# TUTORIAL LL

## Models for Predicting Rapid Intensification of Tropical Storms

**Juan S. Lebrón and José F. Nieves**

*Department of Physics*

*University of Puerto Rico San Juan, PR 00931*

Hurricanes are among the most complex macroscopic systems that are governed, in principle, by the Laws of Physics. One of the purposes of our project is to develop models for making statistical predictions about the intensification of hurricanes without having to deal with the individual physical parameters that characterize such systems. This tutorial gives an example of the kind of *Data Mining* analysis that can be carried out based on the ideas that we have been developing.

## I. INTRODUCTION AND MOTIVATION

The main focus of a hurricane forecast is primarily to predict the trajectory and the intensity. While the method for predicting the trajectories have advanced in recent years together with the advances in computational resources, that has not been the case with respect to the intensity. Recent events of *unexpected* rapid intensification (e.g., Hurricane Wilma in 2005) have exacerbated the situation, and have contributed to the question of what could happen if a tropical cyclone goes through a stage of rapid intensification, unexpectedly, just before touching ground. Fortunately, the recent hurricanes that have suffered rapid intensification just before touching ground (e.g., Humberto and Lorenzo in 2007) have been relatively weak and therefore did not cause much damage. Nevertheless, it is clearly just a question of time before such an event occurs involving a stronger system, with potentially catastrophic consequences.

Most of the information about the position and intensity of hurricanes that appear in public information bulletins are estimated on the basis of satellite images of various kinds (e.g., visible and infrared images), such as that shown in Fig. 1. The technique for doing this was proposed by Vern Dvorak in 1975, and has being refined over the years. It is outside the scope of this tutorial to review the details of the technique and the methods involved in estimating such parameters from the images. Suffice it to say that, while Dvorak pretended that his technique be used for estimating not just the value of the intensity at a particular instance but also forecast it in the future, the fact is that his technique is indeed used for the former but no for the latter.
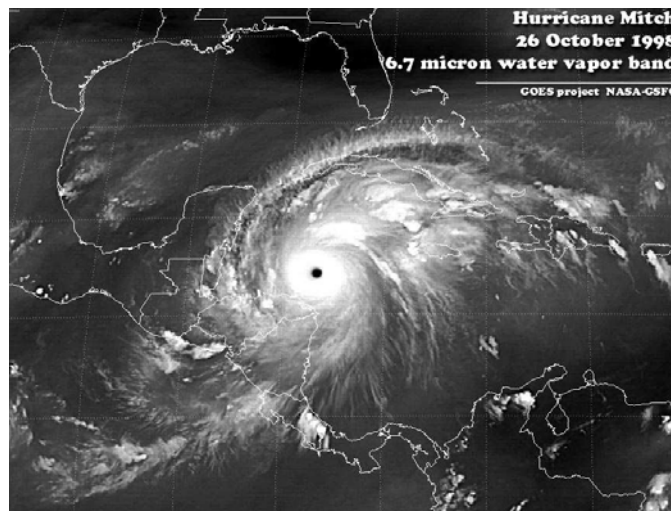


FIGURE 1: Satellite image (water vapor band) of Hurricane Mitch (16 October 1998)

The main purpose of our project is to produce a technique that is similar in spirit to the Dvorak technique for estimating the instantaneous intensity, but which is instead able to forecast the intensity in the future. In particular, the idea is to predict whether a hurricane will undergo rapid intensification during the next 24 hours. Although this can

be expressed as a yes/no result, we also hope to generalize the model to give the expected increase in intensity for the next 24h period, regardless of whether it is rapid or not.

## II. THE METHOD

The method is based on analyzing a satellite image, and determines from it the values of a set of relevant parameters. The salient feature of the technique is that, instead of attempting to determine the precise value of some physical variables (e.g., wind speed, cloud sizes or heights) from the image, the focus is on determining whether such variables may fall within one range of values or another. Thus, for example, we do not bother in trying to estimate the size or height of the clouds, but instead in determining whether or not there are Cirrus clouds at all and, in addition, whether they persist for a given distance from the center, or whether they cover three quadrants, or two.

There are 8 conditions in the model: the one defined as "inner core" and the other ones labeled as 1a, 1b, 2a, 2b, 3, 4a, and 4b. "Inner core" refers to a particular type of organization of the clouds in a hurricane which allows for faster intensification than otherwise. (This is usually equivalent to the hurricane having an "eye" or not). The other 7 conditions characterize the environmental conditions, in a categorical fashion as we have indicated above. It is outside the scope of this tutorial to explain how those values are determined and what they mean in detail. Here we accept that data as given, and the idea is to fit these 7 conditions to the actual observed changes in intensity in order to derive a useful predictive formula. In the dataset, intensity is expressed as a $T$ value which in theory can range from 1 (weak) to 8 (strong hurricane). The model output is

expected to be a $\Delta T$ such that today's $T$ plus the model-obtained $\Delta T$ produces (predicts) tomorrows $T$.

| Inner core | 1a | 1b | 2a | 2b | 3 | 4a | 4b | $\Delta T$ |
|------------|----|----|----|----|---|----|----|------------|
| 0 | 1 | 0 | 0 | 0 | 0 | 1 | 1 | 0 |
| 0 | 1 | 1 | 1 | 0 | 0 | 1 | 0 | 0 |
| 0 | 1 | 1 | 1 | 0 | 0 | 1 | 1 | 0 |
| 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 1 | 1 | 1 | 0 | 0 | 0 | 1 | 1 | 1 |
| 1 | 1 | 1 | 1 | 0 | 0 | 1 | 1 | 1 |
| 0 | 1 | 1 | 0 | 0 | 0 | 1 | 0 | 0 |

TABLE 1: A sample of the data set used in the tutorial.

In the language of statistics, it is clear that we have done is to attempt to characterize such a complex system by a set of categorical variables that presumably carry the most relevant information about physical properties of the system, or at least those properties that are most influential during the periods of rapid intensification.

In this way, we thus end up with a data table that looks like the one shown in Table 1. In short, the problem at hand is to see if the value of $\Delta T$ shown in the last column of that table can be predicted on the basis of the values of the eight parameters labeled "inner core" through 4b.

# III. MODEL BUILDING

For this tutorial we are using the *Data Mining* models of Statistica. The model building process starts by choosing the item labeled *Data Miner Recipes* from the *Data Mining* menu as shown in Fig. 2. In the next dialog appears asking to open an existing recipe of or a new one, we choose New, which brings us to the window shown in Fig. 3, and there we will click on *Open/Connect data file* to open our data file.

NOTE:  Since the DMRecipe is used in many of the other tutorials in this HANDBOOK, we will run through the DMRecipe format rather rapidly in the following example:



FIGURE 2.  Data Mining Menu from STATISTICA Data Miner.

FIGURE 3.  DATA MINING RECIPE initial dialog.

Our original data file was in prepared in a spreadsheet-like comma-separated format. Statistica recognizes the format and asks for confirmation, and after loading the file our workspace looks as shown in Fig. 4, where we will click on *Select variables*. This brings up the window shown in Fig. 5, where we have already selected columns 1-8 as the input categorical parameters, and column 9 as the target categorical variable.



FIGURE 4 . DMRecipe dialog after loading data



FIGURE 5.  Variables selected.

6

Our workspace then looks as shown in Fig. 6 where, instead of clicking on the item labeled *Next Step*, we will click on the arrow that is to it right side to *Run to completion*. That will instruct the Statistica software to go automatically through a series of steps that would have to be carried manually in order to run anyone of the supported models.



