



Table of Contents

Lab 02: Gradient Descent

- 1. Loss landscape
- 2. The “Gradient” in Gradient Descent

3. Forward & Backward

- 3.1. Forward
- 3.2. Backward
- 4. Implementation
 - 4.1. Import library
 - 4.2. Create data
 - 4.3. Training

Lab 02: Gradient Descent

Solution template

Copyright © Department of Computer Science, University of Science, Vietnam National University,
Ho Chi Minh City

- Student name: Phan Huy Đức Tài
- ID: 21127687

How to do your homework

- You will work directly on this notebook; the word **TODO** indicates the parts you need to do.
- You can discuss the ideas as well as refer to the documents, but *the code and work must be yours*.

How to submit your homework

- Before submitting, save this file as `<ID>.jl`. For example, if your ID is 123456, then your file will be `123456.jl`. And export to PDF with name `123456.pdf` then submit zipped source code and pdf into `123456.zip` onto Moodle.

Note

Note that you will get 0 point for the wrong submit.

Content of the assignment:

- Gradient Descent

1. Loss landscape

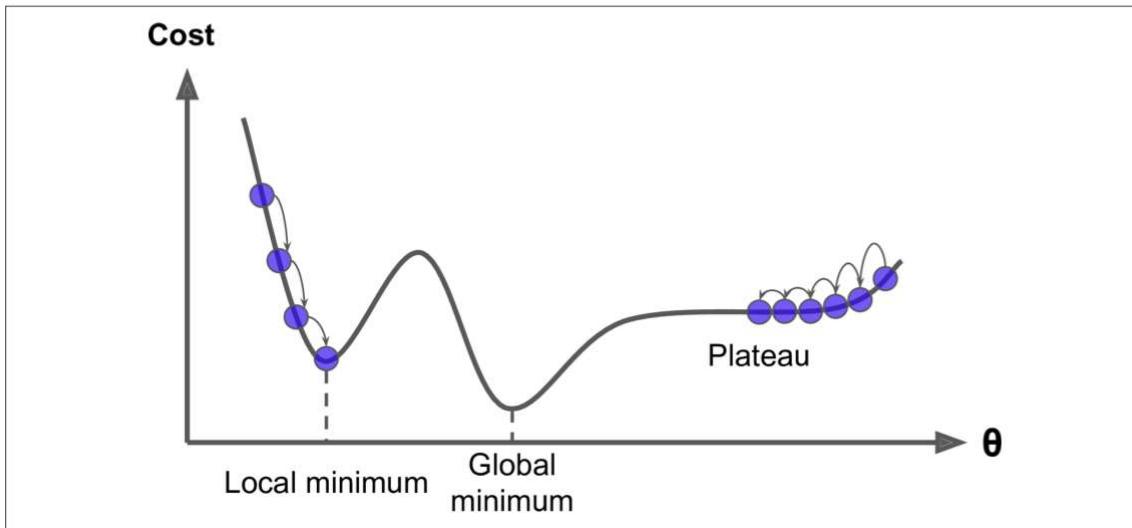


Figure 1. Loss landscape visualized as a 2D plot. Source: codecamp.vn

The gradient descent method is an iterative optimization algorithm that operates over a loss landscape (also called an optimization surface). As we can see, our loss landscape has many peaks and valleys based on which values our parameters take on. Each peak is a local maximum that represents very high regions of loss – the local maximum with the largest loss across the entire loss landscape is the global maximum. Similarly, we also have local minimum which represents many small regions of loss. The local minimum with the smallest loss across the loss landscape is our global minimum. In an ideal world, we would like to find this global minimum, ensuring our parameters take on the most optimal possible values.

Each position along the surface of the corresponds to a particular loss value given a set of parameters **W** (weight matrix) and **b** (bias vector). Our goal is to try different values of **W** and **b**, evaluate their loss, and then take a step towards more optimal values that (ideally) have lower loss.

2. The “Gradient” in Gradient Descent

We can use \mathbf{W} and \mathbf{b} and to compute a loss function L or we are able to find our relative position on the loss landscape, but **which direction** we should take a step to move closer to the minimum.

