

## Homework-1 Answers

Author: Sayantani Bhattacharya

////////////////////////////////////  
**Question1:**

For each loss, answer the following questions. You must provide an explanation for full points.

**1. Perceptron Loss**

a. For this loss, which point(s) have the highest loss value? Why?

Points where  $y \cdot g(x) \leq 0$ , will have the highest loss of 1.

From the table:

- Point (1, 2):  $y=1, g(X)=-0.7 \Rightarrow y \cdot g(X)=-0.7 < 0$
- Point (2, 1):  $y=-1, g(X)=1.5 \Rightarrow y \cdot g(X)=-1.5 < 0$
- **Thus Highest Loss Points:** (1, 2) and (2, 1).

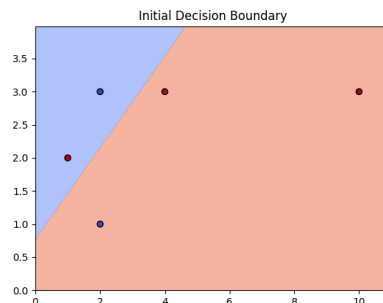
b. For this loss, which point(s) have the lowest loss value? Why?

Points where  $y \cdot g(x) > 0$  have the lowest loss of 0.

From the table:

- Point (2, 3):  $y=-1, g(X)=-1.1 \Rightarrow y \cdot g(X)=1.1 > 0$
- Point (4, 3):  $y=1, g(X)=0.7 \Rightarrow y \cdot g(X)=0.7 > 0$
- Point (10, 3):  $y=1, g(X)=6.1 \Rightarrow y \cdot g(X)=6.1 > 0$
- **Thus, Lowest Loss Points:** (2, 3), (4, 3), and (10, 3).

c. With this dataset and loss, is it possible to find a different set of model parameters that would have a lower total loss? Don't calculate the loss for any specific  $w$  values; just look at the plot and use your intuition. You may answer "Yes", "No", or "Maybe". If you answer "Yes", give a general description of what the new decision boundary would look like and why that would decrease the loss (e.g., "if you rotated the current boundary clockwise, then ..."). If you answer "No", give a general argument for why there is no decision boundary that could decrease the loss. If you answer "Maybe", give a general description of a new decision boundary that might decrease the loss, but it's hard to tell without doing the calculations. Also, provide a general argument for why there's no decision boundary that obviously decreases the loss.



**Answer:** Yes. For instance, observing the above decision boundary, rotating the decision boundary slightly counterclockwise could make (1, 2) correctly classified. A better alignment may increase the number of correct classifications.

d. Is this loss function a good choice for training a multilayer perceptron on a binary classification task? Why or why not?

////////////////////////////////////

- ////////////////////////////////////

////////////////////////////////////

b. Points with  $y \cdot g(x) \geq 1$  have the lowest loss (0).

- Point (10, 3):  $y=1, g(X)=6.1 \Rightarrow y \cdot g(X)=6.1 > 1y = 1$
- **Lowest Loss Point:** (10, 3).

c. Yes. Adjusting  $w$  to push incorrect points like  $(1, 2)$  and  $(2, 1)$  further away from the decision boundary could lower total hinge loss.

d. Yes. Hinge loss enforces a margin between classes, improving the generalization of the classifier.

////////////////////////////////////

### Question2:

This single-layer perceptron focuses on implementing different loss functions, each influencing how weight updates occur. It consists of a simple linear classifier with weights  $w$  and bias  $b$ , trained using gradient descent. With four loss functions, **perceptron loss**, **squared error loss**, **binary cross-entropy**, and **hinge loss**. The **learning rate ( $lr = 0.01$ )** and **100 training epochs** were chosen as reasonable defaults, balancing convergence speed with stability. The **sigmoid function** is used in binary cross-entropy to model probabilistic decision boundaries, while hinge loss encourages margin maximization, similar to SVMs. The perceptron loss updates weights only for misclassified points, while squared error loss can lead to gradient explosion, requiring careful tuning. While other learning methods like the standard **perceptron algorithm** could be used, it is only suitable for **linearly separable data** and does not converge for non-linearly separable datasets. By contrast, **gradient descent allows for smoother updates**, enabling learning even in more complex scenarios. However, the perceptron algorithm is computationally cheaper, updating weights only when a misclassification occurs, while gradient-based methods continuously update weights, leading to potentially better decision boundaries. The choice of the method depends on the dataset characteristics and whether the goal is strict classification (perceptron) or margin-based optimization (hinge loss, cross-entropy).

