

Name: Derrick Kho

PUID: 0037028290

Instructions and Policy: You are allowed to study with others and use online resources for reference, however, the work you turn in must be your own. This means do not copy/paste from Stack Exchange (or from another student.) If you have worked closely with other students, provide their name(s) and a brief (at most one paragraph) description of the interaction; if we feel this oversteps the bounds, we will discuss it with you. Each student should write up their own solutions independently. The requirements below are supposed to be followed in this and further homework assignments.

- For your theoretical submission:
 - **YOU MUST INCLUDE YOUR NAME IN THE HOMEWORK.**
 - The answers **MUST** be submitted via Gradescope.
 - Please write clearly and concisely - clarity and brevity will be rewarded. Refer to known facts as necessary.
 - Theoretical questions **MUST** include the intermediate steps to the final answer.
- For your programming answers:
 - **The Python scripts will be submitted separately via Gradescope**
 - Zero points will be given in any question where the Python code answer doesn't match the answer on Gradescope.
 - If the answer is/includes a plot, it should be added to your theoretical submission unless otherwise specified.

Your code is REQUIRED to run on Python 3.11 in an environment detailed in Homework 0 under the Python Installation section. While your code may run in other environments, any points lost due to environmental issues will not be regraded.

Please make sure you don't use any libraries that are not imported for you. If such a library is used, you will get 0 pt for the coding part of the assignment.

If your code doesn't run on Gradescope, then even if it compiles on another computer, it will still be considered not running and the respective part of the assignment will receive 0 pt.

Theoretical Questions (8+20+14+16+18=76 pts)

Please submit your answers on Gradescope.

Q1 (8 pts): True or False questions

Answer the following as True or False with a justification or example. Points are uniformly distributed within the questions.

- (a) (4 pts) High bias can result in the model learning random noise in the training data.

False. A high bias would result in a model that underfits the training data and is likely to be a simple model. Hence, this could mean that the model would be rather generalised, and factors specific types of solutions and thus random noise would not be able to learn random noise in the training data.

- (b) (4 pts) In k-fold cross-validation, each sample is used to train the model k times.

False. This is because for every iteration i , we train the model on $S_1 \cup S_2 \cup \dots \cup S_{i-1} \cup S_i^c \cup \dots \cup S_k$. Thus, when the model is trained k times, each sample is used for a total of $k-1$ times, and not k times.

Q2 (20 pts): Cross-validation

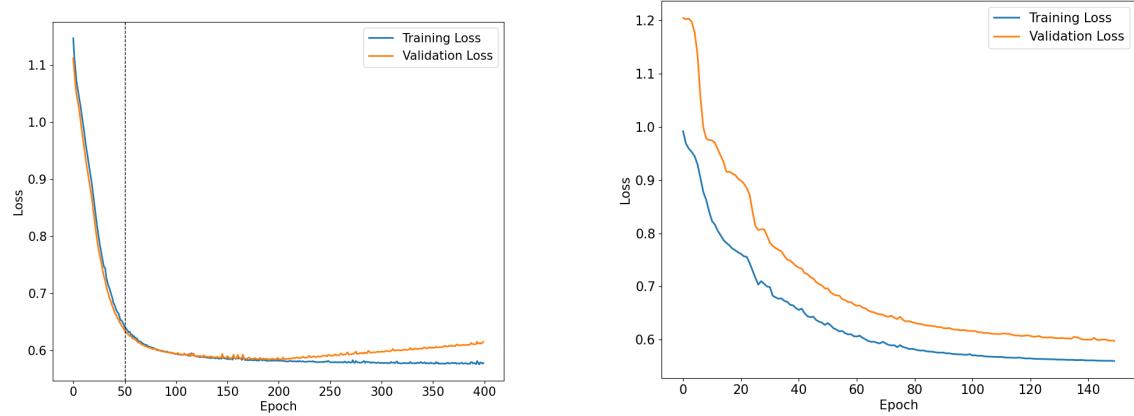


Figure 1: Example of training curves showing training and validation loss for two hypothetical models M_1 and M_2 .

1. (6 pts) Review the training and validation loss curve in Figure 1a (model M_1). The y -axis shows a loss (lower is better) for the training and validation data. The x -axis shows the iterations to update the model parameters (e.g. update steps on Perceptron algorithm or steps of gradient descent for logistic regression).

Hint: Note that these are training curves, not learning curves.

- (3 pts) What are the best parameters for model M_1 ? The parameters at 150 iterations or at 400 iterations? Please justify your answer.

The parameters at 150 iterations. This is because at 150 iterations, the training loss and validation loss is the lowest, at ≈ 0.6 for both, as compared to at 400 iterations where the validation loss is higher at ≈ 0.65 , despite having a lower training loss. This indicates that at 400 iterations, the

- (3 pts) What is the disadvantage of choosing the model with parameters obtained after 50 iterations?

→ model has overfitted to the training data.

The disadvantage is that the model has not fully converged and may be underfitting to the training data. This is because the model may not have learned enough from the training data to capture the underlying pattern of the data and thus, resulting in higher training and validation losses.

2. (4 pts) Review the training and validation loss curve in Figure 1b of model M_2 .

Which of the following explanations and decisions are appropriate for the observed behavior in 1b?
Choose all answers that are correct. You don't need to justify your answer.

- (A) The validation data is scarce and not very representative.
- (B) The training data does not provide sufficient information to learn the problem relative to the validation data used to evaluate it.
- (C) Getting more validation data will not reduce the training loss.
- (D) The training loss has converged, so obtaining more and better training data will not reduce the training loss further.

C

3. (4 pts) Which statement about k -fold cross-validation is incorrect?

Find out the false statements.

- (A) When using k -fold cross-validation, using a smaller k means it will take more time to run the evaluation.
- (B) With larger k , the training set gets larger, while the test set gets smaller. X
- (C) When splitting data into folds, it is desirable to minimize the variance across folds. X
- (D) When $k = N$, where N is the size of the data set, k -fold cross-validation is not the same as leave-one-out cross-validation.

