

Lab03: Logistic Regression.

- Student ID: 21127329
- Student name: Châu Tân Kiệt

```
• md"""
• # Lab03: Logistic Regression.
•
• - Student ID: 21127329
• - Student name: Châu Tân Kiệt
• """
```

How to do your homework

You will work directly on this notebook; the word `TODO` indicate the parts you need to do.

You can discuss ideas with classmates as well as finding information from the internet, book, etc...; but *this homework must be your*.

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.

Danger

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

Contents:

- Logistic Regression.

1. Feature Extraction

Import Library

```
• begin
•   using Distributions , Plots , Images , LinearAlgebra , Statistics ,
•   Random
• end
```

MersenneTwister(2024)

- Random.seed!(2024)

Load data

- *# I DON'T KNOW WHAT YOUR OS, SO I ATTACHED DATA FOR YOU*
 - *# YOU DON'T NEED TO DOWNLOAD AND EXTRACT BY YOURSELF*
 - *# IF YOU WANT TO USE JULIA TO DOWNLOAD DATA, LET'S USE THESE CODE*
 - *# FOR EXTRACTING MNIST .gz FILE, PLEASE USE gzip*
 - *#*
 - *# function download_dataset(save_path::String="data")*
 - *# # setup directory*
 - *# mkpath(joinpath(dirname(@__FILE__), save_path))*
 - *# data_dir = joinpath(dirname(@__FILE__), save_path)*
 - *# mkpath(joinpath(data_dir, "train"))*
 - *# train_dir = joinpath(data_dir, "train")*
 - *# mkpath(joinpath(data_dir, "test"))*
 - *# test_dir = joinpath(data_dir, "test")*
 - *#*
 - *# # download dataset*
 - *# mkpath(joinpath(train_dir, "images"))*
 - *# download("http://yann.lecun.com/exdb/mnist/train-images-idx3-ubyte.gz",joinpath(train_dir, "images/train-images-idx3-ubyte.gz"))*
 - *# train_images_file = joinpath(train_dir, "images/train-images-idx3-ubyte.gz")*
 - *#*
 - *# mkpath(joinpath(train_dir, "labels"))*
 - *# download("http://yann.lecun.com/exdb/mnist/train-labels-idx1-ubyte.gz",joinpath(train_dir, "labels/train-labels-idx1-ubyte.gz"))*
 - *# train_labels_file = joinpath(train_dir, "labels/train-labels-idx1-ubyte.gz")*
 - *#*
 - *# mkpath(joinpath(test_dir, "images"))*
 - *# download("http://yann.lecun.com/exdb/mnist/t10k-images-idx3-ubyte.gz",joinpath(test_dir, "images/t10k-images-idx3-ubyte.gz"))*
 - *# test_images_file = joinpath(test_dir, "images/t10k-images-idx3-ubyte.gz")*
 - *#*
 - *# mkpath(joinpath(test_dir, "labels"))*
 - *# download("http://yann.lecun.com/exdb/mnist/t10k-labels-idx1-ubyte.gz",joinpath(test_dir, "labels/t10k-labels-idx1-ubyte.gz"))*
 - *# test_labels_file = joinpath(test_dir, "labels/t10k-labels-idx1-ubyte.gz")*
 - *# end*
-
- *# If you have downloaded the dataset yet, please uncomment this line below and run this cell. Otherwise, keep it in uncomment state.*
 - *# download_dataset()*

