
Final Report: Predicting student's learning performance with KDD Cup 2010

Binbin Xiong
Information Networking Institute
Carnegie Mellon University
Pittsburgh, PA 15213
bxiong@andrew.cmu.edu

Yujing Zhang
Information Networking Institute
Carnegie Mellon University
Pittsburgh, PA 15213
yujingz1@andrew.cmu.edu

Abstract

KDD Cup 2010 is an educational data mining problem of practical importance. It is to learn a model from students' past behavior and then predict their future performance. We extracted five features: *student ID*, *problem name*, *step*, *knowledge component*, and *opportunity count* from the dataset to predict the binary value of *correct first attempt*. Then, we implemented several classifiers including Naive Bayes, Logistic Regression, and SVM as our training models. In addition, we used Adaboost and k-means these two techniques to optimize the training process. At last, we compared each classifier with accuracy and Root Mean Squared Error to analyze the results.

1 Background

What is a student's learning generality and speed? What is the difference and similarity between problems How can we the difficulty degree of a problem? These problems might be solved by analyzing the knowledge concepts and skills each problem needs.

1.1 Problem Description

In this report, we are tackling the problem of student learning process. We predict the student learning performance on mathematical problems given information regarding past performance, with the 2010 Knowledge Discovery and Data Mining competition dataset. This problem is not only a technical topic for researchers, but also of practical importance.

1.2 KDD 2010 Data Set Description

The dataset is divided into a training set and a test set. Four key terms form the building blocks of the dataset. These are problem, step, knowledge component, and opportunity. A problem is a task for a student to perform, and it typically involves multiple steps, which are observable parts of the solution to a problem. A knowledge component is a piece of information that is used to solve a step. It can be concept, principle, fact, skill, etc. Each step can be labeled with one or more hypothesized knowledge components. And every knowledge component is associated with one or more steps. Finally, an opportunity is a chance for a student to demonstrate whether he or she has learned a given knowledge component. In this report, we predict the binary value of "correct first attempt" as student-step performance.

2 Related Works

In[2], the authors points out that feature extraction is playing a vital role in this problem. In addition, as the student is learning during the tutoring process, temporal information is of great importance. Their work can be separated into two parts: feature extraction and classification. For feature extraction, they use two techniques: sparse features generated by expanding a multi-category feature into a set of binary features, and condensed features using feature

combinations to indicate some relationships. For classification part, they tried several machine learning techniques to train the learner. They use Logistic Regression for sparse feature set, and use Random Forest, AdaBoost for dense feature set. Finally, they combined all results from different classifiers by regularized linear regression, and won the first place of the competition.

The final prediction of the 2nd team[3] was a combination of Bayesian Hidden Markov Models (HMMs) and bagged decision tree. Their HMM model learns student specific parameters and then used them to train skill specific models. It was to predict the probability of knowledge for each student at each opportunity and the probability of correctness on each step. Then they used the prior skill data learned to generate features for the test sets, and implemented Random Forest to do the prediction.

The team of 3rd place[4] used Singular Value Decomposition (SVD) technique, and won the third place. Based on the idea of SVD model, they tried several factor models in order to correct the shortcomings of pure SVD. They introduced many biases into the model, like a step dependent bias, a student dependent bias, a KC dependent bias, a student and KC dependent bias, etc. In addition, they used KC dependent N dimensional feature vectors. Both of the bias and feature vector methods significantly lowered the value of RMSE.

3 Learning Model

According to implicit relationship among data fields, we extract several key features and build a feature model. Then we use both discriminative and generative models including Naive Bayes, Logistic Regression, and SVM to learn and test the *Algebra I 2005-2006* data set. In order to improve learning performance, we implement two techniques: AdaBoost and k-means, to obtain a higher classification accuracy as well as a lower rooted mean squared error.

3.1 Feature Extraction

There are 19 columns in total in the dataset, but some of them are irrelevant to our prediction, like row number. We extract some important features to build the learning model. Different problems and knowledge components have different difficulty degrees. Some are easy to solve, but others are hard. In addition, the more times a knowledge component the student meet, the easier the step for the student. And different students have different learning rate. So we decide to use *student ID*, *problem name*, *step*, *knowledge component*, and *opportunity count* as five features of our learning model. In addition, we use the value of *corret first attempt* as student-step performance predictor.

