

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

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

- Student name:
- ID:

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

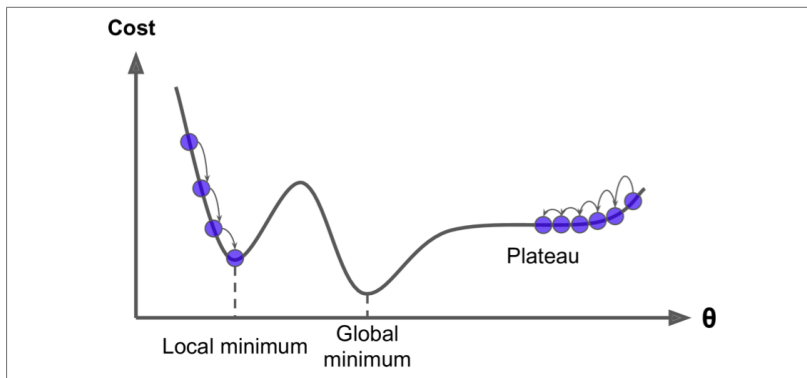


Figure 1. Loss landscape visualized as a 2D plot. Source: codecamp.vn

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 corresponds to a particular loss value given a set of parameters \mathbf{W} (weight matrix) and \mathbf{b} (bias vector). Our goal is to try different values of \mathbf{W} and \mathbf{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} 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 blank 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

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

$$\hat{y}_i = \sigma(\mathbf{h}_i) \Rightarrow \frac{\partial \hat{y}_i}{\partial \mathbf{h}_i} = \dots$$

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

```
1 md"""
2 **TODO**: Consider one sample  $\mathbf{x}_i$ . Fill in the blank
3
4  $\mathbf{h}_i = \mathbf{x}_i^T \mathbf{W} \Rightarrow \frac{\partial \mathbf{h}_i}{\partial \mathbf{W}} = \dots$ 
5
6  $\hat{y}_i = \sigma(\mathbf{h}_i) \Rightarrow \frac{\partial \hat{y}_i}{\partial \mathbf{h}_i} = \dots$ 
7
8  $\text{loss}_i = (\hat{y}_i - y_i)^2 \Rightarrow \frac{\partial \text{loss}_i}{\partial \hat{y}_i} = \dots$ 
9
10 """
```

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

$$\nabla Loss = \frac{\partial Loss(W)}{\partial W} = \frac{1}{n} \sum_{i=1}^n \dots$$

$= \dots$

$= \dots$

4. Implementation

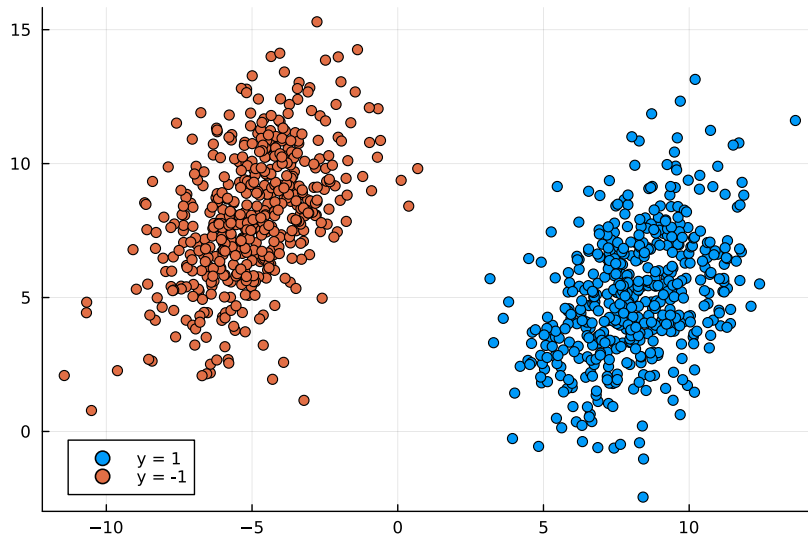
4.1. Import library

```
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     # number of data points
6     n = 1000
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
26    # Visualization
27    plt = scatter(positive[1, :], positive[2, :], label="y = 1")
28    scatter!(plt, negative[1, :], negative[2, :], label="y = -1")
29    # DOT NOT MODIFY THIS CODE
30 end

```

```

4×1000 Matrix{Float64}:
 9.49479  8.0947  8.76046  8.89715  ...  -0.966978  -4.6834  -6.78691  -2.24275
10.4304  4.07179  3.81086  7.02982  ...  10.7883   9.84328  4.07816  9.94067
 1.0      1.0      1.0      1.0      ...  1.0      1.0      1.0      1.0
 0.0      0.0      0.0      0.0      ...  1.0      1.0      1.0      1.0

```

```

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 end

```

```
► (700x3 Matrix{Float64}): , [1.0, 0.0, 1.0, 1.0, 1.0, 1.0, 0.0, 0.0, 0.0, ... more ,1.0], 300x  
-7.42717 5.31941 1.0 -7.
```

```
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],
17    test_data[:,1:3], test_data[:,4]
18    # DOT NOT MODIFY THIS CODE
19 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         """compute the sigmoid activation value for a given input"""
5         #return?
6
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        #return?
17
18    end
19 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 weight
6     matrix 'W'
7     """
8     # return?
9
10  end
11
12  function predict(W, X)
13    #TODO
14    """
15    Take the inner product between our features and weight matrix,
16    then pass this value through our sigmoid activation
17    """
18    # preds = ...
19
20
21    # apply a step function to threshold the outputs to binary
22    # class labels
23    preds[preds <= 0.5] .= 0
24    preds[preds > 0] .= 1
25
26    return preds
27  end
28 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     # return?
10
11  end
12 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

```

0.1
 1 begin
 2     #initialize our weight matrix and necessary hyperparameters
 3     W = rand(Normal(), (size(X_train)[2], 1))
 4     num_epochs=100
 5     learning_rate=0.1
 6 end

```

Train our model

```

MethodError: no method matching -(::String, ::Vector{Float64})
Closest candidates are:
-(!Matched::FillArrays.Zeros{T, N, Axes} where Axes, ::AbstractArray{V, N}) where {T, V,
-(!Matched::SparseArrays.AbstractSparseMatrixCSC, ::Array) at /buildworker/worker/package
-(!Matched::Distributions.MvNormal, ::AbstractVector{T} where T) at /home/lnhutnam/.julia:
...

```

```

1. train @ {Other: 9} [inlined]
2. top-level scope @ {Local: 3} [inlined]

```

```

 1 begin
 2     #training model
 3     θ, losses = train(W, X_train, y_train, learning_rate, num_epochs)
 4     #visualiza training process
 5     plot(1:num_epochs, losses, legend=false)
 6 end

```

Evaluate result

UnDefVarError: preds not defined

1. `predict(::Matrix{Float64}, ::Matrix{Float64})` @ `Other: 23`
2. `top-level scope` @ `Local: 2` [inlined]

```
1 begin
2   y_pred = predict(W, 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 + false_positives
23 + 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:
30 $f1_score\n")
31 end
```

Another cell defining θ contains errors.

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

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

- Accuracy:
- Recall:
- Precision:
- F1:

TODO: Try out different learning rates. Give me your observations