# Bagging, Boosting, and Stacking of Ensemble Methods for Machine Learning

Minghe Ren
Boston University
sawyermh@bu.edu

Sijie Xiang
Boston University
xiangs18@bu.edu

Xinqiao Wei
Boston University
weixq95@bu.edu

## Abstract

*Ensemble learning is a supervised machine learning method which combines multiple machine learning methods as base learners to obtain prediction rather than a single method to do the prediction. There are several algorithms in ensemble learning: bagging, boosting, and stacking. The goals of them are to decrease the variance and bias to improve predictive force [1].Ensemble methods have already showed the improvement of accuracy for certain basic classifiers for both artificial dataset and real world dataset [2]. The aims of our study are to explore that how the ensemble methods could improve the accuracy of classification by combining basic classifier and learn its characteristics.*

## 1. Introduction

In class, we have already explored some algorithms such as Quadratic Discriminant Analysis, K-Nearest Neighbors, Linear Regression, Naive Bayes, SVM and etc.. Each classifier would have different performance on the same dataset. Is it possible to combine them go get a better performance? The ensemble learning method is the solution. Hence, our goal is to explore the most popular ensemble learning methods, learn their characteristics and evaluate their performance against base learning algorithm. We focused on three ensemble learning methods: bagging, boosting, and stacking. All of them are also called "meta algorithm".

For bagging, we randomly select samples from the training set to generate sub-training sets with equally weights for all samples Then we train each sub-training set and test the classifiers on the test set and generate the prediction by voting. The model predictive force would not increase as training set increase. For boosting, we would generate sub-training sets randomly selected with weights and train those datasets. The sample weights would be updating by mis-classification rate from previous classification. We are focusing on one of boosting algorithms, called Adaboost. For stacking, multiple base learning algorithms or models are used to build the final learning algorithm. The outputs of base learners will be combined with original data

as new data, which would become a new dataset for next level learning algorithm.

In order to show their performances, firstly, we use MATLAB build-in ensemble learning functions and compare the performance between using ensemble learning algorithms and base learners. Then we wrote our own ensemble learning algorithms. In this section, we compared our algorithms' performances against MATLAB building functions. For base learners, we use Discriminant Analysis, Decision Tree and SVM.

We use two datasets, Adult Dataset [3] and Ecoli Dataset [4], for our study. Adult Dataset is for binary classification. Ecoli Dataset is for multi-class classification.

## 2. Literature Review

In class, all the algorithms we learned have different properties. Some of those algorithm perform classification on certain assumption. Discriminant Analysis and Linear Regression are both based on Gaussian assumption and Discriminant Analysis have more constraints on the assumption [5].

Decision tree Algorithm is using a predictive model, which is a tree to go from observations about an item to conclusions about the item's target value (represented in the leaves). Tree models where the target variable can take a discrete set of values are called classification trees; in these tree structures, leaves represent class labels and branches represent conjunctions of features that lead to those class labels. Decision trees where the target variable can take continuous values (typically real numbers) are called regression trees [6].

Support vector machine would generate an optimal hyperplane as output and use this hyperplane to classify.

Ensemble Learning refers to learning a weighted combination of base models of the form

$$f(y|x, \pi) = \sum_{m \in M} w_m f_m(y|x) \qquad (1)$$

where $w_m$ are tunable parameters. Ensemble learning can also be called committee method, because each base classifier $f_m$ gets a weighted "vote" [7].

## 3. Dataset

We have use two datasets. The first one is for multi-class classification and the second one is for binary classification.

### 3.0.1 Ecoli Dataset

This is a relatively small dataset that aims to predict protein localization sites. A summary statistics was conducted to check missing values and uniqueness of all instances. While no missing attributes values were detected and first column was just accession number, we decided to delete the first column. The last column is protein class distribution labels that we were trying to classify. We converted this column from string to categorical variables(1-8) so that all attributes are numerical. All other columns are intact. Last, we set seed and split the entire dataset into 70% training and 30% testing so that the whole process is reproducible.