FIGURE 6.  Click the RUN TO COMPLETION selection on the upper-right corner to run the complete DMRecipe set of computations.

The end result of that process is shown in Fig. 7, which indicates that the best model is the Boosted trees model, and the software has automatically selected it as the model for deployment. The various menu items that are shown on the left hand side in that window, under *Evaluation*, can be chose to display various tables and graphs related to the evaluation of the models as shown, for example, in Figs. 8 and 9.

FIGURE 7.  Summary Results spreadsheet result of the DMRecipe computations.

Figure 7 results show that the Boosted Trees have the lowest Error Rate at 17.4%, meaning that the ACCURACY for predicting Hurricanes is  100-Error rate, or  82.6% using the BT modeling.  If we decide to accept this BT Model, we can deploy it on new data coming in during the hurricane season in future years to get predictions at this 82.6% accuracy.



FIGURE 8.  Crosstabulation Table for BOOSTED TREES spreadsheet result.

Lift Chart - Lift value
Cumulative
Selected category of deltaT: 1

FIGURE 9. Lift Chart Graph comparing the 3 data mining algorithms used in this modeling process. All are quite similar, however the Boosted Trees did have a slight higher value at the left side of the graph, which is the important place to look with Life Charts.

9

## IV. DEPLOYMENT

Since the Boosted Trees gave the highest accuracy, we will use it for deployment. We find this easier to do using the Interactive Boosted Trees (IBT). Thus we will nor re-run the analysis just using the IBT module, as follows.

NOTE:   Since the INTERACTIVE modules of data mining algorithms have not been extensively used in other tutorials in this HANDBOOK, we will run through this Interactive BOOSTED TREES format thoroughly, step-by-step,   in the following example:
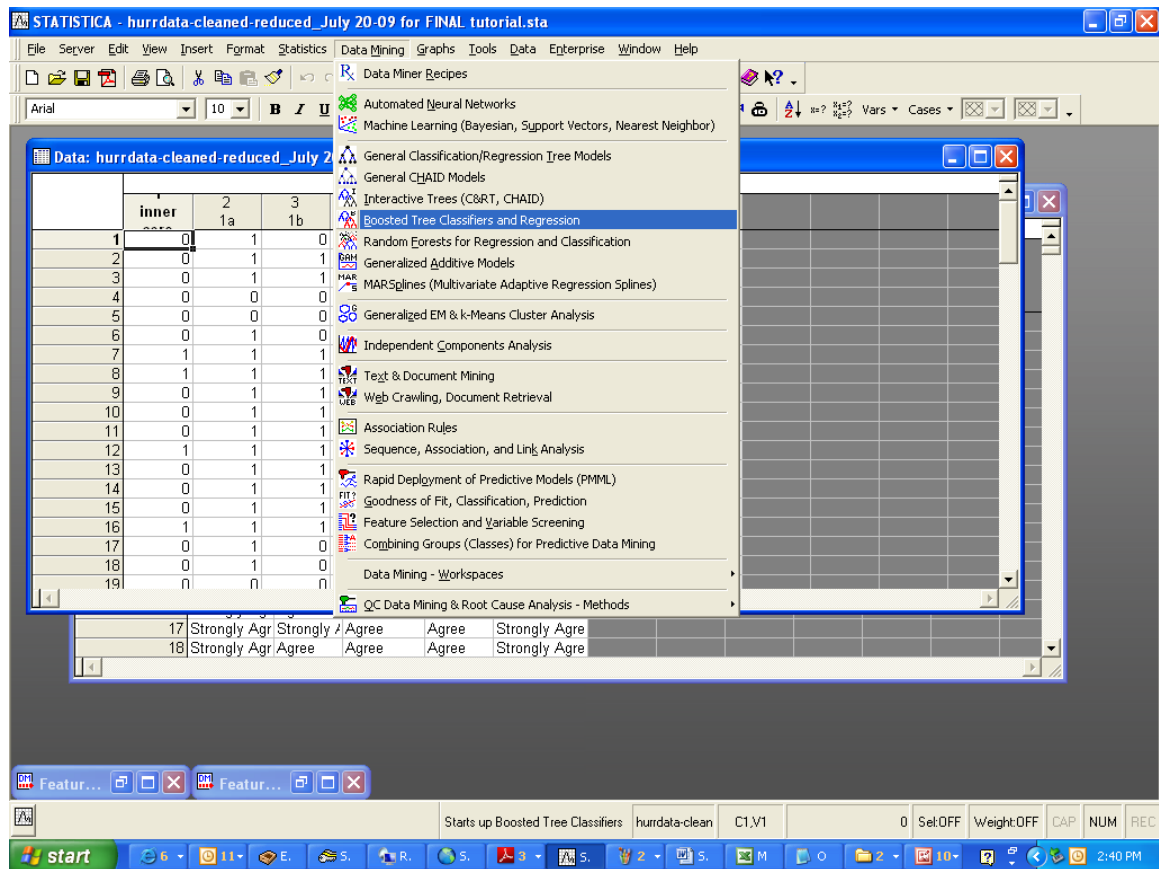


FIGURE 10.  Selecting Interactive BOOSTED TREES algorithm from the menu.

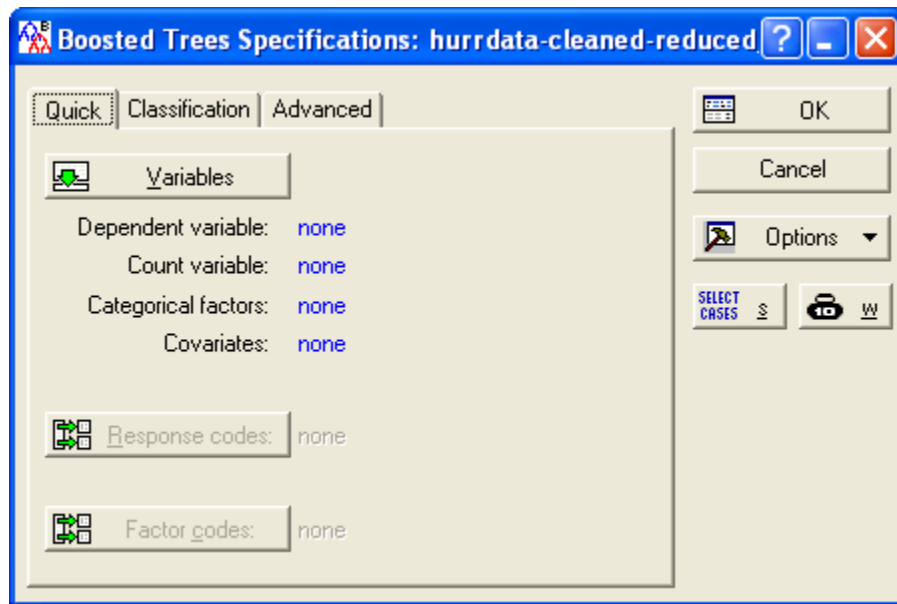FIGURE 11.  Selecting CLASSIFICATION ANALYSIS as our "type of analysis" from the BOOSTED TREES initial dialog.


