`Ctrl` + `S`

# Lab03: Logistic Regression.

- Student ID:
- Student name:

**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

```
1  begin
2      using  Distributions, Plots, Images, LinearAlgebra, Random
3  end
```

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

# Load data

```julia
1   # I DON'T KNOW WHAT YOUR OS, SO I ATTACHED DATA FOR YOU
2   # YOU DON'T NEED TO DOWNLOAD AND EXTRACT BY YOURSELF
3
4   # IF YOU WANT TO USE JULIA TO DOWNLOAD DATA, LET'S USE THESE CODE
5   # FOR EXTRACTING MNIST .gz FILE, PLEASE USE gzip
6
7   # function download_dataset(save_path::String="data")
8   #    # setup directory
9   #    mkpath(joinpath(dirname(@__FILE__), save_path))
10  #    data_dir = joinpath(dirname(@__FILE__), save_path)
11  #    mkpath(joinpath(data_dir, "train"))
12  #    train_dir =  joinpath(data_dir, "train")
13  #    mkpath(joinpath(data_dir, "test"))
14  #    test_dir = joinpath(data_dir, "test")
15
16  #    # download dataset
17  #    mkpath(joinpath(train_dir, "images"))
18  #    download("http://yann.lecun.com/exdb/mnist/train-images-idx3-
        ubyte.gz",joinpath(train_dir, "images/train-images-idx3-ubyte.gz"))
19  #    train_images_file = joinpath(train_dir, "images/train-images-idx3-ubyte.gz")
20
21  #    mkpath(joinpath(train_dir, "labels"))
22  #    download("http://yann.lecun.com/exdb/mnist/train-labels-idx1-
        ubyte.gz",joinpath(train_dir, "labels/train-labels-idx1-ubyte.gz"))
23  #    train_labels_file = joinpath(train_dir, "labels/train-labels-idx1-ubyte.gz")
24
25  #    mkpath(joinpath(test_dir, "images"))
26  #    download("http://yann.lecun.com/exdb/mnist/t10k-images-idx3-
        ubyte.gz",joinpath(test_dir, "images/t10k-images-idx3-ubyte.gz"))
27  #    test_images_file = joinpath(test_dir, "images/t10k-images-idx3-ubyte.gz")
28
29  #    mkpath(joinpath(test_dir, "labels"))
30  #    download("http://yann.lecun.com/exdb/mnist/t10k-labels-idx1-
        ubyte.gz",joinpath(test_dir, "labels/t10k-labels-idx1-ubyte.gz"))
31  #    test_labels_file = joinpath(test_dir, "labels/t10k-labels-idx1-ubyte.gz")
32  # end
```

```julia
1   # If you have downloaded the dataset yet, please uncomment this line below and run
      this cell. Otherwise, keep it in uncomment state.
2   # download_dataset()
```

```
10000
```

```julia
1  begin
2      data_dir = joinpath(dirname(@__FILE__), "data")
3      train_x_dir = joinpath(data_dir, "train/images/train-images.idx3-ubyte")
4      train_y_dir = joinpath(data_dir, "train/labels/train-labels.idx1-ubyte")
5
6      test_x_dir = joinpath(data_dir, "test/images/t10k-images.idx3-ubyte")
7      test_y_dir = joinpath(data_dir, "test/labels/t10k-labels.idx1-ubyte")
8
9      NUMBER_TRAIN_SAMPLES = 60000
10     NUMBER_TEST_SAMPLES = 10000
11 end
```

```
60000×784 adjoint(::Matrix{Float64}) with eltype Float64:
 0.0  0.0  0.0  0.0  0.0  0.0  0.0  0.0  …  0.0  0.0  0.0  0.0  0.0  0.0  0.0  0.0
 0.0  0.0  0.0  0.0  0.0  0.0  0.0  0.0     0.0  0.0  0.0  0.0  0.0  0.0  0.0  0.0
 0.0  0.0  0.0  0.0  0.0  0.0  0.0  0.0     0.0  0.0  0.0  0.0  0.0  0.0  0.0  0.0
 0.0  0.0  0.0  0.0  0.0  0.0  0.0  0.0     0.0  0.0  0.0  0.0  0.0  0.0  0.0  0.0
 0.0  0.0  0.0  0.0  0.0  0.0  0.0  0.0     0.0  0.0  0.0  0.0  0.0  0.0  0.0  0.0
 0.0  0.0  0.0  0.0  0.0  0.0  0.0  0.0  …  0.0  0.0  0.0  0.0  0.0  0.0  0.0  0.0
 0.0  0.0  0.0  0.0  0.0  0.0  0.0  0.0     0.0  0.0  0.0  0.0  0.0  0.0  0.0  0.0
