

## Policies

- Due 11:59 PM, January 26<sup>th</sup>, via Gradescope.
- You are free to collaborate on all of the problems, subject to the collaboration policy stated in the syllabus.
- In this course, we will be using Google Colab for code submissions. You will need a Google account.
- You are allowed to use up to 48 late hours across the entire term. Late hours must be used in units of whole hours. Specify the total number of hours you have used when submitting the assignment.
- **No use of large language models is allowed.** Students are expected to complete homework assignments based on their understanding of the course material.

## Submission Instructions

- Submit your report as a single .pdf file to Gradescope (entry code 2P8P28), under "Problem Set 3".
- In the report, **include any images generated by your code** along with your answers to the questions.
- Submit your code by **sharing a link in your report** to your Google Colab notebook for each problem (see naming instructions below). Make sure to set sharing permissions to at least "Anyone with the link can view". **Links that can not be run by TAs will not be counted as turned in.** Check your links in an incognito window before submitting to be sure.
- For instructions specifically pertaining to the Gradescope submission process, see [https://www.gradescope.com/get\\_started#student-submission](https://www.gradescope.com/get_started#student-submission).

## Google Colab Instructions

For each notebook, you need to save a copy to your drive.

1. Open the github preview of the notebook, and click the icon to open the colab preview.
2. On the colab preview, go to File → Save a copy in Drive.
3. Edit your file name to "lastname\_firstname\_set\_problem", e.g."yue\_yisong\_set3.prob1.ipynb"

## 1 Decision Trees [30 Points]

*Relevant materials: Lecture 5*

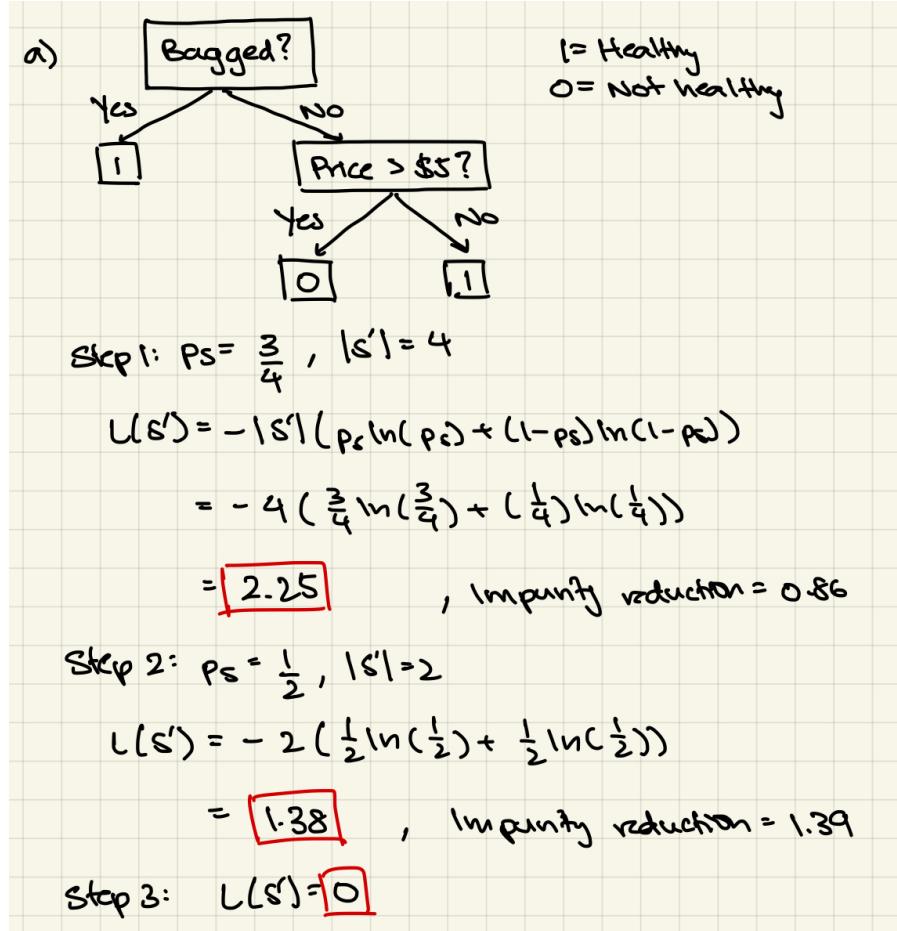
**Problem A [7 points]:** Consider the following data, where given information about some food you must predict whether it is healthy:

No.	Package Type	Unit Price > \$5	Contains > 5 grams of fat	Healthy?
1	Canned	Yes	Yes	No
2	Bagged	Yes	No	Yes
3	Bagged	No	Yes	Yes
4	Canned	No	No	Yes

Train a decision tree by hand using top-down greedy induction. Use *entropy* (with natural log) as the impurity measure. Since the data can be classified without error, the stopping criterion will be no impurity in the leaves.

