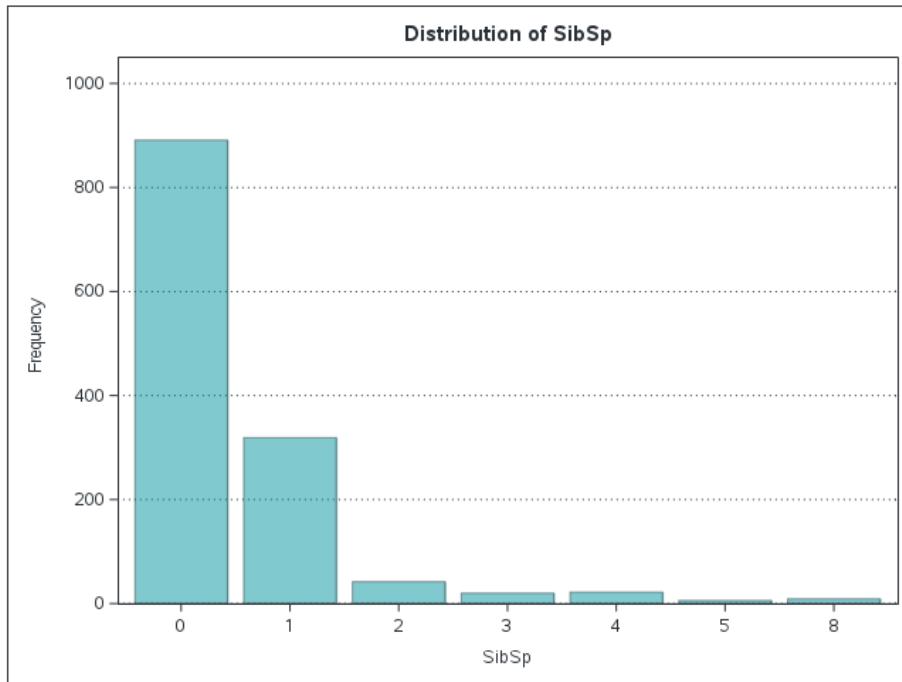
  

Pclass	Frequency	Percent	Cumulative Frequency	Cumulative Percent
Lower	709	54.16	709	54.16
Middle	277	21.16	986	75.32
...	...	...	...	...

- TMI? Yeah, I think so! Let's make the following changes to the SAS Studio task:
  - Drop *PassengerID*, *Name*, *Ticket* from **Automatic Characterization**
  - Change *SibSp* and *parch* to **Categorical Variables**
    - They will display better as discrete variables
  - Click on the **Options** tab. Then under **Categorical Variables** options, select:
    - Treat missing values as valid level
    - Limit categorical values to a max of 30 variables
      - This is not required but useful to know it exists

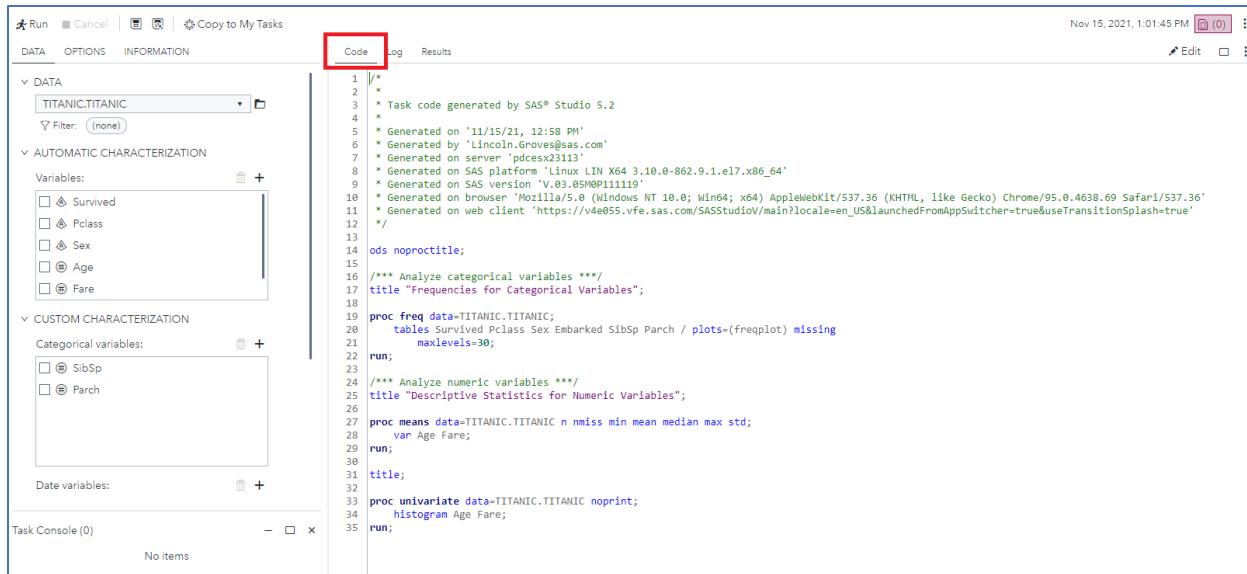


- Rerun the SAS Studio Task and examine the output
  - In particular, examine the distribution of Sibling/Spouse:



- That's a bit easier to read...
- As a last step, examine the SAS Code. We'll return to the SAS Code eventually.