 ⋮                        ⋮             ⋱                   ⋮
 0.0  0.0  0.0  0.0  0.0  0.0  0.0  0.0     0.0  0.0  0.0  0.0  0.0  0.0  0.0  0.0
 0.0  0.0  0.0  0.0  0.0  0.0  0.0  0.0  …  0.0  0.0  0.0  0.0  0.0  0.0  0.0  0.0
 0.0  0.0  0.0  0.0  0.0  0.0  0.0  0.0     0.0  0.0  0.0  0.0  0.0  0.0  0.0  0.0
 0.0  0.0  0.0  0.0  0.0  0.0  0.0  0.0     0.0  0.0  0.0  0.0  0.0  0.0  0.0  0.0
 0.0  0.0  0.0  0.0  0.0  0.0  0.0  0.0     0.0  0.0  0.0  0.0  0.0  0.0  0.0  0.0
 0.0  0.0  0.0  0.0  0.0  0.0  0.0  0.0     0.0  0.0  0.0  0.0  0.0  0.0  0.0  0.0
```

```julia
1  begin
2      # Init arrays
3      train_x = Array{Float64}(undef, 28^2, NUMBER_TRAIN_SAMPLES)
4      train_y = Array{Int64}(undef, NUMBER_TRAIN_SAMPLES)
5
6      # Init io streams
7      io_images = open(train_x_dir)
8      io_labels = open(train_y_dir)
9
10     # Iterating through sample length
11     for i ∈ 1:NUMBER_TRAIN_SAMPLES
12         seek(io_images, (i-1)*28^2 + 16) # offset 16 to skip header
13         seek(io_labels, (i-1)*1 + 8) # offset 8 to skip header
14         train_x[:,i] = convert(Array{Float64}, read(io_images, 28^2))
15         train_y[i] = convert(Int, read(io_labels, UInt8))
16     end
17
18     # Close io streams
19     close(io_images)
20     close(io_labels)
21
22     # Transpose features
23     train_x = train_x'
24 end
```

```
10000×784 adjoint(::Matrix{Float64}) with eltype Float64:
 0.0  0.0  0.0  0.0  0.0  0.0  0.0  0.0  …  0.0  0.0  0.0  0.0  0.0  0.0  0.0  0.0
 0.0  0.0  0.0  0.0  0.0  0.0  0.0  0.0     0.0  0.0  0.0  0.0  0.0  0.0  0.0  0.0
 0.0  0.0  0.0  0.0  0.0  0.0  0.0  0.0     0.0  0.0  0.0  0.0  0.0  0.0  0.0  0.0
 0.0  0.0  0.0  0.0  0.0  0.0  0.0  0.0     0.0  0.0  0.0  0.0  0.0  0.0  0.0  0.0
 0.0  0.0  0.0  0.0  0.0  0.0  0.0  0.0     0.0  0.0  0.0  0.0  0.0  0.0  0.0  0.0
 0.0  0.0  0.0  0.0  0.0  0.0  0.0  0.0  …  0.0  0.0  0.0  0.0  0.0  0.0  0.0  0.0
 0.0  0.0  0.0  0.0  0.0  0.0  0.0  0.0     0.0  0.0  0.0  0.0  0.0  0.0  0.0  0.0
 ⋮                             ⋮        ⋱                      ⋮
 0.0  0.0  0.0  0.0  0.0  0.0  0.0  0.0     0.0  0.0  0.0  0.0  0.0  0.0  0.0  0.0
 0.0  0.0  0.0  0.0  0.0  0.0  0.0  0.0  …  0.0  0.0  0.0  0.0  0.0  0.0  0.0  0.0
 0.0  0.0  0.0  0.0  0.0  0.0  0.0  0.0     0.0  0.0  0.0  0.0  0.0  0.0  0.0  0.0
 0.0  0.0  0.0  0.0  0.0  0.0  0.0  0.0     0.0  0.0  0.0  0.0  0.0  0.0  0.0  0.0
 0.0  0.0  0.0  0.0  0.0  0.0  0.0  0.0     0.0  0.0  0.0  0.0  0.0  0.0  0.0  0.0
 0.0  0.0  0.0  0.0  0.0  0.0  0.0  0.0     0.0  0.0  0.0  0.0  0.0  0.0  0.0  0.0
```

```julia
1  begin
2      # Init arrays
3      test_x = Array{Float64}(undef, 28^2, NUMBER_TEST_SAMPLES)
4      test_y = Array{Int64}(undef, NUMBER_TEST_SAMPLES)
5
6      # Init io streams
7      io_images_test = open(test_x_dir)
8      io_labels_test = open(test_y_dir)
9
10     # Iterating through sample length
11     for i ∈ 1:NUMBER_TEST_SAMPLES
12         seek(io_images_test, (i-1)*28^2 + 16) # offset 16 to skip header
13         seek(io_labels_test, (i-1)*1 + 8) # offset 8 to skip header
14         test_x[:,i] = convert(Array{Float64}, read(io_images_test, 28^2))
15         test_y[i] = convert(Int, read(io_labels_test, UInt8))
16     end
17
18     # Close io streams
19     close(io_images)
20     close(io_labels)
21
22     # Transpose features
23     test_x = test_x'
24 end
```

```
((60000, 784), (60000), (10000, 784), (10000))
```

