

Homework 4

CSCI 5525 S20: Machine Learning

Due on Apr. 15th 11:59 pm (Midnight)

Homework Policy. (1) You are encouraged to collaborate with your classmates on homework problems, but each person must write up the final solutions individually. You need to fill in above to specify which problems were a collaborative effort and with whom. (2) Regarding online resources, you should **not**:

- Google around for solutions to homework problems,
- Ask for help on online.
- Look up things/post on sites like Quora, StackExchange, etc.

Submission. Submit a PDF using this LaTeX template for written assignment part and submit .py files for all programming part. You should upload all the files on Canvas.

Written Assignment

Instruction. For each problem, you are required to write down a full mathematical proof to establish the claim.

Problem 1. Boosting Derivation.(Total 5 points)

We discussed in class that AdaBoost minimizes the exponential loss greedily. In particular, the derivation of α_t is by finding the minimizer of

$$\epsilon_t(\exp(\alpha_t) - \exp(-\alpha_t)) + \exp((- \alpha_t))$$

where ϵ_t is the weighted classification error of h_t . Show that $\alpha_t = \frac{1}{2} \ln(\frac{1-\epsilon_t}{\epsilon_t})$ is the minimizer.

In this question, let

$$l = \epsilon_t(\exp(\alpha_t) - \exp(-\alpha_t)) + \exp((- \alpha_t))$$

We take the derivation of $\frac{\partial l}{\partial \alpha_t} = 0$:

$$\frac{\partial l}{\partial \alpha_t} = \epsilon_t(\exp(\alpha_t) + \exp(-\alpha_t)) - \exp(-\alpha_t) = 0$$

We have:

$$\begin{aligned}\epsilon_t(\exp(\alpha_t) + \exp(-\alpha_t)) &= \exp(-\alpha_t) \\ \epsilon_t(\exp(2\alpha_t) + 1) &= 1 \\ \exp(2\alpha_t) + 1 &= \frac{1}{\epsilon_t} \\ \exp(2\alpha_t) &= \frac{1 - \epsilon_t}{\epsilon_t} \\ \alpha_t &= \frac{1}{2} \ln\left(\frac{1 - \epsilon_t}{\epsilon_t}\right)\end{aligned}$$

Problem 2. Decision Tree.(Total 5 points)

Consider a binary dataset with 400 examples, where half of them belongs to class A and another half belongs to class B.

Next consider two decision stumps (i.e. trees with depth 1) T1 and T2, each with two children. For T1, its left child has 150 examples in class A and 50 examples in class B; for T2, its left child has 0 example in class A and 100 examples in class B. (You can infer what are in the right child.)

- **(3 point)** For each leaf of T1 and T2, compute the corresponding classification error, entropy (base e) and Gini impurity. (Note: the value/prediction of each leaf is the majority class among all examples that belong to that leaf.)
- **(2 points)** Compare the quality of T1 and T2 (that is, the two different splits of the root) based on classification error, weighted entropy (base e), and weighted Gini impurity respectively.

(a) For T1:

the classification error for the left leaf is:

$$u(S) = \min\left\{\frac{150}{200}, \frac{50}{200}\right\} = \frac{1}{4}$$

the classification error for the right leaf is:

$$u(S) = \min\left\{\frac{50}{200}, \frac{150}{200}\right\} = \frac{1}{4}$$

the entropy for the left leaf is:

$$u(S) = p \log\left(\frac{1}{p}\right) + (1 - p) \log\left(\frac{1}{1 - p}\right) = 0.75 \log 0.75 + 0.25 \log 0.25 = 0.56233$$

the entropy for the right leaf is:

$$u(S) = p \log\left(\frac{1}{p}\right) + (1 - p) \log\left(\frac{1}{1 - p}\right) = 0.75 \log 0.75 + 0.25 \log 0.25 = 0.56233$$

the Gini impurity for the left leaf is:

$$u(S) = 2p(1 - p) = 2 * 0.75 * 0.25 = 0.375$$

the Gini impurity for the right leaf is:

$$u(S) = 2p(1 - p) = 2 * 0.75 * 0.25 = 0.375$$

For T2:

the classification error for the left leaf is:

$$u(S) = \min\{\frac{0}{100}, \frac{100}{100}\} = 0$$

the classification error for the right leaf is:

$$u(S) = \min\{\frac{200}{300}, \frac{100}{300}\} = \frac{1}{3}$$

the entropy for the left leaf is:

$$u(S) = p \log(\frac{1}{p}) + (1 - p) \log(\frac{1}{1 - p}) = 0 + 1 \log 1 = 0$$

the entropy for the right leaf is:

$$u(S) = p \log(\frac{1}{p}) + (1 - p) \log(\frac{1}{1 - p}) = \frac{1}{3} \log(\frac{1}{3}) + \frac{2}{3} \log(\frac{2}{3}) = 0.6365$$

the Gini impurity for the left leaf is:

$$u(S) = 2p(1 - p) = 2 * 0 * 1 = 0$$

the Gini impurity for the right leaf is:

$$u(S) = 2p(1 - p) = 2 * \frac{1}{3} * \frac{2}{3} = \frac{4}{9}$$

(b) The weighted classification error for T1 is:

$$\sum_l w(l)u(S_l) = \frac{1}{2} * \frac{1}{4} + \frac{1}{2} * \frac{1}{4} = \frac{1}{4}$$