10000

```
• begin
•     data_dir = joinpath(dirname(@__FILE__), "data")
•     train_x_dir = joinpath.(data_dir, "train\\images\\train-images.idx3-ubyte")
•     train_y_dir = joinpath.(data_dir, "train\\labels\\train-labels.idx1-ubyte")
•
•     test_x_dir = joinpath.(data_dir, "test\\images\\t10k-images.idx3-ubyte")
•     test_y_dir = joinpath.(data_dir, "test\\labels\\t10k-labels.idx1-ubyte")
•     NUMBER_TRAIN_SAMPLES = 60000
•     NUMBER_TEST_SAMPLES = 10000
• end
```

```
begin
    # Init arrays
    train_x = Array{Float64}(undef, 28^2, NUMBER_TRAIN_SAMPLES)
    train_y = Array{Int64}(undef, NUMBER_TRAIN_SAMPLES)

    # Init io streams
    io_images = open(train_x_dir)
    io_labels = open(train_y_dir)

    # Iterating through sample length
    for i ∈ 1:NUMBER_TRAIN_SAMPLES
        seek(io_images, (i-1)*28^2 + 16) # offset 16 to skip header
        seek(io_labels, (i-1)*1 + 8) # offset 8 to skip header
        train_x[:,i] = convert(Array{Float64}, read(io_images, 28^2))
        train_y[i] = convert(Int, read(io_labels, UInt8))
    end

    # Close io streams
    close(io_images)
    close(io_labels)

    # Transpose features
    train_x = train_x'
end
```

```
((60000, 784), (60000), (10000, 784), (10000))  
• size(train_x), size(train_y), size(test_x), size(test_y)
```

Extract Features

So we basically have 70000 samples with each sample having 784 features - pixels in this case and a label - the digit the image represent.

Let's play around and see if we can extract any features from the pixels that can be more informative. First I'd like to know more about average intensity - that is the average value of a pixel in an image for the different digits

```
compute_average_intensity (generic function with 1 method)
```

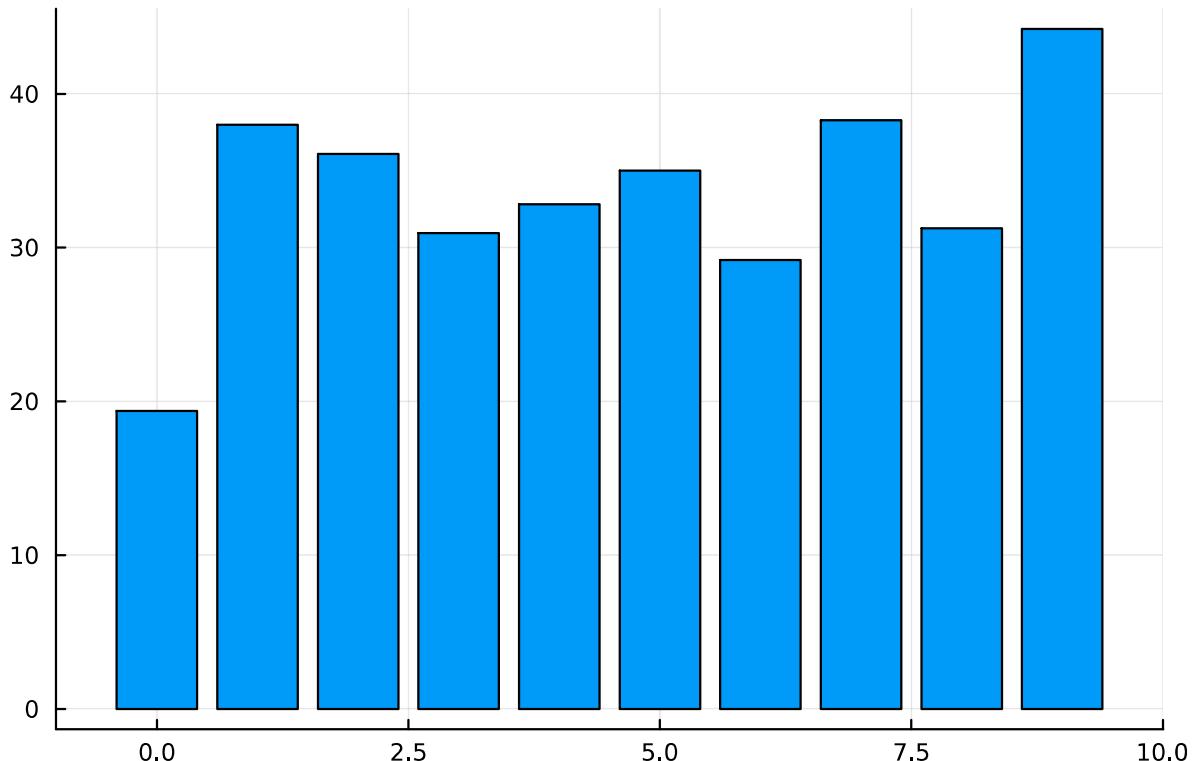
- *#TODO compute average intensity for each label*
- **function** **compute_average_intensity**(**x**, **y**)
- **mean_** = **zeros**(**Float64**, 10) # 10 is number of labels
- *#TODO compute average intensity for each label*
- **for** **i** in 1:10
- **mean_[i]** = **mean**(**x**[**y**.== **i** % 10, :])
- **end**
- **return** **mean_**
- **end**

```
l_mean =
```

```
[19.3797, 37.9887, 36.0902, 30.9482, 32.8311, 35.012, 29.2046, 38.2898, 31.2604, 44.2168]
```