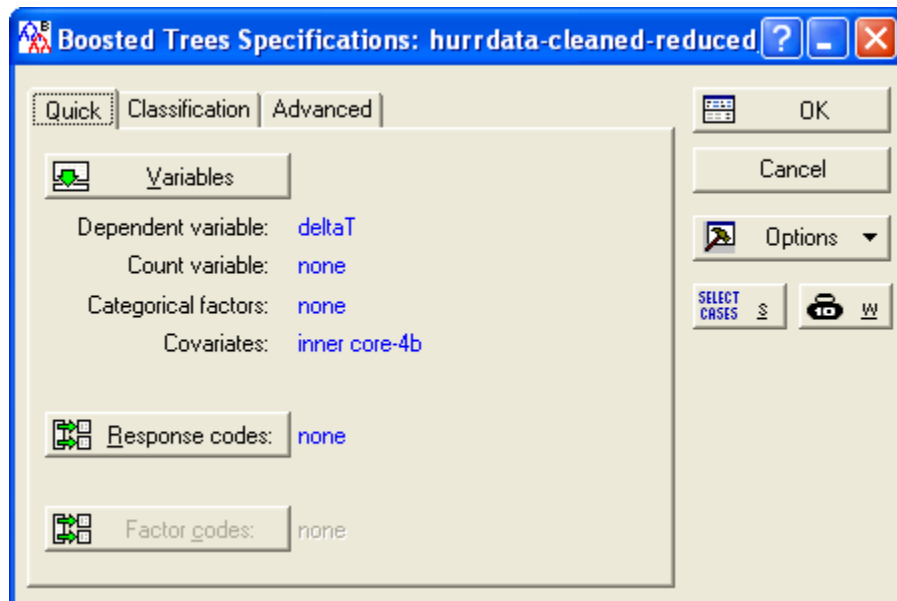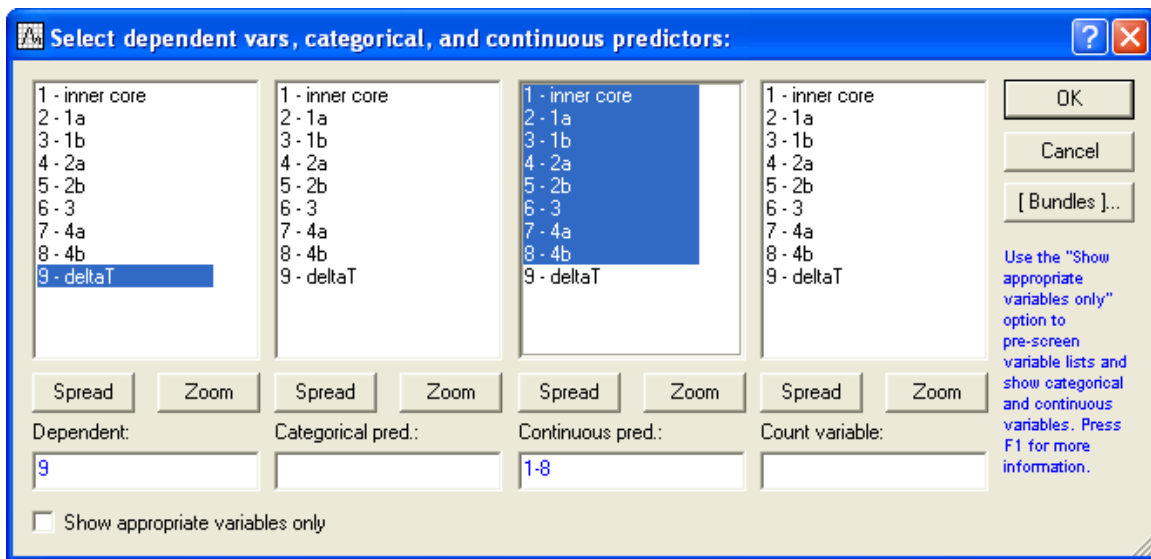
FIGURE 12.  The INTERACTIVE BOOSTED TREES initial dialog.

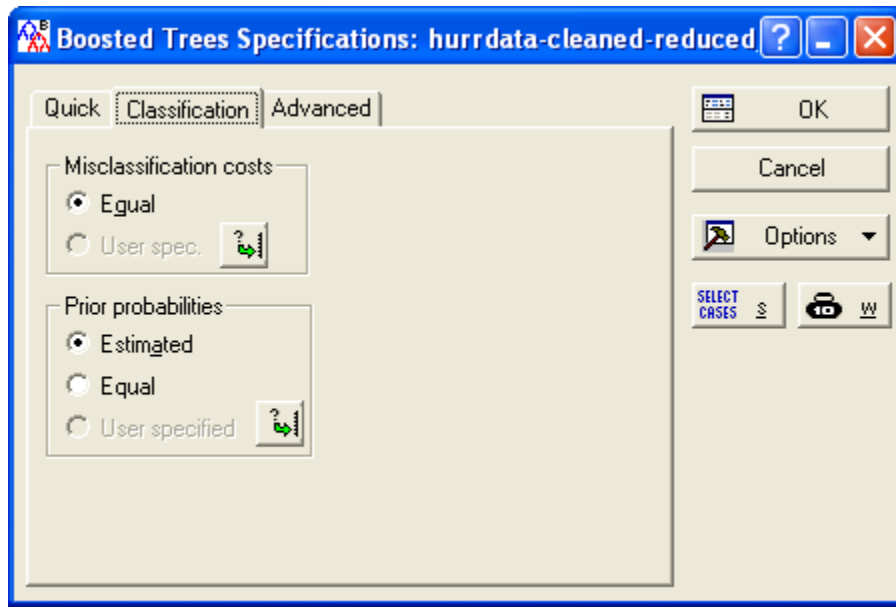FIGURE 13. Variables selected are "showing as entered" on the BT DIALOG.

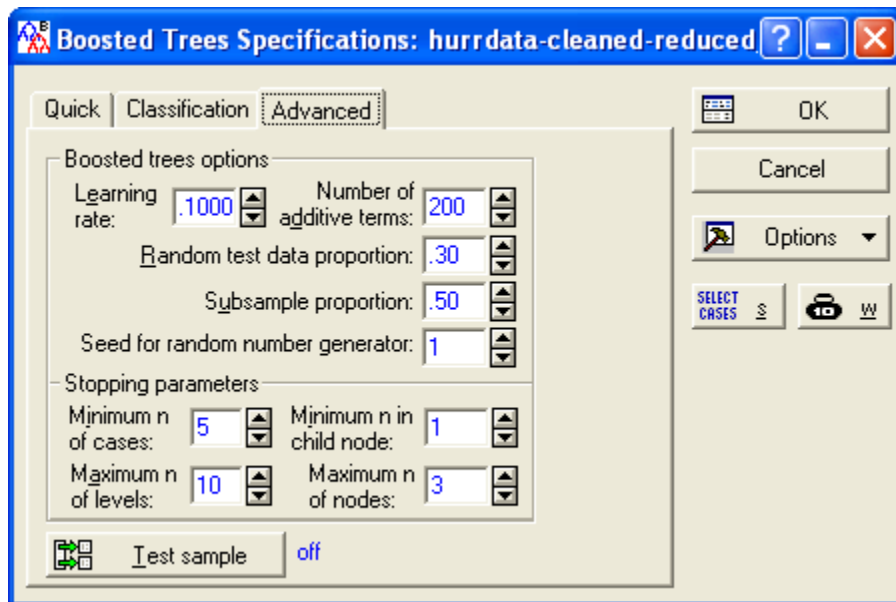FIGURE 14.  Classification Tab is left at it's default values.



FIGURE 15. The Advanced Tab is also left at it's default values for this analysis.