The weighted entropy for T1 is:

$$\sum_l w(l)u(S_l) = 2 * \frac{1}{2} * 0.56233 = 0.56233$$

The weighted Gini impurity for T1 is:

$$\sum_l w(l)u(S_l) = 2 * \frac{1}{2} * 0.375 = 0.375$$

The weighted classification error for T2 is:

$$\sum_l w(l)u(S_l) = \frac{1}{3} * \frac{3}{4} + 0 = \frac{1}{4}$$

The weighted entropy for T2 is:

$$\sum_l w(l)u(S_l) = \frac{3}{4} * 0.6365 + 0 = 0.477$$

The weighted Gini impurity for T2 is:

$$\sum_l w(l)u(S_l) = \frac{3}{4} * \frac{4}{9} + 0 = \frac{1}{3}$$

We can see that the classification errors for T1 and T2 are the same, but when it comes to entropy or Gini impurity, $T2 < T1$. So we can say T2 has a better split than T1.

Problem 3. Game Theory for Boosting

(7 points) Boosting explored from a game-theoretic perspective.

Some Background: A two-player zero-sum game is specified by matrix M and the two players are denoted the row player and the column player. The game is played by the row player choosing a row i and the column player choosing a column j , and the *loss* for the row player is $M(i, j)$ which is also the *reward* for the column player. Instead of choosing individual rows/columns, we will allow both players to choose distributions over rows/columns, so if row player choose P and column player choose Q , the loss/reward is

$$\sum_i \sum_j P(i) M(i, j) Q(j) \triangleq M(P, Q)$$

If we are acting as row player, we would like to choose a distribution P that achieves low loss, no matter what column player does. In other words, we would like to choose P that minimizes $\max_Q M(P, Q)$. On the other hand, if we were column player, we would like to choose Q that maximizes $\min_P M(P, Q)$. Von Neumann's celebrated minimax theorem states that in fact both these values are same, or in some sense it does not matter which player goes first in the game:

$$\max_Q \min_P M(P, Q) = \min_P \max_Q M(P, Q) \triangleq V$$

where V is called the value of the game. In this problem, we will use boosting to compute the optimal strategy in a particular game.

Let \mathcal{H} be finite, and fix a target concept $c : \mathcal{X} \rightarrow \{+1, -1\}$ (not necessarily in \mathcal{H}) and sample $S = \{X_i, c(X_i)\}_{i=1}^n$ of size n . (The concept c is just a formulation on how the labels are generated.)

We will form a matrix $M \in \{0, 1\}^{n \times |\mathcal{H}|}$ where $M(i, h) = \mathbb{1}\{h(X_i) = c(X_i)\}$. Here, the row player specifies distributions over samples, just as in boosting, and the column player chooses distributions over hypotheses.

- a) (4 points) Assume that the empirical γ -weak learning assumption holds so that for every distribution P over examples, there exists a hypothesis $h \in \mathcal{H}$ such that $\mathbb{P}_{x \sim P}[h(x) \neq c(x)] \leq 1/2 - \gamma$. What does this mean about the value of the game? (Hint: First, think about what $M(P, Q)$ means. Then, find out what the relationship between V and γ is.)
- b) (3 points) Let Q^* be distribution achieving value of this game, i.e. $\min_P M(P, Q^*) = V$. Since Q^* is distribution over hypotheses, what is the empirical error over Q^* ? (Hint: consider $Q^*(h)$ as the weight α of weak learner h and then show what is the final predictor and what is the empirical error.)
- c) (Hurray!!! Extra credits: 10 points) Boosting can be viewed as an iterative algorithm to compute Q^* using a weak learner. At every round, we choose P_t , a distribution over samples, and then compute

$$Q_t = \max_Q M(P_t, Q)$$

which is actually a single hypothesis h_t due to linearity. Then we update P_t to be

$$P_{t+1}(x) = \frac{P_t(x)}{Z_t} \times (\exp(-\eta \mathbb{1}\{h_t(x) = c(x)\}))$$

After T rounds, we output $\bar{Q} = \frac{1}{T} \sum_{t=1}^T Q_t$ which is a distribution over hypothesis and the final predictor is $H(x) = \text{sign}(\bar{Q}(x))$.

Prove that once $T = \Omega(\log(n)/\gamma^2)$ rounds and for appropriate choice of η , this variant of boosting guarantees (**Hint:** you may find it helpful to look at the Taylor expansion of $\exp(-x)$.)

$$\forall x \in X, \frac{1}{T} \sum_{i=1}^T M(x, h_i) > 1/2,$$

which implies that $H(x)$ has zero training error.

- d) (**Let's have more!!! Extra credits: 5 points**) Informally, what does this mean about \bar{Q} , what is it converging to as $T \rightarrow \infty$?

Your answer. (a) The value of the game is trying to find a maxmin strategy of M , the column player makes the V close to 1 by picking different hypothesis which satisfy:

$$\forall P, \exists h_t, s.t. M(P, h_t) = \Pr_{x \sim P}[h_t(x) = c(x)] \geq \frac{1}{2} + \gamma$$

It's easy to find that:

$$V = \min_P \max_Q M(P, Q) \geq \min_P M(P, h_t) \geq \frac{1}{2} + \gamma$$

The game is a boosting algorithm.