- **l_mean = compute_average_intensity(train_x, train_y)**

Plot the average intensity using matplotlib



- **bar(0:9, l_mean, legend=false)**

```
(60000)
```

- **begin**
- *#TODO compute average intensity for each data sample*
- **intensity** = []
- **for** **i** in 1: (**NUMBER_TRAIN_SAMPLES**)
- **intensity** = **append!**(**intensity**, **vec**([**mean**(**train_x**[**i**, :])]))
- **end**
- **size(intensity)**
- **end**

Some digits are symmetric (1, 3, 8, 0) some are not (2, 4, 5, 6, 9). Creating a new feature capturing this could be useful. Specifically, we calculate $s = -\frac{s_1 + s_2}{2}$ for each image:

- s_1

: flip the image along y-axis and compute the mean value of result

- s_2

: flip the image along x-axis and compute the mean value of result

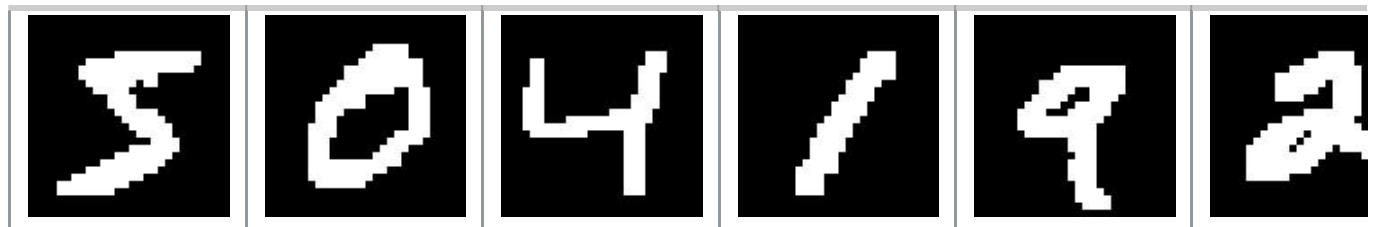
`compute_symmetry` (generic function with 1 method)

```
• function compute_symmetry(train_x)
•     symmetry = []
•     for i in 1:size(train_x)[1]
•         img = reshape(train_x[i,:], (28,28))
•         s1 = mean(abs.(img - reverse(img, dims=1)))
•         s2 = mean(abs.(img - reverse(img, dims=2)))
•         s = -0.5 .* (s1 + s2)
•         append!(symmetry, s)
•     end
•     return symmetry
• end
```

(60000)