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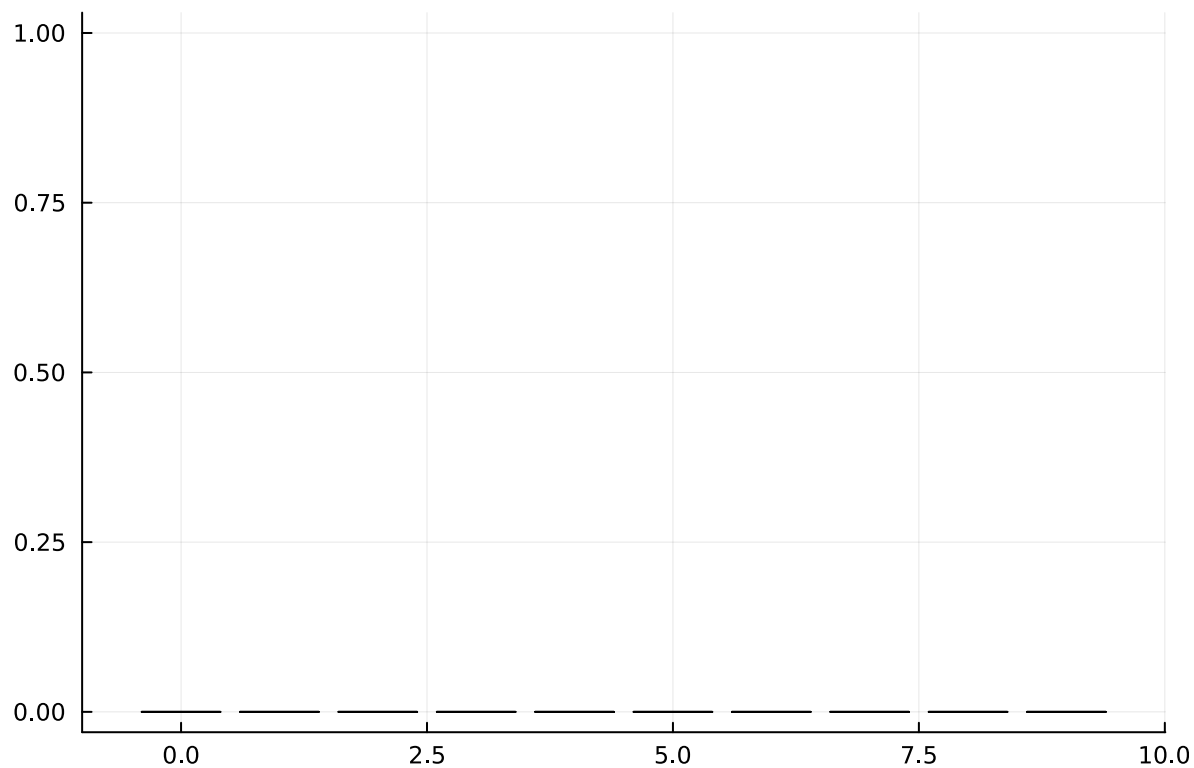
```
compute_average_intensity (generic function with 1 method)
```

```
1  #TODO compute average intensity for each label
2  function compute_average_intensity(x, y)
3      mean_ = zeros(10) # 10 is number of labels
4      #TODO compute average intensity for each label
5
6      return mean_
7  end
```