A, D

4. (6 pts) Given the supervised learning dataset with 3 training examples $(X_1, Y_1) = (-2, -1)$, $(X_2, Y_2) = (-1, -1)$, $(X_3, Y_3) = (0, 1)$, use the Perceptron algorithm to learn a classifier. If needed, use $\text{sign}(0) = 1$. The Perceptron classifier will have two parameters $w = (w_0, w_1)$ (the first parameter, w_0 acts as the bias). What is the average 0-1 loss of the Perceptron algorithm evaluated by 3-fold cross-validation? (Please also describe, for each fold/iteration of the cross-validation, what is the training data used and the loss obtained).

Using Perceptron algorithm on these 3 training examples:

We initialise $w_0 = 0, w_1 = 0$ first. Then, we run the algorithm;

On first iteration, we use $(x_1, y_1), (x_2, y_2)$ as training data and (x_3, y_3) as validation set:

$$\text{sign}(w_1^T x_1 + w_0) = \text{sign}(0) = 1$$

$$\therefore w_1 = 0 + (-1)(-2) = 2$$

$$w_0 = 0 + (-1) = -1$$

$$\text{sign}(w_1^T x_2 + w_0) = \text{sign}(2(-1) + (-1)) = \text{sign}(-3) = -1$$

$$\text{sign}(w_1^T x_1 + w_0) = \text{sign}(2(-2) + (-1)) = \text{sign}(-5) = -1$$

On validation set, $\text{sign}(w_1^T x_3 + w_0) = \text{sign}(2(0) + 1) = \text{sign}(1) = 1$
 $\therefore 0-1 \text{ loss obtained} = 0 \text{ as it is correctly classified.}$

On second iteration, we use (x_1, y_1) and (x_3, y_3) as training data and (x_2, y_2) as validation set and $w_0 = w_1 = 0$ initially.

$$\text{sign}(w_1^T x_1 + w_0) = \text{sign}(0) = 1$$

$$\therefore w_1 = 0 + (-2)(-1) = 2$$

$$w_0 = 0 + (-1) = -1$$

$$\text{sign}(w_1^T x_2 + w_0) = \text{sign}(2(0) - 1) = \text{sign}(-1) = -1$$

$$\therefore w_1 = 2 + (1)(0) = 2$$

$$w_0 = -1 + (1) = 0$$

$$\text{sign}(w_1^T x_1 + w_0) = \text{sign}(2(-2) + 0) = \text{sign}(-4) = -1$$

$$\text{sign}(w_1^T x_3 + w_0) = \text{sign}(2(0) + 0) = \text{sign}(0) = 1$$

Q2.4 On validation set, $\text{sign}(w_1^T x_2 + w_0) = \text{sign}(2(-1) + 0) = \text{sign}(-2) = -1$
 (continued) $\therefore 0-1 \text{ loss obtained} = 0$ as it is correctly classified.

On third iteration, we use (x_2, y_2) and (x_3, y_3) as training data and (x_1, y_1) as validation set, and $w_r = w_0 = 0$ initially.

$$\text{sign}(w_1^T x_2 + w_0) = \text{sign}(0) = 1$$

$$\therefore w_1 = 0 + (-1)(-1) \\ = 1$$

$$w_0 = 0 + (-1) \\ = -1$$

$$\text{sign}(w_1^T x_3 + w_0) = \text{sign}(1(0) + (-1)) \\ = \text{sign}(-1) = -1$$

$$\therefore w_1 = 1 + (0)(1) \\ = 1$$

$$w_0 = -1 + 1 \\ = 0$$

$$\text{sign}(w_1^T x_2 + w_0) = \text{sign}((1(-1) + 0)) = \text{sign}(-1) = -1$$

$$\text{sign}(w_1^T x_3 + w_0) = \text{sign}((1(0) + 0)) = \text{sign}(0) = 1$$

On validation set, $\text{sign}(w_1^T x_1 + w_0) = \text{sign}(1(-2) + 0) = \text{sign}(-2) = -1$
 $\therefore 0-1 \text{ loss obtained} = 0$ as it is correctly classified.

$$\therefore \text{Average } 0-1 \text{ loss of the Perceptron algorithm} = \frac{1}{3} \sum_{i=1}^3 E_{0-1} = \frac{0+0+0}{3} = 0$$

Q3 (14 pts): Ensemble methods (AdaBoost)

Consider the dataset in the following table:

Index (i)	Feature 1 (X_{i1})	Feature 2 (X_{i2})	Class label (y_i)
1	-1	1	-1
2	1	1	+1
3	1	-1	+1
4	-1	-1	+1

For this question, you have to show the first few steps of the AdaBoost algorithm. The weak learner can be one of the four possible decision stumps illustrated in Fig. 2 (i.e., at each step you must choose one of the classifiers in the figure, to minimize the weighted empirical error rate, defined as $\varepsilon_t = \sum_{i=1}^n D_t(i) \mathbb{I}[h_t(x_i) \neq y_i]$,

where t means the number of iterations and $D_t(i)$ means the reweighting of data at round t , breaking any ties by choosing the classifier with lower number). You may refer to this link https://drive.google.com/file/d/10WABd0_gvk0Cgt3WsiNx_b-hoOUUvXM/view?usp=sharing for details of AdaBoost.

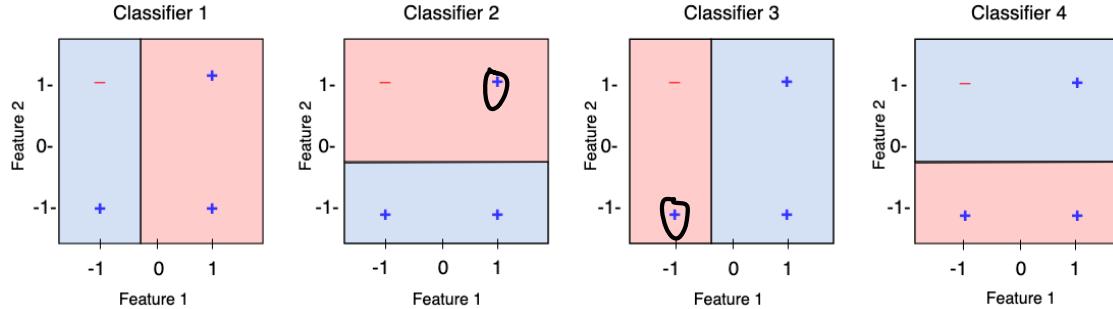


Figure 2: Choices of classifiers for the dataset. The solid red areas are assigned label -1 and the hatched blue areas are assigned label $+1$