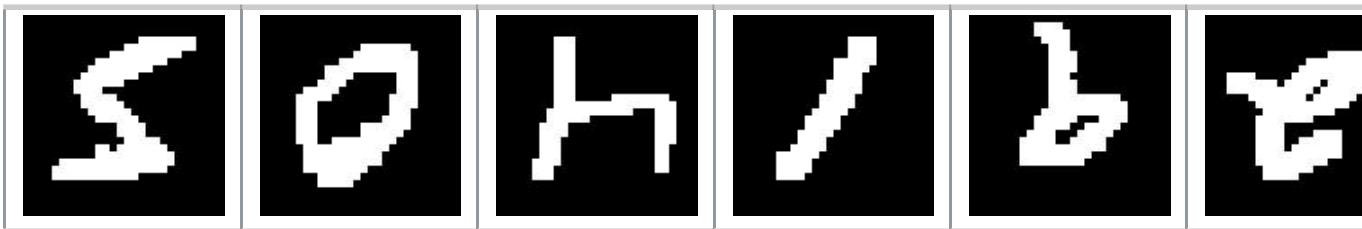
```
• begin
•     symmetry = compute_symmetry(train_x)
•     size(symmetry)
• end
```

Visualize 10 samples in order to illustrate symmetry



(a vector displayed as a row to save space)

```
• begin
•     num_img = 10
•     img_flat = train_x[1:num_img,:]
•     img = [reshape(img_flat[i,:], (28,28))' for i in 1:num_img]
•     [colorview(Gray, Float32.(img[i])) for i in 1:num_img]
• end
```



(a vector displayed as a row to save space)

```
• begin
•     img_reverse_flat = reverse(img_flat, dims=2)
•     img_reverse = [reshape(img_reverse_flat[i,:], (28,28))' for i in 1:num_img]
•     [colorview(Gray, Float32.(img_reverse[i])) for i in 1:num_img]
• end
```

Our new data will have 70000 samples and 2 features: intensity, symmetry.

(60000, 2)

```
• begin
•     #TODO create X_new by horizontal stack intensity and symmetry
•     train_x_new = [intensity,symmetry]
•     train_x_new = hcat(train_x_new...)
•
•     size(train_x_new)
• end
```

2. Training

Usually logistic regression is a good first choice for classification. In this homework we use logistic regression for classifying digit 1 images and not digit 1's images.

Normalize data

First normalize data using Z-score normalization

- **TODO: Study about Z-score normalization**
- We transform v into v' by formula $\hat{v} = \frac{v - \bar{A}}{\sigma}$

$$\bar{A}$$

: mean of A

- σ

: standard deviation

- **TODO: Why should we normalize data?**

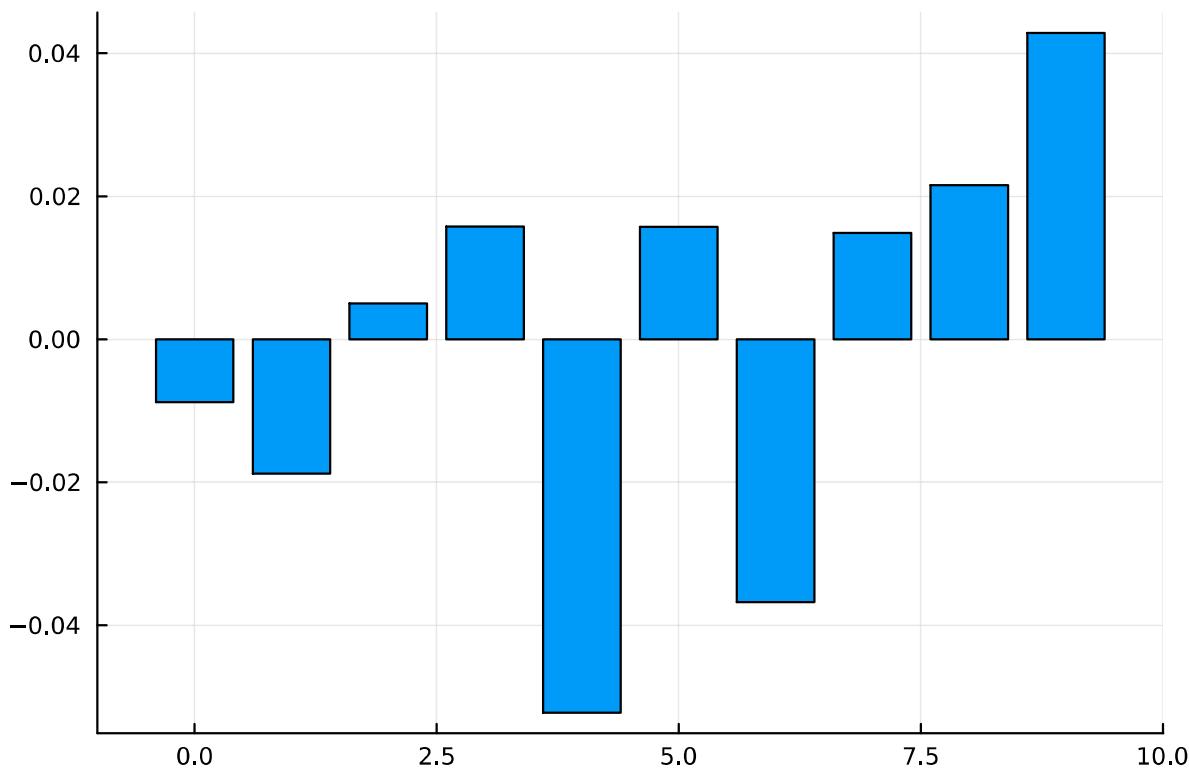
- Makes features comparable: Normalizing data ensures that features with larger ranges don't dominate the analysis.
- Better model convergence: Normalized data helps machine learning algorithms find optimal parameters faster.
- Better interpretation: Normalized data makes it easier to understand the importance of different features.
- Reduces overfitting: Normalization can help machine learning models generalize better to unseen data.