Hit **OK** to run, which will produce a series of windows as shown in the figures that
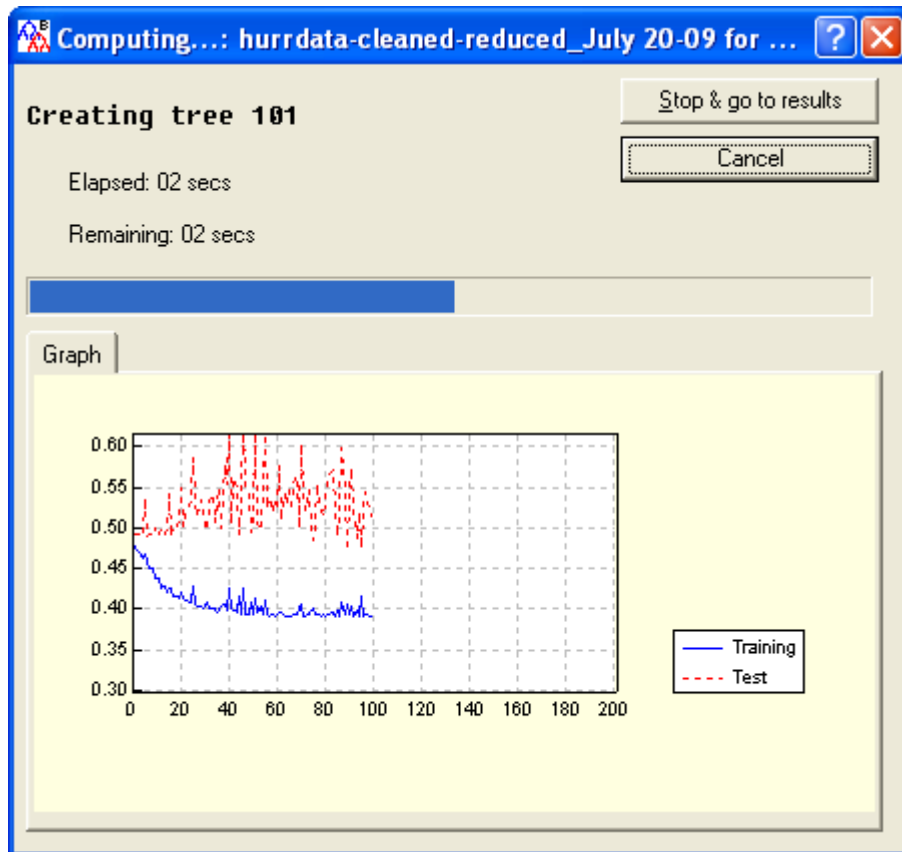
follow.

FIGURE 16. Active Dialog that is visible for a few seconds [or minutes, depending on data size] of the BOOSTED TREES growing process where both the TRAINING data [blue color] and TEST data [red color] "models" are being graphed as per number of trees being grown; when it hits a point where things "level off" for both this will be selected as the "End Point" for an optimal Boosted Trees model.
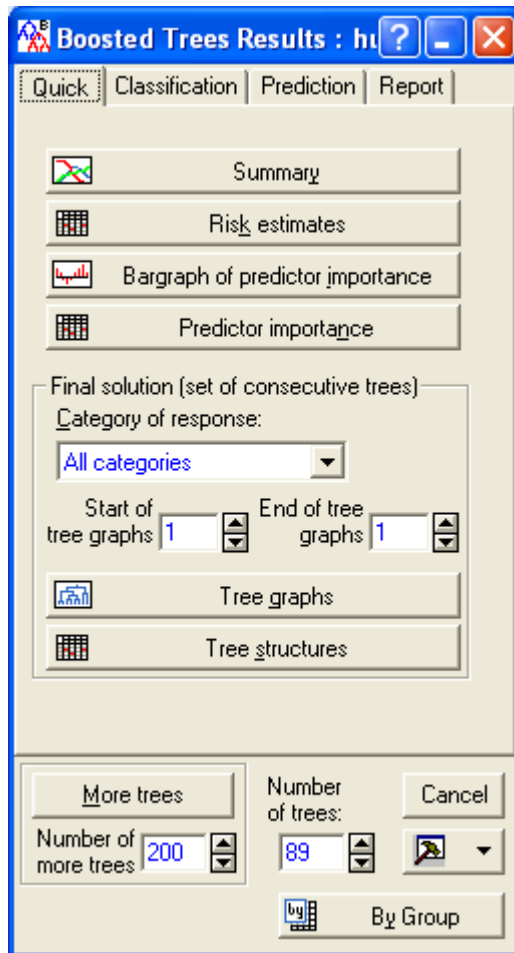
FIGURE 17.  When the Boosted Trees computations concludes, the above RESULTS DIALOG will pop on the computer screen.
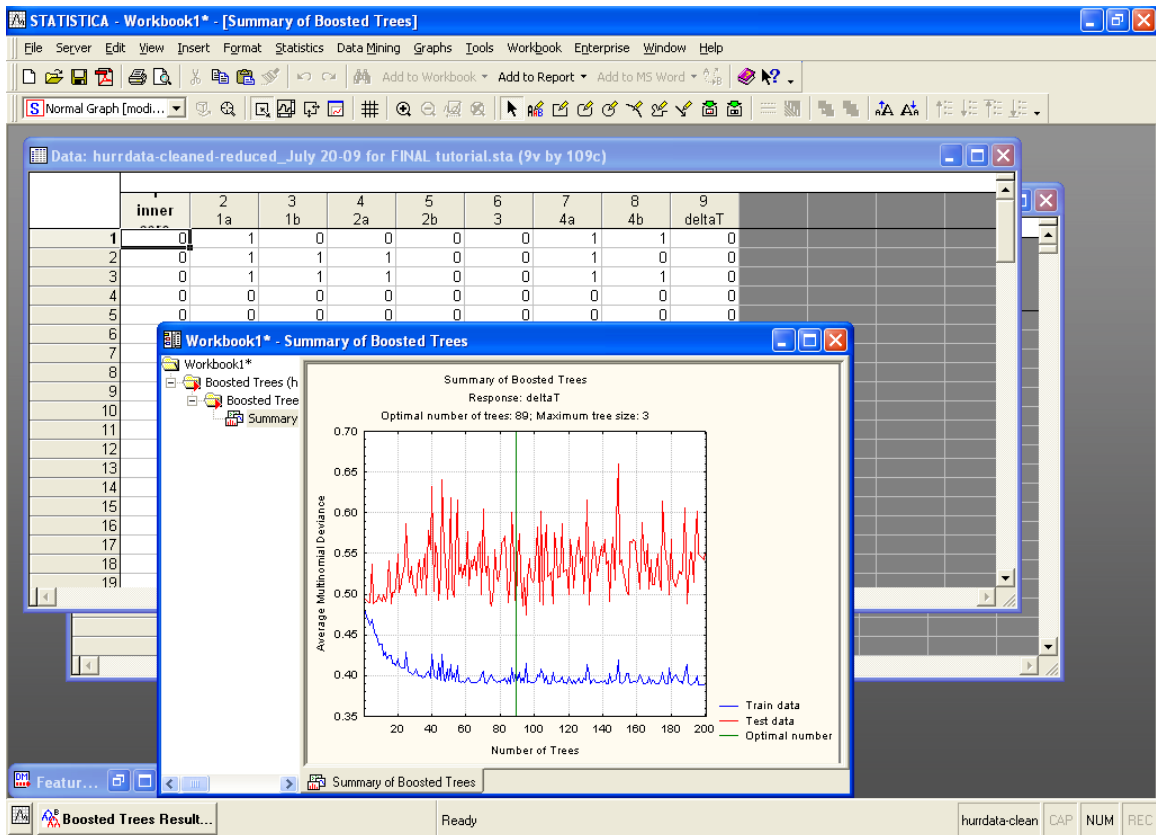
FIGURE 18. Clicking on the SUMMARY button on the RESULTS DIALOG will give a workbook with the result in the left panel of the RESULTS DIALOG, as seen above, will show the final Boosted Trees modeling with the Green Vertical Line being the point of optimum growth for this dataset.
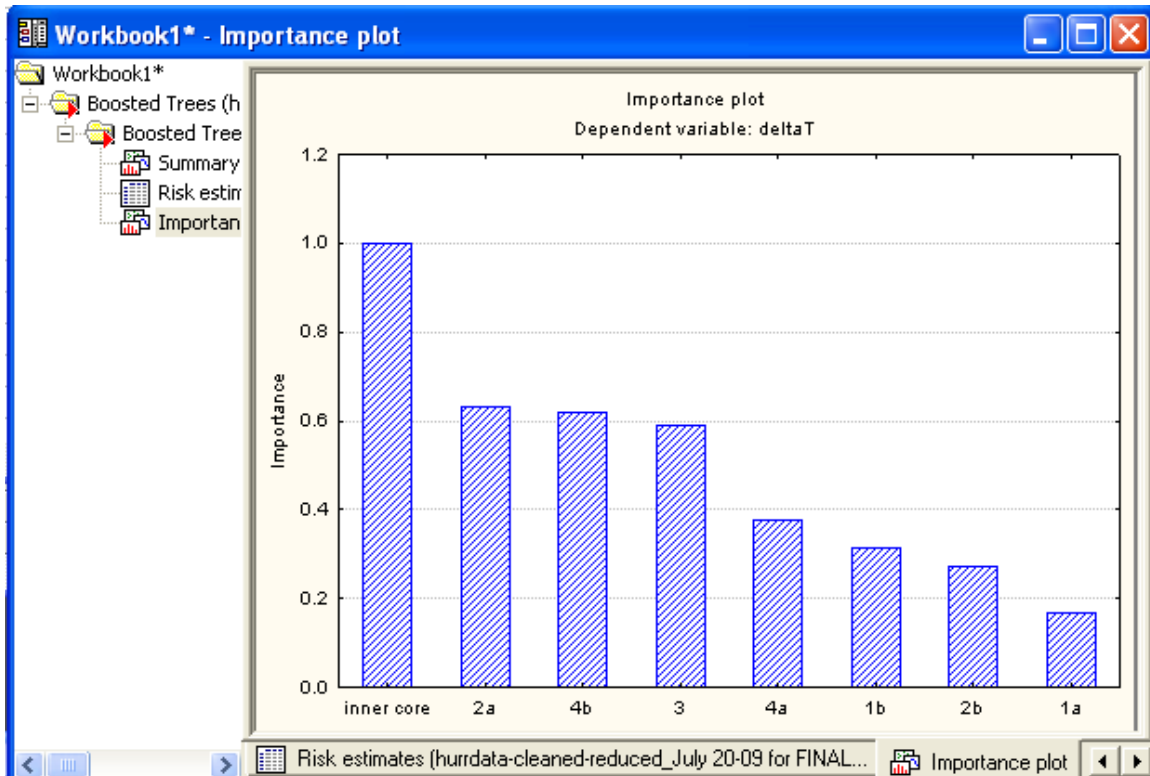
FIGURE 19.  Clicking on the PREDICTOR IMPORTANCE  button on the RESULTS DIALOG will give the Importance Plot as seen above;  this indicates that the "Inner Core" is the most important variable in predicting a hurricane, with the '2a', '4b', and '3' variables being the next in importance, almost equal as they have about the same heights; then the remaining variables trail off in importance.  From this we might try using just the top 4 predictor variables in further modeling, as this might give us the highest Accuracy Score in  Hurricane Prediction.