- But it is nice to know that SAS is doing a lot of work for you behind the scenes.

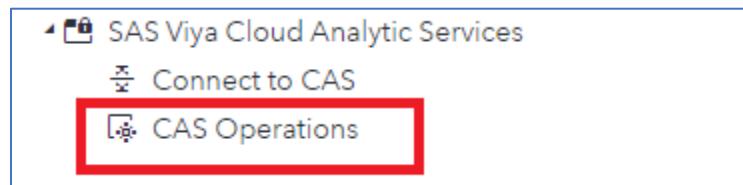


```

1 /*
2 * Task code generated by SAS® Studio 5.2
3 *
4 * Generated on '11/15/21, 12:58 PM'
5 * Generated by 'Lincoln.Grove@sas.com'
6 * Generated on server 'pdcess23113'
7 * Generated on SAS platform 'Linux LIN X64 3.10.0-862.9.1.el7.x86_64'
8 * Generated on SAS version 'V8.0.89M0P111119'
9 * Generated on browser 'Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/95.0.4638.69 Safari/537.36'
10 * Generated on web client 'https://v4e055.vfe.sas.com/SASStudioV/main?locale=en_US&launchedFromAppSwitcher=true&useTransitionSplash=true'
11 */
12 /*
13 *
14 ods noproctitle;
15 *
16 /** Analyze categorical variables ***
17 title "Frequencies for Categorical Variables";
18 *
19 proc freq data=TITANIC.TITANIC;
20   tables Survived Pclass Sex Embarked SibSp Parch / plots=(freqplot) missing
21   maxlevels=30;
22 run;
23 *
24 /** Analyze numeric variables ***
25 title "Descriptive Statistics for Numeric Variables";
26 *
27 proc means data=TITANIC.TITANIC n nmiss min mean median max std;
28   var Age Fare;
29 run;
30 *
31 title;
32 *
33 proc univariate data=TITANIC.TITANIC noprint;
34   histogram Age Fare;
35 run;

```

- Our next step in the data prep process is to use **SAS Viya Prepare and Explore Data**
  - To start, we need to push data up into CAS
    - Wait, **what is CAS?**
      - CAS stands for Cloud Analytics Services. It's a multinet cloud environment founded upon microservices that can run models on a distributed system.
      - The shorter answer is that data run faster in the cloud. And CAS is SAS' tool to run models optimally on the cloud.
    - Why push it?
      - By pushing data up to CAS, we can use a suite of tools specifically designed to leverage machine learning in the cloud.
  - The process:
    - Start with a **SAS Studio Tasks** and navigate to the following:



- Select the following settings for the **CAS Operations Task** in SAS Studio:

```

1 /*
2 *
3 * Task code generated by SAS® Studio 5.2
4 *
5 * Generated on '11/15/21, 2:09 PM'
6 * Generated by 'Lincoln.Groves@sas.com'
7 * Generated on server 'pdcesx23113'
8 * Generated on SAS platform 'Linux LIN X64 3.10.0-862.9.1.el7.x86_64'
9 * Generated on SAS version 'V.03.05M0P111119'
10 * Generated on browser 'Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/95.0.463.54 Safari/537.36'
11 * Generated on web client 'https://v4e055.vfe.sas.com/SASStudioV/main?locale=en_US&launchedFromAppSwitcher=true&useTrans'
12 */
13
14 proc casutil;
15   load data=TITANIC.TITANIC casout='Titanic_CAS' outcaslib='casuser';
16 run;