- All We need to do is follow the slope of the gradient $\nabla_{\mathbf{W}}$. We can compute the gradient $\nabla_{\mathbf{W}}$ across all dimensions using the following equation:

$$\frac{df(x)}{dx} = \lim_{h \rightarrow 0} \frac{f(x + h) - f(x)}{h}$$

- But, this equation has 2 problems:
 - 1. It's an **approximation** to the gradient.
 - 2. It's painfully slow.

In practice, we use the **analytic gradient** instead.

3. Forward & Backward

~~In this section, you will be asked to fill in the black to form the forward process and backward process with the data defined as follows:~~

- Feature: \mathbf{X} (shape: $n \times d$, be already used bias trick)
- Label: \mathbf{y} (shape: $n \times 1$)
- Weight: \mathbf{W} (shape: $d \times 1$)

3.1. Forward

TODO: Consider one sample \mathbf{x}_i . Fill in the blank

$$h_i = \mathbf{x}_i^T \mathbf{W} \Rightarrow \frac{\partial h_i}{\partial \mathbf{W}} = \mathbf{x}_i^T$$

$$\hat{y}_i = \sigma(h_i) \Rightarrow \frac{\partial \hat{y}_i}{\partial h_i} = \hat{y}_i * (1 - \hat{y}_i)$$

$$loss_i = (\hat{y}_i - y_i)^2 \Rightarrow \frac{\partial loss_i}{\partial \hat{y}_i} = 2(\hat{y}_i - y_i) = 2\hat{y}_i - 2y_i$$

3.2. Backward

Our loss function is MSE:

$$Loss = \frac{1}{n} \sum_{i=1}^n loss_i = \frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i)^2$$

Goal: Compute $\nabla Loss = \frac{\partial loss(W)}{\partial W}$

How to compute $\nabla Loss$?: Use Chain-rule. Your work is to fill in the blank

TODO: Fill in the blank

$$\begin{aligned}\nabla Loss &= \frac{\partial Loss(W)}{\partial W} = \frac{1}{n} \sum_{i=1}^n \frac{\partial loss(W)}{\partial \hat{y}_i} \frac{\partial \hat{y}_i}{\partial h_i} \frac{\partial h_i}{\partial W} \\ &= \frac{2}{n} \sum_{i=1}^n (\hat{y}_i - y_i)(\hat{y}_i - \hat{y}_i^2)x_i^T \\ &= \frac{2}{n} \sum_{i=1}^n \hat{y}_i^2 x_i^T - \hat{y}_i^3 x_i^T - \hat{y}_i y_i x_i^T + \hat{y}_i^2 y_i x_i^T\end{aligned}$$

4. Implementation

4.1. Import library

```
1 using Distributions, Plots, LinearAlgebra, Random
```

```
MersenneTwister(2024)
```

```
1 Random.seed!(2024)
```