| | Predictor importance (hurrdata-cleaned-reduced_July 20-09 for FINAL tutorial.s Response: deltaT | |
|---|---|---|
| | Variable Rank | Importance |
| inner core | 100 | 1.000000 |
| 4b | 63 | 0.632256 |
| 3 | 62 | 0.617178 |
| 1a | 59 | 0.587712 |
| 4a | 38 | 0.378208 |
| 2a | 31 | 0.313140 |
| 2b | 27 | 0.271388 |
| 1b | 17 | 0.166068 |

FIGURE 20.  Another way of looking at the Predictor Importance is in the table where the predictor variables are "ranked" on basis of 100 being highest.
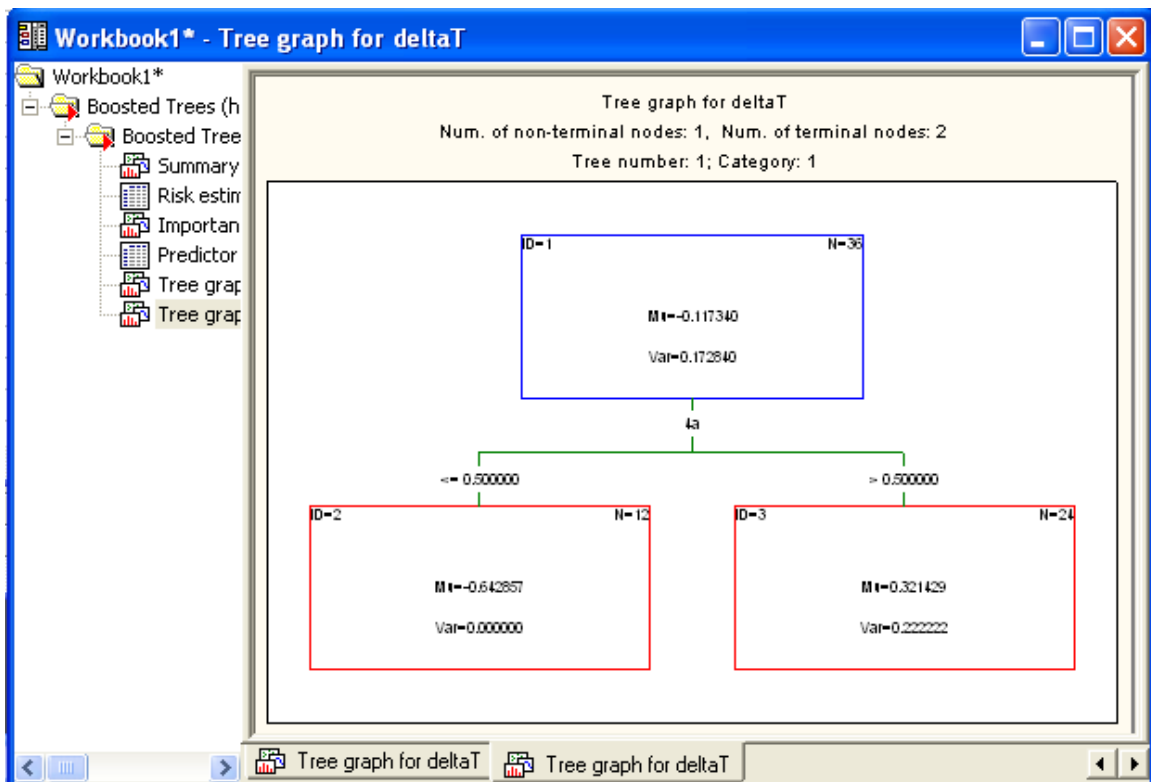
FIGURE 21.  Clicking on the TREE GRAPH  on the RESULTS DIALOG gives the pictorial representation of the "boosted trees" model.