1. (6 pts) Complete the following table, with the correct values computed during the execution of the AdaBoost algorithm, for $t = 1, \dots, 5$. Please refer to lecture slides for details of AdaBoost.

t	ε_t	α_t	$D_t(i)$				h_t	$\mathbb{I}(h_t(x_i) \neq y_i)$			
			$i = 1$	$i = 2$	$i = 3$	$i = 4$		$i = 1$	$i = 2$	$i = 3$	$i = 4$
1	0.25	0.549	0.25	0.25	0.25	0.25	2	0	1	0	0
2	0.166	0.807	0.166	0.501	0.166	0.166	3	0	6	0	1
3	0.301	0.474	0.0916	0.301	0.0996	0.500	2	0	1	0	0
4	0.357	0.294	0.0913	0.500	0.0713	0.357	3	0	0	0	1
5	0.406	0.190	0.0578	0.373	0.0578	0.440	2	0	1	0	0

2. (4 pts) Complete the following table, with the predictions computed by the AdaBoost algorithm for the training data, at the end of training (i.e. after $T = 5$ rounds).

	$\hat{h}_t(x_i)$			
	$i = 1$	$i = 2$	$i = 3$	$i = 4$
$\sum_{t=1}^T \alpha_t h(x_i)$	-2.261	-0.059	2.261	0.059
$\hat{h}(x_i)$	-1	-1	+1	+1

Q3.1

Calculations

$$\varphi = 1: \quad g_t = 0.25(1) \\ = 0.25$$

$$x_t = \frac{1}{2} \ln \left(\frac{1-0.25}{0.25} \right) \\ = 0.549 \text{ (3sf)}$$

$$t=2: \quad D_t(1) = D_t(3) = D_t(4) \\ = 0.25 \cdot e^{-0.549} \\ = 0.444 \text{ (3sf)} \\ D_t(2) = 0.25 e^{0.549} \\ = 0.433 \text{ (3sf)}$$

Renormalizing:

$$D_t(1) = D_t(3) = D_t(4) \\ = \frac{0.444}{0.444(3) + 0.433} \\ \approx 0.66 \text{ (3sf)}$$

$$D_t(2) = \frac{0.433}{0.444(3) + 0.433} \\ \approx 0.501 \text{ (3sf)}$$

$$g_t = 0.66(1) \\ \approx 0.66$$

$$x_t = \frac{1}{2} \ln \left(\frac{1-0.66}{0.66} \right) \\ = 0.807 \text{ (3sf)}$$

$\ell = 3:$

$$D_t(1) = D_t(3) \\ = 0.66 e^{-0.807} \\ \approx 0.0741 \text{ (3sf)}$$

$$D_t(2) = 0.501 e^{-0.807} \\ \approx 0.224 \text{ (3sf)}$$

$$D_t(4) = 0.444 e^{0.807} \\ \approx 0.372 \text{ (3sf)}$$

Renormalizing:

$$D_t(1) = \frac{0.0741}{(2) 0.0741 + 0.224 + 0.372} \\ \approx 0.0996 = D_t(3)$$

$$D_t(2) = \frac{0.224}{(2) 0.0741 + 0.224 + 0.372} \\ \approx 0.301$$

$$D_t(4) = \frac{0.372}{(2) 0.0741 + 0.224 + 0.372} \\ \approx 0.500 \text{ (3sf)}$$

$$g_t = 0.301 \quad (1) = 0.301$$

$$x_t = \frac{1}{2} \ln \left(\frac{1-0.301}{0.301} \right) \\ = 0.421$$

$$\ell = 4:$$

$$D_\ell(1) = D_\ell(3)$$

$$= 0.0996 e^{-0.421}$$

$$= 0.0659 \text{ (3st)}$$

$$D_\ell(2) = 0.301 e^{0.421}$$

$$= 0.459 \text{ (3st)}$$

$$P_\ell(4) = 0.500 e^{-0.421}$$

$$= 0.328 \text{ (3st)}$$

Renormalizing:

$$D_\ell(1) = \frac{0.0659}{2(0.0659) + 0.459 + 0.328}$$

$$= 0.07(3) = D_\ell(3)$$

$$L_\ell(2) = \frac{0.459}{2(0.0659) + 0.459 + 0.328}$$

$$= 0.500 \text{ (3st)}$$

$$D_\ell(4) = \frac{0.328}{2(0.0659) + 0.459 + 0.328}$$

$$= 0.357 \text{ (3st)}$$

$$g_\ell = 0.357 \quad (1 = 0.357)$$

$$\alpha_\ell = \frac{1}{2} \ln \left(\frac{1 - 0.357}{0.357} \right)$$

$$= 0.294$$

$$\ell = 5:$$

$$D_\ell(1) = D_\ell(3)$$

$$= 0.093 e^{-0.294}$$

$$= 0.0531 \text{ (3st)}$$

$$D_\ell(2) = 0.500 e^{-0.294}$$

$$= 0.373 \text{ (3st)}$$

$$P_\ell(4) = 0.328 e^{0.294}$$

$$= 0.440 \text{ (3st)}$$

Renormalizing:

$$D_\ell(1) = \frac{0.0531}{2(0.0531) + 0.373 + 0.440}$$

$$= 0.0578 = D_\ell(3)$$

$$L_\ell(2) = \frac{0.373}{2(0.0531) + 0.373 + 0.440}$$

$$= 0.406 \text{ (3st)}$$

$$D_\ell(4) = \frac{0.440}{2(0.0531) + 0.373 + 0.440}$$

$$= 0.449 \text{ (3st)}$$

$$g_\ell = 0.406 \quad (1 = 0.406)$$

$$\alpha_\ell = \frac{1}{2} \ln \left(\frac{1 - 0.406}{0.406} \right)$$

$$= 0.190 \text{ (3st)}$$

Q3.2

(calculations)

$i = 1$	$i = 2$	$i = 3$	$i = 4$
$\hat{h}_t(x_i)$			$\sum_{t=1}^T \alpha_t h(x_i)$

