

## 1. Classification vs Regression

This is a classification problem, as we are aiming to predict what is the likelihood of the student to fail so that we can intervene at correct time so as to help him out. Moreover, this is a binary classification problem with two classes fails / pass.

## 2. Exploring the Data

Can you find out the following facts about the dataset?

- Total number of students : **395**
- Number of students who passed : **265**
- Number of students who failed : **130**
- Graduation rate of the class (%) : **67.09**
- Number of features (excluding the label/target column) : **30**

Use the code block provided in the template to compute these values.

## 3. Preparing the Data

Execute the following steps to prepare the data for modeling, training and testing:

- Identify feature and target columns
- Preprocess feature columns
- Split data into training and test sets

Starter code snippets for these steps have been provided in the template.

The code has been written and added in the ipython notebook

## 4. Training and Evaluating Models

Choose 3 supervised learning models that are available in scikit-learn, and appropriate for this problem. For each model:

### Naive Bayes

This is very quick to train and predicting because it is the simplest model which can be used in order to predict thing very fast as its complexity is  $O(cd + nd)$  where  $c$  is the number of classes and the  $n$  is the number of data points and  $d$  is the number of features. It is best used in text classification problems like spam email detection where we need predictions very fast, face recognition etc. But being very simplistic it assumes that all the features are

independent of each other thus the prediction results is not very accurate. But it even it does not overfit the data. F1- Score increases with increase in number of training data points.

Training Data	100	200	300
Training Time (sec)	0.001	0.001	0.001
Prediction Time (Test)(sec)	0.001	0.001	0.001
F1 - Score (Training)	0.494117647059	0.810606060606	0.785542168675
F1 - Score (Testing)	0.408602150538	0.75	0.791044776119

### Random Forest Classifier

With fairly small number of data points random forest can easily be used irrespective of its high time complexity. This is best suited for the problems where features are very likely to be dependent on each other or can be explained through trees thus using an ensemble of tree is a good approach. As it is evident from the following table that the training is very overfitted with less number of data points thus, as it assumes that all the features are mutually exclusive from each other which is not the case all the time, thus resulting in high F1- Score with less training data where it is more likely to overfit the data because of mutually exclusive assumption.

Training Data	100	200	300
Training Time (sec)	0.045	0.022	0.037
Prediction Time (Test)	0.001	0.001	0.001
F1 - Score (Training)	1.0	0.992307692308	0.992405063291
F1 - Score (Testing)	0.629921259843	0.755905511811	0.832116788321

### Support Vector Machines

Though this is the most computationally expensive algorithm as it has cubic time complexity, but because of its good classification results where we separable problems can easily be classified with high accuracy such as bioinformatics classification and image classification. As it comes with a kernel trick which can be used to convert non linear problem to linearly

separable problem and thus finding the optimal hyperplane between the classes, but it seems to overfit if too many features were given, as in our case we have just 30 features therefore it performed really well. With the increase in the number of records the F1- Score for training and testing and training and testing time is increasing because the time complexity is a function of  $O(n^3)$  for training.

Training Data	100	200	300
Training Time	0.064	0.003	0.010
Prediction Time (Test)	0.001	0.001	0.002
F1 - Score (Training)	0.882352941176	0.881355932203	0.845986984816
F1 - Score (Testing)	0.774647887324	0.783783783784	0.833333333333

## 5. Choosing the Best Model

In our case as we can see there is not much difference in time complexity as we have just 30 features and 300 training data points thus, we can use SVC for the particular problem. As the F1- Score training and testing quite impressive and quite evidently increasing too which is better than random forest (testing) though random forest has good training F1- Score which indicates overfitting thus lead poor F1 - Score (testing), whereas even Naive Bayes even do not have such a good F1 - Score as compared to SVC. Therefore, SVC is the best choice.

As SVC have 2 hyperparameters that can be tuned as well i.e. gamma and c for which GridCV Search and Stratified ShuffleSplit has been used.

What is the model's final F1 score? : 0.864864864865

```
In [37]: # Test algorithm's performance
print "F1 score for training set: {}".format(predict_labels(best_clf, X_train, y_train))
print "F1 score for test set: {}".format(predict_labels(best_clf, X_test, y_test))

Predicting labels using SVC...
Done!
Prediction time (secs): 0.005
F1 score for training set: 0.819956616052
Predicting labels using SVC...
Done!
Prediction time (secs): 0.001
F1 score for test set: 0.864864864865
```