# Gradient-Boosted Trees with MLlib

Use the MLlib implementation of Gradient-Boosted Trees to generate a predictive model



Product: IBM® SPSS® Modeler

Extension type: Analysis

**Table of Contents**

**Description:**

Gradient-Boosted Trees (GBTs) are a type of ensemble classification algorithm that uses decision trees to build a predictive model.  GBTs can take numeric or categorical input variables and can classify a binary target variable or predict a numeric target with regression.

**Requirements:**

- SPSS Modeler v18.0 or later
- Python 2.7 Anaconda distribution

**Installation:**

Initial one-time set-up for PySpark Extensions

If using v18.0 of SPSS Modeler, navigate to the options.cfg file (Windows default path: C:\Program Files\IBM\SPSS\Modeler\18.0\config).  Open this file in a text editor and paste the following text at the bottom of the document:

eas_pyspark_python_path, "C:/Users/IBM_ADMIN/Anaconda/python.exe"

The underlined path should be replaced with the path to your python.exe from your Anaconda installation.

Extension Hub Installation
1. Go to the Extension menu Modeler and click "Extension Hub"
2. In the search bar, type the name of this extension and press enter
3. Check the box next to "Get extension" and click OK at the bottom of the screen
4. The extension will install and a pop-up will show what palette it was installed

Manual Installation
1. Save the .mpe file to your computer
2. In Modeler, click the Extensions menu, then click Install Local Extension Bundle
3. Navigate to where the .mpe was saved and click open
4. The extension will install and a pop-up will show what palette it was installed
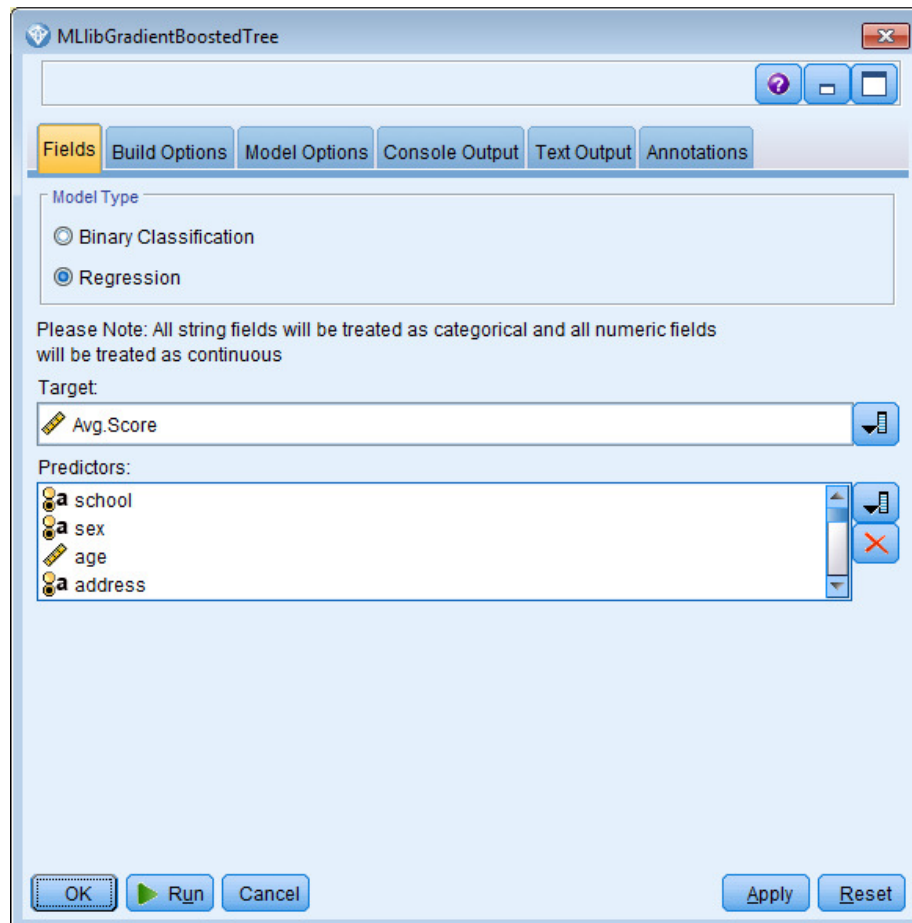
**Python Libraries used:**

- Spark MLlib – Gradient-Boosted Trees [3]

## User Interface

This extension only has two tabs which allow you to control the features used by the model as well as additional parameters that can be tuned to optimize performance.

