# Predictive Modeling Fundamentals I

*Lab 3: Build the first predictive model*

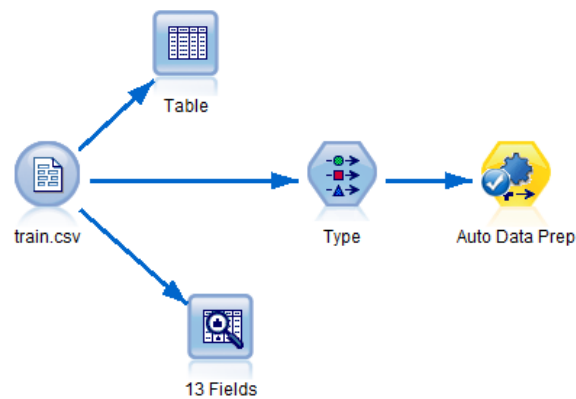# Contents
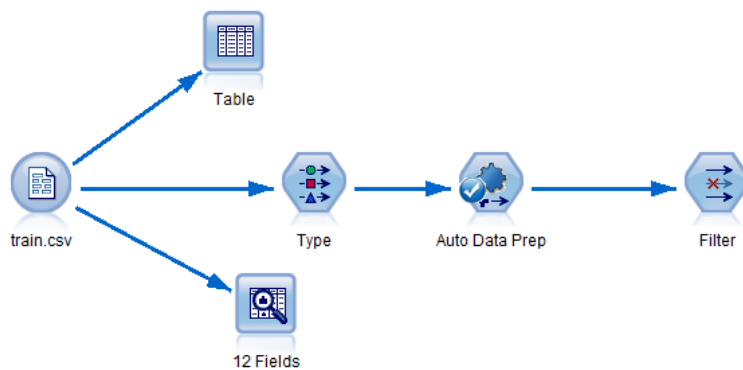
# *Build the first predictive model*

## 1.1 Predicting using Logistic regression

There are many algorithms available in IBM SPSS Modeler. To start simple we will create a Model using Logistic regression, also known as nominal regression. It is a statistical technique for classifying records based on values of input fields. In our use case, the input fields will be the characteristics of each of the passengers of the Titanic and the target will be the field Survived.
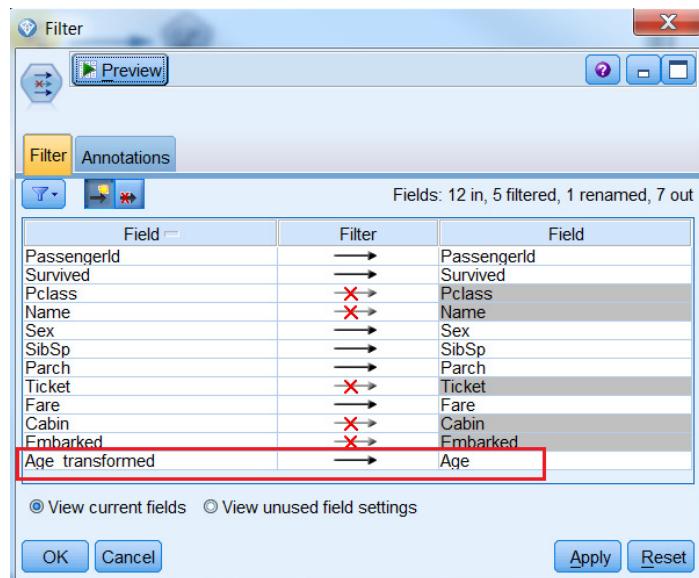
_1. Open the stream you prepared in Lab 2. In this stream you should have something similar to this:
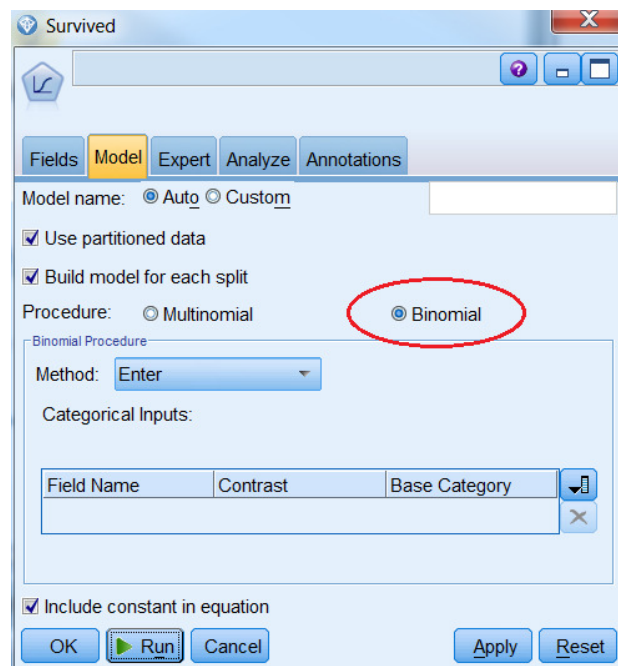
_2. Add a **Filter** node from the **Field Ops** palette and connect it to the **Auto Data Prep** node.
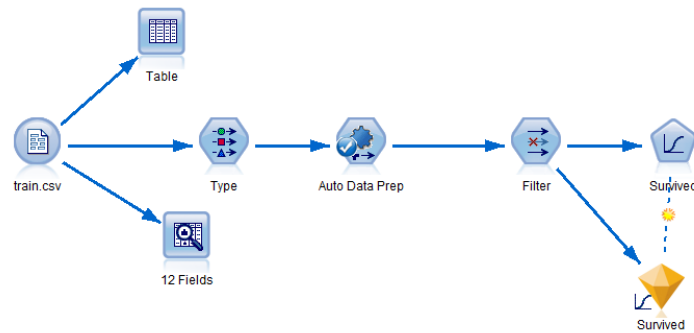
_3. In the **Filter** node, click the filter icon in the center column for the variables that should not be included in the model  and double click the Field name on the right side  for the Age_transformed variable (that was imputed) and rename "Age" to return to original form as in below screenshot.  Click Apply then OK.