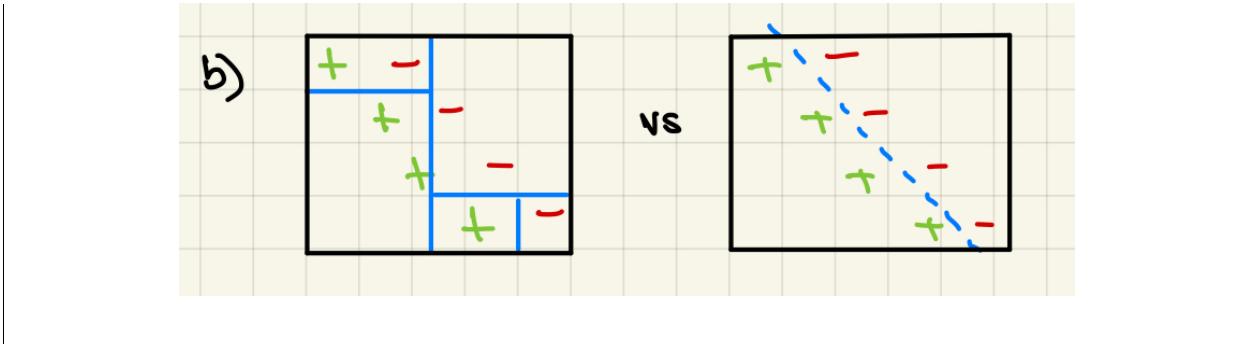
Submit a drawing of your tree showing the impurity reduction yielded by each split (including root) in your decision tree.

**Solution A:**

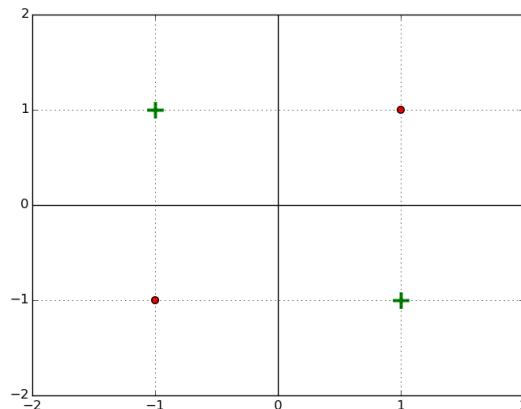


**Problem B [4 points]:** Compared to a linear classifier, is a decision tree always preferred for classification problems? If not, draw a simple 2-D dataset that can be perfectly classified by a simple linear classifier but which requires an overly complex decision tree to perfectly classify.

**Solution B:** No, a decision tree is not always preferred for classification problems. See the figure below for an example:



**Problem C [15 points]:** Consider the following 2D data set:



i. [5 points]: Suppose we train a decision tree on this dataset using top-down greedy induction, with the Gini index as the impurity measure. We define our stopping condition to be if no split of a node results in any reduction in impurity. Submit a drawing of the resulting tree. What is its classification error ((number of misclassified points) / (number of total points))?

ii. [5 points]: Submit a drawing of a two-level decision tree that classifies the above dataset with zero classification error. (You don't need to use any particular training algorithm to produce the tree.)

Is there any impurity measure (i.e. any function that maps the data points under a particular node in a tree to a real number) that would have led top-down greedy induction with the same stopping condition to produce the tree you drew? If so, give an example of one, and briefly describe its pros and cons as an impurity measure for training decision trees in general (on arbitrary datasets).

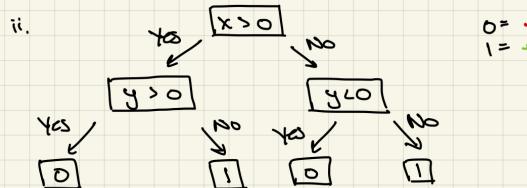
iii. [5 points]: Suppose there are 100 data points in some 2-D dataset. What is the largest number of unique thresholds (i.e., internal nodes) you might need in order to achieve zero classification training error (on the training set)? Please justify your answer.

**Solution C:**

c) i. ~~root~~:  $L(s') = |s'| (1 - p_s^2 - (1-p_s)^2); p_s = 0.5, |s'| = 4$   
 $= 4(1 - 0.5^2 - 0.5^2)$   
 $= 2$

Note (i):  $L(s') = |s'| (1 - p_s^2 - (1-p_s)^2) + |s'| (1 - p_s^2 - (1-p_s)^2); p_s = 0.5, |s'| = 2$   
 $= 2(1 - 0.5^2 - 0.5^2) + 2(1 - 0.5^2 - 0.5^2)$   
 $= 2 \rightarrow \text{no reduction in impurity}$

Class Error =  $2/4 = 1/2$       Tree:  $\boxed{\square}$



An impurity measure that would lead to the same stopping condition is one where the  $|s'|$  term becomes  $|s'|^2$ .

$$\rightarrow L(s') = |s'|^2 (1 - p_s^2 - (1-p_s)^2)$$

Advantages to this is that it can correctly classify all points in the dataset. However, it may lead to overfitting if other datasets are not also structured the same way.