### 3.0.2 Adult Dataset

This dataset is to predict the income of people, whether is greater and equal to or less than 50k, based on 14 features. Some of the features are numeral and some are categorical. We transfer all the categorical features into numeral and transform the table to a matrix. This dataset has two classification labels which are strings. Hence we transfer them into 0 and 1, 0 for $< 50$ and 1 for $>= 50$. In this dataset, there are 2399 out of 32561 training examples contain missing features. We attempt to remove all those examples, since there are only 7.4% of examples. However, if we remove those examples, some class in features will get missing. So we replace the missing data by the previous non missing data in the same feature for both training and testing sets. Due the the dataset is in random order itself, the replacement for missing data would also be random.

## 4. Bagging

- Bagging = Boostrap AGGregatING

- $B$ is the number of "bags" or base hypotheses

- L is the base learning algorithm

### 4.1. Problem Formulation and Solution Approaches

#### 4.1.1 Math Behind Bagging

Let $f(x)$ be the target value for $x$.
　Let $h_1, ..., h_B$ be the base hypotheses.
　Let $h_{avg}$ be the average prediction of $h_1, ..., h_B$.
　Let $e(h, x) = (f(x) - h(x))^2$

$$e(h_{avg}, x) = \frac{\sum_{i=1}^{B} e(h_i, x)}{B} - \frac{\sum_{i=1}^{B}(h_i(x) - h_{avg}(x))^2}{B} \tag{2}$$

### 4.1.2 Bagging Algorithm

---
**Algorithm 1** Bagging
---
1: **for** $i \leftarrow 1 : B$ **do**
2:　　$examples_i \leftarrow$ a bootstrap sample of examples
3:　　$h_i \leftarrow$ apply $L$ to $examples_i$
4: **end for**
5: **return** $h_1, h_2, ..., h_b$
---

In this algorithm, we randomly selected data points from original dataset with replacement into small datasets. And then we train each of the base learners by using the small datasets. In the end, the final predicted labels are decided by vote algorithm among all the base classifiers.

### 4.1.3 Intuition and Key Ideas

The key ideas behind the bagging is not complicated. A bootstrap sample is formed by sampling with replacement. Now idea is to use every bootstrap sample for training our model and we note down the variances of these model outputs. Now if we aggregate all these variances, we will get the aggregated model output which as very less variance. In other words If we take the variance of the average of these outputs, it will be lower than the average of the variances by a factor of $1/N$[8]. The mathematical intuition is as the following:

$$\begin{aligned} var(\bar{X}) &= var(\frac{1}{N}\sum_i X_i) \\ &= \frac{1}{N^2}(\sum_i var(X_i)) \\ &= \frac{1}{N}\frac{\sum_i var(X_i)}{N} \end{aligned} \tag{3}$$

### 4.1.4 Over-Fitting Issues

It shows in figure1. Please note that Bagging is used typically when we want to reduce the variance while retaining the bias. Overfitting occurs when the model has high variance and low bias and vice versa for underfitting. Again, We need to understand that Bagging decreases the variance, while boosting decreases the bias. Therefore, bagging is less vulnerable to overfitting than that of boosting. However, overfitting did occur when we implemented bagging. This can be shown in Ecoli dataset as test CCR of Decision Tree dropped almost 7% compared to the train CCR while CCR level remained approximately the same in SVM and Discriminant Analysis.

### 4.2. Implementation
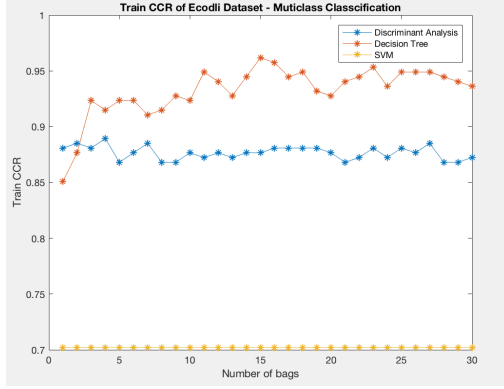
$$function[test\_set\_predicted\_Y] =$$
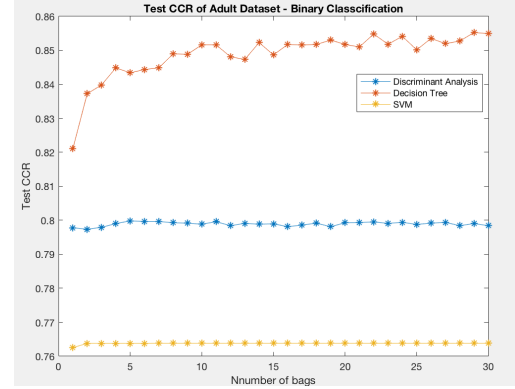
Figure 1: Train CCR of Ecoli Dataset