_4. Add the **Logistic** node from the **Modeling** palette and connect it to the **Filter** node. Select **Binomial** procedure in the **Model** tab as you are predicting a binary outcome – survive/not survive.



_5. Run the stream. A yellow nugget will appear in your canvas. This is the predictive model.

## 1.2 Validating the model of Logistic Regression

_1. Double click in the yellow Nugget called **Survived**. There is some useful information to understand better the model that you generated:

- **Summary Tab**: Information about what's been used to create the model (inputs, target, configuration, etc.)

- **Advanced Tab**: For those looking for advanced information, here there is all the summary of the model. Scroll down until you see the Classification Table. This is the misclassification or confusion matrix, which is used to calculate and understand accuracy of the model. We can see how many cases the model classifies correctly and incorrectly.

**Classification Table(a)**

| | | | Predicted | | |
|---|---|---|---|---|---|
| | | | Survived | | Percentage Correct |
| | Observed | | 0 | 1 | |
| Step 1 | Survived | 0 | 475 | 74 | 86.5 |
| | | 1 | 112 | 230 | 67.3 |
| | Overall Percentage | | | | 79.1 |
| a. The cut value is .500 | | | | | |

_2. Connect an **Analysis Node** from the **Output** palette. The Analysis Node enables you to evaluate the ability of a model to generate accurate predictions. It performs various comparisons between predicted values and actual values (your target field). **Run** the stream.

In this case, our model is 81.48% accurate. Not bad, but can we do better?

_3. Attach a **Table** node from the **Output** palette to the Model nugget, double click it and hit run.



_4. You will see that fields have been created. $L-Survived shows the predicted outcome and $LP-Survived shows predicted probability for that outcome. For example, passenger in row 9 is predicted to survive (1) with a probability of ~67%.
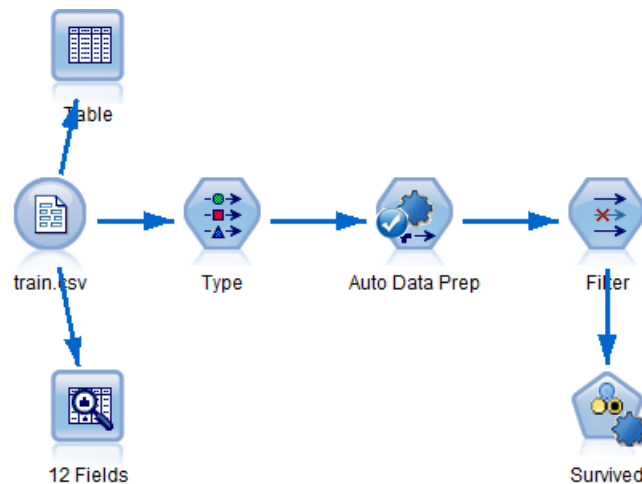
## 1.3 Create a Model using Auto-Modeling tools