- (b) The final predictor is:

$$\hat{h} = \text{sign}\left(\sum_t Q^*(h_t) h_t(x)\right)$$

So the empirical error should be:

$$\begin{aligned} E &= \\ &= \sum_x P(x) \mathbb{1}\{\text{sign}\left(\sum_t Q^*(h_t) h_t(x)\right) \neq c(x)\} \\ &= \Pr_{x \sim P}[\hat{h}(x) \neq c(x)] \\ &\leq \frac{1}{2} - \gamma \end{aligned}$$

In this case, if T is big enough and Q selects the best hypothesis, the empirical error will reach 0.

- (c) Let $t = T$, $T > 1$, we have:

$$\begin{aligned} P_{T+1}(x) &= \frac{P_T(x)}{Z_T} \times (\exp(-\eta \mathbb{1}\{h_T(x) = c(x)\})) \\ &= \frac{P_T(x)(\exp(-\eta M(x, h_T)))}{Z_T} \\ &= \frac{P_{T-1}(x)(\exp(-\eta(M(x, h_T) + M(x, h_{T-1}))))}{Z_T Z_{T-1}} \\ &= \frac{P_1(x)(\exp(-\eta \sum_{t=1}^T M(x, h_t)))}{\prod_{t=1}^T Z_t} \end{aligned}$$

Let $n = |x|$. In boosting problem, $P_1(x) = \frac{1}{n}$, and $\forall(T, x), P_T(x) \leq 1$, which means $\ln(P_T(x)) \leq 0$.

We have:

$$\begin{aligned}
P_{T+1}(x) &= \frac{P_1(x)(\exp(-\eta \sum_{t=1}^T M(x, h_t)))}{\prod_{t=1}^T Z_t} \\
&= \frac{\exp(-\eta \sum_{t=1}^T M(x, h_t))}{n \prod_{t=1}^T Z_t} \\
&\iff \\
\frac{1}{T} \sum_{t=1}^T M(x, h_t) &= -\frac{\ln(P_{T+1}(x)n \prod_{t=1}^T Z_t)}{\eta T} \\
&= -\frac{\ln(P_{T+1}(x))}{\eta T} - \frac{\ln n}{\eta T} - \frac{\sum_{t=1}^T \ln Z_t}{\eta T} \\
&\geq -\frac{\ln n}{\eta T} - \frac{\sum_{t=1}^T \ln Z_t}{\eta T}
\end{aligned}$$

We want to make Z_t larger. We notice that $M(x, h_t)$ can only be 0 or 1 and that is to say:

$$\exp(-\eta M(x, h_t)) = \begin{cases} \exp(-\eta), & M(x, h_t) = 1 \\ 1, & M(x, h_t) = 0 \end{cases}$$

So in this case:

$$\exp(-\eta M(x, h_t)) = \begin{cases} \exp(-\eta), & M(x, h_t) = 1 \\ 1, & M(x, h_t) = 0 \end{cases} = \exp(-\eta)M(x, h_t) + 1 - M(x, h_t)$$

Thus:

$$\begin{aligned}
Z_t &= \sum_x P_t(x)(\exp(-\eta M(x, h_t))) \\
&= \sum_x P_t(x)[\exp(-\eta)M(x, h_t) + 1 - M(x, h_t)] \\
&= \sum_x P_t(x) + (\exp(-\eta) - 1) \sum_x P_t(x)M(x, h_t) \\
&= 1 + (\exp(-\eta) - 1)M(P_t, h_t)
\end{aligned}$$

Besides, we know that: $\ln(1+x) < x$:

$$\begin{aligned}
\frac{1}{T} \sum_{t=1}^T M(x, h_t) &\geq -\frac{\ln n}{\eta T} - \frac{\sum_{t=1}^T \ln Z_t}{\eta T} \\
&\geq -\frac{\ln n}{\eta T} - \frac{\sum_{t=1}^T \ln(1 + (\exp(-\eta) - 1)M(P_t, h_t))}{\eta T} \\
&\geq -\frac{\ln n}{\eta T} - \frac{\sum_{t=1}^T (\exp(-\eta) - 1)M(P_t, h_t)}{\eta T} \\
&= -\frac{\ln n}{\eta T} - (\exp(-\eta) - 1) \frac{\sum_{t=1}^T M(P_t, h_t)}{\eta T}
\end{aligned}$$

Also, from the problem (a), we know that:

$$M(P_t, h_t) = \Pr_{x \sim P}[h_t(x) = c(x)] = \sum_x P_t(x) M(x, h_t) \geq \frac{1}{2} + \gamma$$