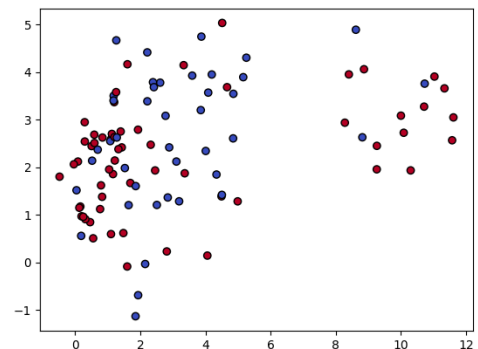
### Results with Augmented data (Type1: gaussian):

Dataset representative of the initial 5 samples provided.

It is generated by adding gaussian noise to

randomly chosen points from the initial sample points.

(Attached as augmented\_data\_gaussian.csv)

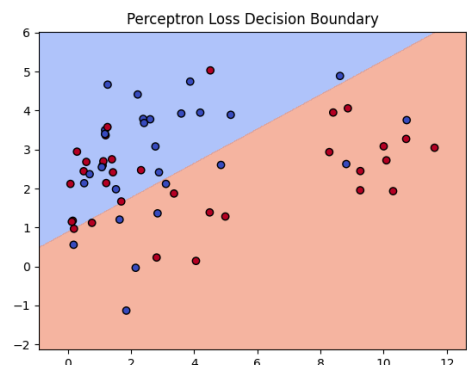


### Training with Perceptron Loss...

Weights: [ 0.39676684 -0.90201702], Bias: 0.7999999999999999

Validation Accuracy: 30.00%

Test Accuracy: 60.00%



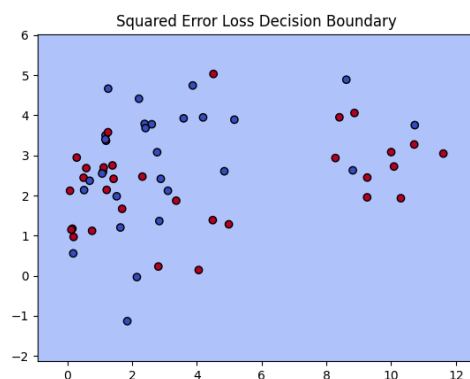
### Training with Squared Error Loss...

Weights: [nan nan], Bias: nan

Validation Accuracy: 0.00%

Test Accuracy: 0.00%

Squared error loss is **not ideal** for classification since it can amplify large errors too aggressively. If the value of  $y - g(x)$  is too large, as in this case further operations with these numbers often result in nan, leading to complete training failure.