4.2. Create data

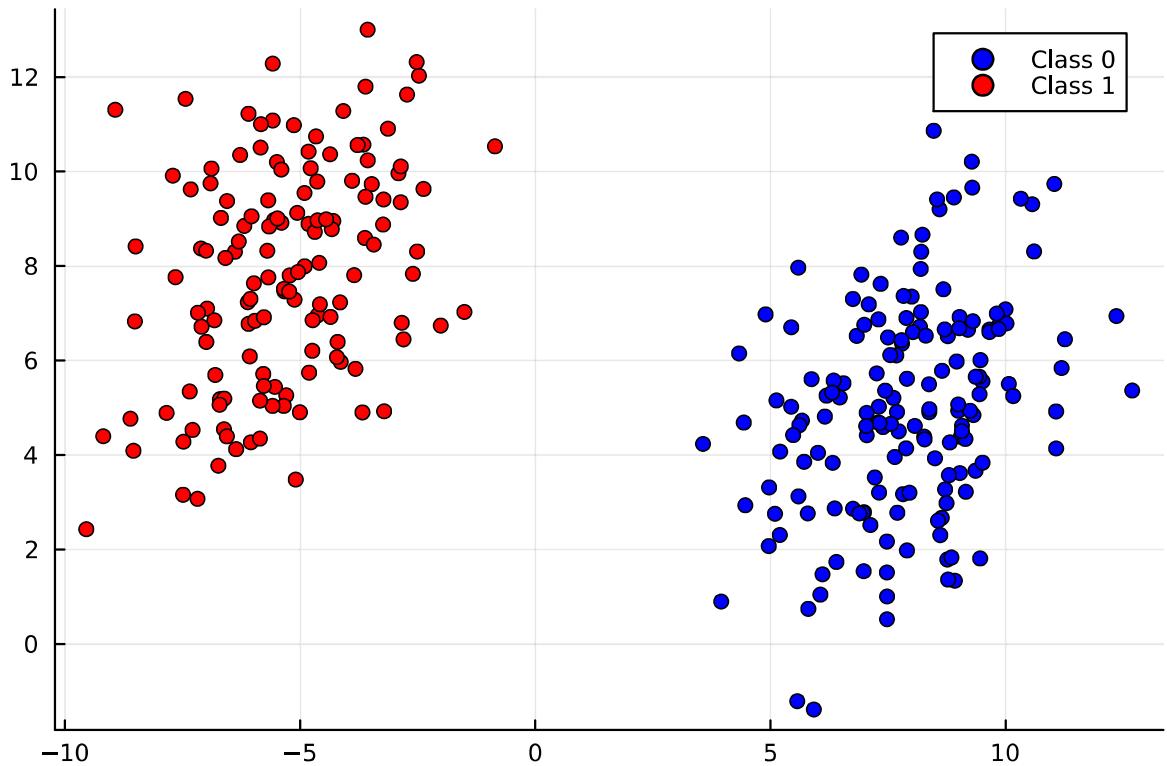
```

1 begin
2 # DOT NOT MODIFY THIS CODE
3 # generate a 2-class classification problem with 1,000 data points, each data
4 # point is a 2D feature vector
5 n = 1000
6
7 # dimensionality of data
8 d = 2
9
10 # mean
11 mu = 5
12
13 # variance
14 Sigma = 8
15
16 # Generate two class for synthesized data
17 positive = rand(MvNormal([Sigma, mu], 3 .* [1 (mu - d)/mu; (mu - d)/mu d]), n / 2)
18 negative = rand(MvNormal([-mu, Sigma], 3 .* [1 (mu - d)/mu; (mu - d)/mu d]), n / 2)
19
20 # Combine two class of generated data.
21 # X = features
22 # y = label
23 X = hcat(positive, negative)
24 y = vcat(ones(n / 2) .- 1, ones(n / 2))'
25 print(size(X))
26 print(size(y))
27
28 # Visualization
29 #plt = scatter(positive[1, :], positive[2, :], label="y = 1")
30 #scatter!(plt, negative[1, :], negative[2, :], label="y = -1")
31 # DOT NOT MODIFY THIS CODE
32 end

```

(2, 1000)(1, 1000)

?



```

1 begin
2   plt = scatter(X_test'[1,:][y_test .== 0], X_test'[2,:][y_test .== 0],
3                 label="Class 0", color=:blue, legend=:topright, markersize=4) # assign to plt
4   var
5   scatter!(X_test'[1,:][y_test .== 1], X_test'[2,:][y_test .== 1], label="Class
6   1", color=:red, markersize=4)
7 end

```

```

1 begin
2   # insert a column of 1's as the last entry in the feature matrix
3   # -- allows us to treat the bias as a trainable parameter
4
5   X_aug = vcat(X, ones(n)')
6   data = vcat(X_aug, y)
7   print(size(X_aug))
8 end

```

(3, 1000)

?

```
(700x3 Matrix{Float64}: , [1.0, 1.0, 0.0, 1.0, 1.0, 0.0, 0.0, 0.0, 1.0, more ,0.0], 3  
-6.19381 6.61483 1.0
```

```
1 begin  
2 # DOT NOT MODIFY THIS CODE  
3 # Split data, use 50% of the data for training and the remaining 50% for testing  
4 # Prepare data  
5 D = data'[shuffle(1:end), :]  
6  
7 # Calculate the number of samples for each split  
8 n_train = Int(n * 0.7)  
9  
10 # Split the samples into train, and test sets  
11 train_data = D[begin:n_train, :]  
12 test_data = D[n_train + 1: end, :]  
13 println(size(train_data), size(test_data))  
14  
15 # Move samples to train-test features and labels  
16 X_train, y_train, X_test, y_test = train_data[:,1:3], train_data[:,4],  
test_data[:,1:3], test_data[:,4]  
17 # DOT NOT MODIFY THIS CODE  
18 end
```