**Now, if we want to deploy this model to new data that has been subsequently collected, we need to follow the steps below:**
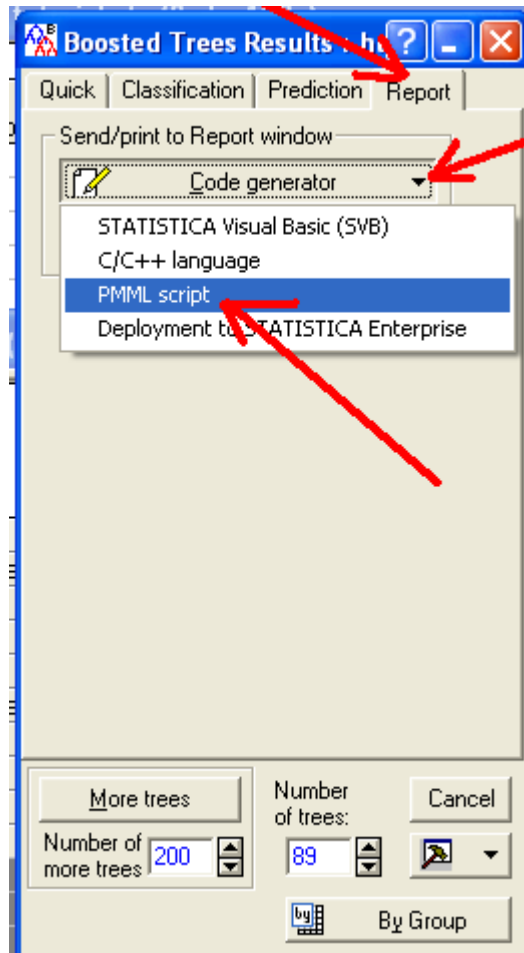
FIGURE 22. Using the REPORT TAB on the Boosted Trees RESULTS DIALOG, select PMML script for the method of creating the model for use in future deployment to new data.

Using the **Report** tab, click on the **Code Generator** button above, then, release on

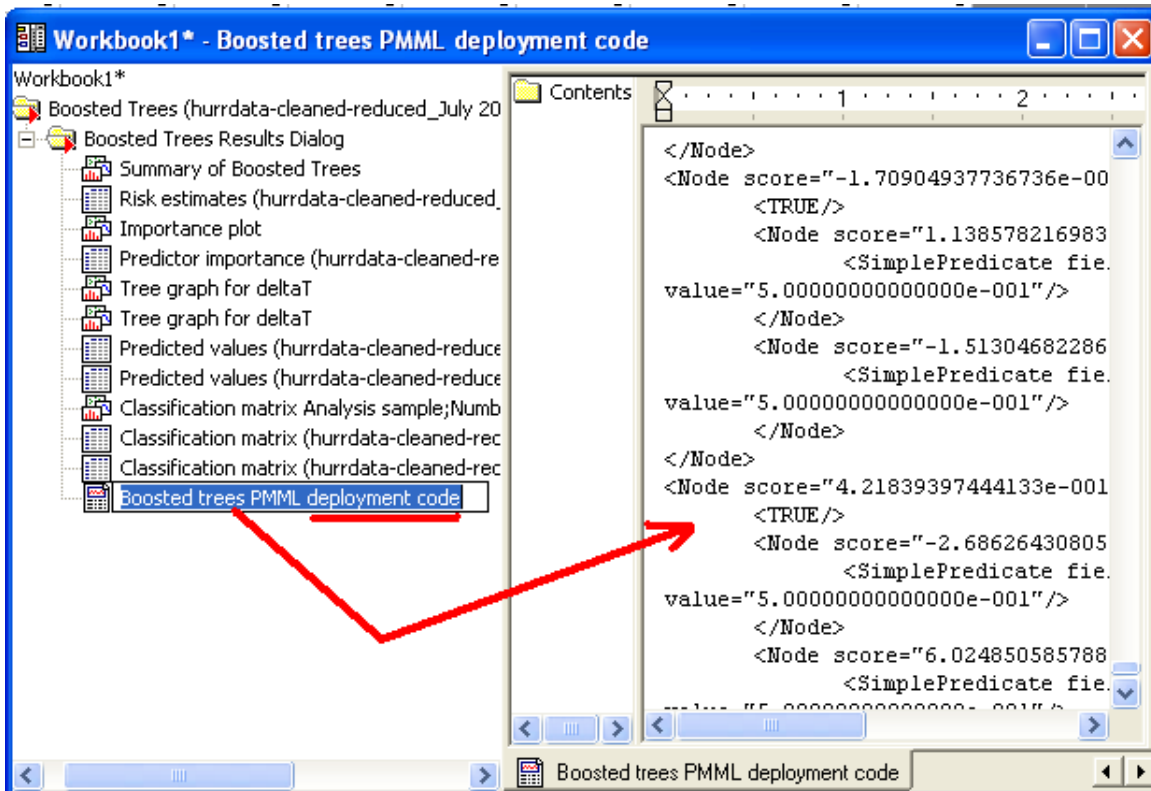**PMML script**, to produce it in a workbook window.

FIGURE 23. Workbook results from selecting the PMML deployment; here the name BOOSTED TREES PMML DEPLOYMENT CODE is in the left panel, and the actual code is in the right panel.