$i=1 : \sum_{t=1}^T \alpha_t h(x_i) = 0.549(-1) + 0.807(-1) + 0.421(1) + 0.294(1) + 0.190(-1)$

$$= -2.261$$

$$\hat{h}(x_i) = \text{sign}(-2.261)$$

$$= -1$$

$i=2 : \sum_{t=1}^T \alpha_t h(x_i) = 0.549(-1) + 0.807(1) + 0.421(-1) + 0.294(1) + 0.190(-1)$

$$= -0.059$$

$$\hat{h}(x_i) = \text{sign}(-0.059)$$

$$= -1$$

$i=3 : \sum_{t=1}^T \alpha_t h(x_i) = 0.549(1) + 0.807(1) + 0.421(1) + 0.294(1) + 0.190(1)$

$$= 2.261$$

$$\hat{h}(x_i) = \text{sign}(2.261)$$

$$= 1$$

$i=4 : \sum_{t=1}^T \alpha_t h(x_i) = 0.549(1) + 0.807(-1) + 0.421 + 0.294(-1) + 0.190(1)$

$$= 0.059$$

$$\hat{h}(x_i) = \text{sign}(0.059)$$

$$= 1$$

3. (4 pts) Assuming we stopped the algorithm after finishing the step $t = 4$, what is the predicted label for a new sample $\underline{x} = (0, 1)$? If the true label of this new test example is +1, what is the 0-1 loss for this example?

Predicted label = +1

\because True label = +1, and predicted label = +1,

0-1 loss for this example = 0.

Q4 (20 pts): SVM

Consider an SVM that obtains a decision boundary through the maximization optimization in Equation (1), which is equivalent to the minimization optimization in Equation (2), where β and β_0 denote the parameters, x_i the i -th training example, $y_i \in \{-1, 1\}$ the class label for the i -th training example, M the margin, N the number of instances.

$$\begin{aligned} & \max_{\beta, \beta_0, \|\beta\|_2=1} M \\ \text{subject to } & y_i (x_i^T \beta + \beta_0) \geq M, \quad i = 1, \dots, N. \end{aligned} \tag{1}$$

which is equivalent to

$$\begin{aligned} & \min_{\beta, \beta_0} \frac{1}{2} \|\beta\|_2^2 \\ \text{subject to } & y_i (x_i^T \beta + \beta_0) \geq 1, \quad i = 1, \dots, N. \end{aligned} \tag{2}$$

Mark true or false for the following assertions (Q4.1, Q4.2, Q4.3) and justify your answer. Questions with no justification will receive 0 points.

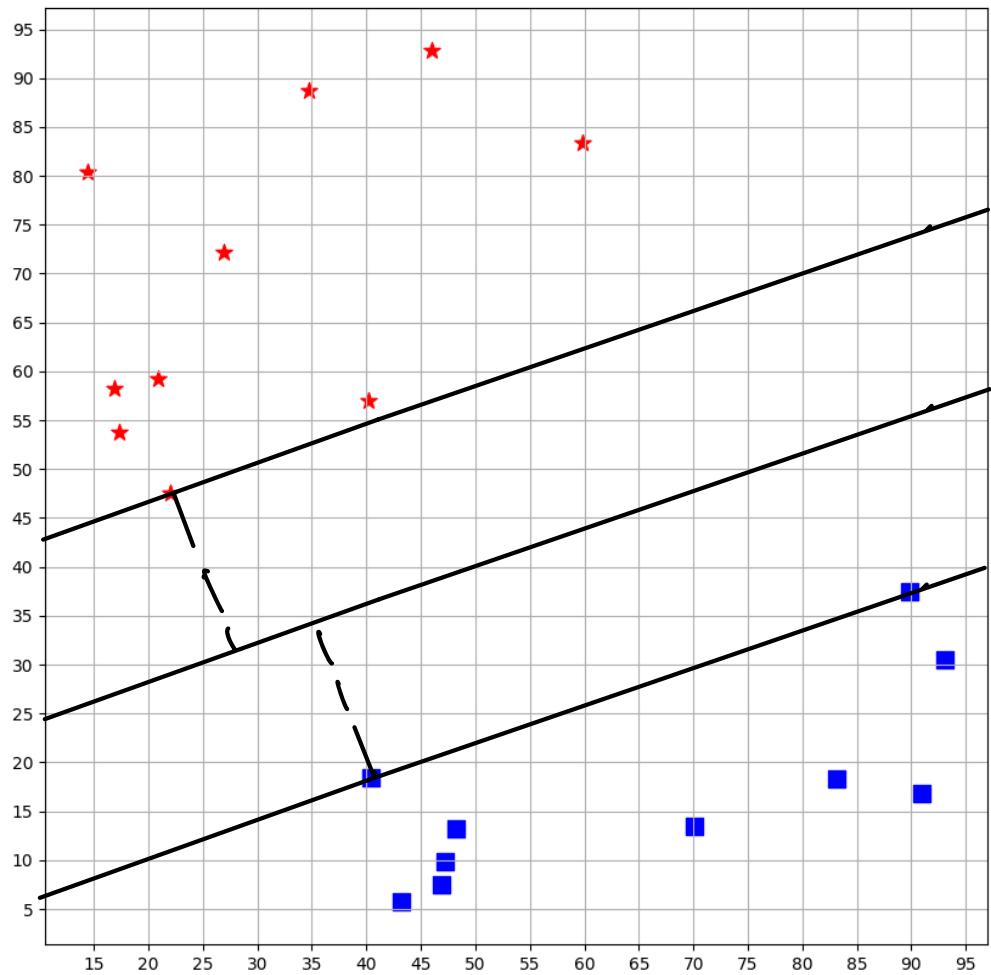
1. (4 pts) This SVM objective works for both linearly separable and non-linearly-separable data.