- iii. For a 2D dataset with 100 points, the largest number of unique thresholds you might need for 0 classification error is 99. This is because in the worst case scenario, the 100 points may all have different labels which would require 99 thresholds to separate them.

**Problem D [4 points]:** Suppose in top-down greedy induction we want to split a leaf node that contains N data points composed of D continuous features. What is the worst-case complexity (big-O in terms of N and D) of the number of possible splits we must consider in order to find the one that most reduces impurity? Please justify your answer.

Note: Recall that at each node-splitting step in training a DT, you must consider all possible splits that you can make. While there are an infinite number of possible decision boundaries since we are using continuous features, there are not an infinite number of boundaries that result in unique child sets (which is what we mean by "split").

**Solution D:** *The worst-case complexity of the number of possible splits is  $O(ND)$ . From lecture 5, we know that the # of possible queries = # of possible splits =  $D * N$  where  $D$  is the # of features and  $N$  is the # of training points. Thus,  $O(ND)$  is the worst-case complexity.*

## 2 Overfitting Decision Trees [30 Points]

*Relevant materials: Lecture 5*

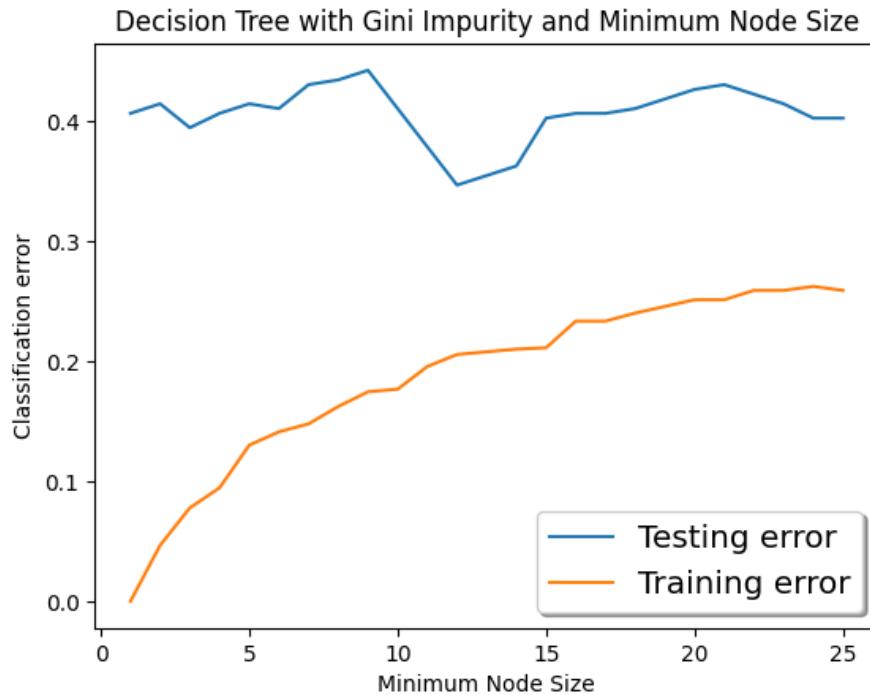
In this problem, you will use the Diabetic Retinopathy Debrecen Data Set, which contains features extracted from images to determine whether or not the images contain signs of diabetic retinopathy. Additional information about this dataset can be found at the link below:

<https://archive.ics.uci.edu/ml/datasets/Diabetic+Retinopathy+Debrecen+Data+Set>

In the following question, your goal is to predict the diagnosis of diabetic retinopathy, which is the final column in the data matrix. Use the first 900 rows as training data, and the last 251 rows as validation data. Please feel free to use additional packages such as Scikit-Learn. Include your code in your submission.