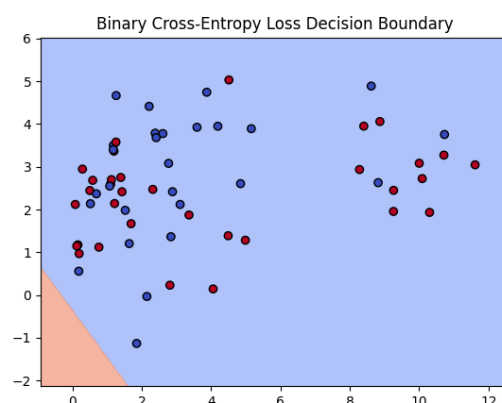


### Training with Binary Cross-Entropy Loss...

Weights: [-420.36778968 -379.99504201], Bias: -134.53133908768777

Validation Accuracy: 35.00%

Test Accuracy: 35.00%

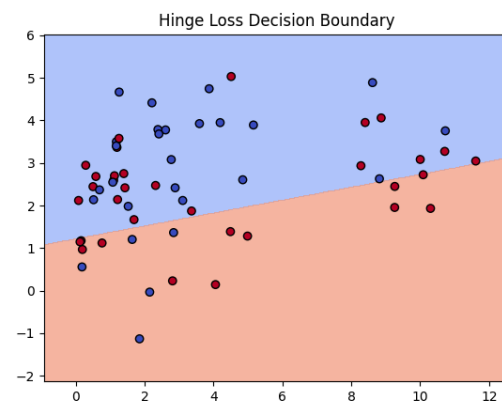


### Training with Hinge Loss...

Weights: [ 0.23647438 -1.56055594], Bias: 1.9000000000000004

Validation Accuracy: 35.00%

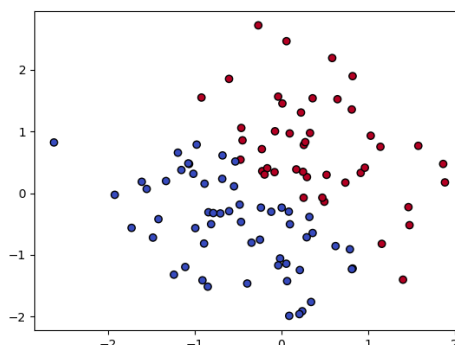
Test Accuracy: 50.00%



## Results with Augmented data (Type2: uniform):

Uniform separable dataset:

(Attached as augmented\_data\_uniform.csv file)

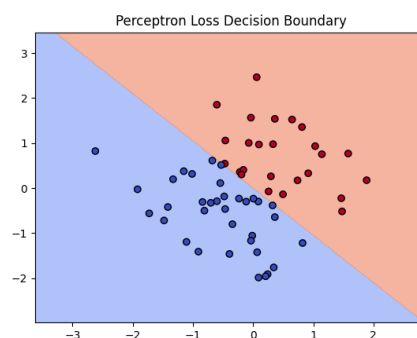


### Training with Perceptron Loss...

Weights: [0.26332301 0.2506396 ], Bias: 0.0

Validation Accuracy: 100.00%

Test Accuracy: 100.00%

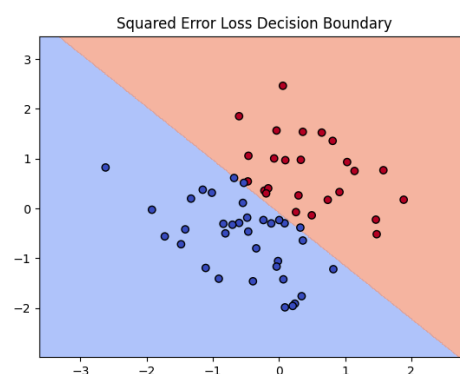


### Training with Squared Error Loss...

Weights: [0.53567998 0.50313879], Bias: 0.0442054554956561

Validation Accuracy: 95.00%

Test Accuracy: 100.00%

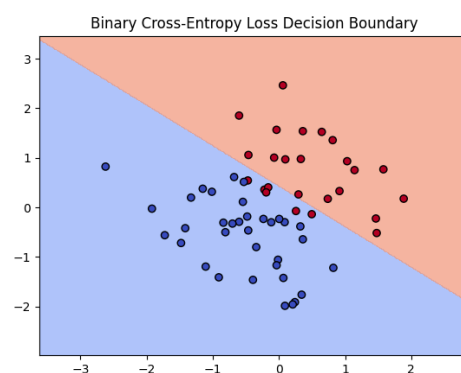


### Training with Binary Cross-Entropy Loss...

Weights: [151.85025306 185.67790979], Bias: -79.51072587217323

Validation Accuracy: 95.00%

Test Accuracy: 85.00%

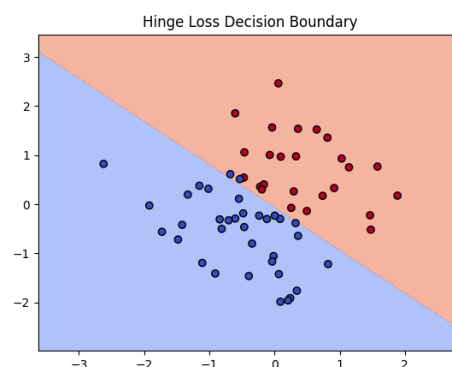


### Training with Hinge Loss...

Weights: [4.39791677 5.02188899], Bias: 0.30000000000000004

Validation Accuracy: 95.00%

Test Accuracy: 100.00%



////////////////////////////////////

### Question 3:

a. In the `run_one_epoch` function, the code uses a combination of `where`, `argmax`, and `torch.nn.CrossEntropyLoss`. How do these calculations work with the model's output to compute the loss and accuracy? Why does the model output a real-valued tensor of shape  $(N, 2)$ ? You may want to look at the documentation for the loss function.

