# Session 2

# Tutorial 2A.  Estimating the Naïve Bayes Model

As the starting point for this tutorial, we will retrieve model settings that have been saved in a file.  Alternatively, we could open the data file 'leukemia.Tr&Val.sav' and apply the model settings using the GUI.

**Opening a Previously Saved Project:**
> ➢ File → Load Project…
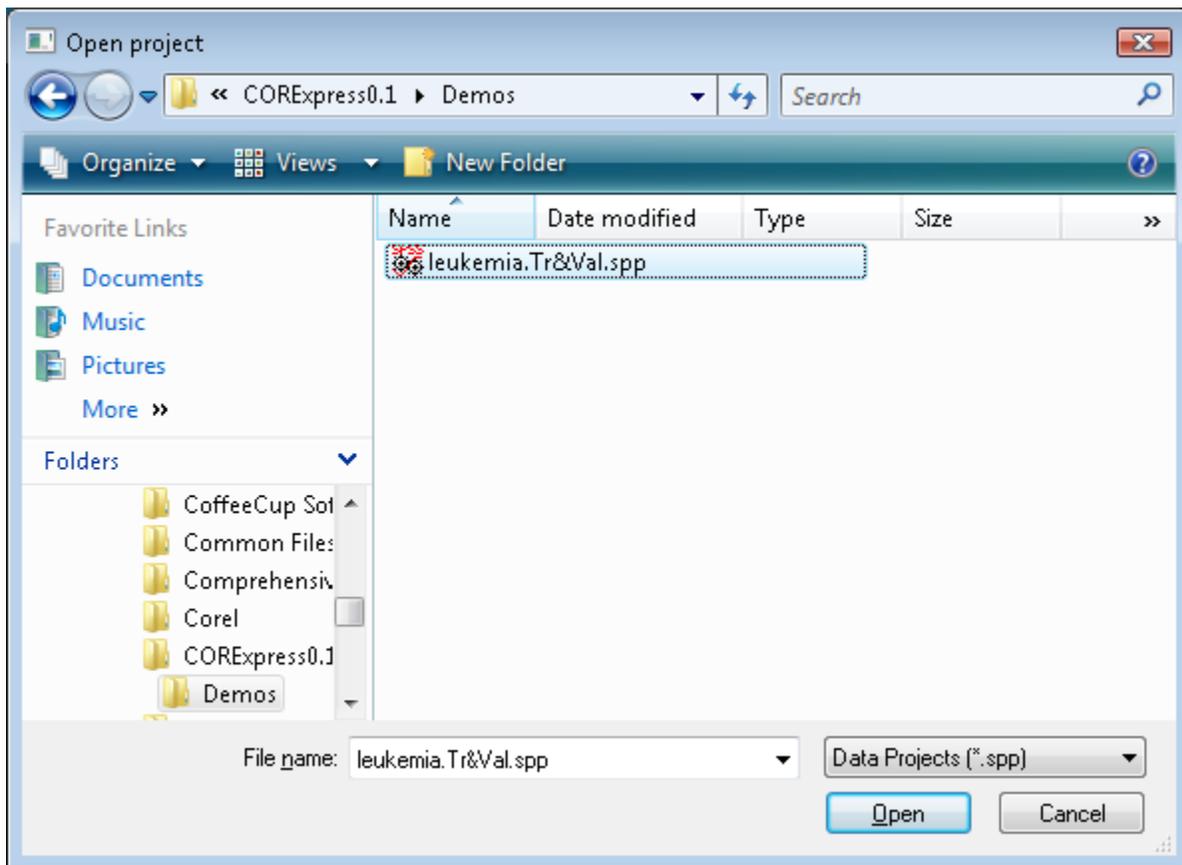> ➢ Select 'leukemia.Tr&Val.spp' and click Open to load the project