Tab 1 - Fields

- Model Type
    - Binary Classification – Classify each instance into a binary class (e.g. 0 or 1)
    - Regression – Predict a numeric target variable using GBTs to perform regression
- Target  - This is the flag variable (for Binary Classification) or numeric variable (for Regression) that the model will classify or predict respectively.
- Predictors – These are the input features or independent variables used to generate the model. Since decision trees are created for this model, input features can be a mix of numeric and categorical variables.  Number variables also do not need to be normalized for this algorithm.



Tab 2 – Build Options

For more details on the parameters of this function visit the MLlib documentation [3]

- Loss Function – Read more here about the Loss Functions used.
- Number of Iterations (Trees) – The number of trees to be used for the ensemble
- Max Tree Depth – Increase depth of trees to improve model on training data however this can lead to over fitting the training data and also increases time required to train.
- Learning Rate – Decreasing the learning rate will increase training time and will also help improve stability.
- Max Bins -  Number of bins used for finding splits at each node

**Example:**

For this example, we will attempt predicting the average grades for students based on a number of variables.  The dataset used for this example can be found at the UCI MachineLearning Repository [1]. This dataset can be found at http://archive.ics.uci.edu/ml/datasets/Student+Performance.  The data was originally used by Paulo Cortez [2].

This is a toy example to demonstrate how this extension can be used locally to use Gradient Boosted Trees for smaller datasets, as well as in a Spark environment using Analytic Server.

1.  Download data and open in Modeler using a Var File node.  This file is ';' (semi-colon) separated, so you will need to mark other for type of delimiter and type in a semi-colon.
2.  This dataset provides 3 different grade response variables.  This is more than we are interested in for now, so rather than just predicting one, let's average the scores as the target we want to predict.
    a.  To do this, add a derive node following the var.file node.
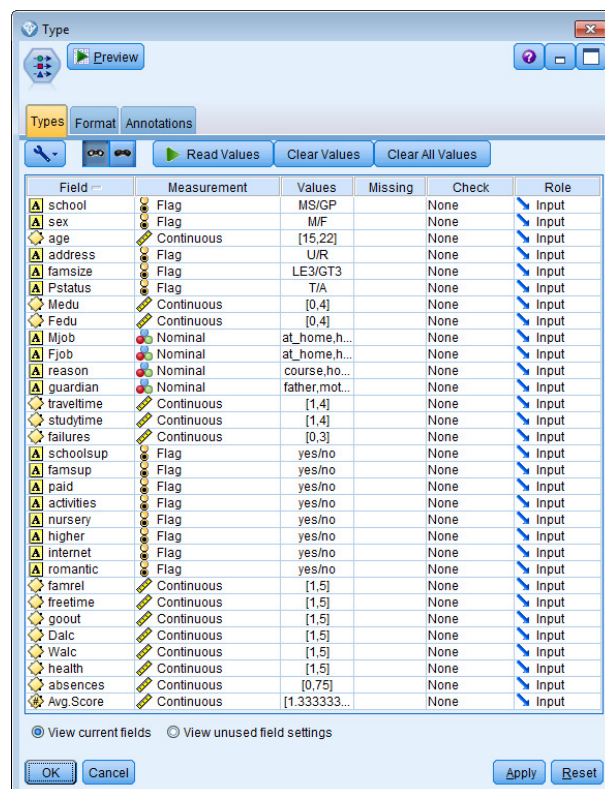    b.  Add the following formula to average the scores.