Then, selecting the **Boosted trees PMML deployment code** in the left panel above produces the following window where we choose **Save Item(s) As…,** and save the file as indicated, remembering the directory and file name.
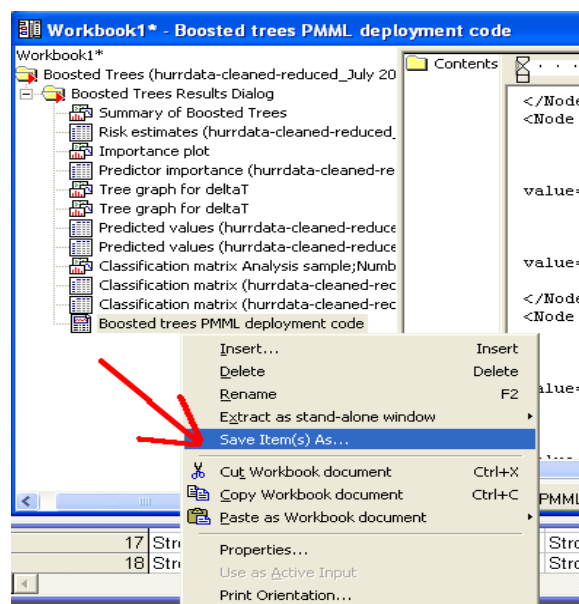
FIGURE 24.  By right-clicking on the BOOSTED TREES PMML DEPLOYMENT CODE file name a "flying menu" appears;  on this flying menu select SAVE ITEM AS, and release, to bring up a folder / file dialog, as shown below [you may have to browse for the right one….].
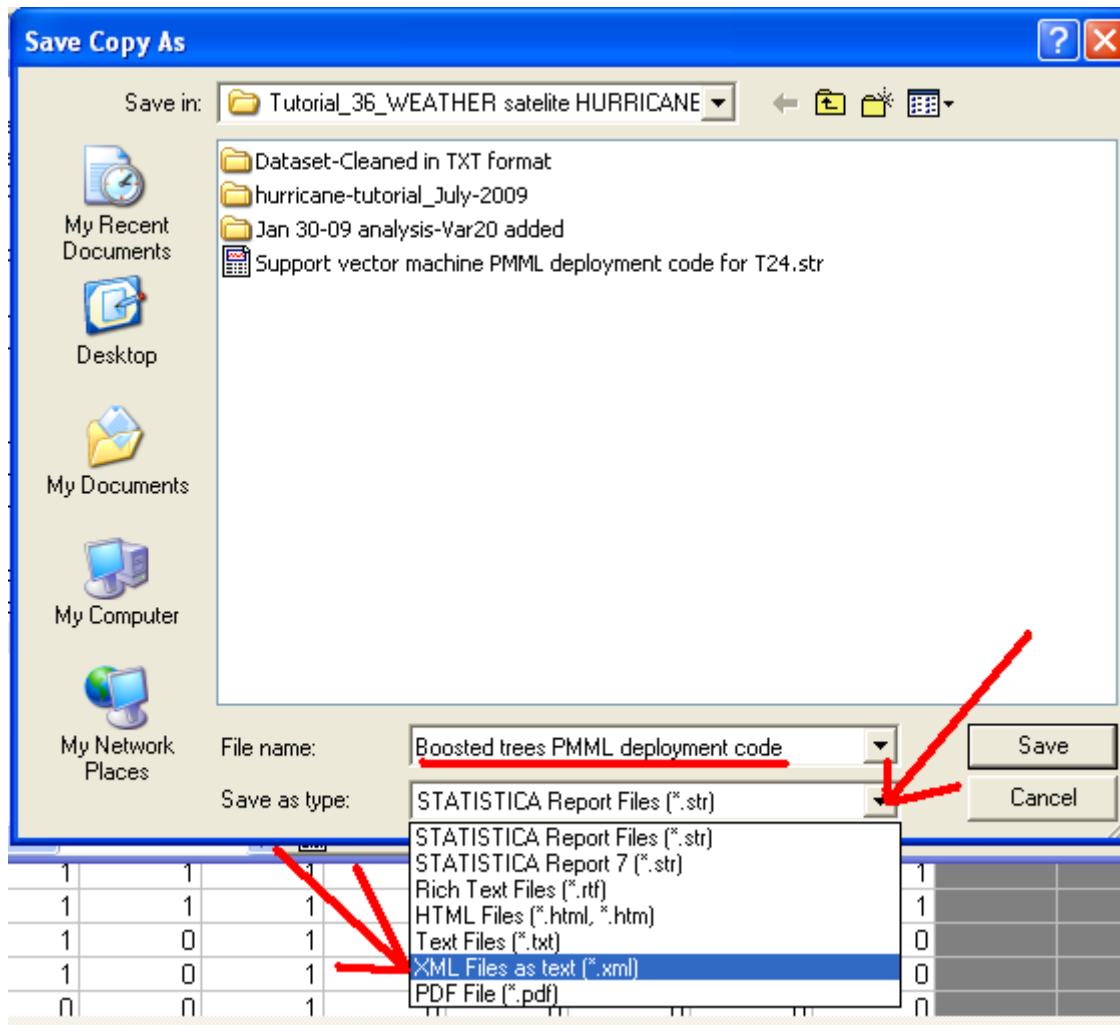


FIGURE 25.  Select the SML FILES AS TEXT (".xml") to save the deployment code into this type of xml file.
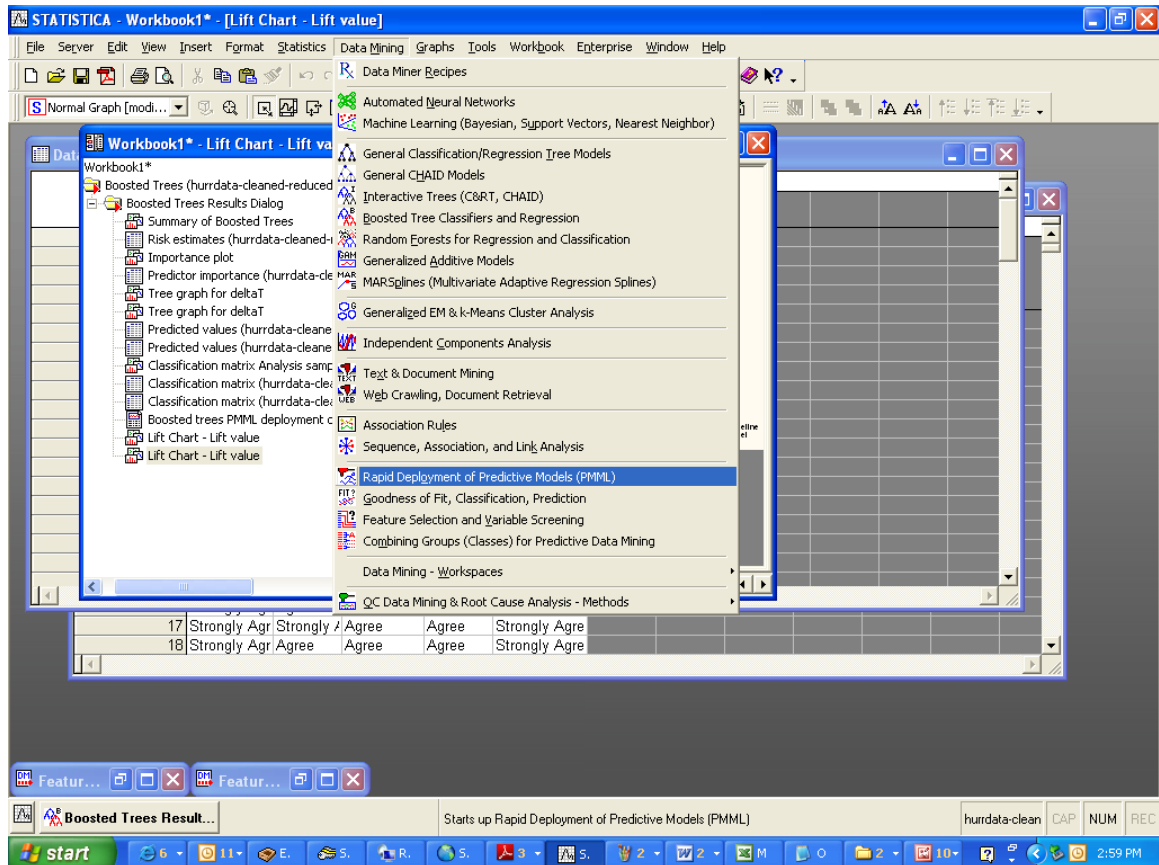
Now we open the **<u>Rapid Deployment Module</u>**