Figure 2: Bagging for Adult dataset

$$Ensemble\_bagging\_modified($$
$$train\_set\_X, train\_set\_Y,$$
$$test\_set\_X, n\_prime\_percent,$$
$$class\_total, iterations, algotype)$$

The implementation of Bagging is very straightforward. Inputs $train\_set\_X$, $train\_set\_Y$ and $test\_set\_X$ are mostly given after partitioned the original dataset into 70% training and 30% testing. $Class\_total$ is used to help me decide the type of classification I am dealing with, whether is binary classification or multiclass classification, in which OVO is utilized for SVM multiclass classification purpose. Algotype is used to indicate the type of weak learner I can choose from$((1 --> QDA, 2 --> DT, 3 --> SVM)$. The key inputs here are $n\_prime\_percent$ and iterations. The former helps me to control how many train data points I can sample from the training set with replacement while the latter suggests how many bags for bootstrap samples. For each weak learner, I trained B classifiers that are later used to classify testing set. Since our datasets are categorical, our final predicted labels are decided by taking majority vote. Moreover, for three learners, test CCR is calculated and plotted for each bag.

### 4.3. Experimental Results

It shows in figure 2 and 3. Even though bagging is less vulnerable to overfitting than that of boosting, overfitting did occur in our datasets as I had demonstrated in the sections above. My own implementation version of Bagging test CCR are comparable to Base algorithm and Bagging build-in function in Matlab in terms of Discriminant Analysis and Decision Tree while SVM is doing a slightly worse than base algorithm.

### 5. Adaboost

Boosting is a machine learning algorithm where trying to decrease the bias of data. Each of the base learning
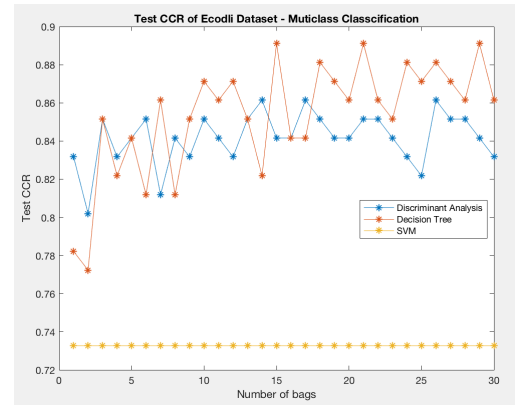


Figure 3: Bagging for Ecoli dataset

method used in boosting is called a weak learner. The performance of one weak learner might be unsatisfied. However, Combining these weak learners together to obtain a strong learner is the key idea of boosting. Adaboost is short for Adaptive Boosting and its one of the most popular method in Boosting.

The final classifier is combined by a weighted sum of sequentially trained weak classifier. During each iteration of the training process, a weight has been assigned to the sample which is misclassified by previous classifier. So the sample that has been wrongly classified has higher weight than correctly classified sample and is more likely to be considered in next training iteration.

There are two ways of using weak learners:

- Using single weak learner multiple times

- Using different weak learners in each iteration

## 5.1. Problem Formulation and Solution Approaches

### 5.1.1 Math Behind Adaboost

We have to consider a binary classification problem with exponential loss:

$$L_m(\phi) = \sum_{i=1}^{N} w_{i,m} exp(-\beta \tilde{y}_i \phi(x_i)) \qquad (4)$$

where $w_{i,m} \triangleq exp(-\tilde{y}_i f_{m-1}(x_i))$ is the weight applied to sample $i$, and $\tilde{y}_i \in \{-1, +1\}$. We can rearrange the equation to

$$L_m = e^{-\beta} \sum_{\tilde{y}_i = \phi(x_i)} w_{i,m} + e^{\beta} \sum_{\tilde{y}_i \neq \phi(x_i)} w_{i,m}$$
$$= (e^{\beta} - e^{-\beta}) \sum_{i=1}^{N} w_{i,m} \mathbb{I}(\tilde{y}_i \neq \phi(x_i)) + e^{-\beta} \sum_{i=1}^{N} w_{i,m} \qquad (5)$$