Answer:

The `run_one_epoch` function executes one training or evaluation step (based on the value of the `train` variable) for the neural network model. It takes in a model, optimiser, input data ( $X$ ) and label ( $y$ ). The function executes a forward pass on the input, using the model, and stores the output.

Since  $y$  contains values that may not be binary, `torch.where(y > 0, 1, 0)` maps all positive values to 1 and all non-positive values to 0. This ensures the labels are in the correct format for `CrossEntropyLoss`, which expects class indices (0 or 1). The `argmax` is used to pick the class with the highest score as the model's prediction. Which is then compared against the label, to calculate the accuracy of the model. Given logits  $(N, 2)$ , the `torch.nn.CrossEntropyLoss` applies the softmax function to get probabilities and computes the negative log-likelihood of the correct class.

In binary classification, the model outputs two logits per sample (one for each class). PyTorch's `CrossEntropyLoss` requires logits of shape  $(N, C)$ , where  $C=2$  for two classes. The two outputs correspond to logits before softmax, and `CrossEntropyLoss` automatically applies softmax internally. This formulation allows the model to generalize easily to multi-class problems where  $C$  is greater than 2.

b. Read through the `run_experiment`, `pretrain_and_train`, and `plot_results` functions to understand how the figure is being built. In your own words, describe in general what the six panels show. What data is being used in which panels?

Answer:

The six panels:

1. Pretrained model on pretrain data

This panel visualizes the decision boundary of the model after pretraining on the pretraining dataset. Input:  $X_{\text{pretrain}}$ , Label:  $y_{\text{pretrain}}$ . It shows whether the model has learned any structure from the pre-training phase before being trained on the main dataset.

2. Pretrained model on train data

This panel visualizes how the pretrained model performs on the training dataset (before any training on this dataset). Input:  $X_{\text{train}}$ , Label:  $y_{\text{train}}$ . This gives insight into how well the pretrained model generalizes to a different dataset before fine-tuning.

3. Trained model on train data

This panel visualizes the decision boundary of the fully trained model on the training dataset. Input:  $X_{\text{train}}$ , Label:  $y_{\text{train}}$ . It demonstrates how well the model has learned the training data, i.e., whether it has effectively captured the patterns in the dataset.

4. Trained Model on Test Data

This panel visualizes the decision boundary of the trained model on test data. Input:  $X_{\text{test}}$ , Label:  $y_{\text{test}}$ . This helps determine the generalization ability of the model. A good model should perform well on unseen test data.

5. Loss per epoch

A line plot of the training loss and test loss over the training epochs. Loss should decrease over time. Train loss always decreases over time, while testset might not. If the test loss increases while training loss continues decreasing, the model may be overfitting.

6. Accuracy per Epoch

////////////////////////////////////

a. Overfitting Experiments: Include at least two experiments that show the model overfitting. Describe what those experiments show and how you know the model is overfitting. What arguments had the most effect on whether your model overfits?

Experiment 2 is the classic example for overfitting. Here I increased the model layer size to [500,100]. This increased the model's ability to memorize training data but led to poor generalization on the test set. The results showed that training accuracy approached **100%**, while test accuracy stagnated around **40-50%**, indicating overfitting. Sample Experiment and Experiment 1 are also examples of overfitting, observing the last two plots using the above statement.

Answer: Pretraining was tested by varying the amount of pretraining data, epochs, and pretraining quality to observe its impact on test accuracy. In Experiment 5, when I set `n_pretrain_epochs = 0`, the model trained directly on the spiral dataset. It achieved a final mean test set accuracy of 38.5%, suggesting the model struggled to learn from scratch. Especially when compared to Experiment 7, when I set `n_pretrain_epochs = 100`. The model started with helpful knowledge but remained flexible enough to adapt. This setup resulted in the **highest final mean test accuracy of 61.5%**. And finally in Experiment 6, I assigned `n_pretrain_epochs` to 5000 (a very high value.) This caused the model to over-specialize on the ring dataset, making it harder to adapt to the spiral data. Test accuracy **dropped to 30.8%**.

c. Best Experiment: Across all the experiments you ran, with which arguments did you achieve the highest Final Mean Test Accuracy? What patterns led you to find these arguments? Which arguments had the largest impact on your