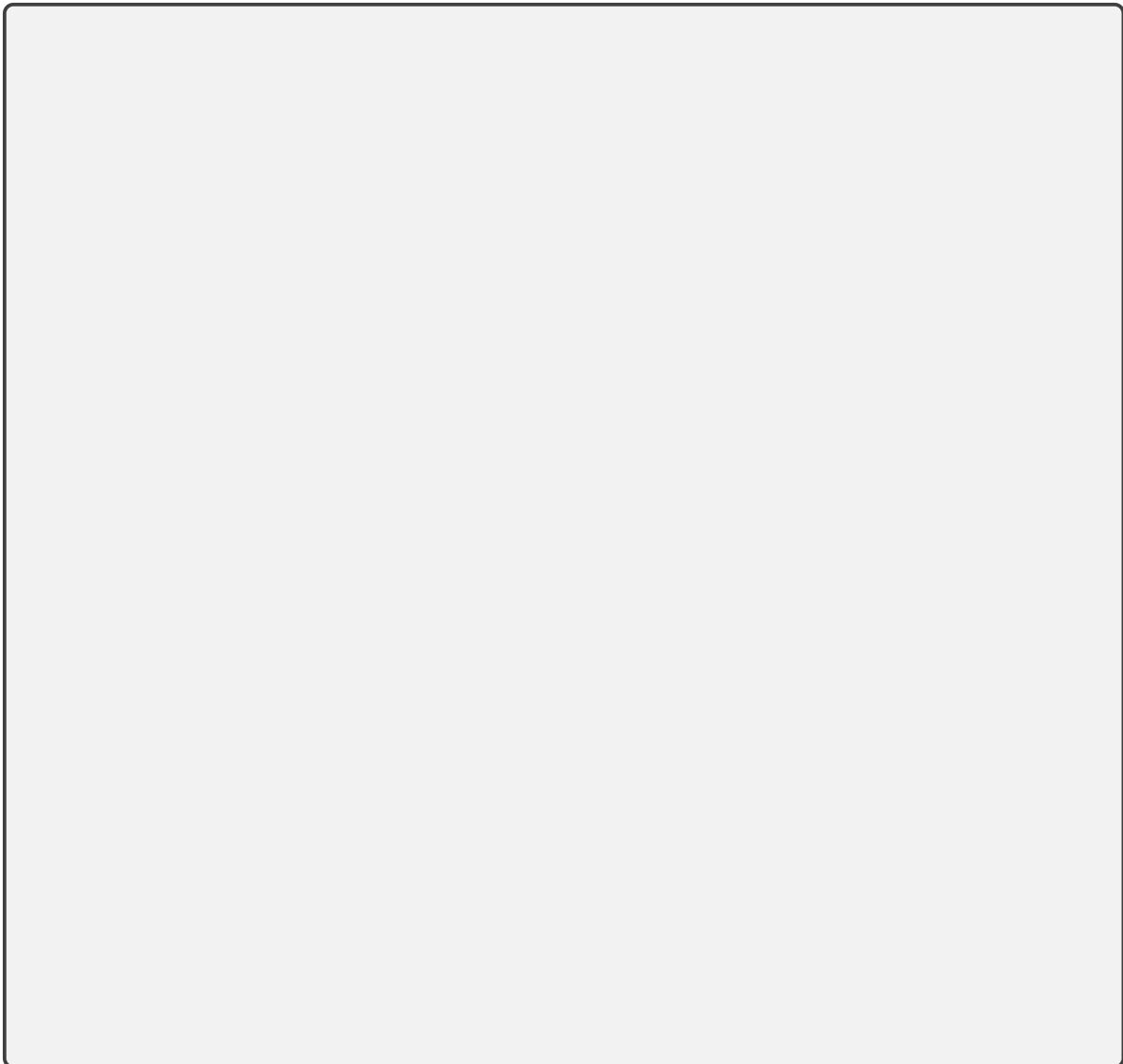
False. This SVM objective is for hard margin SVMs, which works only for linearly separable data and not non-linearly separable data.

2. (4 pts) After obtaining the decision boundary in Equation 1, if we modify β while maintaining the restriction $\|\beta\|_2^2 = 1$, this would modify the minimum distance from the decision boundary to the origin.

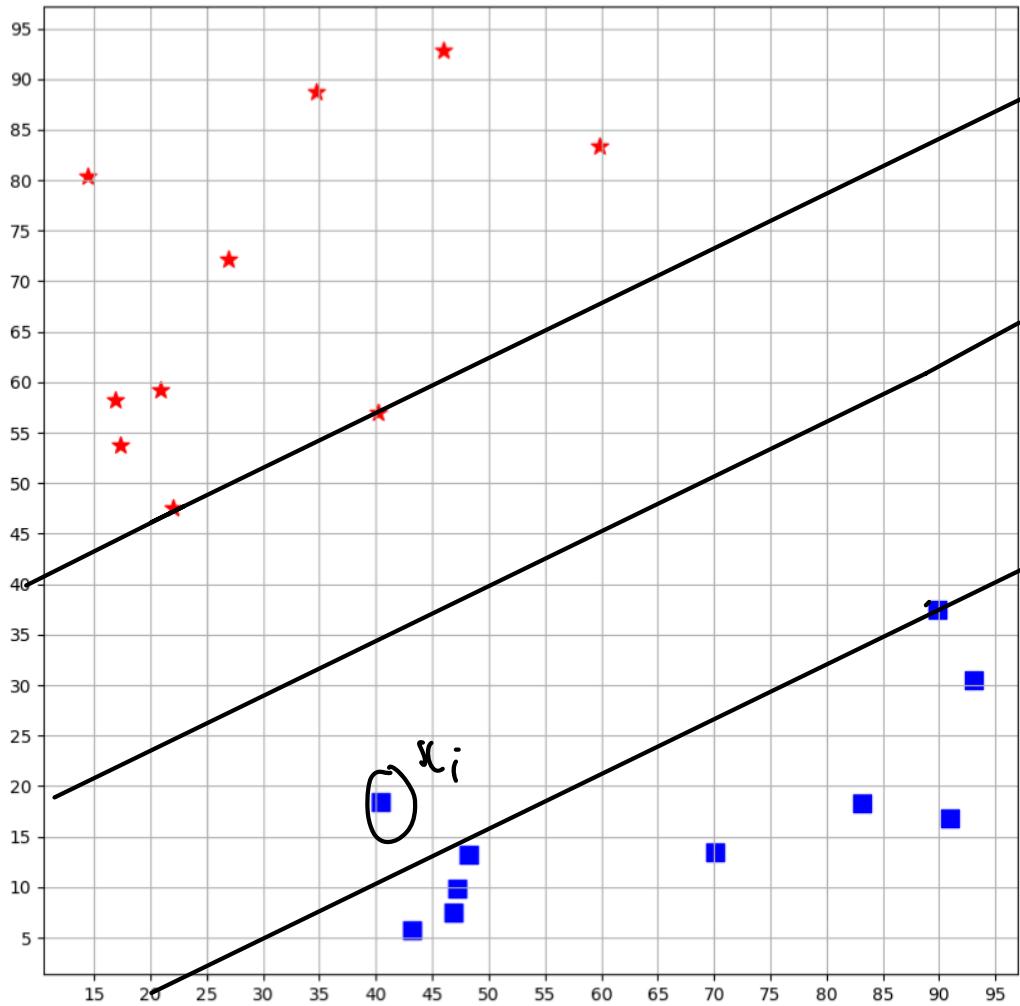
False. This is because as the norm $(\|\beta\|_2^2 = 1)$, this means that the distance from the decision boundary to the origin remains unchanged and thus, would not modify the minimum distance.