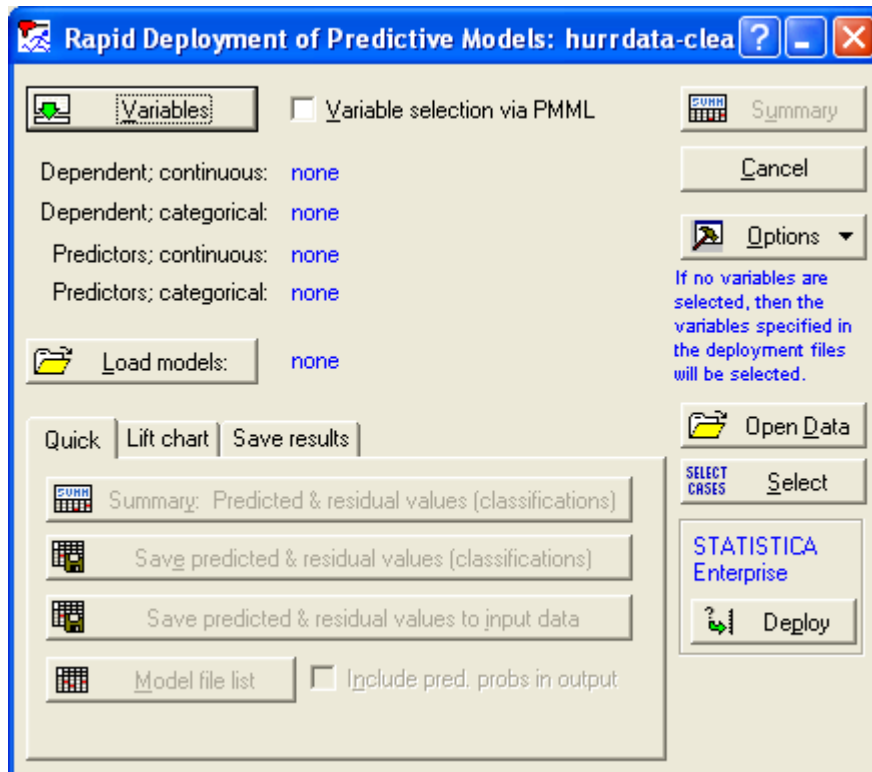FIGURE 26. Now select the RAPID DEPLOYMENT of PREDICTIVE MODELS (PMML) selection, as show above.

FIGURE 27.  The RAPID DEPLOYMENT OF PREDICTIVE MODELS dialog.

Then click on the **Load Models** button, and choose the items as indicated below.
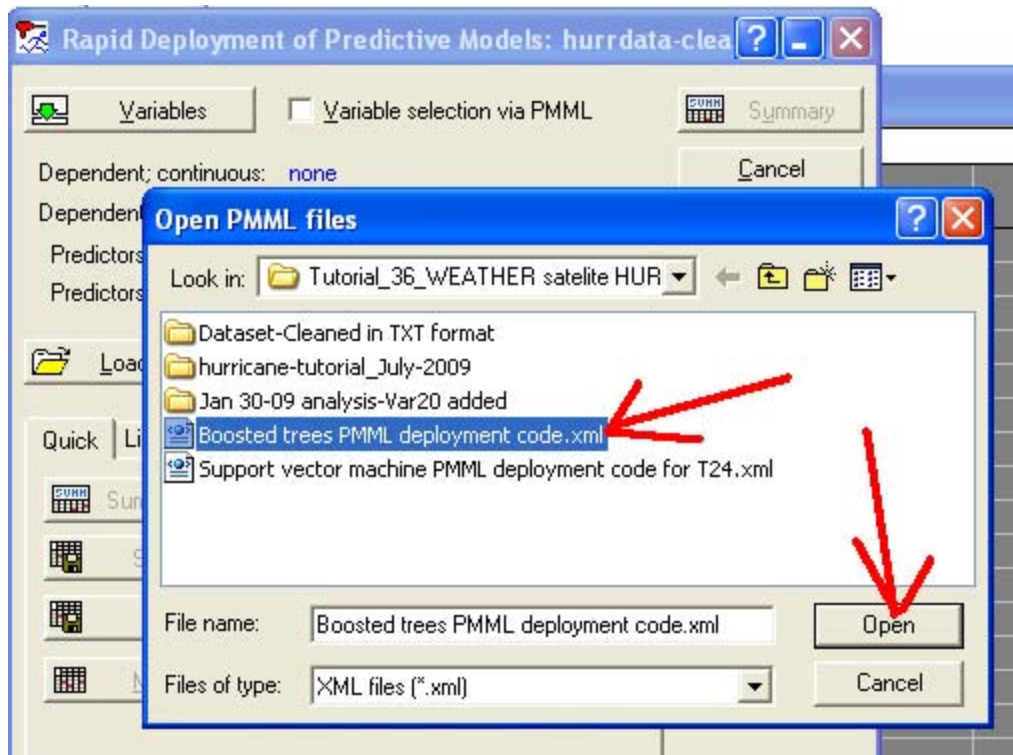
FIGURE 28. Select the BOOSTED TREES PMML DEPLOYMENT CODE.XML that you have saved away from your BT Modeling, and click OPEN to put into the RAPID DEPLOYMENT process.
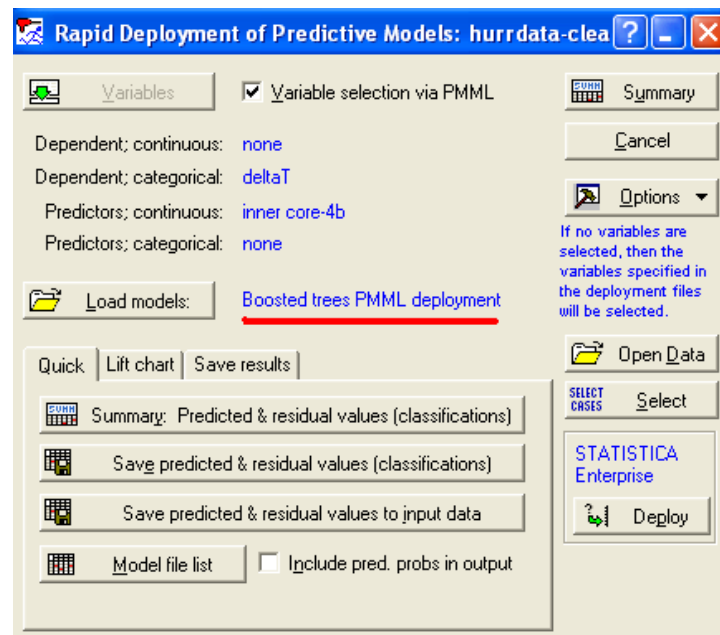


FIGURE 29.  The RAPID DEPLOYMENT DIALOG after the model has been loaded.

24

Note that the variables have already been selected (this is because our dataset was sitting open, behind the Rapid Deployment Screen, and selected as the "active spreadsheet input". We could have brought up a  2$^{nd}$ dataset, with the same variable names and same basic format, and it would have selected those variables from this new dataset.

Now hit the **<u>Summary</u>** button to get the results.



FIGURE 30.   The new data has been "scored" by the Boosted Tree model as to whether it is a DELTA-T  "yes" or "no" to a hurricane, and then giving in the 3$^{rd}$ column whether this prediction was "correct" or "incorrect".  [e.g. if  85% of the cases are CORRECT then we'd say that we have an 85% Accuracy Rate in predicting "Delta-T"………..].

## V. CONCLUSIONS

Of course, summarizing the data in the manner that we are suggesting is where the physical insight into the problem comes into play. What this tutorial aims to show is that once the data has been prepared properly, (i) the process of building a model is systematic and straightforward with the help of the software, and (ii) the statistical (data mining) approach that we are pursuing can be fruitful.

**Bibliography:**
1) "Tropical Cyclone Intensity Analysis and Forecasting from Satellite Imagery"
   Vernon F Dvorak, Monthly Weather Review 103, 420 (1975)

2) "The Dvorak Technique Explained", NOAA,
   http://www.ssd.noaa.gov/PS/TROP/dvorak.html

3) "Development of an Objective Scheme to Estimate Tropical Cyclone Intensity
   from Digital Geostationary Satellite Infrared Imagery"
   Christopher S. Velden, Timothy L. Olander and Raymond M. Zehr,
   http://cimss.ssec.wisc.edu/tropic/research/products/dvorak/odt.html