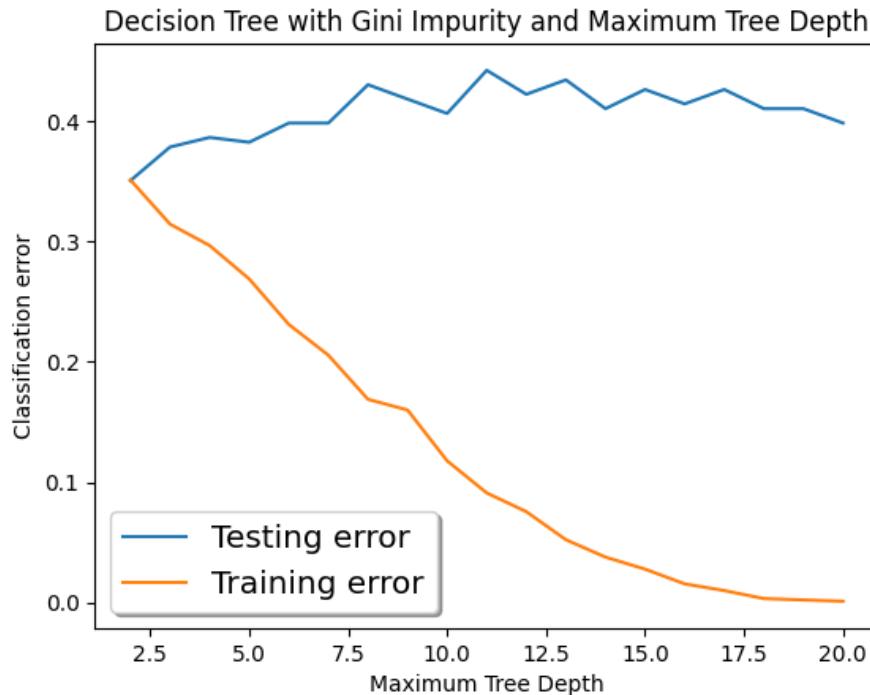
**Problem A [7 points]:** Train a decision tree classifier using Gini as the impurity measure and minimal leaf node size as early stopping criterion. Try different minimal leaf node sizes from 1 to 25 in increments of 1. Then, on a single plot, plot both training and test classification error versus leaf node size. To do this, fill in the `classification_err` and `eval_tree_based_model_min_samples` functions in the code template for this problem.

**Solution A:** Code: [mantripragada\\_ishaan\\_set3\\_prob2.ipynb](#)



**Problem B [7 points]:** Train a decision tree classifier using Gini as the impurity measure and maximal tree depth as early stopping criterion. Try different tree depths from 2 to 20 in increments of 1. Then, on a single plot, plot both training and test classification error versus tree depth. To do this, fill in the `eval_tree_based_model_max_depth` function in the code template for this problem.

**Solution B:**

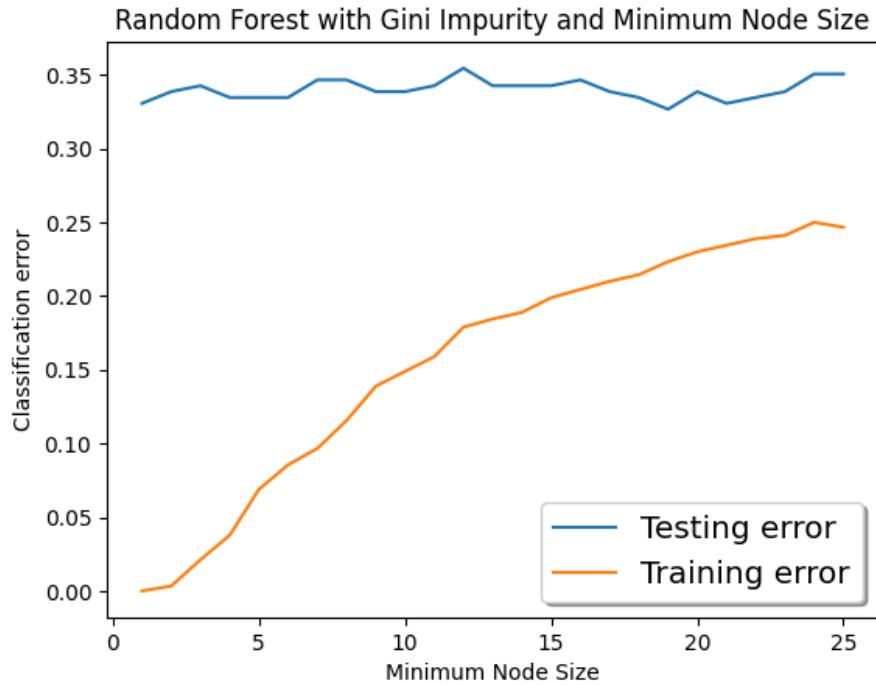


**Problem C [4 points]:** For both the minimal leaf node size and maximum depth parameters tested in the last two questions, which parameter value minimizes the test error? What effects does early stopping have on the performance of a decision tree model? Please justify your answer based on the two plots you derived.

**Solution C:** For the minimal leaf node size, the parameter value that minimizes the test error is 12. After that, we see that the test error increases which are indications of overfitting. Thus, early stopping around 12-15 maximum tree depth would result in a better model performance. For the maximum depth, the test error seems to remain relatively constant, indicating no sign of learning. Early stopping in this case would not seem to affect the model's performance.