(700, 4)(300, 4)



4.3. Training

Sigmoid function and derivative of the sigmoid function

sigmoid_deriv (generic function with 1 method)

```
1 begin  
2     function sigmoid_activation(x)  
3         #TODO  
4         #reshape(x, (size(x, 1),))  
5         """compute the sigmoid activation value for a given input"""  
6         return 1. ./ (1. .+ exp.(-x))  
7     end  
8  
9     function sigmoid_deriv(x)  
10        #TODO  
11        """  
12            Compute the derivative of the sigmoid function ASSUMING  
13            that the input 'x' has already been passed through the sigmoid  
14            activation function  
15        """  
16        temp = x' * (1. .- x)  
17        #print("X dimension: ", string(size(x)))  
18        #print("Dimension: ", string(size(temp)))  
19        return temp  
20    end  
21  
22    #print(sigmoid_activation([3, 4, 5]))  
23    #print(sigmoid_deriv(sigmoid_activation([3, 4, 5])))  
24 end
```

Compute output

```
predict (generic function with 1 method)
1 begin
2     function compute_h(W, X)
3         #TODO
4         """
5             Compute output: Take the inner product between our features 'X' and the
6             weight
7             matrix 'W'
8             """
9             return X * W
10
11 end
12
13
14 function predict(W, X)
15     #TODO
16     """
17         Take the inner product between our features and weight matrix,
18         then pass this value through our sigmoid activation
19         """
20         preds = sigmoid_activation(compute_h(W, X))
21
22         # apply a step function to threshold the outputs to binary
23         # class labels
24         preds[preds .<= 0.5] .= 0
25         preds[preds .> 0] .= 1
26
27         return preds
28     end
end
```

Compute gradient

```
compute_gradient (generic function with 1 method)
1 begin
2     function compute_gradient(error, y_hat, trainX)
3         #TODO
4         """
5             the gradient descent update is the dot product between our
6             features and the error of the sigmoid derivative of
7             our predictions
8             """
9
10     deriv_error = sigmoid_deriv(y_hat)
11     # features: trainX (n x d)
12
13     return trainX' * (error * deriv_error)
14 end
15 end
```

Training function

```

train (generic function with 1 method)
1 begin
2     function train(W, trainX, trainY, learning_rate, num_epochs)
3         losses = []
4         for epoch in 1:num_epochs
5             y_hat = sigmoid_activation(compute_h(W, trainX))
6             # now that we have our predictions, we need to determine the
7             # 'error', which is the difference between our predictions and
8             # the true values
9             error = y_hat .- trainY
10            append!(losses, 0.5 * sum(error .^ 2))
11            grad = compute_gradient(error, y_hat, trainX)
12            W -= learning_rate * grad
13
14            if epoch == 1 || epoch % 5 == 0
15                println("Epoch=$epoch; Loss=$(losses[end])")
16            end
17        end
18        return W, losses
19    end
20 end