3. (4 pts) **SVM decision boundary.** Consider again an SVM whose decision boundary is obtained by Equations 1 and 2 like in the previous questions. Now, consider the training dataset with two training classes (signaled as squares and triangles) given by the scatter plot in the following figure. Draw the decision boundary for this dataset obtained by our SVM on it. The drawing needs to correctly separate the examples but there are multiple possible solutions. Each point of the line you draw must match the points of the true decision boundary within the error margin of 10 unit intervals along the x -axis and y -axis.





4. (4 pts) SVM decision boundary with removed example. Now pick a single data point \mathbf{x}_i that, if removed, would yield the maximum possible increase in the margin of a new SVM decision boundary over the dataset without \mathbf{x}_i . Indicate which point is \mathbf{x}_i in the figure below (reproduction of last question). Draw the new decision boundary after \mathbf{x}_i 's removal on it. The drawing needs to correctly separate the remaining examples but there are multiple possible solutions. Each point of the line you draw must match the points of the true new decision boundary within the error margin of 10 unit intervals along the x -axis and y -axis.
-



Q5 (18 pts): Bagging and Boosting

Bagging is generally used to help stabilize classifiers with unstable learning algorithms (optimal score searching algorithms). A classifier has a stable learning algorithm if, by changing the training data, the predicted class labels in the test data don't change. For instance, the predictions of a decision tree might significantly change with a small change in the training data. This definition depends on the amount of data, of course. Classifiers that are unstable with 10^3 training examples may be stable with 10^9 examples. Bagging works by aggregating the answers of unstable classifiers trained over multiple training datasets. These multiple datasets are often not independent, generally sampled with replacement from the same training data.

Boosting works by converting weak classifier (very simple models) to strong ones (models that can describe complex relationships between the inputs and the class labels). A weak learner is a classifier whose output of an test example attributes x_i is only slightly correlated with its true class t_i . That is, the weak learner classifies the data better than random, but not much better than random. In boosting, weak learners are trained sequentially in a way that the current learner gives more emphasis to the examples that past learns made mistakes on.

1. (8 pts) Suppose we decide to use an SVM classifier with a small training dataset. Assume the network can perfectly fit the training data but we want to make sure it is accurate in our test data (without having access to the test data). Would you use boosting or bagging to help improve the classification accuracy? Describe what would be the problem of using the other approach.

I would use bagging to help improve the classification accuracy. This is because using boosting may lead to overfitting. This is because the network can perfectly fit the training data, there would not be misclassified points for boosting to focus on. Thus, using boosting may result in placing more emphasis on outliers in the test data and eventually, cause the model to overfit. However, by using bagging, where an aggregate of answers from unstable classifiers is used, overfitting is less likely and thus, I would use bagging over boosting to help improve the classification accuracy.

2. (10 pts) Read the following code and answer the following questions:

```
1 def fit(self, X, Y):
2
3     if self.loss == "mse":
4         loss = MSELoss()
5     elif self.loss == "crossentropy":
6         loss = CrossEntropyLoss()
7
8
9     if self.classifier:
10        Y = to_one_hot(Y.flatten())
11    else:
12        Y = Y.reshape(-1, 1) if len(Y.shape) == 1 else Y
13
14    N, M = X.shape
15    self.out_dims = Y.shape[1]
16    self.learners = np.empty((self.n_iter, self.out_dims), dtype=object)
17    self.weights = np.ones((self.n_iter, self.out_dims))
18    self.weights[1:, :] *= self.learning_rate
19
20    Y_pred = np.zeros((N, self.out_dims))
21    for k in range(self.out_dims):
22        t = loss.base_estimator()
23        t.fit(X, Y[:, k])
24        Y_pred[:, k] += t.predict(X)
25        self.learners[0, k] = t
26
27    for i in range(1, self.n_iter):
28        for k in range(self.out_dims):
29            y, y_pred = Y[:, k], Y_pred[:, k]
30            neg_grad = -1 * loss.grad(y, y_pred)
31
32            t = DecisionTree(
33                classifier=False, max_depth=self.max_depth, criterion="mse"
34            )
35
36            t.fit(X, neg_grad)
37            self.learners[i, k] = t
38
39            step = 1.0
40            h_pred = t.predict(X)
41            if self.step_size == "adaptive":
42                step = loss.line_search(y, y_pred, h_pred)
43
44            self.weights[i, k] *= step
45            Y_pred[:, k] += self.weights[i, k] * h_pred
46
47    def predict(self, X):
48        Y_pred = np.zeros((X.shape[0], self.out_dims))
49        for i in range(self.n_iter):
50            for k in range(self.out_dims):
51                Y_pred[:, k] += self.weights[i, k] * self.learners[i, k].predict(X)
52
53        if self.classifier:
54            Y_pred = Y_pred.argmax(axis=1)
55
56    return Y_pred
```

- (a) (5 pts) Which classifier is this? Why are the trees used in this classifier so shallow? Describe how the classifier works using pseudo-code.

Hint: It uses decision trees but this is not a decision tree classifier.

This is an AdaBoost classifier. The trees used in this classifier are shallow because they are decision stumps which are weak learners. This classifier works by the following pseudo-code:

AdaBoost():

Input: Data set $S = \{(x_1, y_1), \dots, (x_n, y_n)\}$ (where $y_i \in \{-1, 1\}$)
and weak learning algorithm A

For $t = 1, 2, \dots, T$
Configure reweighting D_t for data
Run $A(D_t)$ to get weak classifier h_t ,
Carry over weight $w_t \in \mathbb{R}$

- (b) (5 pts) What happens if we have mislabeled data? Why mislabeled data could be a problem?
Hint: We are looking for an answer that uses the training weights of the examples at each iteration as a justification.