3.2 Naive Bayes

Naive Bayes classifier is a simple probabilistic classifier based on applying Bayes Theorem with conditional independent given class label. Assume $X = \{x_1, x_2, \dots, x_n\}$ is a collection of features, and $Y = \{1, 0\}$. Each x_i is a discrete feature, and $x_i \in X$ is independent given Y. Then:

$$P(Y = y \mid x_1, x_2, \dots, x_n) = \frac{P(Y = y) \prod_{i=1}^n P(x_i \mid Y = y)}{\sum_{k=0}^1 P(Y = k) \prod_{i=1}^n P(x_i \mid Y = k)} \quad (1)$$

$$Y \leftarrow \arg \max_{y_k} P(Y = y_k) \prod_{i=1}^n P(x_i \mid Y = y_k)$$

Naive Bayes classifier can be trained in linear time, rather than by expensive iterative approximation as used for many other types of classifiers. Applying MAP estimation in Naive Bayes is necessary to avoid overfitting and extreme probability values.

3.3 Logistic Regression

Logistic Regression is also a probabilistic statistical classification model. For binary classification:

$$P(Y = 1 \mid X) = \frac{1}{1 + \exp(w_0 + \sum_{i=1}^n w_i X_i)} \quad P(Y = 0 \mid X) = \frac{\exp(w_0 + \sum_{i=1}^n w_i X_i)}{1 + \exp(w_0 + \sum_{i=1}^n w_i X_i)} \quad (2)$$

This problem can be solved through gradient ascent, which is to update the vector of parameters iteratively, until it satisfies some termination condition. In addition, to reduce overfitting, Logistic Regression is allowed to add a

regularization term $P(W)$. In our implementation, we plan to use both L1 and L2 norm, which can be represented as:

$$L1 : \ln P(W) \propto \sum_i |w_i| \quad L2 : \ln P(W) \propto \sum_i w_i^2 \quad (3)$$

Compared to Naive Bayes, Logistic Regression has no assumption that features are independent with each other conditional on label, but it can only produce a linear decision boundary.

3.4 Support Vector Machine

Support Vector Machine (SVM) is a discriminative classifier formally defined by a separating hyperplane. In other words, given labeled training data (supervised learning), the algorithm outputs an optimal hyperplane that gives the maximum margin γ to the training examples.

$$\begin{aligned} \min_w \quad & ||w||_2^2 + C \sum_{i=1}^m \xi_i \\ \text{subject to} \quad & y_i(w^T x_i) \geq 1 - \xi_i \quad \forall i \\ & \xi_i \geq 0 \quad \forall i \end{aligned} \quad (4)$$

Here the C is the slack variant that we need to assign before starting the training process.

The goal of SVM is to find a linear boundary with the maximum margin of the different labels. Therefore, if the data is linearly separable, the *linear* boundary that the SVM finds should be with a rather high accuracy. With the same cost function shown above, in addition to performing linear classification, SVMs can efficiently perform a non-linear classification using the *kernel trick*, which can implicitly mapping their inputs into high-dimensional feature spaces. A typical kernel function that can map the feature into a infinite space is the Gaussian kernel (which is also known as RBF), which can be represented as:

$$K(X_i, X_j) = \exp(-\gamma ||X_i - X_j||^2), \text{ for } \gamma > 0 \quad (5)$$

Thus if we train a non-linear classifier, even if the training data is not linear separable, it should still give a rather decent accuracy, which means that by applying both linear and non-linear SVM classifiers, the change of the accuracy of them *should* tell us something of the distribution of training data. Therefore, we may implement both linear and non-linear with RBF kernel L1 SVM classifiers in our system.

3.5 AdaBoost

AdaBoost is adaptive boosting algorithm. It can turn weak algorithm into a strong learner by concentrating on examples misclassified by previous iteration and combining hypothesis of each iteration as the final hypothesis.

For dataset $S = \{(x_1, y_1), \dots, (x_m, y_m)\}; x_i \in X, y_i \in Y = \{-1, 1\}$ and weak algorithm A, construct example distribution $D_t (t = 1, 2, \dots, T)$.

$$D_{t+1} = \frac{D_t(i)}{Z_t} \exp\{-\alpha_i y_i h_t(x_i)\} \quad (D_1 = \frac{1}{m})$$

Run A on D_t producing $h_t: X \rightarrow \{-1, 1\}, \epsilon_t = P_{x_i}(h_t(x_i) \neq y_i)$

$$\alpha_t = \frac{1}{2} \ln \frac{1 - \epsilon_t}{\epsilon_t}$$

$$H_{final}(x) = \text{sign}\left(\sum_t \alpha_t h_t(x)\right)$$

AdaBoost can be used to improve the learning performance for any weak algorithms, like Naive bayes, decision stumps, etc. It converges vary fast, and has no parameters to tune. In addition, when training accuracy converges, the test accuracy could be further improved as the number of rounds increasing.