```

Initialize our weight matrix and list of losses

0.01

```

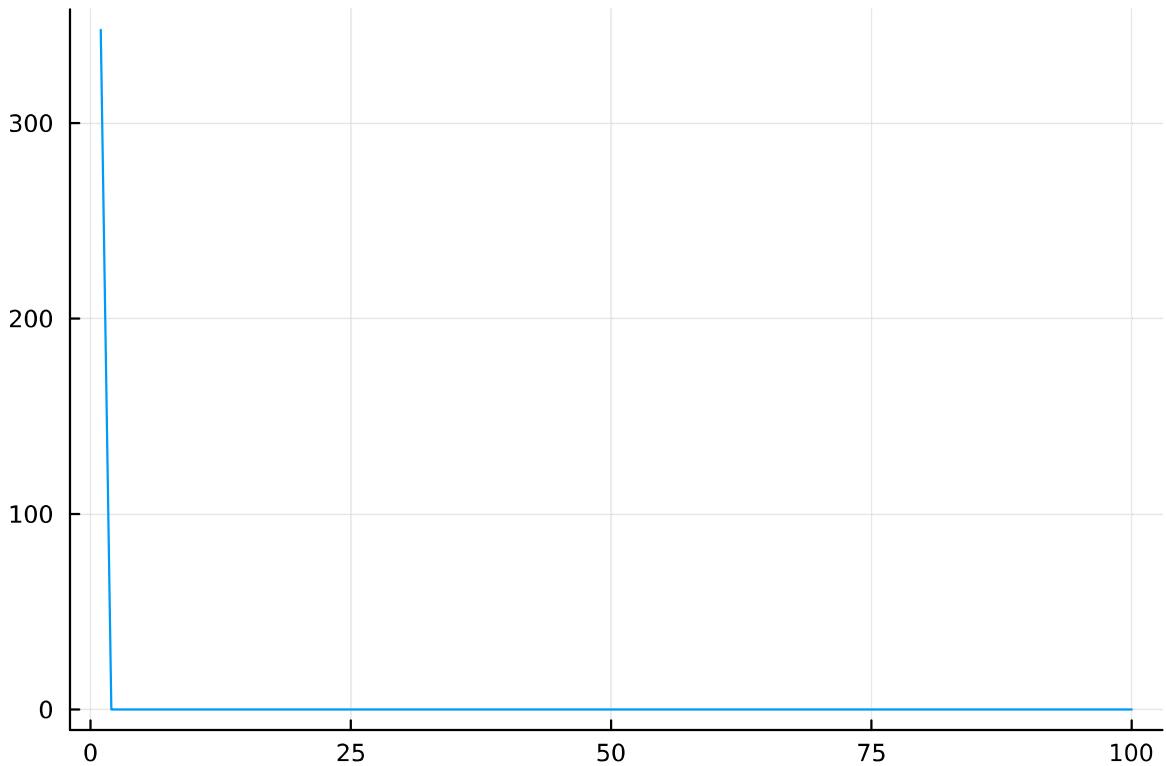
1 begin
2     #initialize our weight matrix and necessary hyperparameters
3     W = rand(Normal(), (size(X_train)[2], 1))
4     print(size(W))
5     #print(size(compute_h(W, X_train)),size(X_train), size(W))
6     num_epochs=100
7     #learning_rate=0.08
8     #learning_rate = 0.1
9     learning_rate = 0.01
10 end

```

(3, 1)

?

Train our model



```

1 begin
2     #training model
3     #print(sigmoid_activation(compute_h(W, X_train)))
4     #print(sigmoid_activation(compute_h(W, X_train)) - y_train)
5     theta, losses = train(W, X_train, y_train, learning_rate, num_epochs)
6     #visualiza training process
7     plot(1:num_epochs, losses, legend=false)
8 end

```

Epoch=1; Loss=347.7942690062178
 Epoch=5; Loss=5.1017344600080736e-194
 Epoch=10; Loss=5.1017344600080736e-194
 Epoch=15; Loss=5.1017344600080736e-194
 Epoch=20; Loss=5.1017344600080736e-194
 Epoch=25; Loss=5.1017344600080736e-194
 Epoch=30; Loss=5.1017344600080736e-194
 Epoch=35; Loss=5.1017344600080736e-194
 Epoch=40; Loss=5.1017344600080736e-194
 Epoch=45; Loss=5.1017344600080736e-194
 Epoch=50; Loss=5.1017344600080736e-194
 Epoch=55; Loss=5.1017344600080736e-194
 Epoch=60; Loss=5.1017344600080736e-194
 Epoch=65; Loss=5.1017344600080736e-194
 Epoch=70; Loss=5.1017344600080736e-194
 Epoch=75; Loss=5.1017344600080736e-194
 Epoch=80; Loss=5.1017344600080736e-194
 Epoch=85; Loss=5.1017344600080736e-194
 Epoch=90; Loss=5.1017344600080736e-194
 Epoch=95; Loss=5.1017344600080736e-194
 Epoch=100; Loss=5.1017344600080736e-194

?