**Fig. 1:** Loading a previously saved project

**Opening the Model Specifications for the Saved Project:**
> ➢ Double click on "Leukemia" in the Datasets window

The control window will now show the saved model specifications and the corresponding model output.
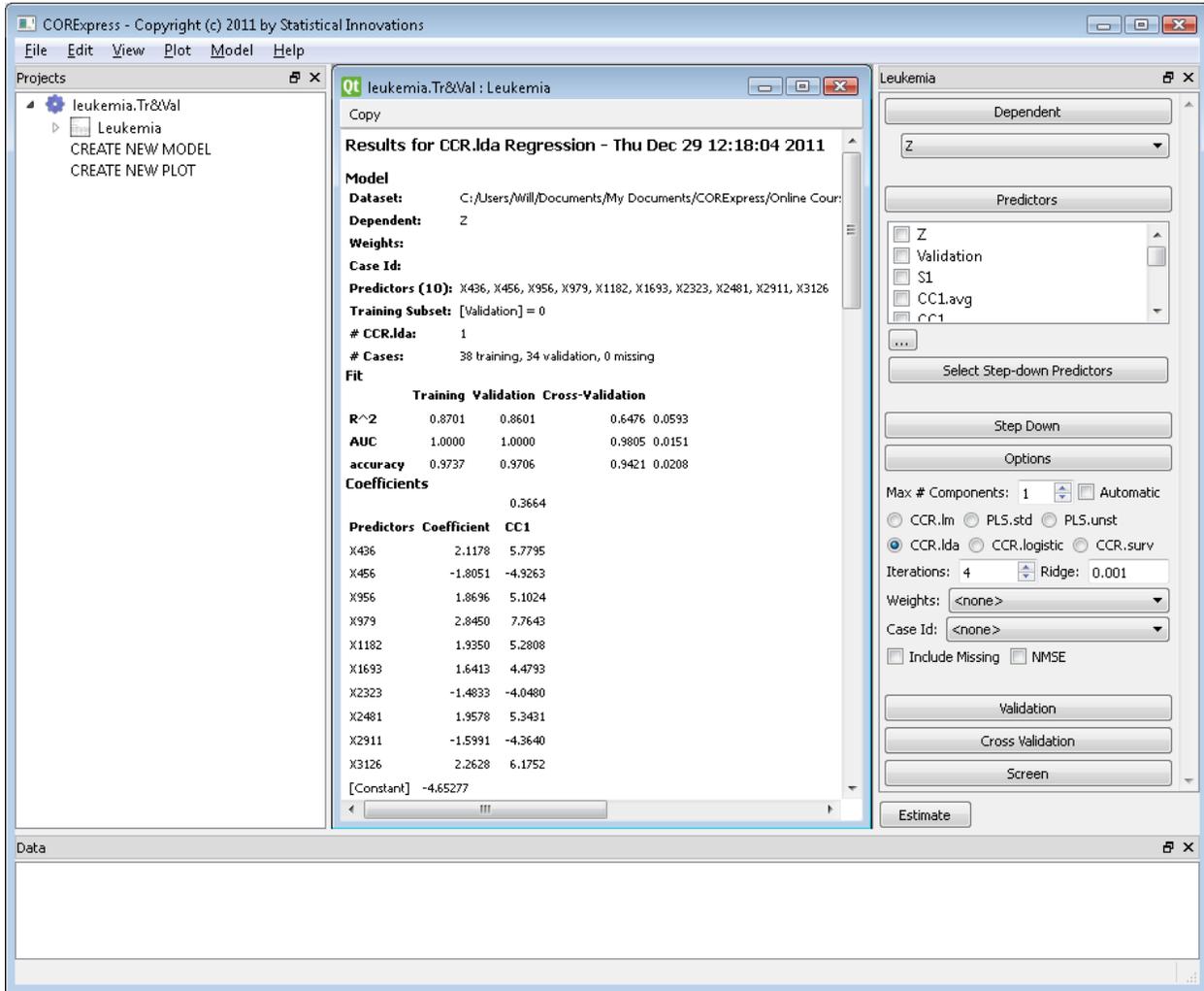


**Fig. 2:** Results from Previously Saved Project

Note the following:

- The dependent variable is Z, coded '1' and '0', to distinguish the two types of leukemia.
- We selected all 3,571 predictors as candidate predictors.
- For step down options, we have selected a minimum of 1 predictor and a maximum of 10 predictors.