Let $\eta > 0$, so $\exp(-\eta) < 1$:

$$\begin{aligned}
\frac{1}{T} \sum_{t=1}^T M(x, h_t) &\geq -\frac{\ln n}{\eta T} - (\exp(-\eta) - 1) \frac{\sum_{t=1}^T M(P_t, h_t)}{\eta T} \\
&= -\frac{\ln n}{\eta T} + (1 - \exp(-\eta)) \frac{\sum_{t=1}^T M(P_t, h_t)}{\eta T} \\
&\geq -\frac{\ln n}{\eta T} + (1 - \exp(-\eta)) \frac{\sum_{t=1}^T (\frac{1}{2} + \gamma)}{\eta T} \\
&= -\frac{\ln n}{\eta T} + (1 - \exp(-\eta)) \frac{T(\frac{1}{2} + \gamma)}{\eta T} \\
&= -\frac{\ln n}{\eta T} + (1 - \exp(-\eta)) \frac{(\frac{1}{2} + \gamma)}{\eta}
\end{aligned}$$

From the above formula, if we want $\frac{1}{T} \sum_{t=1}^T M(x, h_t) \geq \frac{1}{2}$, we need to find:

$$\begin{aligned}
-\frac{\ln n}{\eta T} + (1 - \exp(-\eta)) \frac{(\frac{1}{2} + \gamma)}{\eta} &\geq \frac{1}{2} \\
&\Leftrightarrow \\
T &\geq \frac{\ln n}{(1 - \exp(-\eta))(\frac{1}{2} + \gamma) - \frac{\eta}{2}} \\
&= \frac{\ln n}{(1 - \exp(-\eta))\gamma + \frac{1}{2}(1 - \exp(-\eta) - \eta)}
\end{aligned}$$

For $1 - \exp(-\eta) - \eta$, it's obvious that $\forall \eta > 0, 1 - \exp(-\eta) - \eta < 0$, besides, when $\eta > 0, 1 -$

$\exp(-\eta) \in [0, 1]$, so we assume that $1 - \exp(-\eta) = \alpha\gamma, \alpha > 0$, then we have:

$$\begin{aligned}
T &\geq \frac{\ln n}{(1 - \exp(-\eta))\gamma + \frac{1}{2}(1 - \exp(-\eta) - \eta)} \\
&> \frac{\ln n}{(1 - \exp(-\eta))\gamma} \\
&= \frac{1}{\alpha} \frac{\ln n}{\gamma^2} \\
&= \Omega\left(\frac{\ln n}{\gamma^2}\right)
\end{aligned}$$

Thus, we get the proof that:

$$\begin{aligned}
&\text{When: } T > \Omega\left(\frac{\ln n}{\gamma^2}\right) \\
&\frac{1}{T} \sum_{t=1}^T M(x, h_t) > \frac{1}{2}
\end{aligned}$$

(d) \bar{Q} is the average of all the weak learners selected in the learning process. From the problem c, we know that when $T > \Omega\left(\frac{\ln n}{\gamma^2}\right)$, $\frac{1}{T} \sum_{t=1}^T M(x, h_t) > \frac{1}{2}$, which means that when T is big enough, there are always more than half of the hypothesis that are correct in predicting x . As $c : \mathcal{X} \rightarrow \{+1, -1\}$ is binary classification, when more than half of the hypothesis can correctly predict the x , we always have $H(x) = c(x)$. In other words, \bar{Q} converges to $c(x)$.

Programming Assignment

Instruction. For each problem, you are required to report descriptions and results in the PDF and submit code as python file (.py) (as per the question).

- **Python** version: Python 3.
- Please follow PEP 8 style of writing your Python code for better readability in case you are wondering how to name functions & variables, put comments and indent your code
- **Packages allowed:** numpy, pandas, matplotlib
- **Submission:** Submit report, description and explanation of results in the main PDF and code in .py files.
- Please **PROPERLY COMMENT** your code in order to have utmost readability otherwise **1 MARK** would be deducted

Programming Common

This programming assignment focuses on boosting and bagging approaches.

Problem 4. "Blind" boosting without training data.

Problem 4.1. Blind boosting for regression (9 Points)

Imagine we do not have the training dataset visible but a scoring function is made available to you. This score function gives us an idea about how well your model is doing on the test dataset. Higher score indicates worse performance. The only knowledge you have is that the predict target, $y \in \mathbb{R}$. Can we achieve better performance than a completely random guessing using **just** the knowledge of how well your model does on test dataset. If **yes**, then explain how? You have to implement this model. Show how the score changes as we increase the number of weak learners in the boosting model (show a plot).

(Hint: Please try to use the idea of gradient boosting and also try to randomly guess the gradient.)

Note: For getting the score, use the *score* function in **test_score.py** script by importing it in your **hw4_boosting_regression.py** file you would submit as solution for this question. The function *score* accepts only one parameter which is of type 'numpy.ndarray' and has shape of $(N,1)$ where N is the number of samples predicted passed as argument to this function. Here, we have 21283 samples in testing dataset therefore $N = 21283$.

Keep the **test_score.py** and **true_values_regression.csv** which contains the true values of test dataset in the same location where **hw4_boosting_regression.py** file would be. This *score* basically calculates how much error is made in predictions by comparing with true value in **true_values.csv** testing dataset.

Submission: Submit all plots requested and explanation in latex PDF. Submit your program in a file named (hw4_boosting_regression.py).

I'm confused about this problem so I realized gradient boosting(at least its idea) and gradient guessing(don't know loss and true value)

For the gradient boosting part, I calculate the derivation of the loss function to get the gradient and add the gradient to the final prediction with a learning rate. The gradient can be regarded as a prediction of some hypothesis.