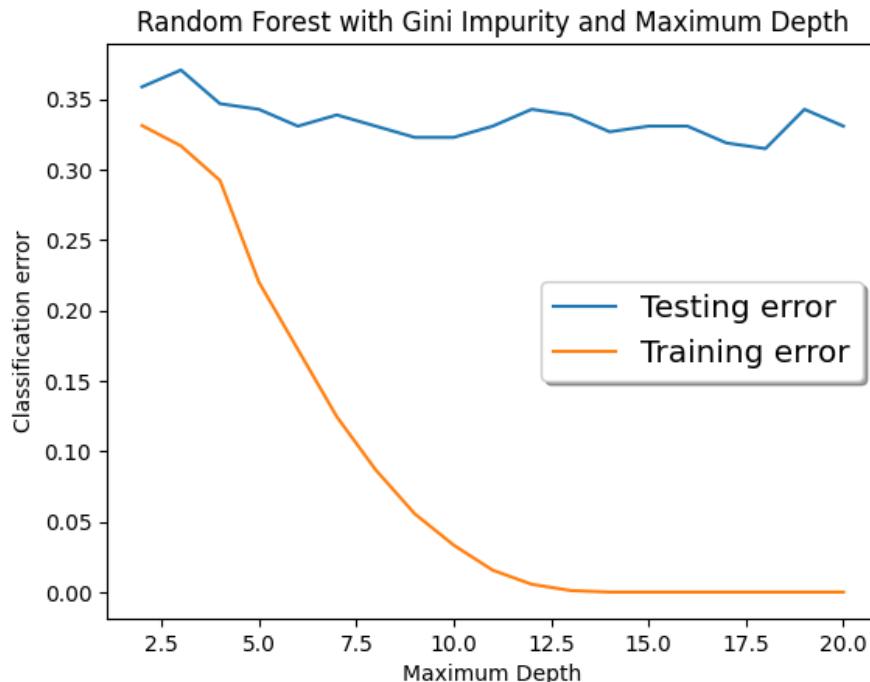
**Problem D [2 points]:** Train a random forest classifier using Gini as the impurity measure, minimal leaf node size as early stopping criterion, and 1,000 trees in the forest. Try different node sizes from 1 to 25 in increments of 1. Then, on a single plot, plot both training and test classification error versus leaf node size.

**Solution D:**



**Problem E [2 points]:** Train a random forest classifier using Gini as the impurity measure, maximal tree depth as early stopping criterion, and 1,000 trees in the forest. Try different tree depths from 2 to 20 in increments of 1. Then, on a single plot, plot both training and test classification error versus tree depth.

**Solution E:**



**Problem F [4 points]:** For both the minimal leaf node size and maximum depth parameters tested, which parameter value minimizes the random forest test error? What effects does early stopping have on the performance of a random forest model? Please justify your answer based on the two plots you derived.

**Solution F:** For the minimal leaf node size, the test error again seems to remain relatively constant around 0.35 which indicates no sign of learning. Early stopping would not have any positive affect on the model's performance. For the maximum depth, the parameter value that minimizes the test error is around 18. We also notice that the training error seems to drop to 0 around a maximum depth of 13, which again indicates that early stopping may not have a positive impact.

**Problem G [4 points]:** Do you observe any differences between the curves for the random forest and decision tree plots? If so, explain what could account for these differences.

**Solution G:** Based on the plots shown, the testing errors for the random forest plots are much less volatile compared to the decision tree plots. Also, early stopping seems to improve the decision tree plot, but not the random forests. The reason for this is because decision trees are low bias, high variance models compared to random forests.

### 3 The AdaBoost Algorithm [40 points]

*Relevant materials: Lecture 6*

In this problem, you will show that the choice of the  $\alpha_t$  parameter in the AdaBoost algorithm corresponds to greedily minimizing an exponential upper bound on the loss term at each iteration.

**Problem A [3 points]:** Let  $h_t : \mathbb{R}^m \rightarrow \{-1, 1\}$  be the weak classifier obtained at step  $t$ , and let  $\alpha_t$  be its weight. Recall that the final classifier is

$$H(x) = \text{sign}(f(x)) = \text{sign}\left(\sum_{i=1}^T \alpha_i h_i(x)\right).$$

Suppose  $\{(x_1, y_1), \dots, (x_N, y_N)\} \subset \mathbb{R}^m \times \{-1, 1\}$  is our training dataset. Show that the training set error of the final classifier can be bounded from above if an exponential loss function is used:

$$E = \frac{1}{N} \sum_{i=1}^N \exp(-y_i f(x_i)) \geq \frac{1}{N} \sum_{i=1}^N \mathbb{1}(H(x_i) \neq y_i),$$

where  $\mathbb{1}$  is the indicator function.