- We request that 95% of predictors be removed at each step, until 100 predictors remain, at which time it resumes with the default step-down of 1 predictor at a time. For the 1-component model, choosing the '95%' option speeds up the process substantially, without altering the solution at all.
- CCR.lda has been selected with 1 component (this specifies the Naïve Bayes model, under the assumptions of Linear Discriminant Analysis).
- The cases to be used to develop the model (analysis sample) is the training data set, defined by 'Validation=0' (N=38 cases). (The remaining 34 cases, 'Validation'=1, are retained for use as a 'test' set).
- We selected the "Use Cross Validation" option, requesting the use of 10 rounds of 6-folds on the training data.
- We selected the "Stratify" option, which attempts to equalize the folds with respect to the distribution of the dependent variable Z.

**Estimate the Model:**
  ➢ Click Estimate

The program quickly reduces the number of predictors in the model from 3,571 to 10.
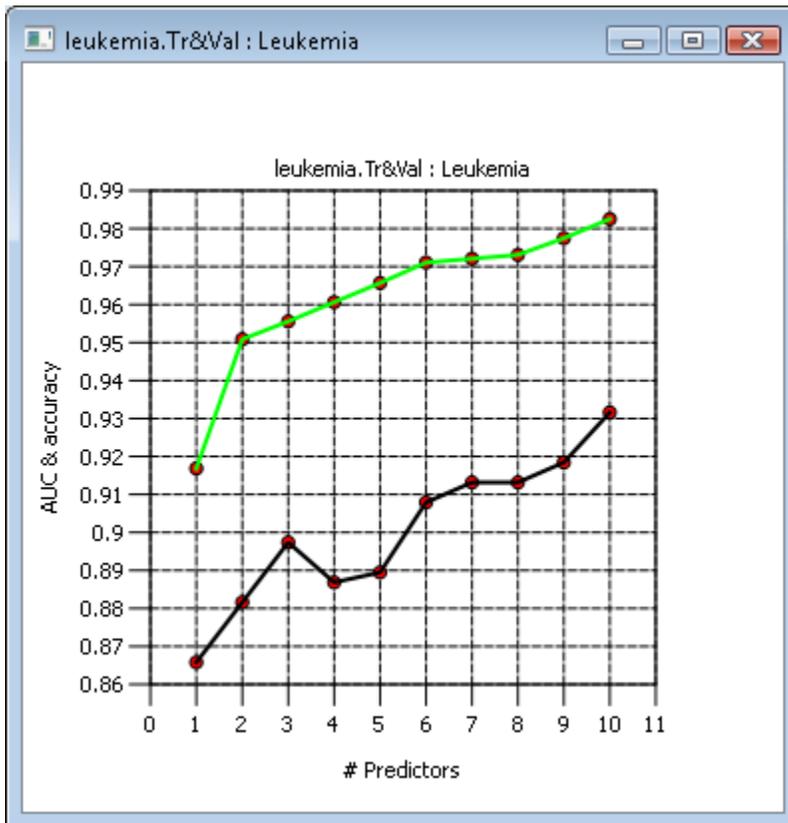


**Fig. 3:** CV-AUC and CV-ACC Plot

3

**Viewing the Training Interactive Plot:**

Two interactive plots are available -- one for the training, the other for the validation data.