3.  Now , let's filter out the 3 columns for the original grades provided:
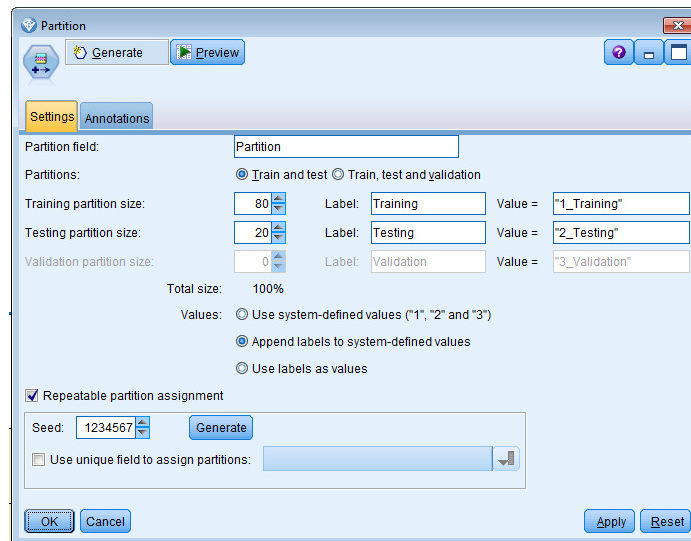


4.  Next we can add a Type node, to read values to make sure all the meta-data is correct.
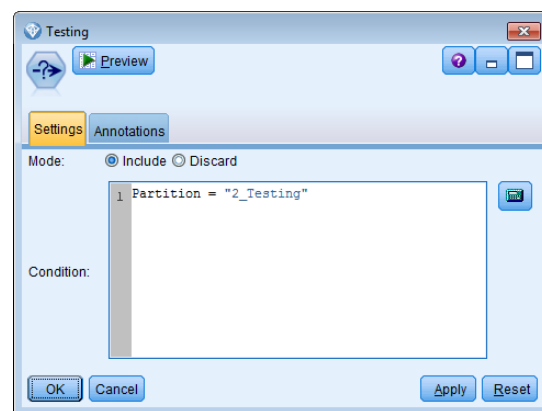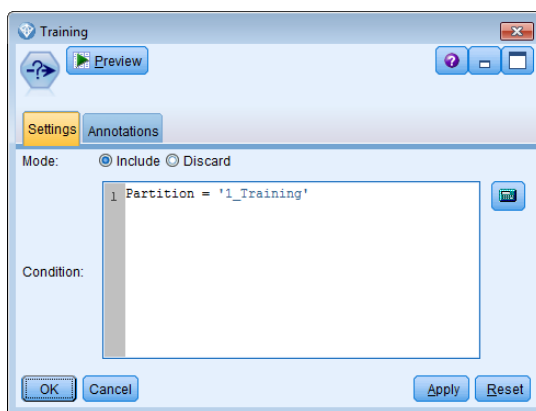
5. With our data prepared, let's partition and split into training and testing
   a. Add a Partition node from the Field Ops Palette
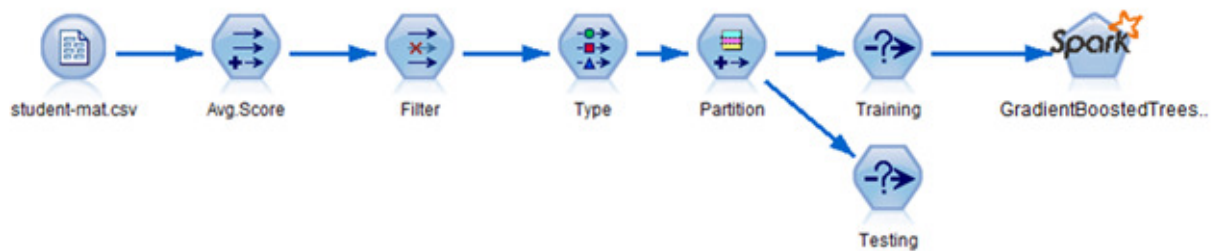   b. I chose to do 80/20 train/test split for this data.