normalize (generic function with 3 methods)

```

• function normalize(x, mean_=nothing, std_=nothing)
•   if (mean_ == nothing && std_ == nothing)
•     #TODO normalize x_train
•     mean_ = mean(x)
•     std_ = stdm(x, mean_)
•     return normalize(x, mean_, std_), mean_, std_
•     #mean_ and std_ will be re-used to pre-process test set
•   end
•   normalized_x = (x.-mean_) ./ std_
•   return normalized_x
• end
•
•

```



```

• begin
•     normalized_train_x, mean_, std_ = normalize(train_x_new)
•     s_mean = compute_average_intensity(normalized_train_x, train_y)
•     bar(0:9, s_mean, legend=false)
• end

```

Construct data

(60000, 1)

```

• begin
•     train_y_new = reshape(deepcopy(train_y), (size(train_y)[1], 1))
•     train_y_new[train_y_new .!= 1] .= 0
•     size(train_y_new)
• end

```

(60000, 3)

```

• begin
•     # construct data by adding ones
•     add_one_train_x = hcat(ones(size(normalized_train_x)[1],), normalized_train_x)
•     size(add_one_train_x)
• end

```

Sigmoid function and derivative of the sigmoid function

sigmoid_activation (generic function with 1 method)

- `function sigmoid_activation(x)`
- `#TODO`
- `"""compute the sigmoid activation value for a given input"""`
- `#return?`
- `return vec([1.0 / (1.0 + exp(-x_i)) for x_i in x])`
- `end`

sigmoid_deriv (generic function with 1 method)

- `function sigmoid_deriv(x)`
- `#TODO`
- `"""`
- `Compute the derivative of the sigmoid function ASSUMING`
- `that the input 'x' has already been passed through the sigmoid`
- `activation function`
- `"""`
- `#return?`
- `return x .* (1.0 - x)`
- `end`

Compute output

compute_h (generic function with 1 method)

- `function compute_h(W, X)`
- `#TODO`
- `"""`
- `Compute output: Take the inner product between our features 'X' and the weight`
- `matrix 'W'`
- `"""`
- `h = X * W`
- `return h`
- `end`

predict (generic function with 1 method)

- `function predict(W, X)`
- `#TODO`
- `"""`
- `Take the inner product between our features and weight matrix,`
- `then pass this value through our sigmoid activation`
- `"""`
- `preds = sigmoid_activation(compute_h(W, X))`
- `# apply a step function to threshold the outputs to binary`
- `# class labels`
- `preds[preds .≤ 0.5] .= 0`
- `preds[preds .> 0] .= 1`
- `return preds`
- `end`