Hence, we need to optimal the function:

$$\phi_m = argmin_\phi w_{i,m} \mathbb{I}(\tilde{y}_i \neq \phi(x_i)) \qquad (6)$$

This could be solved by using weak learner on a weighted dataset with weights $w_{i,m}$. We could find $\beta$ by substituting $\phi_m$ into $L_m$:

$$\beta_m = \frac{1}{2} log \frac{1 - err_m}{err_m} \qquad (7)$$

where $err_m$ can be found by:

$$err_m = \frac{\sum_{i=1}^{N} w_i \mathbb{I}(\tilde{y}_i \neq \phi_m(x_i))}{\sum_{i=1}^{N} w_{i,m}} \qquad (8)$$

The overall prediction would become:

$$f_m(x) = f_{m-1}(x) + \beta_m \phi_m(x) \qquad (9)$$

### 5.1.2 Adaboosting Algorithm

The pseudo code for Adaboost.M1, for binary classification with exponential loss, is shown in Algorithm 2.

Notation: $w_i$ is the weight for each sample points. $N$ is the size of training data. $\phi_m(x)$ is the prediction from the weak learner. $err_m$ is the mis-classification rate by equation 8. $\beta_m$ is the log loss computed by negative logit function in equation 7. Weights for next iteration, in line 7, is calculated by the exponential loss function calculated by exponential loss for each sample by equation 4. The final classifier is combined by a weighted sum $\beta$ of sequentially trained weak classifier by equation 9.

---

**Algorithm 2** Adaboost.M1

1: *Initial:* $w_i = 1/N$
2: **for** $m = 1 : M$ **do**
3:     Select samples with weights $w$ with replacement;
4:     Find weak learner $\phi_m(x)$ on the chosen samples;
5:     Computer $err_m = \frac{\sum_{i=1}^{N} w_i \mathbb{I}(\tilde{y}_i \neq \phi_m(x_i))}{\sum_{i=1}^{N} w_{i,m}}$;
6:     Compute $\beta_m = \frac{1}{2} log \frac{1 - err_m}{err_m}$;
7:     Update $w_i \leftarrow w_i e^{-\tilde{y}_i \beta_m \phi_m(x_i)}$;
8: **end for**
9: **return** $f(x) = sign[\sum_{m=1}^{M} \beta_m \phi_m(x)]$

---

### 5.1.3 Intuition and Key Ideas

The key idea for Adaboosting is to increase the predictive force by combining weak learners to form a strong learner. A single weak learner may not get a ideal result. However, when the weak learners are combining with linear combination, the bias would decrease and such that the predictive force would increase.

### 5.1.4 Over-fitting Issues

As the iteration increases, there shows overfitting problem for each classifier on the temporally dataset. However, our final classification is the combination of classifiers and classifications, overfitting would not imply a large effect in our final classification. The accuracy would decrease a little bit when overfitting occurred.

### 5.1.5 Time Complexities

If we assume the time complexity of a single weak learner is $O(N)$, since different single weak learners have different time complexity. Then the Adaboosting would have the time complexity as $O(MN)$, where $M$ is for the number of iterations. As the number of iterations increase, the time complexity would also increases linearly. The idea for Adaboosting can be seem as trade-off between time complexity and predictive force.

## 5.2. Implementation

In our practice, we use a single weak learner multiply times. Since the weights are increasing for those samples, high probability to be miss-classified, those samples would be more likely to be in more classifier. Before applying the algorithm on the dataset, we partition the training dataset randomly into 70% as new training set and 30% as testing set. Then we start to implement the algorithm.