In the previous section we created a model using Logistic Regression. In IBM SPSS Modeler there are many algorithms to use and you can see in the **Modeling** Palette.
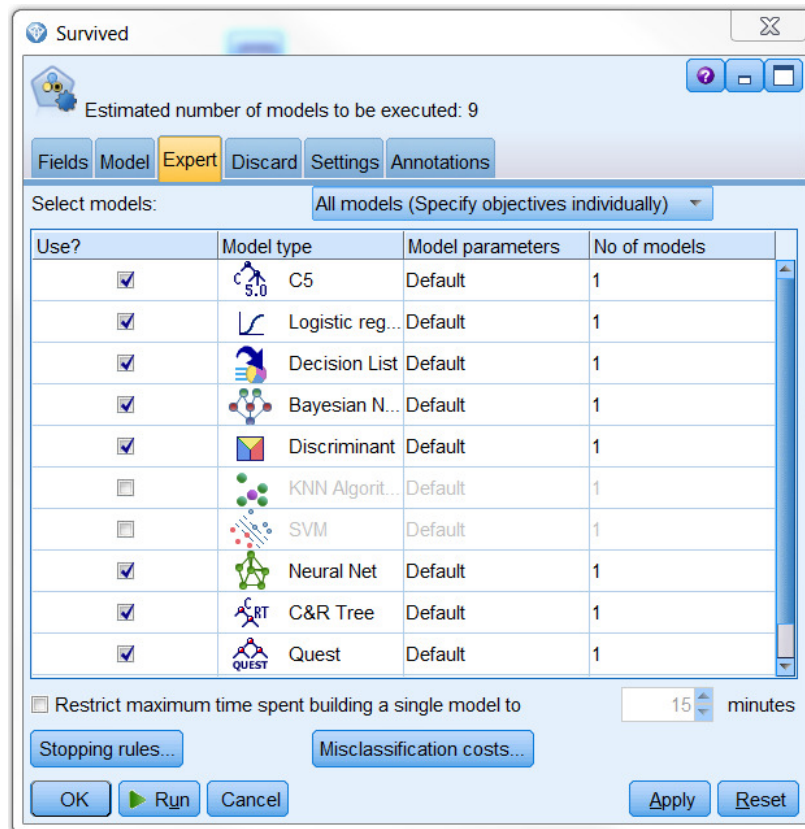


Unless you are an expert, you might feel a bit lost and you might don't know which is the best option. Previously we used Logistic Regression but what if there is a more accurate model? IBM SPSS Modeler provides Auto-Modeling tools that will recommend the best algorithms based on your input data. Let's learn how to use it.
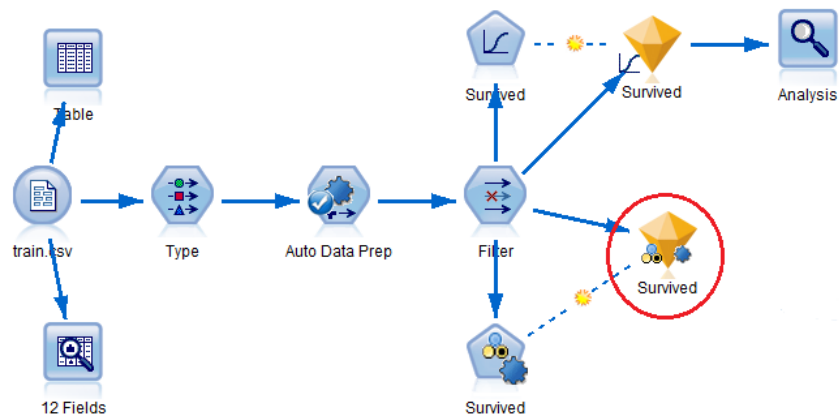
_1.  Add an **Auto Classifier Node** to the canvas and connect it to the **Filter** node.



_2. Double click in the **Auto Classifier Node** (now called Survived).  Here you can change the advanced option of the Auto Classifier. Go to **the Expert Tab** and you will see the list of all the algorithms that will be executed. Take into account that the more you select the longer it will take to be executed. Since our dataset is small everything will move fast.