```

- But wait!
  - The code is not yet ready to be submitted, because we haven't assigned defined a CAS engine libref for our CAS in-memory data tables.
  - We can do this by adding one line to the code to the CAS Operations task.
- To start the coding process, select the **Edit button** on the code generated by SAS:

```

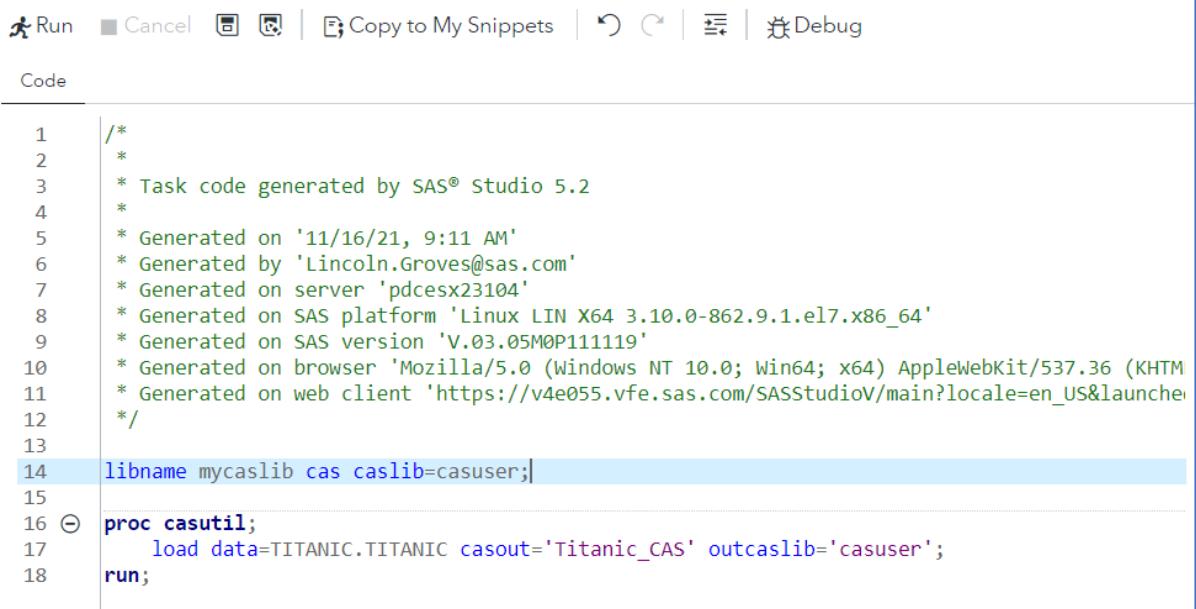
1 /*
2 *
3 * Task code generated by SAS® Studio 5.2
4 *
5 * Generated on '11/16/21, 9:11 AM'
6 * Generated by 'Lincoln.Groves@sas.com'
7 * Generated on server 'pdcesx23104'
8 * Generated on SAS platform 'Linux LIN X64 3.10.0-862.9.1.el7.x86_64'
9 * Generated on SAS version 'V.03.05M0P111119'
10 * Generated on browser 'Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/95.0.463.54 Safari/537.36'
11 * Generated on web client 'https://v4e055.vfe.sas.com/SASStudioV/main?local'
12 */
13
14 proc casutil;
15   load data=TITANIC.TITANIC casout='Titanic_CAS' outcaslib='casuser';
16 run;

```

- This creates a new SAS Program. To get that statement to execute, add the following line of code before the proc casutil statement:

**libname mycaslib cas caslib=casuser;**

- The program now appears as follows:

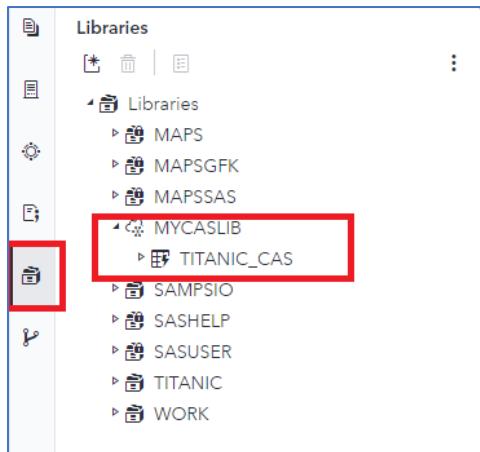


```

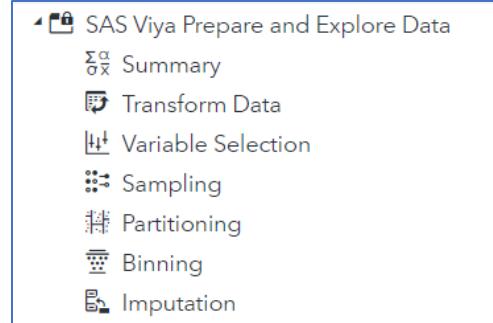
1 /*
2 *
3 * Task code generated by SAS® Studio 5.2
4 *
5 * Generated on '11/16/21, 9:11 AM'
6 * Generated by 'Lincoln.Groves@sas.com'
7 * Generated on server 'pdcesx23104'
8 * Generated on SAS platform 'Linux LIN X64 3.10.0-862.9.1.el7.x86_64'
9 * Generated on SAS version 'V.03.05M0P111119'
10 * Generated on browser 'Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML
11 * Generated on web client 'https://v4e055.vfe.sas.com/SASStudioV/main?locale=en_US&launche
12 */
13
14 libname mycaslib cas caslib=casuser;
15
16 proc casutil;
17   load data=TITANIC.TITANIC casout='Titanic_CAS' outcaslib='casuser';
18 run;

```

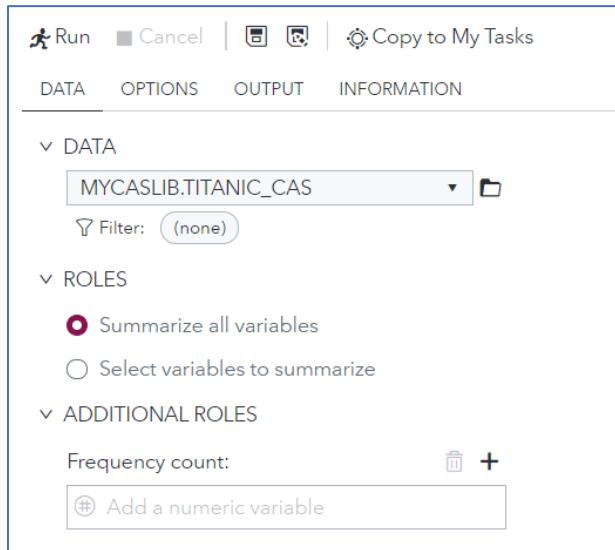
- Submit the code by clicking 
- Again, let's ensure that the file has loaded by checking the library, which can be found under MYCASLIB:



- Summary, One More Time
  - For completeness, let's examine output from the SAS Viya Prepare and Explore Data Tasks once more.
  - Start with the **SAS Viya Prepare and Explore Data** tasks:



- In particular, select **Summary**, with the following basic settings:



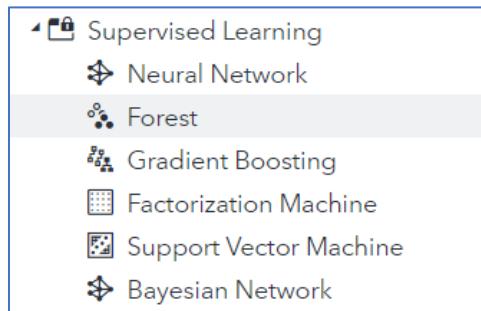
- Here is some of the summary output:

Obs	Variable name	Width of the variable formatted value	Type of the raw values	Recommended level for analytics	Have more unreported levels	N
1	PassengerId		4 N	ID	Y	
2	Survived		16 C	CLASS	N	
3	Pclass		12 C	CLASS	N	
4	Name		164 C	ID	Y	
5	Sex		12 C	CLASS	N	
6	Age		4 N	INTERVAL	Y	
7	SibSp		1 N	CLASS	N	
8	Parch		1 N	CLASS	N	
9	Ticket		36 C	ID	Y	
10	Fare		8 N	INTERVAL	Y	
11	Cabin		30 C	ID	Y	
12	Embarked		14 C	CLASS	N	

Level Details						
Obs	Variable name	Level index	Level frequency	Percentage of the frequency of current level	Percentage of non missing values	Raw numeric value
1	PassengerId	1	1	0.076394194	0.076394194	1
2	PassengerId	2	1	0.076394194	0.076394194	2
3	PassengerId	3	1	0.076394194	0.076394194	3

- What, if any, new information did you learn from this analysis?
- Alright, now on to **SAS Viya Machine Learning**
  - Here comes the good part!
    - One note before we get started: with coding, notice how many more options you have... but how much more time consuming the process can be!
  - The supervised learning tools available:



- Let's do a deeper dive into **Supervised Learning using SAS Studio Tasks**
  - Model 1: **Gradient Boosting Model**
    - Use the following settings:

The screenshot displays the SAS Studio interface for a Gradient Boosting Model. It includes three main panels: DATA, OPTIONS, and INFORMATION.

- DATA Tab:** Shows the dataset 'MYCASLIB.TITANIC\_CAS' selected. Under 'Partition Data', there are checkboxes for 'Validation data' and 'Test data'. Under 'ROLES', the 'Target' section has two options: 'Use a nominal target' (selected) and 'Use an interval target'. The 'Nominal target:' field contains 'Survived'.
- OPTIONS Tab:** Shows the configuration for the model inputs.
  - Interval inputs:** Contains 'Age' and 'Fare'.
  - Nominal inputs:** Contains 'Parch', 'SibSp', 'Embarked', 'Pclass', and 'Sex'.
- INFORMATION Tab:** Not visible in the screenshot.

- Submit the code and examine the **Output**:

Model Information	
Number of Trees	100
Learning Rate	0.1
Subsampling Rate	0.5
Number of Variables Per Split	7
Number of Bins	50
Number of Input Variables	7
Maximum Number of Tree Nodes	31
Minimum Number of Tree Nodes	15
Maximum Number of Branches	2
Minimum Number of Branches	2
Maximum Depth	4
Minimum Depth	4
Maximum Number of Leaves	16
Minimum Number of Leaves	8
Maximum Leaf Size	520
Minimum Leaf Size	5
Seed	106187529
Lasso (L1) penalty	0
Ridge (L2) penalty	1
Actual Number of Trees	100
Average Number of Leaves	13.39

Training	
Number of Observations Read	1309
Number of Observations Used	1309

Variable Importance			
Variable	Importance	Std Dev	Relative Importance
Sex	5.6400	17.6343	1.0000
Fare	4.4698	0.8795	0.7925
Age	3.7470	0.9282	0.6644
Pclass	2.2488	3.3836	0.3987
SibSp	1.2356	0.9545	0.2191
Embarked	0.8325	0.8117	0.1476
Parch	0.7256	0.7450	0.1287

- Is this a good model?
  - Examine which options you can change to refine the results.