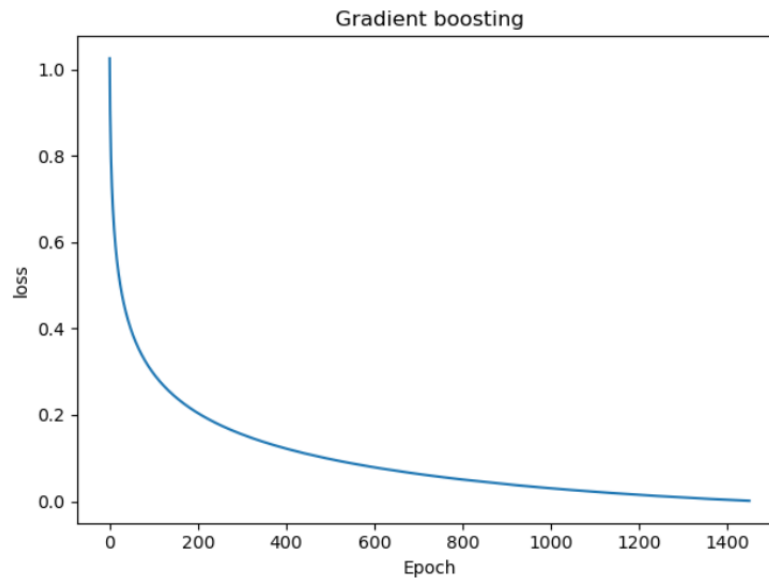


Figure 1: Gradient boosting

If the loss function is unknown, it's possible to guess the gradient. In the code, I initialize output F as an array in which all elements are 0 and guess the gradient for these elements one by one. After that, add the guessing gradient to F and check the score, if the score goes down, this guessing gradient will be kept as an h , otherwise, shrink the upper bound of the guessing gradient.

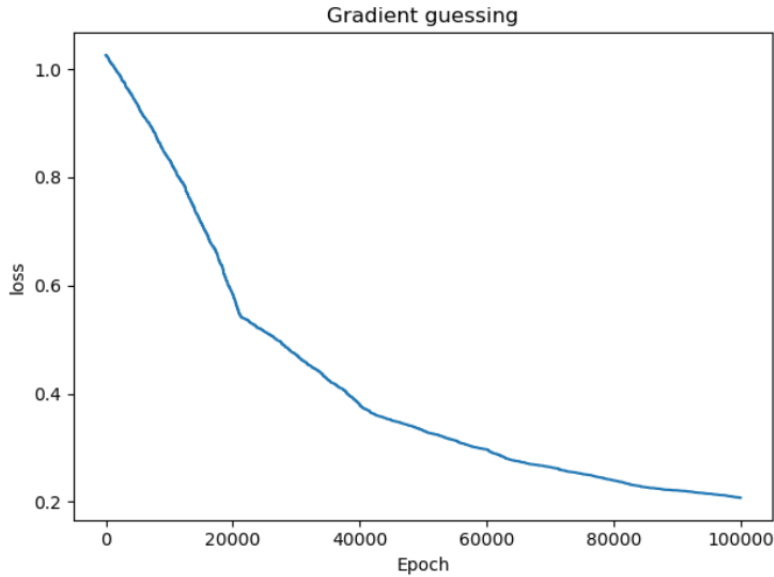


Figure 2: Gradient guessing

Problem 4.2. Blind boosting for classification (9 Points)

Imagine we do not have the training dataset visible but a error rate function is made available to you. This function gives us an idea about how accurate the model is on the test set. The only knowledge you have is that the predict target, $y \in \{0, 1\}$. Can we achieve better performance than a completely random guessing using **just** the knowledge of how well your model does on test dataset. If **yes**, then explain how? You have to implement this model. Show how the score changes as we increase the number of weak learners in the boosting model (show a plot).

Note: For getting the error rate, use the *error_rate* function in **test_error_rate.py** script by importing it in your **hw4_boosting_classification.py** file you would submit as solution for this question. The function *error_rate* accepts only one parameter which is of type 'numpy.ndarray' and has shape of $(N, 1)$ where N is the number of samples predicted passed as argument to this function. Here, we have 21283 samples in testing dataset therefore $N = 21283$.

Keep the **test_error_rate.py** and **true_values_classification.csv** which contains the true values of test dataset in the same location where **hw4_boosting_classification.py** file would be. This *score* basically calculates how much error is made in predictions by comparing with true value in **true_values.csv** testing dataset.

Submission: Submit all plots requested and explanation in latex PDF. Submit your program in a file named (hw4_boosting_classification.py).