3.6 K-Means

K-Means is an unsupervised learning algorithm. It partitions data into k clusters with the nearest distance to the cluster center(mean). Assume that $d(x, y)$ represents the distance between point x and y . The algorithm aims to find center c_1, c_2, \dots, c_k of k clusters, and

$$\text{minimize} \sum_{i=1}^n \min_{j \in \{1, 2, \dots, k\}} d^2(x_i, c_j)$$

At the beginning, initialize k clusters with centers c_i , and partition all data into these k clusters. Then, according to the data distribution, update the center of each cluster. Finally, we can find an optimal local solution through iterations.

4 Experiment Results

4.1 Data Preprocessing

We first extract *student ID*, *problem name*, *step*, *knowledge component*, and *opportunity count* from the dataset to predict the value of *correct first attempt* from the training set. For some classifiers, we map the string fields into positive integers before feeding them to the learning machine.

4.2 Naive Bayes

We use both MLE and MAP to compute the conditional probabilities of each feature. For MLE, the conditional probability is represented as:

$$P(x_i = m \mid Y = y_k) = \frac{\text{count of } x_i = m \text{ among } Y = y_k \text{ cases}}{\text{total count of } Y = y_k \text{ cases}} \quad (6)$$

For MAP, we use Laplace Smoothing. The conditional probability is represented as:

$$P(x_i = m \mid Y = y_k) = \frac{\text{count of } x_i = m \text{ among } Y = y_k \text{ cases} + 1}{\text{total count of } Y = y_k \text{ cases} + |x_i|} \quad (7)$$

The experiment results are shown in Table 1. We set the Naive Bayes with Maximum Likelihood as our benchmark in the following experiment.

Table 1: Classification results of Naive Bayes

CLASSIFIER	TRAIN ACCURACY	TEST ACCURACY	TRAIN RSME	TEST RSME
NB-MLE	67.77%	63.04%	0.5676	0.6079
NB-MAP	67.53%	64.40%	0.5698	0.5966

Where the Root Mean Squared Error is computed by:

$$RMSE(\hat{y}) = \sqrt{E((\hat{y} - y)^2)} = \sqrt{\frac{\sum_{t=1}^n (\hat{y}_t - y)^2}{n}} \quad (8)$$

The prediction accuracy of Naive Bayes is not ideal. The deviation of the differences between predicted values and real values are large. And the result of MAP does not make a great improvement compared to the result of MLE. We assume that the reason of low accuracy may be the conditional independent assumption of Naive Bayes. In our case, if the features are not conditional independent, the result may be bad. For example, if a knowledge component is difficult to understand, the corresponding steps would be hard to solve. In this case, knowledge component and step are not independent.

4.3 Logistic Regression

For the logistic regression learning, we first map each feature's value into $[0, 1]$ space mentioned in previous section, and then implement the Liblinear[3] open source software to train our model. To reduce overfitting, we apply L1 and L2 regularization term in the training process. The resulting accuracy and the RSME of both training set and the test set are given as follow.

Table 2: Classification results of Logistic Regression

CLASSIFIER	TRAIN ACCURACY	TEST ACCURACY	TRAIN RSME	TEST RSME
L1-Logistic Regression	73.03%	74.38%	0.51932	0.50612
L2-Logistic Regression	73.04%	74.38%	0.51891	0.50612

Compared to Naive Bayes, we can see that the classification result of Logistic Regression is much better, which in turn proves our reasoning for the inaccuracy of NB classifier. However, the accuracy of our LR classifier is still not high. In this case, the reason should be either overfitting or that samples are *not* linearly separable. Though Logistic Regression doesn't assume conditional independence during training, it can only produce linear decision surfaces.

Since the training accuracy and the test accuracy are very similar, our regularized model is free from overfitting. Then given the previous result, the bottleneck of the LR classifier might lies in the data distribution rather than overfitting (In fact, one of the experience we learnt from class is that overfitting usually happens when the feature is too complex or the number of training examples are very limited, our LR result gives support to it), that is the data we extracted is not linearly separable.

4.4 Support Vector Machine

For the SVM training, we implement the Liblinear[3] for fast linear SVM learning and LibSVM[4] for the none-linear SVM learning.

In order to increase the accuracy as well as to boost the training speed, we apply a scaling operation to each of the features, mapping the value to $[0, 1]$, where the 0 means the missing value. In addition, we add the regularization term L1 during the training process.