DATA	OPTIONS	OUTPUT	INFORMATION
<div style="border-bottom: 1px solid #ccc; padding-bottom: 5px;"> <b>METHODS</b> </div> <div style="margin-top: 5px;"> <input type="checkbox"/> Automatically tune selected options <small>②</small> </div> <div style="margin-top: 5px;">         Lasso regularization: <small>*</small> <div style="border: 1px solid #ccc; padding: 2px; width: 100px; margin-top: 2px;">0</div> </div> <div style="margin-top: 5px;">         Ridge regularization: <small>*</small> <div style="border: 1px solid #ccc; padding: 2px; width: 100px; margin-top: 2px;">1</div> </div> <div style="margin-top: 5px;">         Sampling proportion of training observations for each iteration: <small>*</small> <div style="border: 1px solid #ccc; padding: 2px; width: 100px; margin-top: 2px;">0.5</div> </div> <div style="margin-top: 5px;">         Learning rate: <small>*</small> <div style="border: 1px solid #ccc; padding: 2px; width: 100px; margin-top: 2px;">0.1</div> </div> <div style="margin-top: 5px;">         Number of trees: <small>*</small> <div style="border: 1px solid #ccc; padding: 2px; width: 100px; margin-top: 2px;">100</div> </div> <div style="margin-top: 5px;">         Number of inputs considered for splitting a node:         <div style="border: 1px solid #ccc; padding: 2px; width: 100px; margin-top: 2px;">All inputs (Default)</div> </div> <div style="margin-top: 5px;">         Maximum depth of a tree: <small>*</small> <div style="border: 1px solid #ccc; padding: 2px; width: 100px; margin-top: 2px;">4</div> </div> <div style="margin-top: 5px;">         Minimum number of observations for a leaf: <small>*</small> <div style="border: 1px solid #ccc; padding: 2px; width: 100px; margin-top: 2px;">5</div> </div> <div style="margin-top: 5px;">         Maximum number of branches for a node: <small>*</small> </div>			

 Code Log Results Edit ⋮   ``` 1 /* 2 * 3 * Task code generated by SAS® Studio 5.2 4 * 5 * Generated on '11/16/21, 9:31 AM' 6 * Generated by 'Lincoln.Groves@sas.com' 7 * Generated on server 'pdcesx23104' 8 * Generated on SAS platform 'Linux LIN X64 3.10.0-862.9.1.el7.x86_64' 9 * Generated on SAS version 'V.03.05MOP111119' 10 * Generated on browser 'Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/91.0.4453.102 Safari/537.36' 11 * Generated on web client 'https://v4e055.vfe.sas.com/SASstudioV/main?loc= 12 */ 13 ods noproctitle; 14 15 proc gradboost data=MYCASLIB.TITANIC_CAS; 16   target Survived / level=nominal; 17   input Age Fare / level=interval; 18   input Parch SibSp Embarked Pclass Sex / level=nominal; 19   run; ``` |

#### ■ Model 2: Neural Network

- Use the following settings:

The screenshot shows the SAS Machine Learning Tools interface. On the left, under 'DATA', the dataset 'MYCASLIB.TITANIC\_CAS' is selected. Under 'ROLES', 'Target' is set to 'Use a nominal target' with 'Survived' as the nominal target. On the right, the 'Inputs' panel is open, divided into 'Interval inputs' and 'Nominal inputs'. 'Interval inputs' include 'Age' and 'Fare'. 'Nominal inputs' include 'Parch', 'SibSp', 'Embarked', 'Pclass', and 'Sex'.

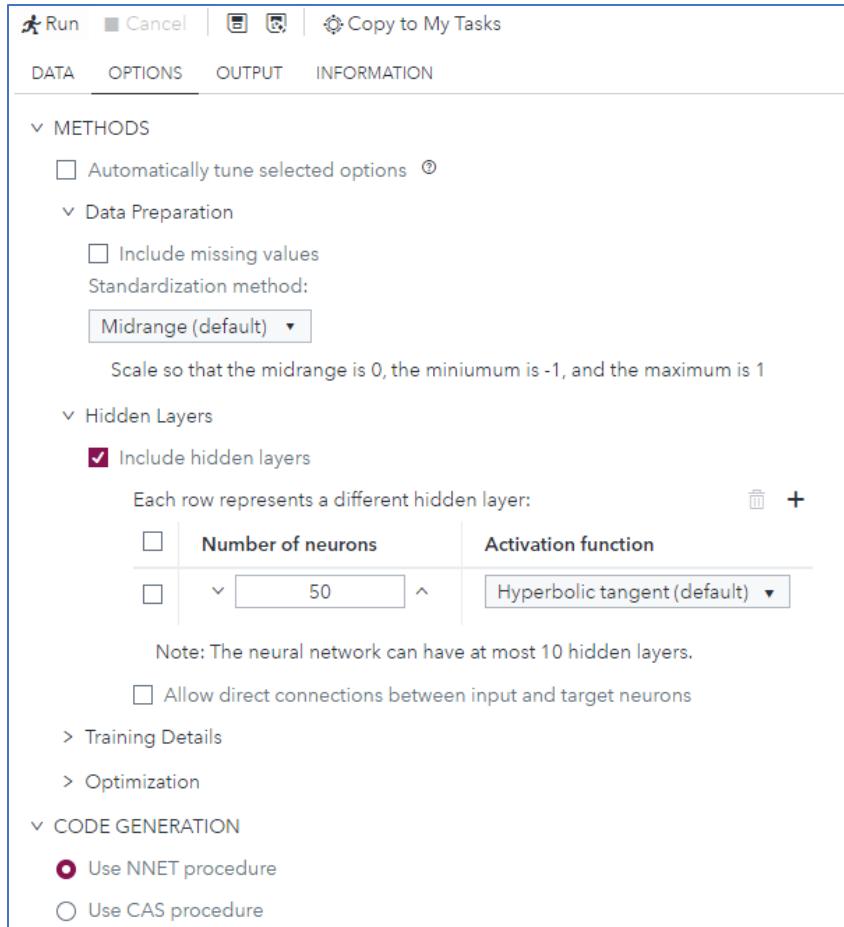
- Submit the code and examine the Output