➢ Click on the drop down arrow next to "Leukemia" in the Datasets window
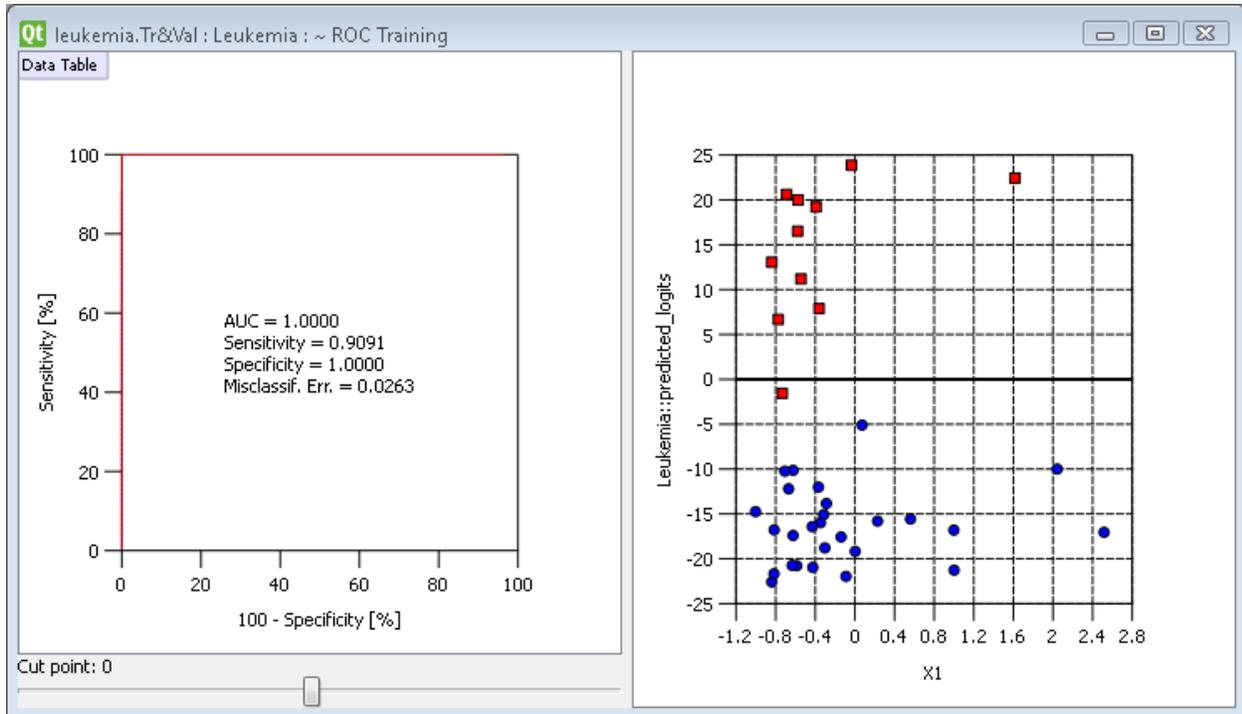➢ Double click on "~ ROC Training



**Fig. 4:** Training Dataset ROC & Scatter Plot

Notice that AUC = 1, meaning that all cases coded Z=1 have higher predicted scores than any case coded Z=0.

The scatterplot shows that 1 of the red square symbols (representing cases in the Z=1 group) is misclassified based on the default logit score cut-point of 0 (logit = 0 corresponds to predicted probability of .5). Reducing the cut-point, lowers the classification line to produce 100% correct classification.
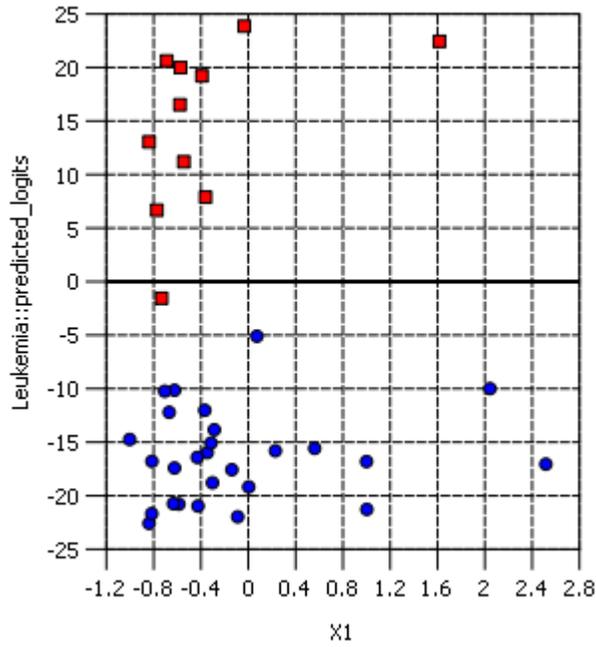
4