**l_mean** =   [0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0]

```
1  l_mean = compute_average_intensity(train_x, train_y)
```

Plot the average intensity using matplotlib



```
1  bar(0:9, l_mean, legend=false)
```

**UndefVarError: intensity not defined**

1. **top-level scope**  @ ⎸ *Local: 3*

```
1  begin
2      #TODO compute average intensity for each data sample
3      intensity =
4      size(intensity)
5  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

- $s_2$

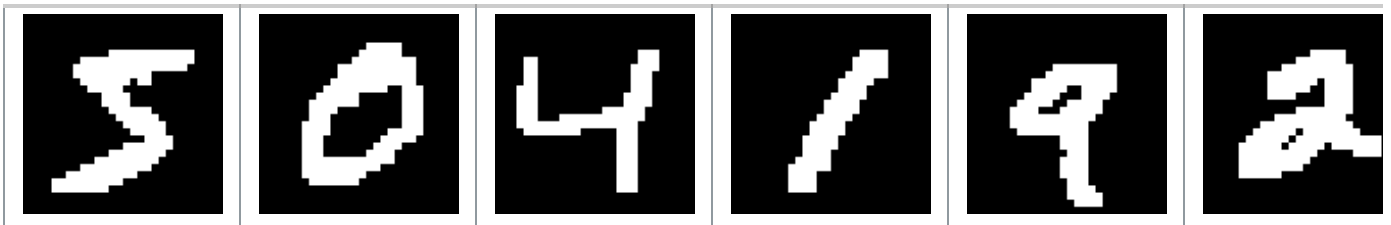  : flip the image along x-axis and compute the mean value of result

compute_symmetry (generic function with 1 method)

```julia
1  function compute_symmetry(train_x)
2      symmetry = []
3      for i in 1:size(train_x)[1]
4          img = reshape(train_x[i,:], (28,28))
5          s1 = mean(abs.(img - reverse(img, dims=1)))
6          s2 = mean(abs.(img - reverse(img, dims=2)))
7          s = -0.5 .* (s1 + s2)
8          append!(symmetry, s)
9      end
10     return symmetry
11 end
```
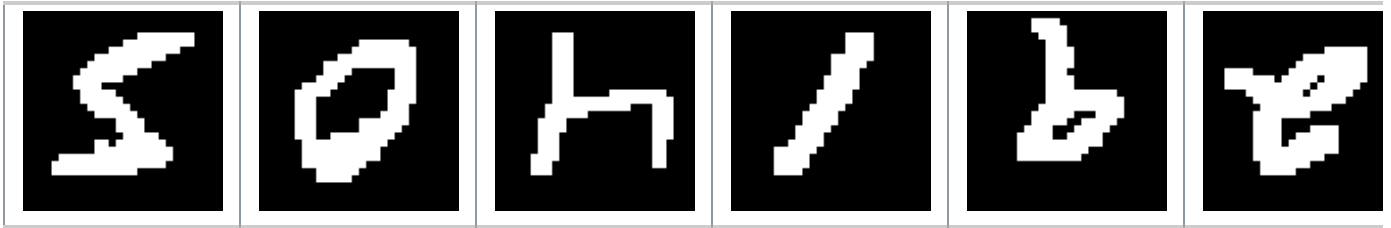
(60000)