Compute gradient

Loss Function: Average negative log likelihood

$$\mathcal{L} = \frac{1}{N} \sum_{i=1}^N - (y^i \ln h_{\mathbf{w}}(\mathbf{x}^i) + (1 - y^i) \ln (1 - h_{\mathbf{w}}(x^i)))$$

Sigmoid Activation: $z = \sigma(h) = \frac{1}{1 + e^{-h}}$

Cross-entropy: $J(w) = -(y \log(z) + (1 - y) \log(1 - z))$

Chain rule: $\frac{\partial J(w)}{\partial w} = \frac{\partial J(w)}{\partial z} \frac{\partial z}{\partial h} \frac{\partial h}{\partial w}$

$$\frac{\partial J(w)}{\partial z} = - \left(\frac{y}{z} - \frac{1-y}{1-z} \right) = \frac{z-y}{z(1-z)}$$

$$\frac{\partial z}{\partial h} = z(1-z)$$

$$\frac{\partial h}{\partial w} = X$$

$$\frac{\partial J(w)}{\partial w} = X^T(z - y)$$

```
compute_gradient (generic function with 1 method)
```

- `function compute_gradient(error, train_x)`
- `#TODO`
- `"""`
- This is the gradient descent update of "average negative loglikelihood" loss function.
- In lab02 our loss function is "sum squared error".
- `"""`
- `gradient = train_x' * error`
- `return gradient`
- `end`

```
train (generic function with 1 method)
```

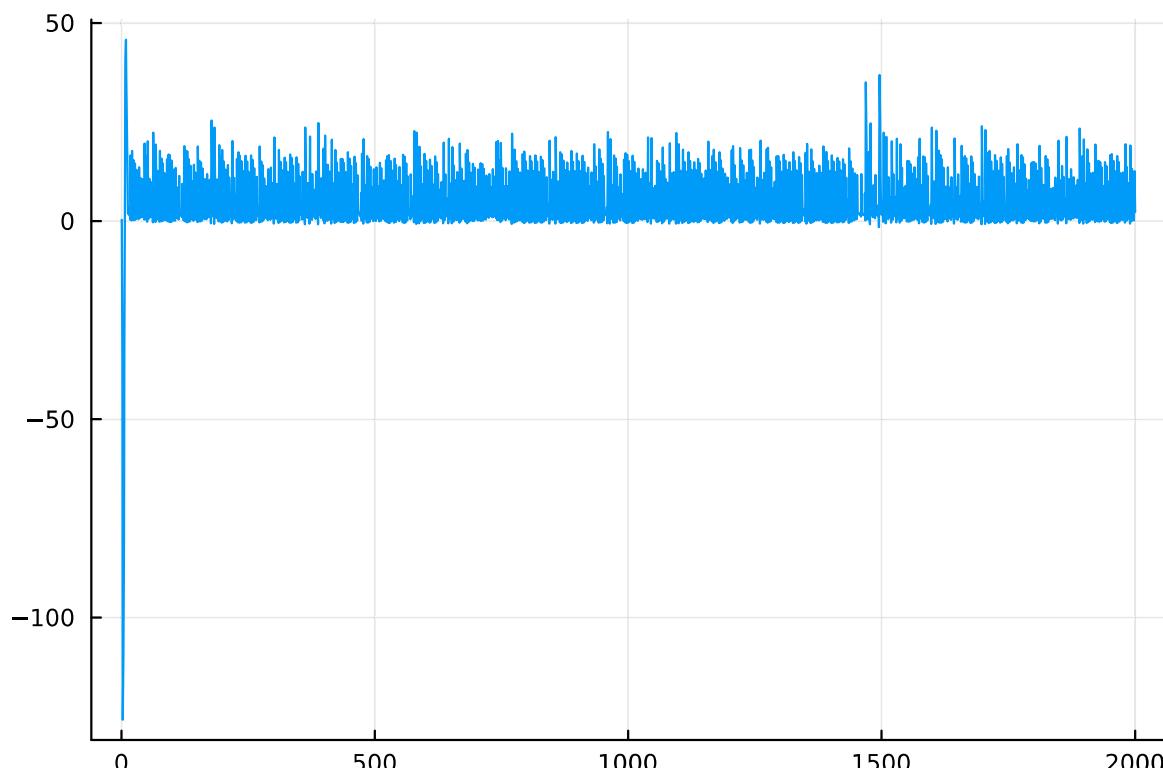
```
• function train(W, train_x, train_y, learning_rate, num_epochs)
•     losses = []
•     W_list = W
•     for epoch in 1:num_epochs
•         h = compute_h(W, train_x)
•         y_hat = sigmoid_activation(h)
•         error = y_hat - train_y
•         loss = -mean([train_y[i] == 1 ? (isinf(log(y_hat[i])) ? abs(h[i]) :
•             log(y_hat[i])) : (isinf(log(1 - y_hat[i]))) ? abs(h[i]) : log(1 - y_hat[i]))])
•         for i = 1 : NUMBER_TRAIN_SAMPLES]
•             # loss = mean(-1 .* train_y .* log.(y_hat) .- (1 .- train_y) .* log.(1 .-
•             ## y_hat)) # wow very nice function, I loved it
•             append!(losses, loss)
•             grad = compute_gradient(error, train_x)
•             W -= learning_rate * grad
•             if epoch == 1 || epoch % 200 == 0
•                 print("Epoch=$epoch; Loss=$(losses[end])\n")
•                 W_list = hcat(W_list, W)
•             end
•         end
•     end
•     return W, losses, W_list'
• end
```