_3. Run the stream. When an automated modeling node is executed, the node estimates candidate models for every possible combination of options, ranks each candidate based on the measure you specify, and saves the best models in a composite automated model nugget.
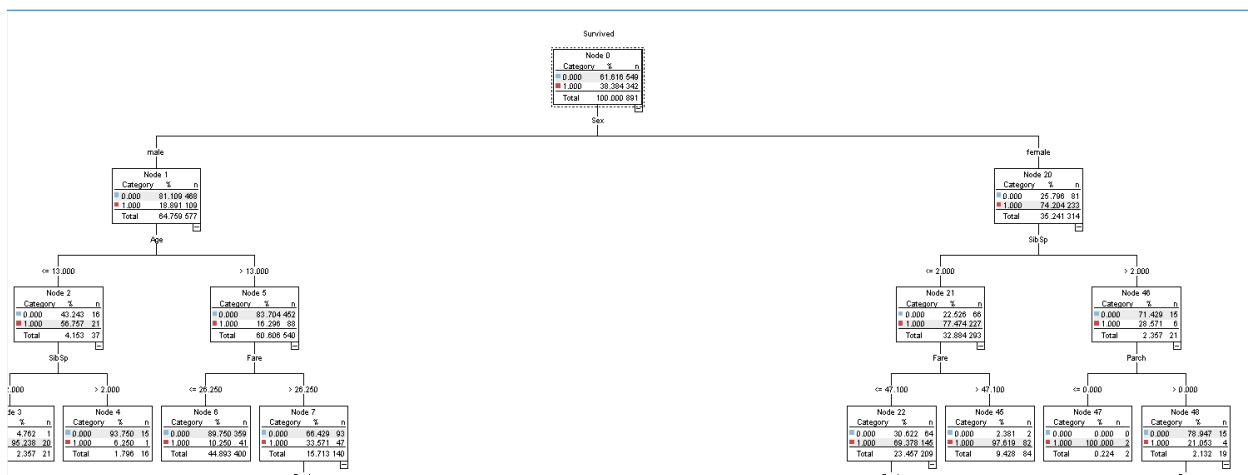


_4. Double click in the new model. The top 3 models are proposed with a list of measures that you can sort: Building Time, Max Profit, Overall Accuracy, etc. In this example the algorithm **C5.1** is the most accurate.

## 1.4 Evaluate the Auto-Model

_1. Double click in the C5.1 Icon:

_2. Here there are all the details of the **C5.1 model** that we just generated.  In the case of the C5.1 algorithms, the results are quite easy to understand since it is a decision tree with a set of conditions. Click on the **Viewer** tab to see the decision tree.



_3. The **Model** tab shows you the same interpretation as the decision tree but in the form of rules. Expand the rules by clicking the ⊞ signs to see them.

Sex = male [ Mode: 0 ]
Sex = female [ Mode: 1 ]
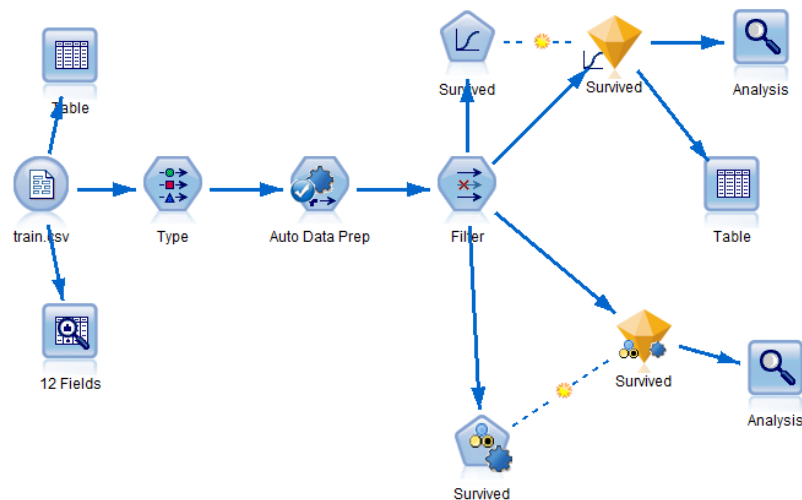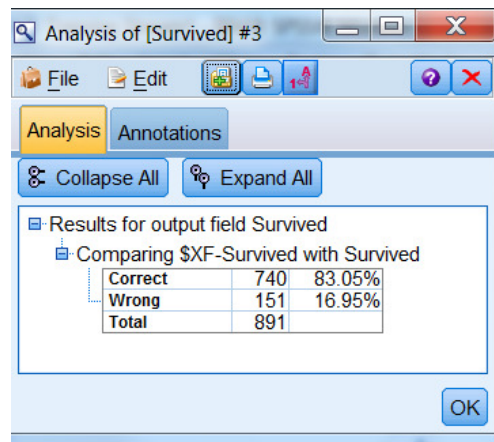  SibSp <= 2 [ Mode: 1 ]
    Fare <= 47.100 [ Mode: 1 ]
    Fare > 47.100 [ Mode: 1 ]  ⇨ **1**
  SibSp > 2 [ Mode: 0 ]
    Parch <= 0 [ Mode: 1 ]  ⇨ **1**
    Parch > 0 [ Mode: 0 ]
      Fare <= 164.867 [ Mode: 0 ]  ⇨ **0**
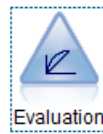      Fare > 164.867 [ Mode: 1 ]  ⇨ **1**