In this part, I initialize F as a 0 array and check the error rate. With the error rate, the numbers of 1 and 0s in the true value are known. After that, I generate h which has the same number of 0 and 1 as the true value and shuffle h and check the error rate, usually the error

rates are around 50%. The boosting algorithm is defined below, err refers to error rate of h .

$$F = F + \sum_t \begin{cases} (1 - err_t) * h_t, err_t \leq 0.5 \\ err_t(1 - h_t), err_t > 0.5 \end{cases}$$

$$prediction = sign(F)$$

With this algorithm, we can always keep an h with an accuracy more than 50%. The loss graph is shown below:

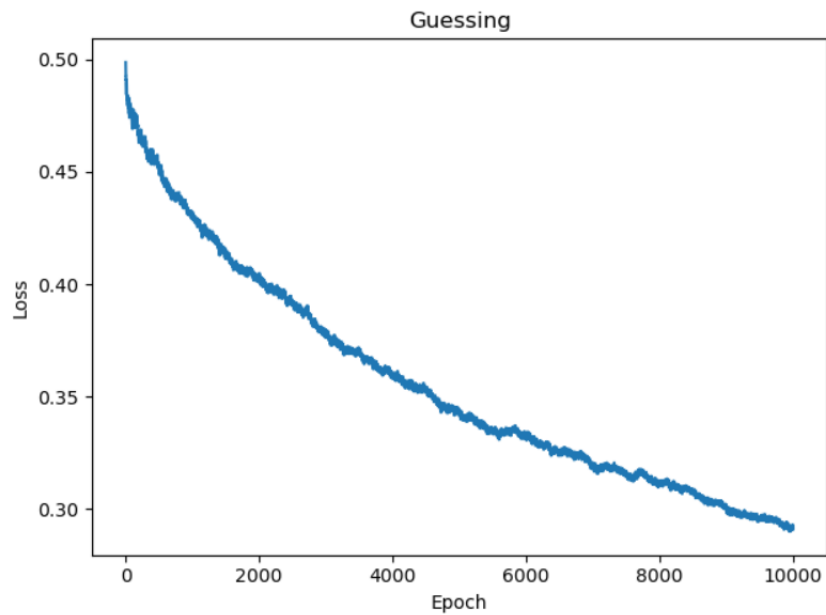


Figure 3: Classification guessing

Problem 5. Adaboost.

(Total 12 Points) For this problem, you will use a cancer dataset as in the file cancer_train.csv and cancer_test.csv. You can find this file in the canvas file folder named hw4. The goal is to predict whether there is a cancer: 0 or 1, included in column y.

- (7 Points) Implement an Adaboost algorithm with 100 weak learners. You are allowed to use decision tree classifier with maximum depth = 1 from scikit standard library (sklearn link). **No other function from scikit learn is allowed.** Use gini as quality measure for the decision tree splits and use the default parameters.
- (2 Points) Plot the misclassification error on both train and test set as the number of weak learners increase.
- (3 Points) plot the margin distribution when the number of weak learners = [25, 50, 75, 100]. In total you will show 4 plots.

Submission: Submit all plots requested and explanation in latex PDF. Submit your program in a file named (hw4_adaboost.py).

(b) The error goes to 0 quickly for training data set when number of the hypothesis increases, but the error goes up a little bit for the testing data set when the number of hypothesis is more than 60.

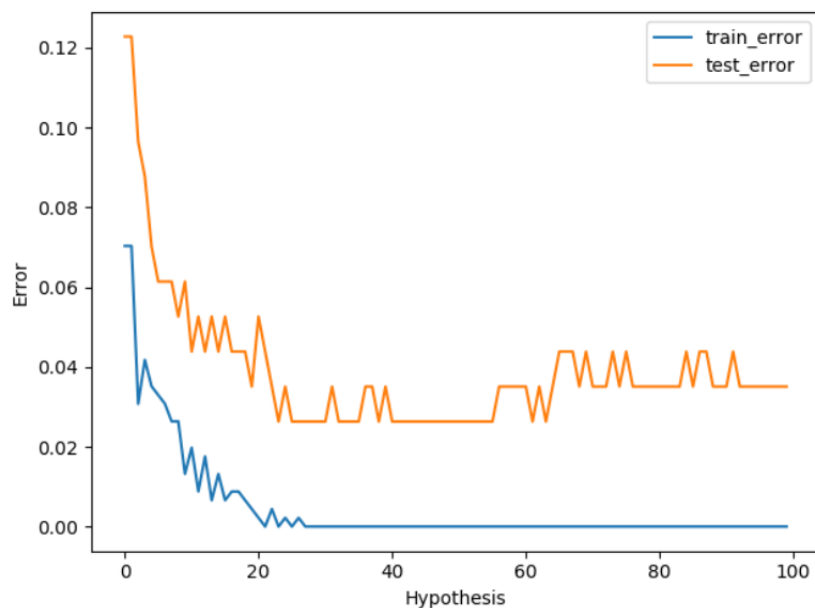


Figure 4: Test and train error

(c) I use the training dataset to plot the margin cumulative distribution, we can see as the number of weaker learner increase, the line is closer to 1.