Train our model

```
(3x1 Matrix{Float64}[:, [0.30099, -36.6037, -125.797, -110.159, -79.6602, -33.1472, 10.68  
173.805  
-203.157  
50.7049
```

```
• begin  
•     W = rand(Normal(), (size(add_one_train_x)[2], 1))  
•  
•     num_epochs=2000  
•     learning_rate=0.01  
•     W, losses, W_list = train(W, add_one_train_x, train_y_new, learning_rate,  
•     num_epochs)  
• end
```

```
Epoch=1; Loss=0.3009902175758327  
Epoch=200; Loss=5.601832058570936  
Epoch=400; Loss=5.693783003145179  
Epoch=600; Loss=-0.3888144251431936  
Epoch=800; Loss=-0.01911694995086421  
Epoch=1000; Loss=-0.39119565465614103  
Epoch=1200; Loss=9.296160704682897  
Epoch=1400; Loss=4.876691763277152  
Epoch=1600; Loss=10.83620779934856  
Epoch=1800; Loss=-0.04033944311027385  
Epoch=2000; Loss=2.2544822081005935
```



```
• plot(1:num_epochs, losses, legend=false)
```

3. Evaluate our model

In this section, you will evaluate your model on train set and test set and make some comment about the result.

Evaluate model on training set

```
tpfptnfn_cal (generic function with 2 methods)
• function tpfptnfn_cal(y_test, y_pred, positive_class=1)
•     true_positives = 0
•     false_positives = 0
•     true_negatives = 0
•     false_negatives = 0
•
•     # Calculate true positives, false positives, false negatives, and true negatives
•     for (true_label, predicted_label) in zip(y_test, y_pred)
•         if true_label == positive_class && predicted_label == positive_class
•             true_positives += 1
•         elseif true_label != positive_class && predicted_label == positive_class
•             false_positives += 1
•         elseif true_label == positive_class && predicted_label != positive_class
•             false_negatives += 1
•         elseif true_label != positive_class && predicted_label != positive_class
•             true_negatives += 1
•         end
•     end
•
•     return true_positives, false_positives, true_negatives, false_negatives
• end
```

```

• begin
•     preds_train = predict(W, add_one_train_x)
•     train_y_n = reshape(train_y_new, length(train_y_new), 1)
•
•     acc = 0
•     precision = 0
•     recall = 0
•     f1 = 0
•
•     for i in 1:10
•         # Calculate true positives, false positives, false negatives, and true
•         # negatives
•         true_positives, false_positives, true_negatives, false_negatives =
•         tpfpfn_fn_cal(train_y_n, preds_train)
•
•         # Calculate precision, recall, and F1-score
•         acc += (true_positives + true_negatives) / (true_positives +
•             false_positives + true_negatives + false_negatives)
•         precision += true_positives / (true_positives + false_positives)
•         recall += true_positives / (true_positives + false_negatives)
•     end
•
•     acc = acc / 10
•     precision = precision / 10
•     recall = recall / 10
•     f1 = 2 * precision * recall / (precision + recall)
•     print(" acc: $acc\n precision: $precision\n recall: $recall\n f1_score: $f1\n")
•
• end