In line 1 of the algorithm 2, we initial the weight for all samples with $1/N$ to ensure that during first iteration, every sample have the equal probability to be selected. Then we starts to do iterations. We have to generate a temporally
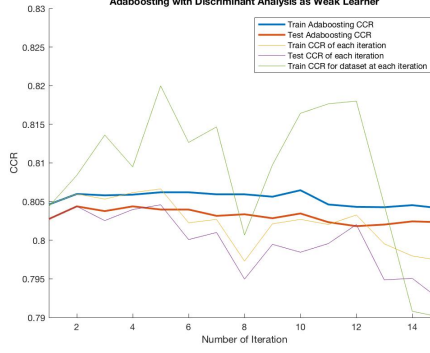
Figure 4: Discriminant Analysis as weak learners for 15 iterations



Figure 5: Discriminant Analysis as weak learners for 100 iterations



Figure 6: SVM and Decision tree as weak learners for 20 iterations

training set for each iteration base on the weights. For selecting data, we use selecting with replacement method, so for each selection, the probability for each sample would not be changed. After generating the temporally dataset, we train it with a weak classifier and test this temporally dataset with this classifier. We could get the error rate for this classifier with equation 7. Then we calculate the classifier weight with negative logit function in line 6 and we update the sample weights with the exponential loss of this classifier. For each iteration, we calculate the CCR for the temporally dataset, the original dataset, and the original test set. Then we do this iteration for $M$ times. The final classification results are the linear combination of the each classifier weights and the classification results from each classifier. We also calculate the CCR for the final classification after each iteration and try to find how many iteration it takes to have the best final classification performance. We use three weak learners, SVM, Linear Discriminant Analysis, and Decision Tree. For all three weak learners, we do 15 iterations for each of them and do 100 iterations for Discriminant Analysis to see how the predict force would change for large number of iteration. For all the weak learners, we used the MATLAB built-in algorithm with default setting for training and testing. We have only implement the algorithm for binary classification, so we have only use the adult dataset for adaboost.

### 5.3. Experimental Results

The Correctly Classified Rate(CCR) for Linear Discriminant Analysis for 15 iterations is the figure 4

The red and blue line represent the **final Classification CCR** of original training set and testing set respectively for each iteration. The green line represents the CCR of the temporally training set for each iteration. The yellow and purple line represent the CCR of original training set and testing set respectively for each iteration.

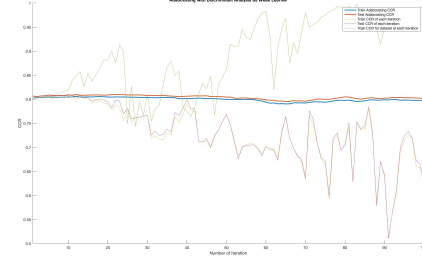Due to the dataset we are using, the CCR for the first

iteration have already reach about $80\%$, which is close to be converge. This is not really a weak learner we expected. So we can only see a slightly increment of CCR from 1st iteration to the 10th iteration.

The figure that shows the CCR of Discriminant Analysis for 100 iterations is in figure 5 In this figure, it's more clear to show that as the CCR for the temporally training set present a increasing trend and the CCR of original training set and testing set show a decreasing trend. This is a overfitting problem. However, the final classification CCR does not get a large effect by the overfitting.

In figure 6, it is more clear the how quickly for CCR of temporally training set to be overfitting. The CCR on the temporally training set is about to be $99\%$ as iteration increases. For SVM as a weak learner, the CCR of final classification CCR clearly increases about $2.5\%$ for testing set but decrease about $2.0\%$ for original training set. For decision tree as weak learner, it shows a similar outcome as SVM as weak learner, although decision tree as weak learner appears overfitting in less iteration.

For each weak learner, the best final classification CCR reached at different iterations. For LDA as weak learner, the best final classification CCR is about at 7 iteration. For SVM as weak learner, the best final classification CCR is about at 11 iteration. For decision tree as weak learner, the best final classification CCR is about at 9 iteration.
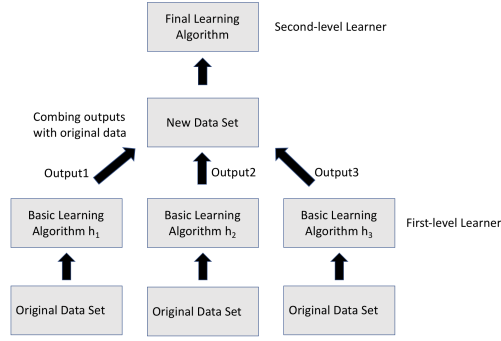
5