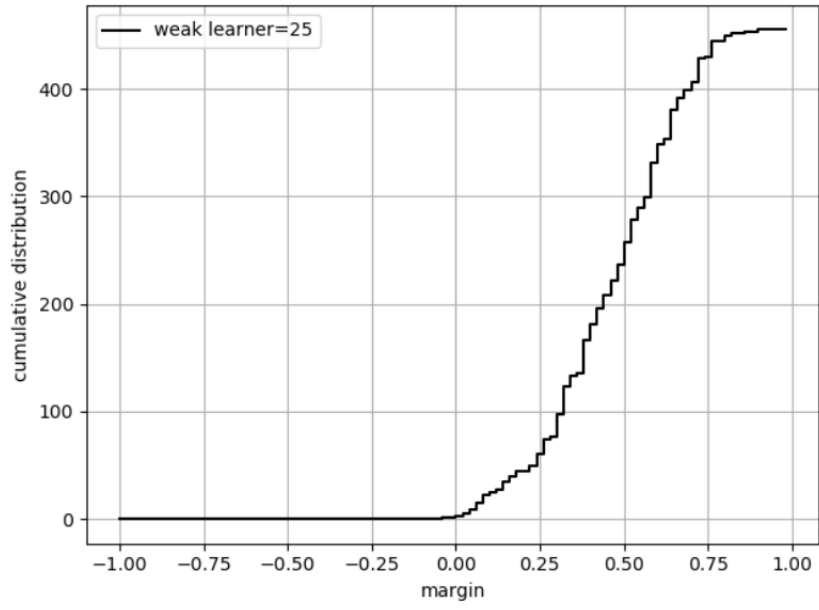


Figure 5: Weak learner=25

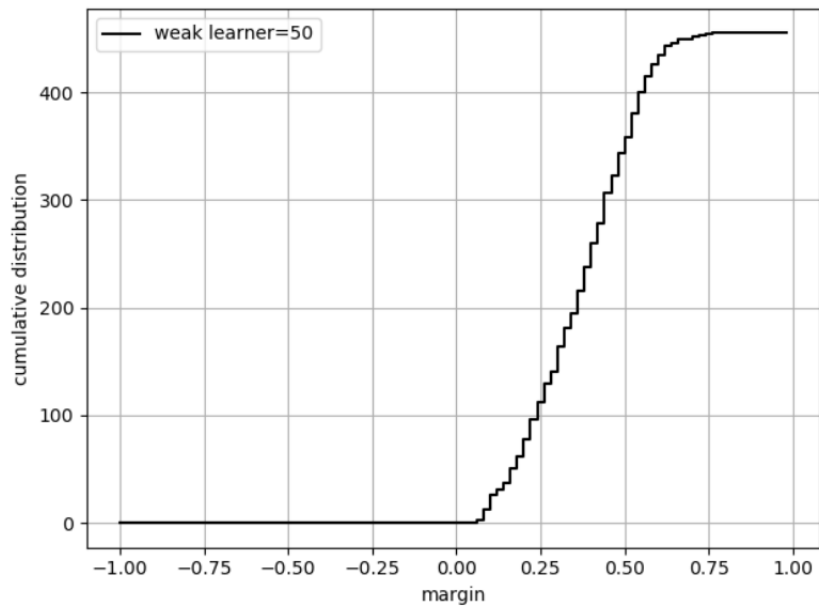


Figure 6: Weak learner=50

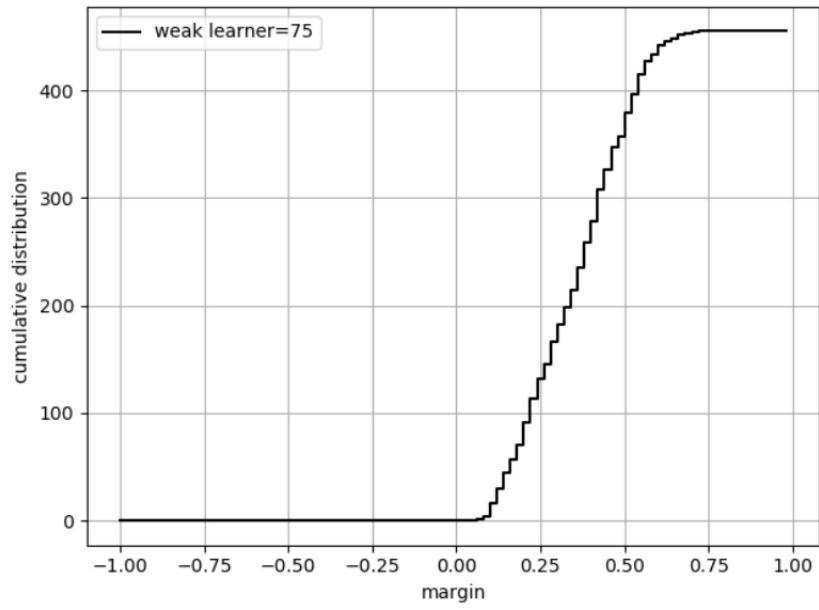


Figure 7: Weak learner=75

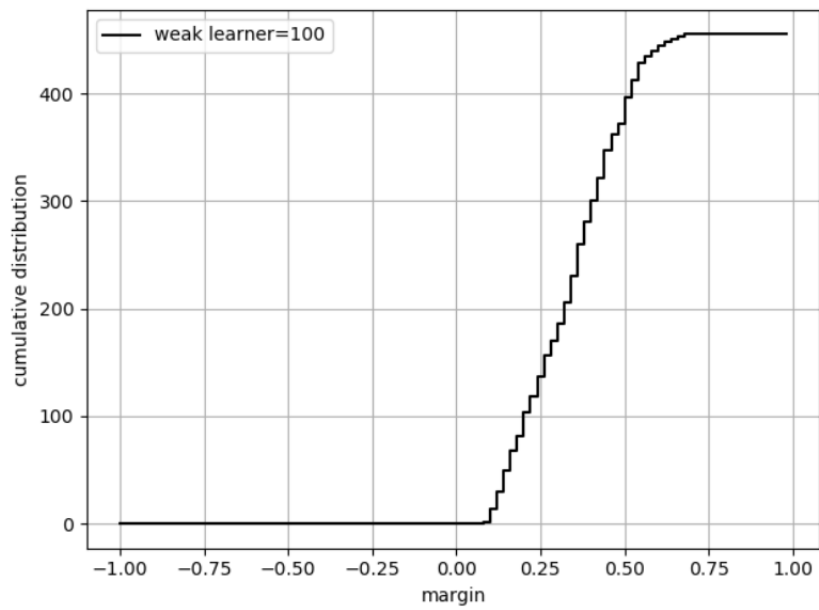


Figure 8: Weak learner=100

Problem 6. Random forest.

(**Total 13 Points**) For this problem, you will use a health dataset as in the file `health_train.csv` and `health_test.csv`. You can find this file in the canvas file folder named `hw4`. The goal is to predict the overall health of an individual: either 0 for poor health or 1 for good health, included in column `y`.

- (**7 Points**) Implement a random forest of 100 decision trees. You are allowed to use decision tree classifier from scikit standard library ([sklearn link](https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestClassifier.html)). **No other function from scikit learn is allowed.** Use gini as quality measure for the decision tree splits and use the default parameters.