Table 3: Classification results of SVM

CLASSIFIER	TRAIN ACCURACY	TEST ACCURACY	TRAIN RSME	TEST RSME
Linear SVM	73.02%	74.38%	0.51921	0.50612
Non-linear SVM	73.56%	74.86%	0.50962	0.50311

4.4.1 Cross Validation

During our implementation, in order to choose the best parameter set for both the linear and the none-linear SVM, we implement the k-fold cross validation method. In k-fold, the original sample is randomly partitioned into k equal size subsamples. Of the k subsamples, a single subsample is retained as the validation data for testing the model, and the remaining k - 1 subsamples are used as training data. The cross-validation process is then repeated k times with each of the k subsamples used exactly once as the validation data. The k results are averaged to produce a single estimation. In our implementation, we set $k = 5$. The potential parameter sets we choose for C in both linear and none-linear SVM is $C \in \{2^{-5}, 2^4, 2^3, \dots, 2^3, 2^4, 2^5\}$ while the variable g in RBF kernel is $g \in \{2^{-5}, 2^4, 2^3, \dots, 2^3, 2^4, 2^5\}$. After the grid search, the best values we get are: for linear SVM, $C = 1.0$, while for none-linear SVM, $C = 1.0$ and $g = 0.25$, while the cross validation accuracy of linear SVM is 73.85%, and for none-linear SVM is 73.96%

It is worth mentioning that during our cross validation process, we find that when we change the parameters, the accuracies of different validation processes are almost the same.

4.4.2 Classification Results

With the validated best parameters, we then train the corresponding classifiers with the processed dataset. In addition to the accuracy, we also compare the RMSE value of these two classifiers both on the training set and on the testing set. the results are shown in Table 3.

From the results shown above, we can see that both kinds of classifiers outperforms the NB classifier. In addition, the result of the linear SVM is very close to the result of Logistic Regression. However, the results we get from the non-linear classifier is surprisingly close to the linear version, which means that given the number of training samples, even when we map the data into a infinite space, the classification still can't be perfectly solved. Given this observation, we deduce that the bottleneck of our solution should lies on the feature selection, rather than the classifier models.

4.5 AdaBoost

We applied AdaBoost to improve the learning performance of Naive Bayes. The training process is shown in the table below:

Table 4: Classification results of AdaBoost

Classifier	Train Accuracy	Test Accuracy	Train RSME	Test RSME
Naive Bayes	67.53%	64.40%	0.5698	0.5966
AdaBosst	73.03%	74.38%	0.5192	0.5061

AdaBoost greatly enhance the accuracy of Naive Bayes. And the convergence rate of AdaBoost is very fast. After 4 iterations, the error rates of both the training dataset and test dataset converge. In addition, the test error rate decreases as the training error rate increases. There is no overfitting problem, because as the training error rate increases, the classifier is more and more confident.

4.6 Feature Extraction

Dimension Reduction: Compared to others features, the dimension of "step" is very high, which makes it a highly

Table 5: Dimension of Each Feature

Feature	Dimension
Student ID	574
Problem	1084
Step	187539
KC	112

sparse vector. However, many of them are very similar, like " $-4x+2=7$ " and " $-4x+2-2=7-2$ ", " $x-500000=2$ " and " $x-500000=40000$ ". We remove all the numbers and spaces in the step names, and then get 26582 features. We use the lower dimensional "step" to feed the Naive Bayes, Logistic Regression, and SVM classifiers again. Here are the results:

5 Conclusion

As a conclusion, we plot all the results from previous sections in the same figure as shown in Figure 1 and Figure 2. Despite of the accuracy, we also compare the relative wall clock time consumption of these different learning models, which can be seen from Table 4.

Table 6: Relative time complexity of different classifiers			
NAIVE BAYES	LR	LINEAR-SVM	KERNEL-SVM
VERY LOW	LOW	VERY LOW	VERY HIGH