The screenshot shows two tables. The 'Model Information' table provides details about the neural network model, such as the number of observations (1309), target variable ('Survived'), and architecture ('MLP'). The 'Iteration History' table details the optimization process, showing 17 iterations with columns for Iteration Number, Objective Function, Norm of Gradient, Loss, Step Size, and various Norm metrics (L1, L2, Maximum, Fit Error).

Model Information							
Model	Neural Net						
Number of Observations Used	1309						
Number of Observations Read	1309						
Target/Response Variable	Survived						
Number of Nodes	77						
Number of Input Nodes	25						
Number of Output Nodes	2						
Number of Hidden Nodes	50						
Number of Hidden Layers	1						
Number of Weight Parameters	1300						
Number of Bias Parameters	52						
Architecture	MLP						
Seed for Initial Weight	1605445419						
Optimization Technique	LBFGS						
Number of Neural Nets	1						
Objective Value	2.0416870809						

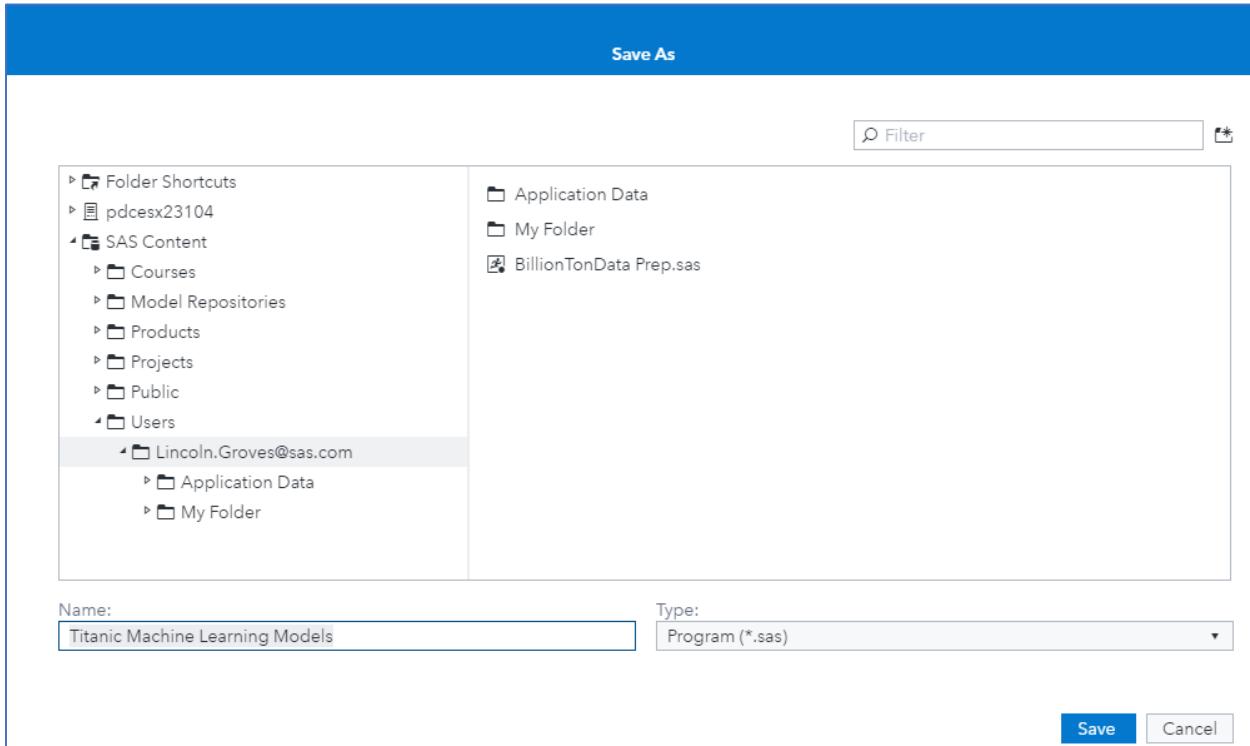
Iteration Number	Objective Function	Norm of Gradient	Loss	Step Size	Iteration History			
					Norm	L1	L2	Maximum
1	2.857615	0.522868	2.774181	0	26.05427	0.834338	0.040000	0.627960
2	2.807446	0.575029	2.730891	1.875958	22.58570	0.765552	0.109539	0.382735
3	2.679234	0.329565	2.605401	1	22.06933	0.738328	0.096917	0.382735
4	2.463666	1.792478	2.280740	5	32.40929	1.829263	0.320387	0.262032
5	2.146350	0.475897	1.922871	1	38.89754	2.234781	0.394969	0.222307
6	2.118857	0.283720	1.918210	0.321574	35.20444	2.006473	0.356874	0.220779
7	2.104127	0.179213	1.896175	1	36.42326	2.079522	0.376426	0.218487
8	2.095734	0.186704	1.887008	1	36.64942	2.087262	0.381073	0.219251
9	2.071317	0.156038	1.859759	1	39.05841	2.115577	0.392531	0.210084
10	2.061105	0.132983	1.848242	1	40.33493	2.128621	0.391608	0.216960
11	2.055460	0.118111	1.835147	1	43.12972	2.203133	0.399375	0.214668
12	2.054750	0.103660	1.833985	1	43.24292	2.207648	0.402831	0.212376
13	2.054034	0.103225	1.831781	1	43.71058	2.222534	0.409630	0.212376
14	2.053518	0.090536	1.830741	1	43.92131	2.227765	0.416256	0.211612
15	2.053091	0.095639	1.829022	1	44.60839	2.240692	0.421766	0.207028
16	2.052352	0.114326	1.828766	1	45.11380	2.235861	0.431612	0.207792
17	2.052128	0.104081	1.828078	1	45.26978	2.240504	0.431587	0.207792