   c. This node will create a new column that gives a field to select training and testing from (the values in this column are "1_Training" and "2_Testing")
   d. To split the data, add two Select nodes from the Record Ops palette where the formula should match the Nodes below
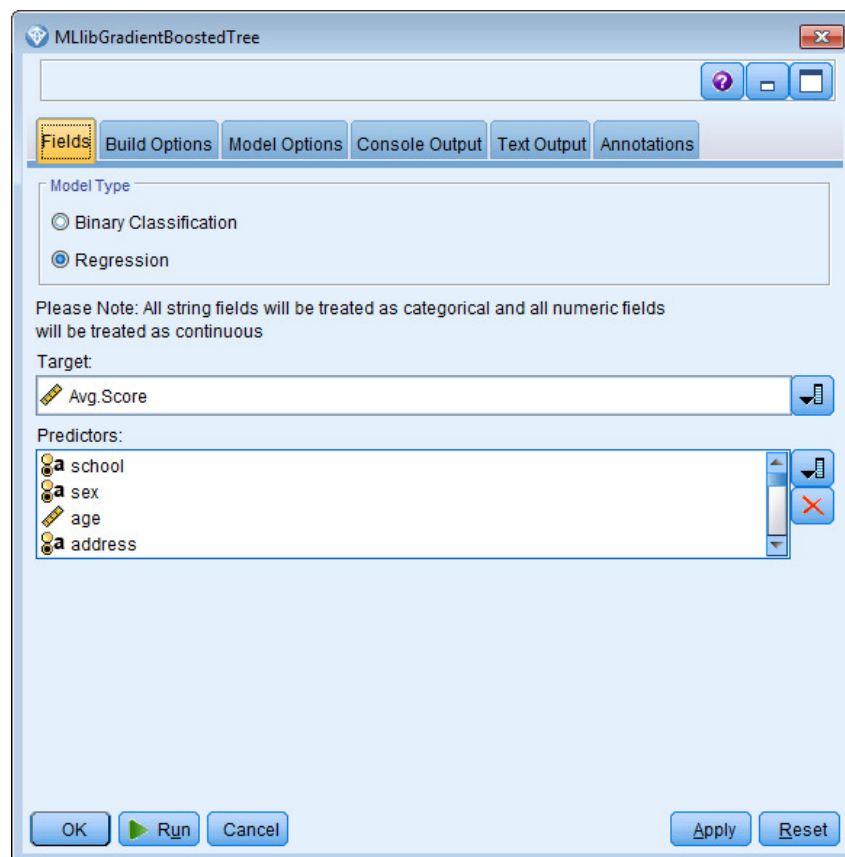
6.  Now, let's add our Gradient-Boosted Tree node so our stream will look like this
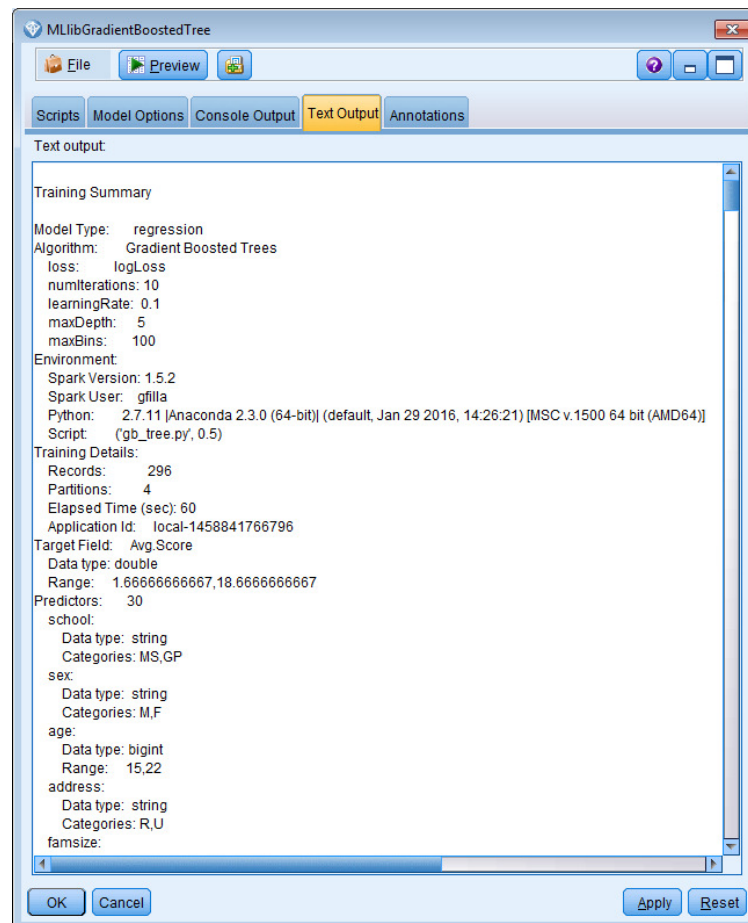