Evaluate result

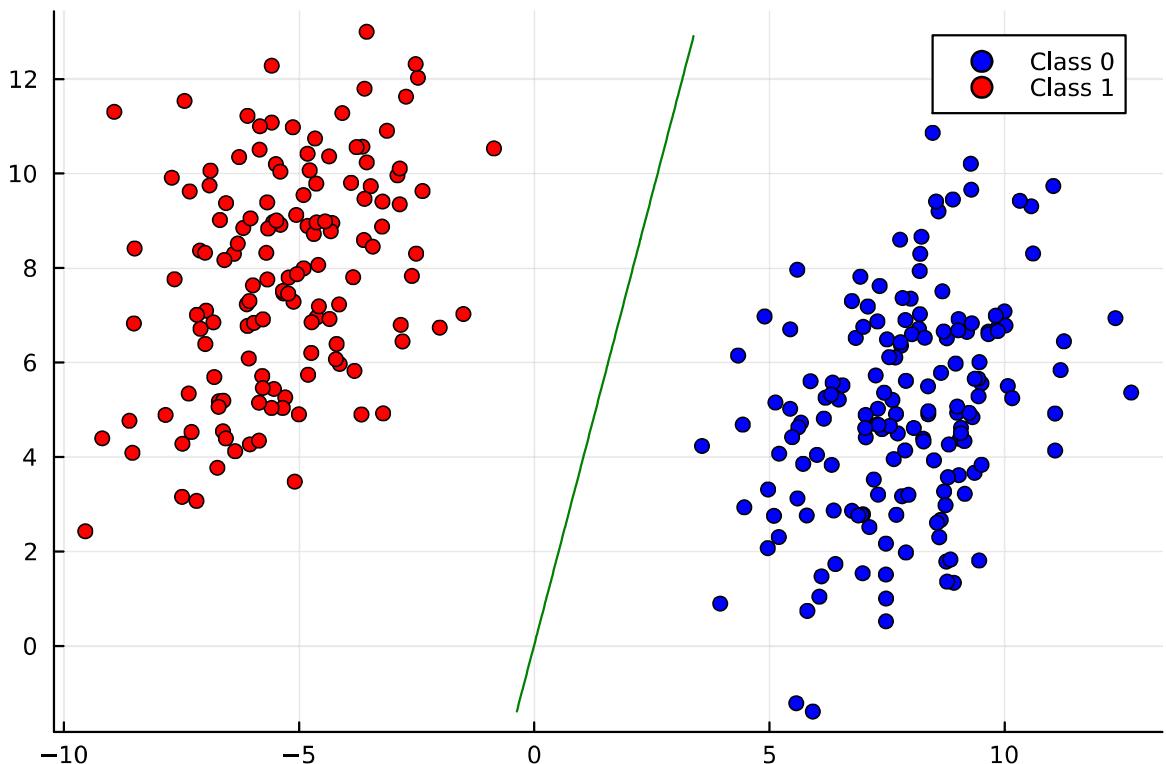
```

1 begin
2     y_pred = predict(θ, X_test)
3     true_positives = 0
4     false_positives = 0
5     true_negatives = 0
6     false_negatives = 0
7
8     # Calculate true positives, false positives, false negatives, and true negatives
9     for (true_label, predicted_label) in zip(y_test, y_pred)
10         if true_label == 1 && predicted_label == 1
11             true_positives += 1
12         elseif true_label == 0 && predicted_label == 1
13             false_positives += 1
14         elseif true_label == 1 && predicted_label == 0
15             false_negatives += 1
16         elseif true_label == 0 && predicted_label == 0
17             true_negatives += 1
18         end
19     end
20
21     # Calculate precision, recall, and F1-score
22     accuracy = (true_positives + true_negatives) / (true_positives +
23     false_positives + true_negatives + false_negatives)
24     precision = true_positives / (true_positives + false_positives)
25     recall = true_positives / (true_positives + false_negatives)
26     f1_score = 2 * precision * recall / (precision + recall)
27
28     # Display
29     print("acc: $accuracy, precision: $precision, recall: $recall, f1_score:
      $f1_score\n")