**Fig. 5:** Training Data

Using the cutpoint slider below the ROC curve, you can adjust the cutpoint so that all the red subjects are above the cutpoint and all blue subjects are below the line.
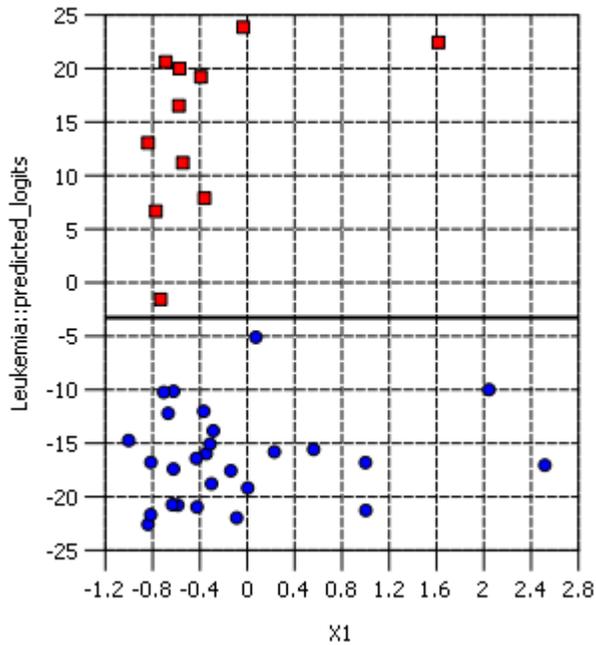


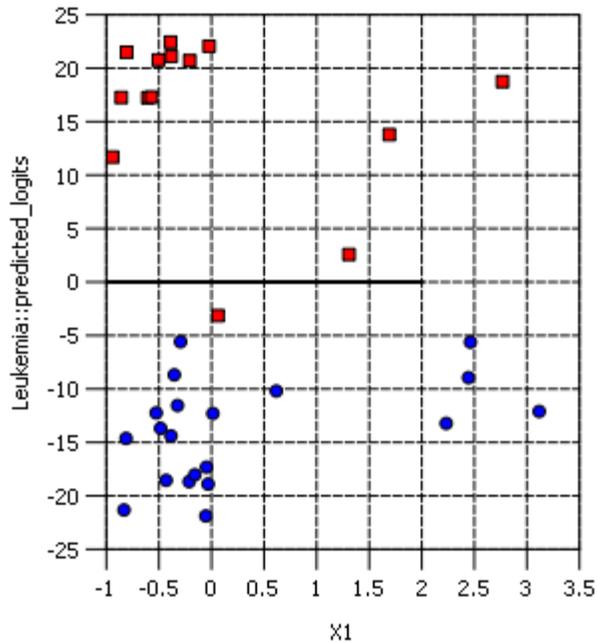**Fig. 6:** Training Data with adjusted cutpoint

**Fig. 7:** Validation Data

For the validation data, we also obtain AUC = 1, which confirms that the results obtained using the training data were not due to chance.

As an exercise, under "Step Down", change the percent from '95' to '50', re-estimate the model and confirm that you get the same results.

**Change Remove by Percent:**
- ➢ Double click on "Leukemia" in the Datasets window to bring up the model specifications in the Control Window.
- ➢ Under Stepdown, delete '95' and type '50'
- ➢ Click Estimate

The reason that the results are the same is that the loadings from the 1-component model are not affected by the number of predictors in the model, eliminating predictors that have the smallest magnitude of the standardized coefficient is equivalent to eliminating predictors that have the smallest magnitude of the loadings. Thus, it does not matter how many predictors are eliminated at the same – the 10 predictors remaining in the model will be same.

To compare the coefficients with those obtained from LDA, we can estimate the LDA model with CORExpress. LDA is equivalent to a saturated CCR model, which with 10 predictors occurs by changing the number of components from 1 to 10.

6

To select these 10 predictors do the following:

**Select Step-down Predictors:**
- ➢ In the Control Window, click the "Select Step-down Predictors" button below the Predictors window

This will check off only the 10 predictors from the previously estimated model. Note that the "Step Down" feature will automatically be unchecked.