Figure 7: Stacking Structure



Figure 8: Stacking Pseudo Code[9]

# 6. Stacking

In stacking (also called stacked generalization), we use multiple classification or regression models combined to build the final learning algorithm. The structure of stacking can be seen in figure 7. Each of the based learners are obtained by training with the original data. What makes the stacking method unique from the other ensemble learning methods is that the outputs of base learner (also called first level learner) will be considered as new features trained at next level. The new generated data will be used as input by the final learner (also called second level learner) and seek to obtain better performance. Usually, Logistic Regression are used as the final learning model.

## 6.1. Problem Formulation and Solution Approaches

### 6.1.1 Math Behind Stacking

To solve the general ensemble learning problem as show in equation 1, the weights can be represented as equation 10:

$$\hat{w} = argmin_w \sum_{i=1}^{N} L(y_i, \sum_{m=1}^{M} w_m f_m(x)) \qquad (10)$$

This will cause overfitting problem as shown in Boosting where the weights become larger and larger. Cross-validation is a popular approach to estimating how well the result learned from a given training data set using unseen new data. We can use method to solve this problem.

$$\hat{w} = argmin_w \sum_{i=1}^{N} L(y_i, \sum_{m=1}^{M} w_m \hat{f}_m^{-i}(x)) \qquad (11)$$

Where $\hat{f}_m^{-i}(x)$ is the classifier trained by data excluding $(xi, yi)$

### 6.1.2 Stacking Algorithm

The algorithm of stacking can be shown in figure 8.
First, it makes $m$ partition of original dataset. The whole

original data has been used to train a different first-level classifier $h(t)$ using different algorithms. And next, we use $h(t)$ to classify the training example $x_i$ and use the outputs of each classifier to build our new train data set. Last, we train a classifier using the new data and second-level algorithm.

### 6.1.3 Intuition and Key Ideas

Generally, the goal of stacking is trying to observe which one of the base learner behaves well or badly and find a proper way as combiner to package these base learners together to build a classifier with better performance. The way of generating new data is that we take all the outputs of base learner, which are all labels. In this way, those individual learners are combined by a second-level learner which is called as meta-learner. Stacking can show more robust performance in the case where the true machine learning model is not in the base model family.

### 6.1.4 Over-Fitting Issues

As mention before, the stacking method is meant to solve the overfitting problem by using cross-validation method and second-level learner uses a new dataset generated by using the outputs of base learner. So, the stacking is trying to reduce the overfitting issue instead of adding more.

## 6.2. Implementation

In our experiment, we chose 5 different algorithms as base learners and we used both data set to conduct the experiment. In Adult dataset for binary classification we used QDA, Knn, Naive Bayes, Decision Tree, SVM as base learners, Logistic Regression and Decision Tree as second-level learner. In Ecoli dataset for multi-class classification, we used QDA, Knn, Decision Tree, SVM as base learners, Decision Tree as second-level learner for comparison. Weve used Matlab built-in functions (fitcdiscr, fitcknn, fitcnb,fitctree,svmtrain) to produce results of base learner as benchmark.
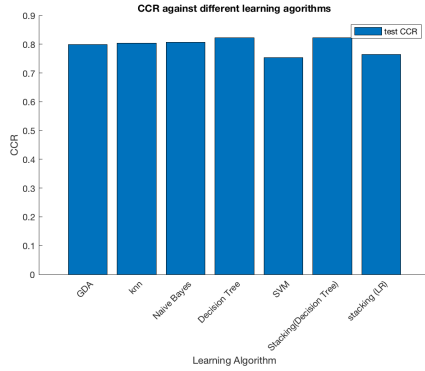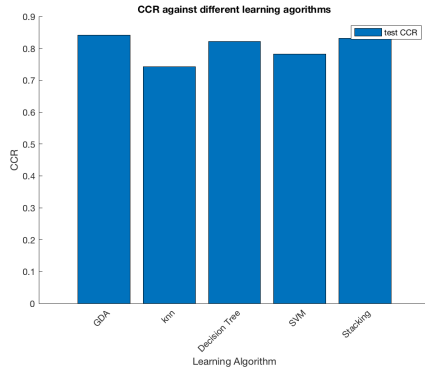
Figure 9: Test CCR for Adult Dataset