```

acc: 1.0, precision: 1.0, recall: 1.0, f1_score: 1.0





```

1 begin
2   # Create a scatter plot
3   scatter!(X_test'[1,:][y_test .== 0], X_test'[2,:][y_test .== 0], label="Class
4   0", color=:blue, legend=:topright, markersize=4)
5   scatter!(X_test'[1,:][y_test .== 1], X_test'[2,:][y_test .== 1], label="Class
6   1", color=:red, markersize=4)
7   # Getting decision boundary configuration
8   b = θ[3]
9   θ_ml = θ[1:2]
10
11   decision(x) = θ_ml' * x + b
12
13   D_test = ([
14     tuple.(eachcol(hcat(X_test'[1,:][y_test .== 0], X_test'[2,:][y_test .==
15     0]))'), 1)
16     tuple.(eachcol(hcat(X_test'[1,:][y_test .== 1], X_test'[2,:][y_test .==
17     1]))'), -1)
18   ])
19
20   # Max, min for visualization decision boundary
21   xmin = minimum(map((p) -> p[1][1], D_test))
22   ymin = minimum(map((p) -> p[1][2], D_test))
23   xmax = maximum(map((p) -> p[1][1], D_test))
24   ymax = maximum(map((p) -> p[1][2], D_test))
25
26   # Display decision boundary
27   contour!(plt, xmin:0.1:xmax, ymin:0.1:ymax,
28     (x, y) -> decision([x, y]),
29     levels=[0], linestyles=:solid, label="Decision boundary",
30     colorbar_entry=false, color=:green)
31 end

```

TODO: Study about accuracy, recall, precision, f1-score.

- Accuracy: Assuming the problem's category is that of classification, the accuracy score represents the percentage of the dataset that the model correctly predicted with no regard to any other statistics. For example, if the model correctly predicted 9 of out 10 datapoints, the accuracy will be 90%. This goes true for whatever the class the model predicted, as long as it is correct.
- Recall: With the above assumption, recall represents the percentage of True Positive (TP) among those that belong to the Positive group (True Positive + False Negative). Hence the formula: $\text{recall} = \text{TP} / (\text{TP} + \text{FN})$
- Precision: Unlike recall, precision represents the percentage of TP among those that are predicted as Positive, that includes the wrongly predicted Negative datapoints. We have the formula: $\text{precision} = \text{TP} / (\text{TP} + \text{FP})$
- F1: This is the harmonic mean of recall and precision. When calculating the recall or precision values, trying to lower one value will result in the other being raised. F1-score is used to prevent this, as the difference in both values will lower the F1-score. Proving that the model is inefficient.

TODO: Try out different learning rates. Give me your observations The data is very sensitive to the learning rate and the initial weights generated. With the learning rate $\alpha = 0.1$, often times the model could not generate a sufficient prediction, reaching 60% to 80% in accuracy but generally leaning towards either the negative or positive samples. There are some cases the model was able to accurately give a decision boundary that separates two groups from one another, but the success rate is quite low. This is the same for many other learning rates such as 0.06, 0.05, ... The submitted version contains the visualization of a good run with learning rate $\alpha = 0.01$, as the model converges almost immediately. However, trying to use the same learning rate for another data generation might not produce the same results. This holds true to the concept of gradient descent, as choosing a learning rate that puts the initial x_0 too far from x^* will result in a longer conversion. And the larger the learning rate, the higher the possibility that it could not converge fully after n epochs.

- 1 `md"""`
- 2 ****TODO: Try out different learning rates. Give me your observations****
- 3 The data is very sensitive to the learning rate and the initial weights generated. With the learning rate $\alpha = 0.1$, often times the model could not generate a sufficient prediction, reaching 60% to 80% in accuracy but generally leaning towards either the negative or positive samples. There are some cases the model was able to accurately give a decision boundary that separates two groups from one another, but the success rate is quite low. This is the same for many other learning rates such as 0.06, 0.05, ...
- 4 The submitted version contains the visualization of a good run with learning rate $\alpha = 0.01$, as the model converges almost immediately. However, trying to use the same learning rate for another data generation might not produce the same results.
- 5 This holds true to the concept of gradient descent, as choosing a learning rate that puts the initial x_0 too far from x^* will result in a longer conversion. And the larger the learning rate, the higher the possibility that it could not converge fully after n epochs.
- 6 `"""`