**Solution A:**

$$\begin{aligned} \text{a)} \quad E &= \frac{1}{N} \sum_{i=1}^N \exp(-y_i f(x_i)) \geq \frac{1}{N} \sum_{i=1}^N \mathbb{1}(H(x_i) \neq y_i) \\ &\rightarrow \exp(-y_i f(x_i)) \geq \mathbb{1}(H(x_i) \neq y_i) \quad \text{for every } (x_i, y_i) \\ &\text{if } (-y_i f(x_i) \geq 0) \rightarrow y_i \text{ and } f(x_i) \text{ have opposite sign} \\ &\quad \hookrightarrow \exp(-y_i f(x_i)) \geq 1 \\ &\quad \hookrightarrow \mathbb{1}(H(x_i) \neq y_i) = 1 \quad \boxed{\exp(-y_i f(x_i)) \geq \mathbb{1}(H(x_i) \neq y_i) \checkmark} \\ &\text{else} \rightarrow y_i \text{ and } f(x_i) \text{ have same sign} \\ &\quad \hookrightarrow \exp(-y_i f(x_i)) \geq 0 \\ &\quad \hookrightarrow \mathbb{1}(H(x_i) \neq y_i) = 0 \quad \boxed{\exp(-y_i f(x_i)) \geq \mathbb{1}(H(x_i) \neq y_i) \checkmark} \end{aligned}$$

**Problem B [3 points]:** Find  $D_{T+1}(i)$  in terms of  $Z_t$ ,  $\alpha_t$ ,  $x_i$ ,  $y_i$ , and the classifier  $h_t$ , where  $T$  is the last

timestep and  $t \in \{1, \dots, T\}$ . Recall that  $Z_t$  is the normalization factor for distribution  $D_{t+1}$ :

$$Z_t = \sum_{i=1}^N D_t(i) \exp(-\alpha_t y_i h_t(x_i)).$$

**Solution B:**

$$\begin{aligned} b) \quad Z_t &= \sum_{i=1}^N D_t(i) \exp(-\alpha_t y_i h_t(x_i)) \\ D_{t+1}(i) &= \frac{D_t(i) \exp(-\alpha_t y_i h_t(x_i))}{Z_t}, \quad D_1(x) = \frac{1}{N} \\ \therefore D_{T+1}(i) &= \left[ \frac{\exp(-\alpha_T y_i h_T(x_i))}{Z_T} \right] \left[ \frac{\exp(-\alpha_{T-1} y_i h_{T-1}(x_i))}{Z_{T-1}} \right] \left[ \dots \right] = \\ &= \boxed{\frac{1}{N} \prod_{t=1}^T \frac{\exp(-\alpha_t y_i h_t(x_i))}{Z_t}} \end{aligned}$$

**Problem C [2 points]:** Show that  $E = \sum_{i=1}^N \frac{1}{N} e^{\sum_{t=1}^T -\alpha_t y_i h_t(x_i)}$ .

**Solution C:**

$$\begin{aligned} c) \quad E &= \frac{1}{N} \sum_{i=1}^N \exp(-y_i f(x_i)), \quad f(x_i) = \sum_{t=1}^T \alpha_t h_t(x_i) \\ &= \sum_{i=1}^N \left( \frac{1}{N} \right) \exp(-y_i \sum_{t=1}^T \alpha_t h_t(x_i)) \\ &= \boxed{\sum_{i=1}^N \left( \frac{1}{N} \right) \exp \left( \sum_{t=1}^T -y_i \alpha_t h_t(x_i) \right)} \end{aligned}$$

**Problem D [5 points]:** Show that

$$E = \prod_{t=1}^T Z_t.$$

**Hint:** Recall that  $\sum_{i=1}^N D_t(i) = 1$  because  $D$  is a distribution.

**Solution D:**

$$\begin{aligned}
 a) \quad D_{T+1}(i) &= \frac{1}{N} \prod_{t=1}^T \frac{\exp(-\alpha_t y_i h_t(x_i))}{Z_t} \\
 &= \frac{1}{N} \prod_{t=1}^T \exp(-\alpha_t y_i h_t(x_i)) \cdot \prod_{t=1}^T \frac{1}{Z_t} \\
 \\
 \left[ D_{T+1}(i) \right] \left[ \prod_{t=1}^T Z_t \right] &= \frac{1}{N} \prod_{t=1}^T \exp(-\alpha_t y_i h_t(x_i)) \\
 (c) \quad E &= \sum_{i=1}^N \left( \frac{1}{N} \exp \left( \sum_{t=1}^T -y_i \alpha_t h_t(x_i) \right) \right) \\
 &= \sum_{i=1}^N \frac{1}{N} \prod_{t=1}^T \exp(-\alpha_t y_i h_t(x_i)) \\
 \\
 \sum_{i=1}^N \left[ D_{T+1}(i) \right] \left[ \prod_{t=1}^T Z_t \right] &= E \\
 \\
 \boxed{E = \prod_{t=1}^T Z_t}
 \end{aligned}$$