- Is this a good model?
  - Examine which options you can change to refine the results.



- In particular, is it better than the gradient boosting model?
- Now let's **Put all the SAS Studio Task Pieces Together**
  - Let's pause and think about everything we did with our no-code options.
  - Going back to the beginning, we:
    - Imported data
    - Checked underlying data to make sure it made sense
    - Created new variables
    - Dropped variables from the analysis
    - Created a data partition
    - Performed descriptive analysis
    - Estimated models
    - Compared models
    - Choose a champion model
  - With coding, you have complete flexibility, but you must code everything!
    - So, we really just got started with our coding work – a complete analysis would take a LOT more time before we would be ready to share results
  - But, let's start now by moving over to SAS Programming
  - One nice element of SAS Studio Tasks is that you can use SAS to write the code for you – and then you can save the code in a single file to be executed at a later date.

- Yes, we've seen this code being generated all along
- Assigning SAS libraries is one of the first thing we do.
  - So, let's use the library assignment program that we created earlier as the starting point for a larger program.
  - Save it to your user-space in the cloud as follows:



- Note: your view will be different
- Add some comments to the top of the program so that you can more easily remember what you were trying to do... months from now:

```

 1 *-----*
 2 |          Titanic Data Analysis
 3 |          Machine Learning in SAS Viya for Learners
 4 *-----*;
 5
 6 **** Load the SAS Libraries + Data ;
 7 libname mycaslib cas caslib=casuser;
 8
 9 proc casutil;
10   load data=TITANIC.TITANIC casout='Titanic_CAS' outcaslib='casuser';
11 run;
  
```

- The next step in the analysis would be to examine the underlying descriptive statistics. Let's grab the code from the SAS Studio Summary Task and add some comments:

The screenshot shows a SAS code editor window with the following code:

```

1  *-----*
2  |          Titanic Data Analysis
3  |          Machine Learning in SAS Viya for Learners
4  *-----*;
5
6 ***** Load the SAS Libraries + Data ;
7 libname mycaslib cas caslib=casuser;
8
9 ⊕ proc casutil;
10    load data=TITANIC.TITANIC casout='Titanic_CAS' outcaslib='casuser';
11 run;
12
13
14 *-----*
15 |          Examine Descriptive Statistics
16 *-----*;
17 ods noproctitle;
18 libname _tmpcas_ cas caslib="CASUSER";
19
20 ⊕ proc cardinality data=MYCASLIB.TITANIC_CAS outcard=_tmpcas_.varSummaryTemp
21     out=_tmpcas_.levelDetailTemp;
22 run;
23
24 ⊕ proc print data=_tmpcas_.varSummaryTemp label;
25   var _varname_ _fmtwidth_ _type_ _rlevel_ _more_ _cardinality_ _nmiss_ _min_
26   _max_ _mean_ _stddev_;
27   title 'Variable Summary';
28 run;
29
30 ⊕ proc print data=_tmpcas_.levelDetailTemp (obs=20) label;
31   title 'Level Details';
32 run;
33
34 ⊕ proc delete data=_tmpcas_.varSummaryTemp _tmpcas_.levelDetailTemp;
35 run;
36
37 libname _tmpcas_|;

```

- Continuing on – lets add the Gradient boosting code now – even though we'd likely spend a LOT of time creating new variables, partitions, and, in general, refining the data before modeling:

The screenshot shows the continuation of the SAS code editor with the following code:

```

33
34 ⊕ proc delete data=_tmpcas_.varSummaryTemp _tmpcas_.levelDetailTemp;
35 run;
36
37 libname _tmpcas_;
38
39
40 *-----*
41 |          Gradient Boosting Model
42 *-----*;
43 ods noproctitle;
44
45 ⊕ proc gradboost data=MYCASLIB.TITANIC_CAS;
46   target Survived / level=nominal;
47   input Age Fare / level=interval;
48   input Parch SibSp Embarked Pclass Sex / level=nominal;
49 run;