_4. As we previously did, add an **Analysis Node** and connect it to the new yellow nugget. This will give us the accuracy.
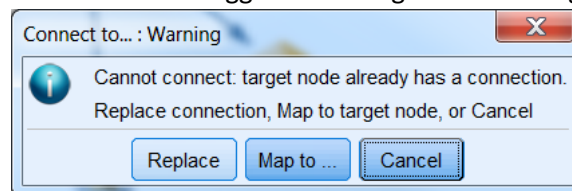


_5. The results of the **Analysis Node** show that our model is 83.05% accurate. Our accuracy just improved!
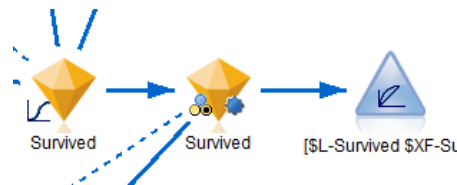
_5. Let us compare the performance of the Decision Tree and Logistic Regression in a visual way. To do so, we can use the ROC (receiver operator characteristic curve) graph. Whichever model has greatest area under this curve, generates more accurate results. Start by adding in the **Evaluation** graph node from **Graphs** palette.
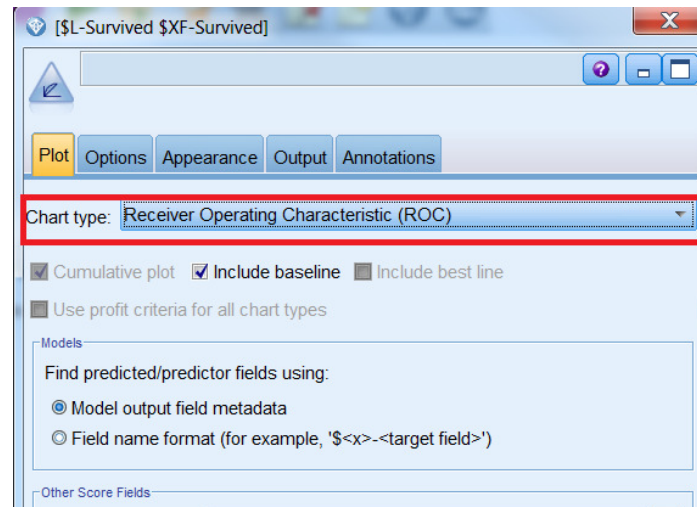


_6. To compare models, we will have to connect the model nuggets to each other. Connect the **Logistic** nugget to the **Auto-Classifier** nugget. You will get this warning. Click replace.
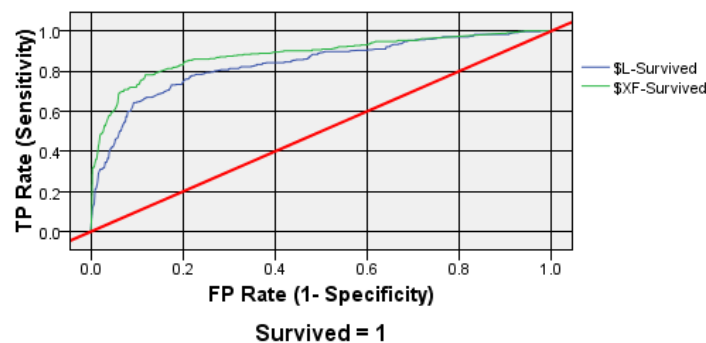


_7. Connect the **Auto-Classifier** nugget to the **Evaluation** Graph node. The caption will change to [$L-Survived….



_8. Double click the Graph node and select Receiver Operator Characteristic (ROC) as Chart type. Click Run.

_9. We see that the Decision Tree (marked by blue line and $XF-Survived) has a greater area under the curve than Logistic Regression (green line and $L-Survived).



## Summary

*Congratulations! You built your first predictive model using IBM SPSS Modeler!*

In this lab you learnt how to create new models using IBM SPSS Modeler. There are many different options that can satisfy the basic need of people without good knowledge on the algorithms but also very advanced options for the advanced users.

We create a model using logistic regression, a very popular algorithm, and we've got an accuracy of 79%. Then we asked IBM SPSS Modeler to propose us which algorithms we should use based on our dataset. The Auto-Modeling tools suggested to use C5.1 instead, and our accuracy improved to 83%. Improving the accuracy of a model can be a very time consuming task!