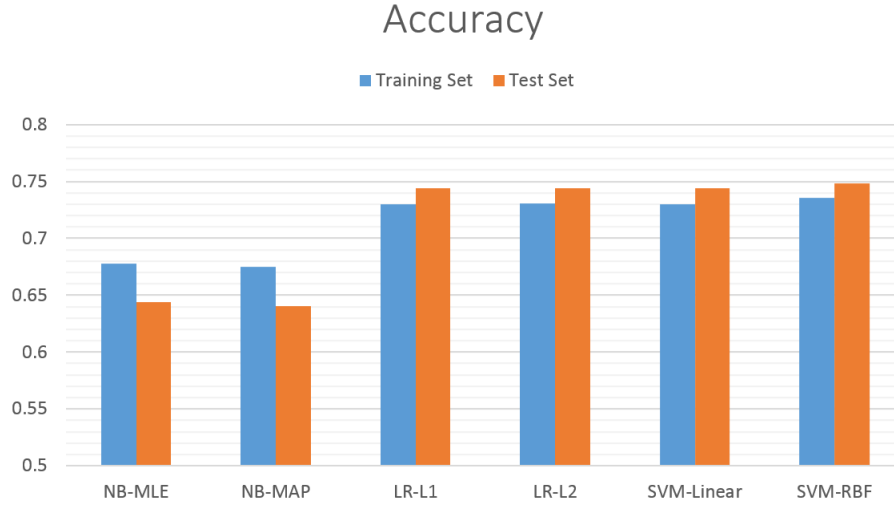


Figure 1: Accuracy of different classifiers

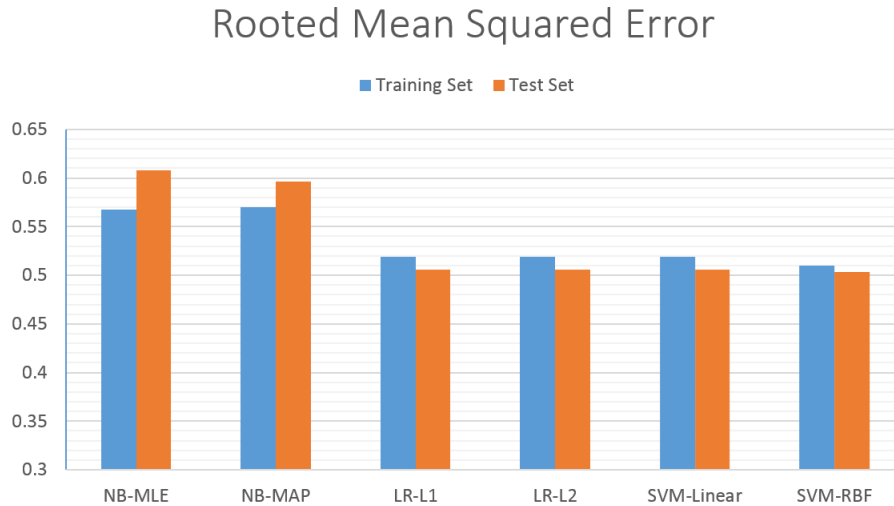


Figure 2: RMSE of different classifiers

Naive Bayes classifier can be done in linear time, but it generates very high RMSE. The features in this problem are not conditional independent, however, Naive Bayes has a strong requirement that each feature is conditional dependent with each other. When applying AdaBoost to Naive Bayes, the result is greatly improved, but still not high enough. So Naive Bayes may not be an ideal classifier for this problem.

For Logistic Regression and linear SVM, the time complexity is also relatively low. Compared to Naive Bayes, the classification result of Logistic Regression is much better. However, the accuracy of our LR and linear SVM classifier is still not high.

For non-linear SVM, the time complexity is very high. The result is similar to the result of linear SVM. Expanding our features to a higher dimension can still not improve the performance.

6 Next Step

On the next step, we will focus more on feature extraction. The best result we get for now will be chosen as the new benchmark in the next step. Our target will still be to improve the predicting accuracy as well as decreasing the RSME.

Specifically, here are the steps that we plan to try:

- Use larger dataset to train each classifier.
- Assemble more features not limited to the given columns to increase feature diversity.
- Add student temporal learning information to the feature vector. For example, for each step, add step name and Knowledge Component from the previous few steps. In addition, we can also record whether a student has met a problem or knowledge component before.

References

- [1] Pardos Z A, Gowda S M, Baker R S J, et al. The sum is greater than the parts: ensembling models of student knowledge in educational software[J]. ACM SIGKDD explorations newsletter, 2012, 13(2): 37-44.
- [2] Yu H F, Lo H Y, Hsieh H P, et al. Feature Engineering and Classifier Ensemble for KDD Cup 2010[J].
- [3] Zach A. Pardos, neil T.Heffernan. Using HMMs and bagged decision trees to leverage rich features of user and schoolroom an intelligent tutoring system dataset.
- [4] Andreas Töschel, Michael Jahrer. Collaborative Filtering Applied to Educational Data Mining.
- [5] Fan, Rong-En, et al. "LIBLINEAR: A library for large linear classification." The Journal of Machine Learning Research 9 (2008): 1871-1874.
- [6] Chang, Chih-Chung, and Chih-Jen Lin. "LIBSVM: a library for support vector machines." ACM Transactions on Intelligent Systems and Technology (TIST) 2.3 (2011): 27.