Figure 10: Test CCR of Ecoli Dataset

### Adult Dataset

| CCR | Base algorithm | Bagging | Bagging (Built-in) | Ada-boosting | Ada-boosting (Built in) | Stacking |
|---|---|---|---|---|---|---|
| Discriminant Analysis | 78.28% | 79.84% | 79.86% | 80.22% | 79.88% | Decision Tree: 81.5% |
| Decision Tree | 78.34% | 85.49% | 85.91% | 80.18% | 85.63% | Logistic Regression: 77.0% |
| SVM | 76.90% | 76.38% | Nah | 80.11% | Nah | |

Figure 11: Adult dataset

### Ecoli Dataset

| CCR | Base algorithm | Bagging | Bagging (Built-in) | Stacking |
|---|---|---|---|---|
| Discriminant Analysis | 84.16% | 83.42% | 84.16% | Decision Tree: 83.19% |
| Decision Tree | 82.18% | 86.13% | 89.11% | |
| SVM | 78.22% | 73.80% | Nah | |

Figure 12: Ecoli dataset

## 7. Conclusion

Here is the conclusion for all three ensemble methods:

| | Bagging | Adaboost | Stacking |
|---|---|---|---|
| Partitioning Dataset | Random | Weights | None |
| Goal Achieved | Minimize variance | 1Increasing predictive force | Both |
| Combining Classifiers | Average weights | Weights majority vote | Logistic regression / Decision Tree |

The comparison of performance for MATLAB built-in function for base learners and ensemble methods are in figure 11 and figure 12. In those figures, we can see that the performance of ensemble methods are all better than a single base algorithm. For some base learner, the algorithms we implemented has a better performance than the built in functions. Since the datasets, we chosen, are well construct, we could not get the obvious increasing trend for CCR.

In the future, we could try some dataset does not have a good performance for a single base algorithm. We also can try to implement the Adaboost for multi-class classification.

## 8. Individual Effort

Data Preprocessing: Sijie Xiang, Xinqiao Wei.
Coding:
Bagging: Sijie Xiang.

### 6.3. Experimental Results

The performances of base learners and stacking are shown in figure 9 and figure 10 testing both on two dataset. X coordinate represents the name of learning algorithm. In Adult dataset, the stacking with DT as final combiner and using DT alone have the highest test CCR $82.20\%$. In Ecoli dataset, Stacking with DT as final combiner has the second highest test CCR $83.19\%$ and the highest CCR $84.16\%$ is using GDA alone without stacking method. Both datasets show that the stacking method has better performance than most of the single learner alone. The reason why the stacking method shows a little bit CCR improvement is that the base learners are already strong enough. We used the outputs of each base classifier equally without adding weights to build new dataset for second-level model which might cause the reason why sometimes the performance of stacking gets lower than one of the base learner. Besides that, the choice of second-level learner also plays an important role on the performance.

Adaboost: Xinqiao Wei, Minhe Ren.

Stacking: Minhe Ren.

Report:

Xinqiao Wei: Introduction, Literature Review, Dataset, Adaboost, Conclusion.

Minhe Ren: Abstract, Stacking.

Sijie Xiang: Bagging.

**Everyone is trying very hard for this project.**

## References

[1] V. Smolyakov, "Ensemble learning to improve machine learning results," Aug. 2017.

[2] E. BAUER and R. KOHAVI, "An empirical comparison of voting classification algorithms: Bagging, boosting, and variants," Sept. 1998.

[3] R. Kohavi and B. Becker, "Adult data set," 1996.

[4] K. Nakai, "Ecoli data set," 1996.

[5] A. Liu, "Stat 697f st - topics in regression (umass)," 2015.

[6] O. Rokach, Lior; Maimon, "Data mining with decision trees: theory and applications.," 2008.

[7] K. P. Murphy, "Machine learning: A probabilistic perspective," 2012.

[8] A. Joshi, "What's the difference between boosting and bagging?."

[9] Z.-H. Zhou, "Ensemble learning."