```

```

acc: 0.8812499999999999 ②
precision: 0.4819798626140962
recall: 0.7597152180361909
f1_score: 0.5897864010593586

```

Evaluate model on test set

In order to predict the result on test set, you have to perform data pre-process first. The pre-process is done exactly what we have done on train set. That means, you have to:

- Change the label in `test_y` to 0 and 1 and store in a new variable named `test_y_new`
- Calculate `test_intensity` and `test_symmetry` to form `test_x_new` (the shape should be `(10000,2)`)
- Normalized `test_x_new` by z-score. Note the you will re-use variable `mean_` and `std_` to calculate `test_x_new` instead of compute new ones. You will store the result in `normalized_test_x`
- Add a column that's full of one to `test_x_new` and store in `add_one_test_x` (the shape should be `(10000,3)`)

```

• begin
• #TODO
• # compute test_y_new
• test_y_new = reshape(deepcopy(test_y), NUMBER_TEST_SAMPLES, 1)
•
• # compute test_intensity and test_symmetry to form test_x_new
• test_intensity = vec([mean(test_x[i, :]) for i in (1 : NUMBER_TEST_SAMPLES)])
• test_symmetry = compute_symmetry(test_x)
• test_x_new = hcat(test_intensity, test_symmetry)
• println(size(test_x_new))
•
• # normalize test_x_new to form normalized_test_x
• normalized_test_x = normalize(test_x_new, mean_, std_)
•
• # add column 'ones' to test_x_new
• add_one_test_x = hcat(ones(Float64, NUMBER_TEST_SAMPLES), normalized_test_x)
• println(size(add_one_test_x))
• end

```

(10000, 2)
(10000, 3)



After doing all these stuffs, you now can predict and evaluate your model

```

• begin
•     preds_test = predict(W, add_one_test_x)
•     test_y_n = reshape(test_y_new, length(test_y_new), 1)
•
•     _acc = 0
•     _p = 0
•     _r = 0
•     _f1 = 0
•
•     for i ∈ 1:10
•         # Calculate true positives, false positives, false negatives, and true
•         # negatives
•         tp, fp, tn, fn = tpfptnfn_cal(test_y_n, preds_test)
•
•         # Calculate precision, recall, and F1-score
•         _acc += (tp + tn) / (tp + fp + tn + fn)
•         _p += tp / (tp + fp)
•         _r += tp / (tp + fn)
•     end
•
•     _acc = _acc / 10
•     _p = _p / 10
•     _r = _r / 10
•     _f1 = 2 * _p * _r / (_p + _r)
•
•     print(" acc: $_acc\n precision: $_p\n recall: $_r\n f1_score: $_f1\n")
• end

```

acc: 0.8949999999999998
precision: 0.526283240568954/
recall: 0.7497797356828195
f1_score: 0.6184593023255813

TODO: Comment on the result

- The runs are to separate one value label and no value label data
- The accuracy score is pretty high (consistently at nearly 90%)
- The f1 score is around mid-range (high 50% - low 60%)
- The contrary result of the precision score and recall score affected the f1 score