```julia
1  begin
2      symmetry = compute_symmetry(train_x)
3      size(symmetry)
4  end
```

Visualize 10 samples in order to illustrate symmetry



(a vector displayed as a row to save space)

```julia
1  begin
2      num_img = 10
3      img_flat = train_x[1:num_img,:]
4      img = [reshape(img_flat[i,:], (28,28))' for i in 1:num_img]
5      [colorview(Gray, Float32.(img[i])) for i in 1:num_img]
6  end
```

(a vector displayed as a row to save space)

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

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

**UndefVarError: train_x_new not defined**

1. **top-level scope**  @  | *Local: 3*

```
1  begin
2      #TODO create X_new by horizontal stack intensity and symmetry
3      train_x_new =
4      size(train_x_new)
5  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**
- **TODO: Why should we normalize data?**

normalize (generic function with 3 methods)

```
1  function normalize(x, mean_=nothing, std_=nothing)
2    if mean_ == nothing && std_ == nothing
3       #TODO normalize x_train
4       #return 'normalized_train_x', 'mean_', 'std_'
5       #mean_ and std_ will be re-used to pre-process test set
6    end
7    return #return 'normalized_train_x' calculated by using mean_ and std_ (both of
   them are passed and not null)
8  end
9
10
```

*Another cell defining train_x_new contains errors.*

```
1  begin
2    normalized_train_x, mean_, std_ = normalize(train_x_new)
3
4    s_mean = compute_average_intensity(normalized_train_x, train_y)
5    bar(0:9, s_mean, legend=false)
6  end
```

# Construct data

(60000, 1)

```
1  begin
2    train_y_new = reshape(deepcopy(train_y), (size(train_y)[1], 1))
3    train_y_new[train_y_new .!= 1] .= 0
4    size(train_y_new)
5  end
```

*Another cell defining train_x_new contains errors.*

```
1  begin
2    # contruct data by adding ones
3    add_one_train_x = hcat(ones(size(normalized_train_x)[1],), normalized_train_x)
4    size(add_one_train_x)
5  end
```

# Sigmoid function and derivative of the sigmoid function

sigmoid_activation (generic function with 1 method)

```
1  function sigmoid_activation(x)
2    #TODO
3    """"compute the sigmoid activation value for a given input"""
4    #return?
5    return
6  end
```

sigmoid_deriv (generic function with 1 method)

```julia
1  function sigmoid_deriv(x)
2      #TODO
3      """
4      Compute the derivative of the sigmoid function ASSUMING
5      that the input 'x' has already been passed through the sigmoid
6      activation function
7      """
8      #return?
9      return
10 end
```

# Compute output

compute_h (generic function with 1 method)

```julia
1  function compute_h(W, X)
2      #TODO
3      """
4      Compute output: Take the inner product between our features 'X' and the weight
5      matrix 'W'
6      """
7      return
8  end
```

predict (generic function with 1 method)

```julia
1  function predict(W, X)
2      #TODO
3      """
4      Take the inner product between our features and weight matrix,
5      then pass this value through our sigmoid activation
6      """
7      preds =
8
9      # apply a step function to threshold the outputs to binary
10     # class labels
11     preds[preds .<= 0.5] .= 0
12     preds[preds .> 0] .= 1
13
14     return preds
15 end
```