**Change Number of Components:**
- ➢ Change the # components from 1 to 10 to estimate the saturated model

**Turn off the Cross Validation Option:**
- ➢ Scroll down and uncheck the "Use Cross Validation" option under "Cross Validation"
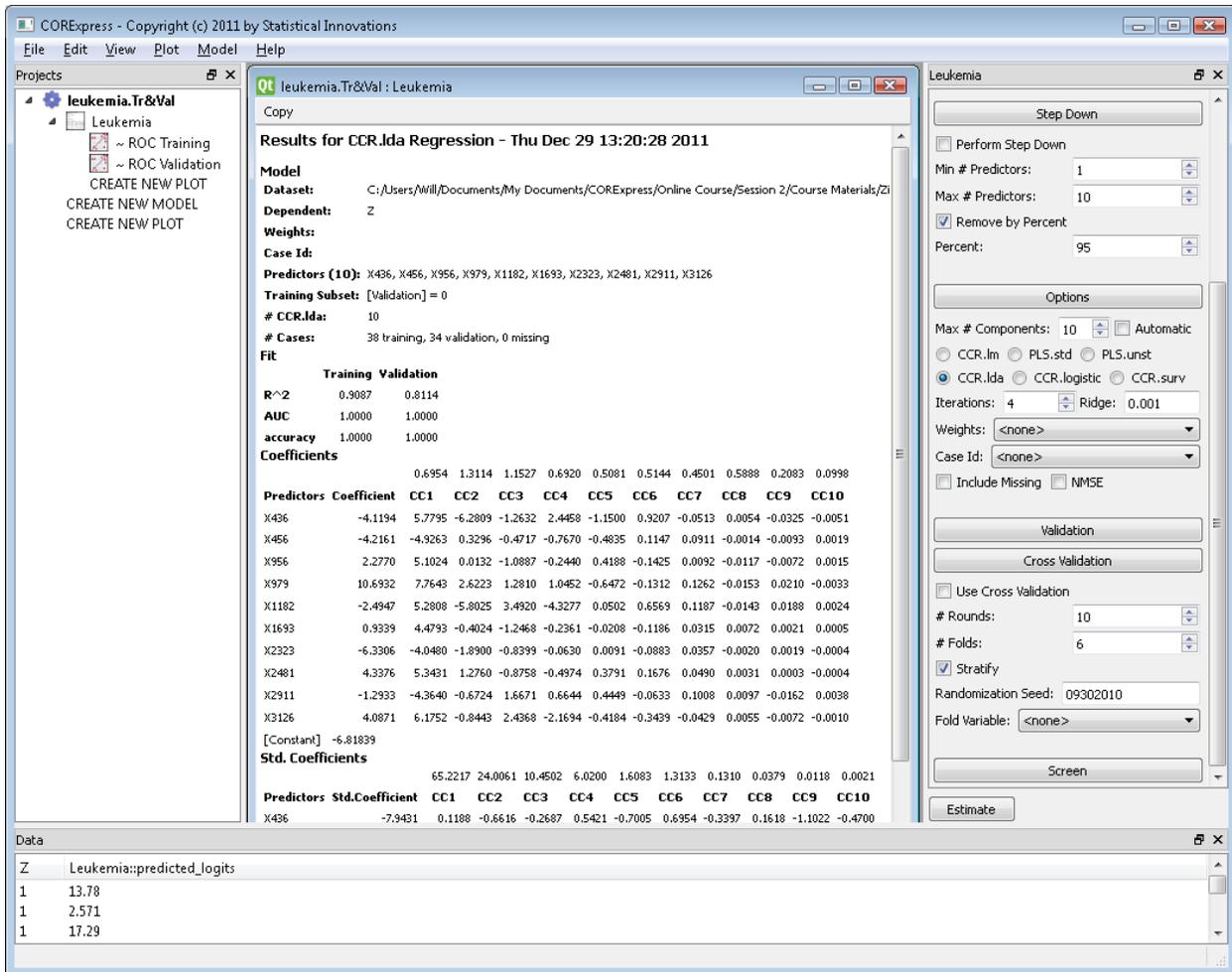
**Estimate the Model:**
- ➢ Click Estimate

**Fig. 8:**

Notice that the coefficients reported under the first component CC1 are unchanged from the Naïve Bayes (NB) model. NB can be viewed as an extreme form of regularization, where all components except for the first component are set to 0.