- (**3 Points**) Vary the size of random feature sets as `[50, 100, 150, 200, 250]` and fit the model. After the model is fit, plot the accuracy on train and test set vs size of the random feature set.
- (**3 Points**) Vary the number of estimators i.e. number of decision trees as `[10, 20, 40, 80, 100]` for feature size 250 and plot the accuracy on train and test set for number of estimators.

Submission: Submit all plots requested and explanation in latex PDF. Submit your program in a file named (`hw4_random_forest.py`).

(b) We can see when the number of feature increases, the tendency of training and test accuracy also increase. However, when the number of features reaches 250, the test accuracy decreases a little bit. I think the model may be overfitting when there are too many features.

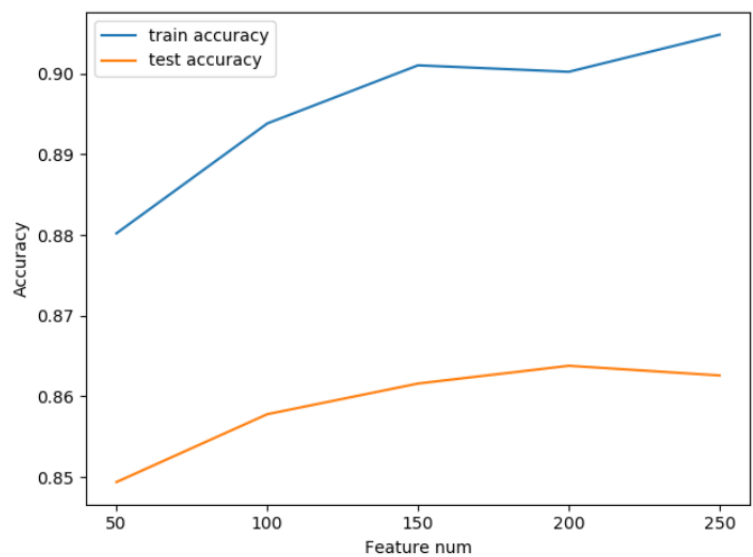


Figure 9: Different feature numbers

(c) When trees increase, the accuracy for training and testing data increase.

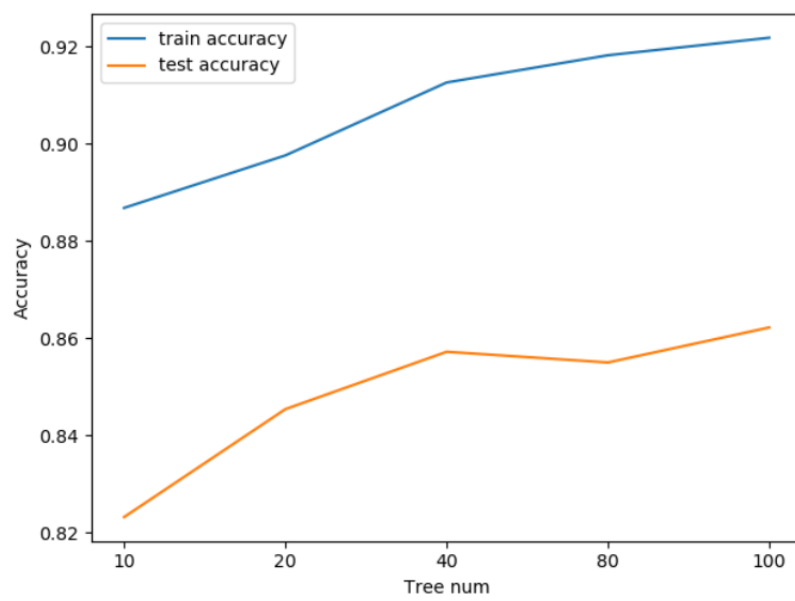


Figure 10: Different tree numbers