**Problem E [5 points]:** Show that the normalizer  $Z_t$  can be written as

$$Z_t = (1 - \epsilon_t) \exp(-\alpha_t) + \epsilon_t \exp(\alpha_t)$$

where  $\epsilon_t$  is the training set error of weak classifier  $h_t$  for the weighted dataset:

$$\epsilon_t = \sum_{i=1}^N D_t(i) \mathbb{1}(h_t(x_i) \neq y_i).$$

**Solution E:**

$$\begin{aligned}
 e) \epsilon_t &= \sum_{i=1}^N D_t(i) \mathbb{1}(h_t(x_i) \neq y_i) \\
 Z_t &= \sum_{i=1}^N D_t(i) \exp(-\alpha_t y_i h_t(x_i)) \\
 &= \sum_{i=1}^N D_t(i) \left[ [1 - \mathbb{1}(h_t(x_i) \neq y_i)] \exp(-\alpha_t) + [\mathbb{1}(h_t(x_i) \neq y_i) \exp(\alpha_t)] \right] \\
 &\Rightarrow \exp(-\alpha_t) \left[ \sum_{i=1}^N D_t(i) - \sum_{i=1}^N D_t(i) \mathbb{1}(h_t(x_i) \neq y_i) \right] + \sum_{i=1}^N D_t(i) \mathbb{1}(h_t(x_i) \neq y_i) \exp(\alpha_t) \\
 &= \boxed{\exp(-\alpha_t) (1 - \epsilon_t) + (\epsilon_t) \exp(\alpha_t)}
 \end{aligned}$$

**Problem F [2 points]:** We derived all of this because it is hard to directly minimize the training set error, but we can greedily minimize the upper bound  $E$  on this error. Show that choosing  $\alpha_t$  greedily to minimize  $Z_t$  at each iteration leads to the choices in AdaBoost:

$$\alpha_t^* = \frac{1}{2} \ln \left( \frac{1 - \epsilon_t}{\epsilon_t} \right).$$

**Solution F:**

f) Minimize  $Z_t$ :

$$Z_t = (1 - \epsilon_t) \exp(-\alpha_t) + \epsilon_t \exp(\alpha_t)$$

$$\frac{\partial Z_t}{\partial \alpha_t} = -(1 - \epsilon_t) \exp(-\alpha_t) + \epsilon_t \exp(\alpha_t) = 0$$

$$\rightarrow (\epsilon_t - 1) \exp(-\alpha_t) + \epsilon_t \exp(\alpha_t) = 0$$

$$\exp(\alpha_t) [(\epsilon_t - 1) \exp(-2\alpha_t) + \epsilon_t] = 0$$

$$(\epsilon_t - 1) \exp(-2\alpha_t) + \epsilon_t = 0$$

$$\exp(-2\alpha_t) = \frac{-\epsilon_t}{\epsilon_t - 1}$$

$$-2\alpha_t = \ln\left(\frac{-\epsilon_t}{\epsilon_t - 1}\right)$$

$$\alpha_t = -\frac{1}{2} \ln\left[\frac{-\epsilon_t}{\epsilon_t - 1}\right]$$

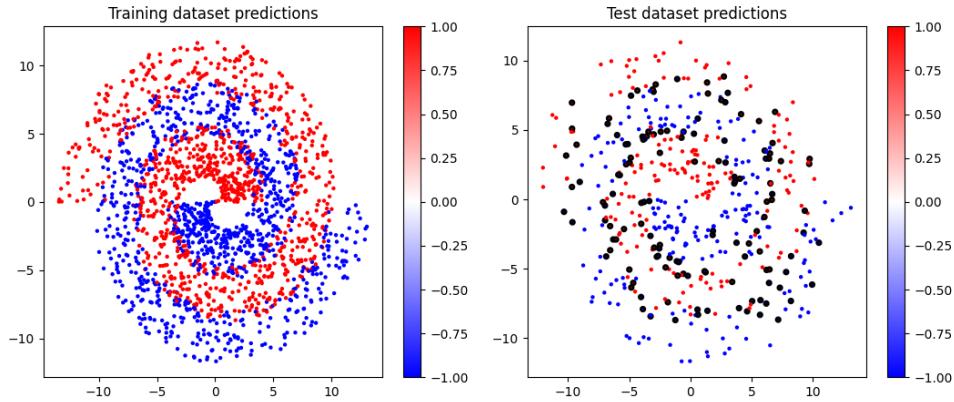
$$= \boxed{\frac{1}{2} \ln\left[\frac{1 - \epsilon_t}{\epsilon_t}\right]}$$

**Problem G [14 points]:** Implement the `GradientBoosting.fit()` and `AdaBoost.fit()` methods in the notebook provided for you. Some important notes and guidelines follow:

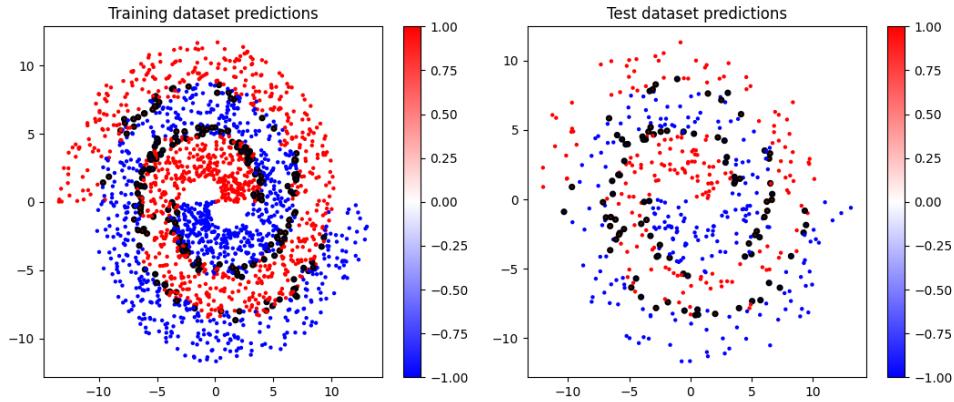
- For both methods, make sure to work with the class attributes provided to you. Namely, after `GradientBoosting.fit()` is called, `self.clfs` should be appropriately filled with the `self.n_clfs` trained weak hypotheses. Similarly, after `AdaBoost.fit()` is called, `self.clfs` and `self.coeffs` should be appropriately filled with the `self.n_clfs` trained weak hypotheses and their coefficients, respectively.
- `AdaBoost.fit()` should additionally return an  $(N, T)$  shaped numpy array `D` such that `D[:, t]` contains  $D_{t+1}$  for each  $t \in \{0, \dots, \text{self.n_clfs}\}$ .
- For the `AdaBoost.fit()` method, **use the 0/1 loss** instead of the exponential loss.
- The only Sklearn classes that you may use in implementing your boosting fit functions are the `DecisionTreeRegressor` and `DecisionTreeClassifier`, not `GradientBoostingRegressor`.

**Solution G:**

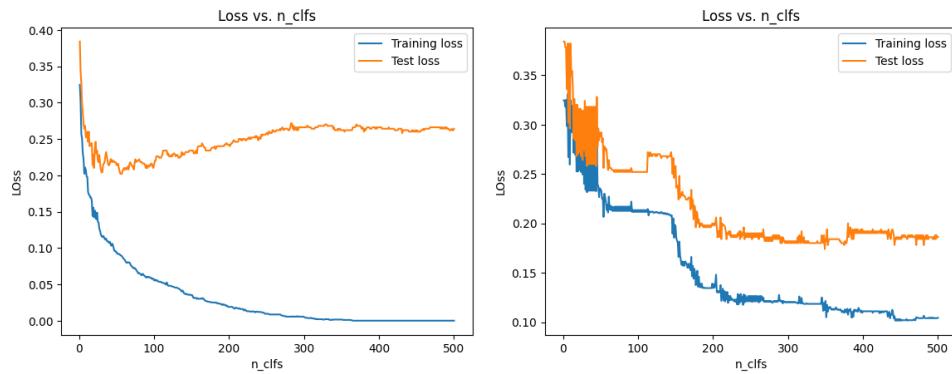
*Gradient Boosting:*



*AdaBoost:*



*Gradient Boosting and AdaBoost loss:*



**Problem H [2 points]:** Describe and explain the behaviour of the loss curves for gradient boosting and for

AdaBoost. You should consider two kinds of behaviours: the smoothness of the curves and the final values that the curves approach.

**Solution H:** Code: [mantripragada\\_ishaan\\_set3\\_prob3.ipynb](#)

*The loss curves for gradient boosting is much smoother than that for AdaBoost. However, the test loss seems to increase when the number of clfs reaches 100. This occurs far before the training loss reaches its minimum which indicates overfitting. For AdaBoost, although there is a slight increase at clfs = 100, the test loss continue to decrease until the clfs reaches 500. There is more variance in the both the training and test loss, but both seem to decrease at similar rates. AdaBoost shows no signs of overfitting as the test loss continue to decrease. The final values for the gradient boosting training loss reaches 0 while the Adaboost training loss only seems to reach 0.1.*

**Problem I [2 points]:** Compare the final loss values of the two models. Which performed better on the classification dataset?

**Solution I:** *For gradient boosting, the final training loss was 0 and the final testing loss was 0.25. For AdaBoost, the final training loss was 0.1 and the final testing loss was around 0.18. Overall, AdaBoost performed better on the classification dataset since its test loss was less than gradient boosting.*

**Problem J [2 points]:** For AdaBoost, where are the dataset weights the largest, and where are they the smallest?

**Hint:** Watch how the dataset weights change across time in the animation.

**Solution J:** *The dataset weights were the largest at the boundary of the blue and red points. The weights were the smallest when the datapoint was not close to any boundary. This makes sense because the weights should be higher when the dataset is harder to classify while less importance can be placed on more clear datapoints.*