When we have mislabeled data, the model would mostly misclassify such data points, leading to an increase in the training weights of such examples at every iteration of AdaBoost. This could lead to the model overfitting to such mislabeled data. Eventually, this would lead to the model being biased to such mislabeled data, and potentially performing poorly on correctly labelled data.

Q5(a) continued

Output: classifier h defined as

$$\hat{h}(x) = \text{sign} \left(\sum_{t=1}^T w_t h_t(x) \right)$$

Thus, D_t increases weight on points where h_{t-1} misclassifies, and decreases weight on points where h_{t-1} correctly classifies.

Programming Part

Overview

In this assignment, you will implement a one-hot encoder, a bagging classifier, and a cross-validation class. You will use these utilities to classify the bank loan dataset.

Files

You will fill in functions in `encoder.py`, `bagging.py`, and `cv.py`. After you finish, you will submit these as well as any other files mentioned in the assignment to Gradescope.

Asides from these three files, you will also find `utils.py` in the same folder. You do not need to modify it. `utils.py` contains some helper functions.

Packages

You will need the following packages for this assignment:

- `numpy`
- `pandas`
- `scikit-learn`

If you want, you can install `tqdm` to see a progress bar when running the autograder. This is not required for this class and might not be heavily focused on in future assignments, but it is a useful package to know about if you are doing any machine learning work in Python.

These packages should be installed if you have installed Anaconda. If you are not using Anaconda, you can install them using `pip` or `conda`. For example, you can run `pip install numpy` in the terminal to install `numpy`.

Note that for this assignment, you should **not** use any other packages. You may see some other packages in the skeleton code, but they are only used for testing and grading. If you use any other packages, you may lose points. If you are not sure whether you can use a package, please ask the course staff.

Evaluation

Unless otherwise specified, your code will be evaluated by a autograder on Gradescope using the same environment as detailed below. You should make sure your code runs without errors in this environment. Otherwise, you may lose points. You can submit your code to Gradescope as many times as you want. We will only grade the latest submission.

There will be two types of test cases in the autograder: public (local) and hidden. Public test cases are visible to you, and you can see the results after you submit your code. You can also run the public test cases locally by running `python <filename>.py` in the same folder as your code. Passing all public test cases **does not** guarantee you will get full credit for the assignment.

Hidden test cases, however, are not visible to you, and you will not be able to see the results until your grade is published. The final score you get for this assignment is the sum of the scores of all public and hidden test cases.

Getting Help

If you have any questions about this assignment, please contact the course staff for help, preferably during office hours. You can also post your questions on the course forum. However, please do not post any code publicly. When asking questions, you should describe your problem in words and post the relevant code snippet. Please avoid showing a screenshot of your code to TAs and ask “why my code is not working”.

1 One-Hot Encoder

In this part, you will implement a one-hot encoder in `encoder.py`. You will use this encoder in Part 2 to encode the categorical features in the bank loan dataset.

In machine learning, we often need to encode categorical features into numerical features so that we can use them in our models. Suppose we have a dataset of students, and each student has a major. The major is a categorical feature. We can encode the major into a numerical feature by assigning a number to each major. For example, we can assign 0 to “Computer Science”, 1 to “Mathematics”, and 2 to “Chemistry”.

However, this encoding is not ideal because it implies that “Computer Science” is closer to “Mathematics” than “Chemistry” if we calculate the “distance” between the two majors. This is not true because the majors are categorical features, and there is no natural ordering between them. To avoid this problem, we can use a one-hot encoder.

A one-hot encoder encodes a categorical feature into a vector of binary features. In our previous example, we can encode the major into a vector of length 3, where the first element is 1 if the major is “Computer Science” and 0 otherwise, the second element is 1 if the major is “Mathematics” and 0 otherwise, and the third element is 1 if the major is “Chemistry” and 0 otherwise. This encoding is called a one-hot encoding because only one element in the vector is 1 and the rest are 0.

1.1 One-Hot Encoder - Constructor

Under `class OneHotEncoder`, fill in the constructor so that it initializes the attribute `self.categories`, `self.data`, and `self.encoded_data` to `None`. We will initialize them later in the `fit` method.

1.2 One-Hot Encoder - Fit

Under `class OneHotEncoder`, fill in the `fit` method so that it takes in a numpy array `data` and save it to `self.data` and initializes `self.categories` to an empty list.

`data` is a 2D numpy array, where each row is a data point and each column is a feature. For example, if we have 3 data points and 2 features, `data` will be a 3 by 2 numpy array. You can assume that `data` only contains categorical features.

`self.categories` is a list of lists (or a list of numpy arrays). You should iterate through each column in `data` and find all the unique values in that column. For example, if the first column in `data` is [“Com-

puter Science”, “Mathematics”, “Mathematics”], then `self.categories[0]` should be [“Computer Science”, “Mathematics”]. You can use `np.unique` to find the unique values in a numpy array.

You may assume that the input data is valid and does not contain any missing values or unknown categories. You do not need to handle these cases.

1.3 One-Hot Encoder - Encode

Under `class OneHotEncoder`, fill in the `encode` method so that it takes in a numpy array `data` and returns a numpy array `encoded_data`. To encode, you should iterate through each column in `data` and create a two dimensional numpy array `encoding_array` of shape `(len(data), len(self.categories[i]))`.

For example, if the first column in `data` is [[“Computer Science”], [“Mathematics”], [“Mathematics”]], then `encoding_array` should be a 3 by 2 numpy array where the first row is [1, 0], the second row is [0, 1], and the third row is [0, 1]. You should use `self.categories` to find the index of the category in each row and set the corresponding element in `encoding_array` to 1.

You may assume that the input data is valid and does not contain any missing values or unknown categories. You do not need to handle these cases.

2 Bagging

In this part, you will implement a bagging classifier in `bagging.py`. You will use this classifier in Part 3 to train a bagging classifier on the bank loan dataset.

Bagging is a machine learning ensemble meta-algorithm designed to improve the stability and accuracy of machine learning algorithms used in statistical classification and regression. It also reduces variance and helps to avoid overfitting. Although it is usually applied to decision tree methods, it can be used with any type of method.

Bagging leads to “improvements for unstable procedures”, which include, for example, artificial neural networks, classification and regression trees, and subset selection in linear regression. On the other hand, it can mildly degrade the performance of stable methods such as K-nearest neighbors.