7.  Open the GBT dialog
    a.  Make sure Model Type is Regression since we are predicting a numeric variable
    b.  Add the Avg. Score as the Target variable
    c.  For predictors you can include all variables except Avg.Score and the Partition field
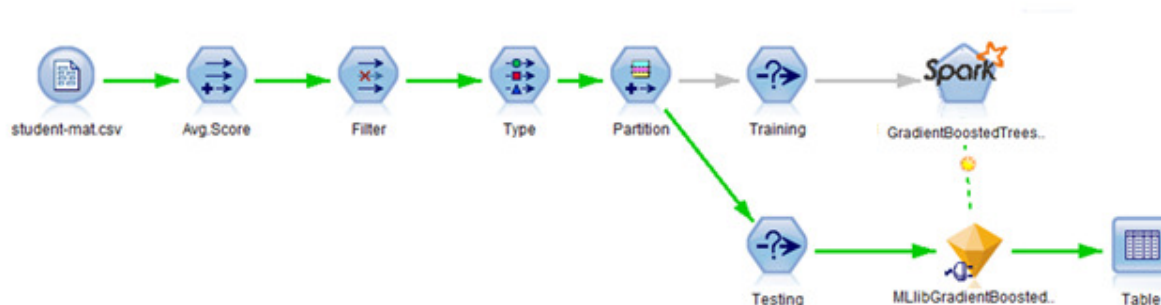


    d.  If you would like to fine tune the model, go to Build Options and tweak the parameters.
8.  Run this stream to create the model.

9.  By double clicking the Model Nugget and clicking on the Text Out tab, you can see a summary of the model training.  This may come in handy when looking into details on the model.



10. Now connect the "Testing" Select node to the model nugget create and add a table as the output
11. Right click the table node and run to get the following stream:

12. The table produced will include the predicted score for each student:



13. That is good, but we want to see how the model performed.  Add an Analysis node from the Output palette and use the setting below, then run the stream.



14. Now you have the Mean Error, MAE, and other metrics that will help evaluate the performance of the model.

15. Finally – our full stream should look like this:



## Important Links

### Learn

- Learn more about SPSS software.
- To learn more about this implementation of GBTs take a look at the Spark documentation
- Read about coding a PySpark Extension for SPSS Modeler
- Visit developerWorks Business analytics for more technical analytics resources for developers.

## Discuss

- Visit the IBM SPSS Community to share tips and experiences with other IBM SPSS developers.
- Follow developerWorks on Twitter to be among the first to hear about new resources.

## References

[1] Lichman, M. (2013). UCI Machine Learning Repository [http://archive.ics.uci.edu/ml]. Irvine, CA: University of California, School of Information and Computer Science.

[2] P. Cortez and A. Silva. Using Data Mining to Predict Secondary School Student Performance. In A. Brito and J. Teixeira Eds., Proceedings of 5th FUture BUsiness TEChnology Conference (FUBUTEC 2008) pp. 5-12, Porto, Portugal, April, 2008, EUROSIS, ISBN 978-9077381-39-7.

[3] Apache Spark version 1.6.1 – Mllib Guide: Gradient-Boosted Trees (http://spark.apache.org/docs/latest/mllib-ensembles.html#gradient-boosted-trees-gbts)