```

- Finally, let's also stick in that Neural Network model:

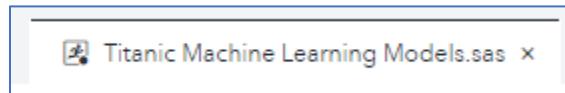
```

50
51
52  *-----*-----*-----*
53 |          Neural Network Modeling |           |
54 *-----*-----*-----*-----*-----*-----*-----*
55 ods noproctitle;
56 libname _tmpcas_ cas caslib="CASUSER";
57
58 ⊕ proc nnet data=MYCASLIB.TITANIC_CAS;
59     target Survived / level=nominal;
60     input Age Fare / level=interval;
61     input Parch SibSp Embarked Pclass Sex / level=nominal;
62     hidden 50;
63     train outmodel=_tmpcas_._Nnet_model_;
64     optimization regL2=0.1;
65 run;
66
67 ⊕ proc delete data=_tmpcas_._Nnet_model_;
68 run;
69
70 libname _tmpcas_;
71

```

- Then, of course, let's save the SAS Program – because we don't want to lose all our hard work!

- A saved file does not have an asterisk before the file name:



- Finally, let's take it from the top and submit the whole program at once.
  - If this is your first time, welcome to your first full-fledged SAS Program!
- Final thoughts on the code-based approach
  - Notice how much longer it takes to produce results with code? And, technically, we didn't come even close to replicating what we created in SAS VA and SAS Model Studio.
  - One big item we omitted was comparing the results of our two ML algorithm. And, to be honest, this was intentional, because it's a LOT of code to process.
  - But for those interested coders, I leave you with one final resource, which is the **Compare Two ML Alogorithm** snippet found here:

The screenshot shows the SAS Snippets interface. On the left, there is a tree view of snippets categorized by type: My Snippets, SAS Snippets, Data, Descriptive, Graph, IML, Macro, SAS Viya Cloud Analytic Services, SAS Viya Machine Learning, SAS Viya Image Processing. The 'SAS Viya Machine Learning' category is highlighted with a red box and contains several sub-snippets. One specific snippet, 'Compare Two ML Algorithms', is selected and highlighted with a red box. The main pane displays the source code for this snippet.

```

1 /*********************************************************************
2  * This example illustrates fitting and comparing two Machine      */
3  * Learning algorithms for predicting the binary target in the   */
4  * HMEQ data set. The steps include:                                */
5  *   */
6  *   (1) PREPARE AND EXPLORE                                     */
7  *       a) Check data is loaded into CAS                         */
8  *   */
9  *   (2) PERFORM SUPERVISED LEARNING                               */
10 *       a) Fit model using Logistic Regression                   */
11 *       b) Fit a model using a Decision Tree                    */
12 *   */
13 *   (3) EVALUATE AND IMPLEMENT                                    */
14 *       a) Score the data                                       */
15 *       b) Assess model performance                            */
16 *       c) Generate ROC and Lift charts                      */
17 *   ****
18 */
19 /* Define a CAS engine libref for CAS in-memory data tables */
20 libname mycaslib cas caslib=casuser;
21 */
22 /*
23 /* Specify the data set names */
24 %let casdata      = mycaslib.hmeq_prep;
25 %let partitioned_data = mycaslib.hmeq_part;
26 */
27 /* Specify the data set inputs and target */
28 %let class_inputs  = reason job;
29 %let interval_inputs = im_clage clno im_debtinc loan mortdue value im_yoj im_ninq derog im_delinq;
30 %let target        = bad;
31 */
32 */
33 */
34 */

```

- You can modify these 182 lines of coding bliss to fit our Titanic data analysis – or simply run the code, as is, and see the output for another data set and research question.
- Regardless, know that SAS has the tools to support you on nearly all your analytical adventures!

Part IV: Learning, continued

- We made it!
- Here are some free SAS resources to continue your learning
  - Free eLearning
    - [https://www.sas.com/en\\_us/learn/academic-programs/resources/free-sas-e-learning.html](https://www.sas.com/en_us/learn/academic-programs/resources/free-sas-e-learning.html)
  - LiveWeb Courses
    - <https://www.sas.com/professor-workshops>
  - Certification series
    - [www.sas.com/certguidedprep](http://www.sas.com/certguidedprep)
- Find these resources here in one location here: [https://www.sas.com/en\\_us/learn/academic-programs.html](https://www.sas.com/en_us/learn/academic-programs.html)
  - SAS Educator Portal
    - [https://www.sas.com/en\\_us/learn/academic-programs/educators.html](https://www.sas.com/en_us/learn/academic-programs/educators.html)
  - SAS Student Skill Builder
    - [https://www.sas.com/en\\_us/learn/academic-programs/students.html](https://www.sas.com/en_us/learn/academic-programs/students.html)
- And check out the free software options here:
  - [https://www.sas.com/en\\_us/software/viya-for-learners.html](https://www.sas.com/en_us/software/viya-for-learners.html)
  - [https://www.sas.com/en\\_us/software/on-demand-for-academics.html](https://www.sas.com/en_us/software/on-demand-for-academics.html)