# Compute gradient

**Loss Function: Average negative log likelihood**

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

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

$$\text{Cross-entropy: } J(w) = - \left( y log(z) + (1 - y) log(1 - z) \right)$$

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

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

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

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

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

compute_gradient (generic function with 1 method)

```
1  function compute_gradient(error, train_x)
2      #TODO
3      """
4      This is the gradient descent update of "average negative loglikelihood" loss
   function.
5      In lab02 our loss function is "sum squared error".
6      """
7      return
8  end
```

train (generic function with 1 method)

```julia
1  function train(W, train_x, train_y, learning_rate, num_epochs)
2      losses = []
3      for epoch in 1:num_epochs
4          y_hat = sigmoid_activation(compute_h(W, train_x))
5          error = y_hat - train_y
6          append!(losses, mean(-1 .* train_y .* log.(y_hat) .- (1 .- train_y) .* log.
   (1 .- y_hat)))
7          grad = compute_gradient(error, train_x)
8          W -= learning_rate * grad
9
10         if epoch == 1 || epoch % 50 == 0
11             print("Epoch=$epoch; Loss=$(losses[end])\n")
12         end
13     end
14     return W, losses
15 end
```

# Train our model

*Another cell defining __train_x_new__ contains errors.*

```julia
1  begin
2      W = rand(Normal(), (size(add_one_train_x)[2], 1))
3
4      num_epochs=2000
5      learning_rate=0.01
6      W, losses = train(W, add_one_train_x, train_y_new, learning_rate, num_epochs)
7  end
```

*Another cell defining __train_x_new__ and __num_epochs__ contains errors.*

```julia
1  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 1 method)

```julia
1  function tpfptnfn_cal(y_test, y_pred, positive_class=1)
2      true_positives = 0
3      false_positives = 0
4      true_negatives = 0
5      false_negatives = 0
6
7      # Calculate true positives, false positives, false negatives, and true negatives
8      for (true_label, predicted_label) in zip(y_test, y_pred)
9          if true_label == positive_class && predicted_label == positive_class
10             true_positives += 1
11         elseif true_label != positive_class && predicted_label == positive_class
12             false_positives += 1
13         elseif true_label == positive_class && predicted_label != positive_class
14             false_negatives += 1
15         elseif true_label != positive_class && predicted_label != positive_class
16             true_negatives += 1
17         end
18     end
19
20     return true_positives, false_positives, true_negatives, false_negatives
21 end
```

*Another cell defining __train_x_new__ contains errors.*

```julia
1  begin
2      preds_train = predict(W, add_one_train_x)
3      train_y_n = reshape(train_y_new, length(train_y_new), 1)
4
5      acc = 0
6      precision = 0
7      recall = 0
8      f1 = 0
9
10     for i ∈ 1:10
11         # Calculate true positives, false positives, false negatives, and true
           negatives
12         true_positives, false_positives, true_negatives, false_negatives =
           tpfptnfn_cal(train_y_n, preds_train)
13
14         # Calculate precision, recall, and F1-score
15         acc += (true_positives + true_negatives) / (true_positives + false_positives
           + true_negatives + false_negatives)
16         precision += true_positives / (true_positives + false_positives)
17         recall += true_positives / (true_positives + false_negatives)
18     end
19
20     acc = acc / 10
21     precision = precision / 10
22     recall = recall / 10
23     f1 = 2 * precision * recall / (precision + recall)
24     print(" acc: $acc\n precision: $precision\n recall: $recall\n f1_score: $f1\n")
25 end
```

# 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)`)

**UndefVarError: add_one_test_x not defined**

1. **top-level scope** @ ▌*Local: 4*

```
1  begin
2      #TODO
3      # compute test_y_new
4      test_y_new =
5
6
7      # compute test_intensity and test_symmetry to form test_x_new
8      test_intensity =
9      test_symmetry =
10     test_x_new =
11
12     # normalize test_x_new to form normalized_test_x
13     normalized_test_x =
14
15     # add column `ones` to test_x_new
16     add_one_test_x =
17     size(add_one_test_x)
18 end
```

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

*Another cell defining **train_x_new**, **add_one_test_x** and **test_y_new** contains errors.*

```julia
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
```

**TODO: Comment on the result**