In this assignment, we will use bagging to improve the performance of a decision tree classifier. We will train multiple decision tree classifiers on different subsets of the training set and combine their predictions to make the final prediction.

2.1 Bagging - Constructor

Under `class BaggingClassifier`, fill in the constructor so that it takes in an integer `n_learners`, a string representing the encoding type for the decision tree, and a random seed `random_seed`. It should initialize the attribute `self.encoder` to an encoder object specified by the encoding type and initialize the attribute `self.random_seed` to the random seed. It should also create a list of `DecisionTreeClassifier` objects with `random_state` set to `self.random_seed` and save the list to `self.learners`.

2.2 Bagging - Preparing Samples

Under `class BaggingClassifier`, fill in the `_prepare_samples` method so that it takes in a numpy array `X` and returns a list of numpy arrays containing the indices of the samples randomly selected with replacement. You should use `self.random_state` to initialize the random state before sampling. You may find `np.random.choice` useful.

2.3 Bagging - Fit

Under `class BaggingClassifier`, fill in the `fit` method so that it takes in a numpy array `X` and a numpy array `y` and trains `self.n_learners` decision tree classifiers on different subsets of the training set. You should use `self.encoder` to first fit and encode the training data. For each learner, you should use the `_prepare_samples` method to get the indices of the samples and use the `fit` method of the `DecisionTreeClassifier` class to train each decision tree classifier.

2.4 Bagging - Predict

Under `class BaggingClassifier`, fill in the `predict` method so that it takes in a numpy array `X` and returns a numpy array `y_pred` containing the predictions of the bagging classifier. You should use `self.encoder` to encode the test data first. For each learner, you should use the `predict` method of the `DecisionTreeClassifier` class to get the predictions of each decision tree classifier. You should then combine the predictions of all the decision tree classifiers to get the final prediction of the bagging classifier. You may find `np.bincount` and `np.argmax` useful.

3 Cross Validation

In this part, you will implement cross validation in `cv.py`. You will use this function to find the best hyperparameters for the bagging classifier you implemented in Part 2.

Cross validation is a model validation technique for assessing how the results of a statistical analysis will generalize to an independent data set. It is mainly used in settings where the goal is prediction, and one wants to estimate how accurately a predictive model will perform in practice.

In a prediction problem, a model is usually given a dataset of known data on which training is run (training dataset), and a dataset of unknown data (or first seen data) against which the model is tested (testing dataset). The goal of cross validation is to define a dataset to “test” the model in the training phase (i.e., the validation dataset), in order to limit problems like overfitting, give an insight on how the model will generalize to an independent dataset (i.e., an unknown dataset, for instance from a real problem), etc.

One round of cross-validation involves partitioning a sample of data into complementary subsets, performing the analysis on one subset (called the training set), and validating the analysis on the other subset (called the validation set or testing set). To reduce variability, multiple rounds of cross-validation are typically performed using different partitions, and the validation results are averaged over the rounds.

The same kind of machine learning model can require different constraints, weights or learning rates to generalize different data patterns. These measures are called hyperparameters and have to be tuned so that the model can optimally solve the machine learning problem. Hyperparameters are set before the learning process begins while the values of other parameters are derived via training. The optimization or the tuning of hyperparameters is the problem of choosing a set of optimal hyperparameters for a learning algorithm.

3.1 Cross Validation - Constructor

Under `class CrossValidation`, fill in the constructor so that it takes in an integer `k` and a random seed `seed`. It should initialize the attribute `self.k` to `k` and initialize the attribute `self.random_seed` to the random seed.

3.2 Cross Validation - Split

Under `class CrossValidation`, fill in the `split` method so that it takes in a numpy array `X` and a numpy array `y` and splits the data into `self.k` folds. It should return a list of numpy arrays containing the samples in each fold. You should use `self.random_seed` to initialize the random state before splitting. Use `np.random.permutation` to randomly generate the indices of the samples and use `np.array_split` to split the indices into `self.k` folds.

3.3 Cross Validation - Combine Folds

Under `class CrossValidation`, fill in the `_combine_folds` method so that it takes in two lists of numpy arrays `X_folds` and `y_folds` and an integer `i`. It should combine all the folds except the `i`th fold into a training set and return the training set and the `i`th fold as a validation set. You should use `np.concatenate` to combine the folds.

3.4 Cross Validation - Score

Under `class CrossValidation`, fill in the `score` method so that it takes in two lists of numpy arrays `X` and `y`, the features and labels, and our Bagging `model`. It should use the model to predict the labels of the features and return the accuracy of the predictions.

3.5 Cross Validation - Cross Validate

Under `class CrossValidation`, fill in the `cross_validate` method so that it takes in two lists of numpy arrays `X` and `y`, the features and labels, and our Bagging `model`. It should use the `split` method to split the data into `self.k` folds and use the `_combine_folds` method to combine the folds into training and validation sets for each fold. It should then fit the model on the training set and use the `score` method to get the accuracy of the model on the validation set. Return the average accuracy of the model on the validation sets and the list of accuracies on the validation sets.

3.6 Cross Validation - Get Best Model

Under `class CrossValidation`, fill in the `get_best_model` method so that it takes in two lists of numpy arrays `X` and `y`, the features and labels, and a list of hyperparameters `params`. For each hyperparameter in `params`, it should create a Bagging model with that hyperparameter and use the `cross_validate` method to get the average accuracy of the model on the validation sets. It should then return the model with the best hyperparameter, the best hyperparameter, and the list of accuracies on the validation sets.

3.7 Cross Validation - Visualize

Under `class CrossValidation`, fill in the `plot_learning_curve` method so that it creates a plot of the learning curve of the model. The x-axis should be the hyperparameters and the y-axis should be the average accuracy of the model on the validation sets using each hyperparameter. You should use `plt.plot` to plot the learning curve. You should also use `plt.xlabel` and `plt.ylabel` to label the x-axis and y-axis. You should use `plt.title` to give the plot a proper title. You should use `plt.savefig` to save the plot as `learning_curve.png`. Return the plot.

Attach the plot to the Gradescope submission entry for plots. If plots are submitted as a part of the code submission and not submitted through the gradescope link for plots, then they won't be graded..

Submission

You will submit the following files to Gradescope:

- `encoder.py`
- `bagging.py`
- `cv.py`
- `learning_curve.png`