

Table of Contents

Homework 5: Support Vector Networks

- Instructions for homework and submission

 - Instructions for doing homework

 - Instructions for submission

 - Content of the assignment

 - Others

- Problem statement

- Recall: Perceptron & Geometry Margin (Maximum 2.5 points)

 - Linear classifiers through origin

 - Perceptron Learning Algorithms

 - Convergence Proof

 - Geometric margin & SVM Motivation

- Linear Support Vector Machine (Maximum 6 points)

 - Hard-margin

 - Soft-margin

- Computing the SVM classifier (To get beyond 8.5 points)

 - SMO algorithm

 - Dual SVM - Hard-margin

 - Dual SVM - Soft-margin

- Multi-classes classification problem with SVMs (To get beyond 10.0 points)

 - Load MNIST dataset

 - Training SVMs

 - Evaluation

- References

Submission by: **Lê Nhật Nam** (20xx2023@student.hcmus.edu.vn)

```
student = ▶ (name = "Lê Nhật Nam", id = "20xx2023")
```

Homework 5: Support Vector Networks

CSC14005, Introduction to Machine Learning

This notebook was built for FIT@HCMUS student to learn about Support Vector Machines/or Support Vector Networks in the course CSC14005 - Introduction to Machine Learning.

Instructions for homework and submission

It's important to keep in mind that the teaching assistants will use a grading support application, so you must strictly adhere to the guidelines outlined in the instructions. If you are unsure, please ask the teaching assistants or the lab instructors as soon as you can. **Do not follow your personal preferences at stochastically**

Instructions for doing 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*.

Instructions for submission

- Before submitting, save this file as `<ID>.jl`. For example, if your ID is 123456, then your file will be `123456.jl`. Submit that file on Moodle.

Danger

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

Content of the assignment

- Recall: Perceptron & Geometrical Margin
- Linear support vector machine (Hard-margin, soft-margin)
- Popular non-linear kernels
- Computing SVM: Primal, Dual
- Multi-class SVM

Others

Other advice for you includes:

- Starting early and not waiting until the last minute
- Proceed with caution and gentleness.

"Living 'Slow' just means doing everything at the right speed – quickly, slowly, or at whatever pace delivers the best results." Carl Honoré.

- Avoid sources of interference, such as social networks, games, etc.

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

```
MersenneTwister(0)
```

```
1 Random.seed!(0)
```



Problem statement

Let $\mathcal{D} = \{(x_i, y_i) | x_i \in \mathbb{R}^d, y_i \in \{-1, 1\}\}_{i=1}^n$ be a dataset which is a set of pairs where $x_i \in \mathbb{R}^d$ is *data point* in some d -dimension vector space, and $y_i \in \{-1, 1\}$ is a *label* of the correspondent x_i data point classifying it to one of the two classes.

The model is trained on \mathcal{D} after which it is present with x_{i+1} , and is asked to predict the label of this previously unseen data point.

The prediction function is donated by $f(x) : \mathbb{R}^d \rightarrow \{-1, 1\}$

Recall: Perceptron & Geometry Margin (Maximum 2.5 points)

In fact, it is always possible to come up with such a "perfect" binary function if training samples are distinct. However, it is unclear whether such rules are applicable to data that does not exist in the training set. We don't need "learn-by-heart" learners; we need "intelligent" learners. More especially, such trivial rules do not suffice because our task is not to correctly classify the training set. Our task is to find a rule that works well for all new samples we would encounter in the access control setting; the training set is merely a helpful source of information to find such a function. We would like to find a classifier that "generalizes" well.

The key to finding a generalized classifier is to constrain the set of possible binary functions we can entertain. In other words, we would like to find a class of classifier functions such that if a function in this class works well on the training set, it is also likely to work well on the unseen images. This problem is considered a key problem named "model selection" in machine learning.

```
1 md"""
2 ## Recall: Perceptron & Geometry Margin (Maximum 2.5 points)
3
4 In fact, it is always possible to come up with such a "perfect" binary function if
  training samples are distinct. However, it is unclear whether such rules are
  applicable to data that does not exist in the training set. We don't need "learn-by-
  heart" learners; we need "intelligent" learners. More especially, such trivial rules
  do not suffice because our task is not to correctly classify the training set. Our
  task is to find a rule that works well for all new samples we would encounter in the
  access control setting; the training set is merely a helpful source of information to
  find such a function. We would like to find a classifier that "generalizes" well.
5
6 The key to finding a generalized classifier is to constrain the set of possible
  binary functions we can entertain. In other words, we would like to find a class of
  classifier functions such that if a function in this class works well on the training
  set, it is also likely to work well on the unseen images. This problem is considered
  a key problem named "model selection" in machine learning.
7 """
```

Linear classifiers through origin

For simplicity, we will just fix the function class for now. We will only consider a type of *linear classifiers*. For more formally, we consider the function of the form:

$$f(\mathbf{x}, \theta) = \text{sign}(\theta_1 \mathbf{x}_1 + \theta_2 \mathbf{x}_2 + \dots + \theta_d \mathbf{x}_d) = \text{sign}(\theta^T \mathbf{x})$$

where $\theta = [\theta_1, \theta_2, \dots, \theta_d]^T$ is a column vector of real valued parameters.

Different settings of the parameters give different functions in this class, i.e., functions whose value or output in $\{-1, 1\}$ could be different for some input \mathbf{x} .

Perceptron Learning Algorithms

After chosen a class of functions, we still have to find a specific function in this class that works well on the training set. This task often refers to estimation problem in machine learning. We would like to find θ that minimize the *training error*, i.e we would like to find a linear classifier that make fewest mistake in the training set.

$$\mathcal{L}(\theta) = \frac{1}{n} \sum_{t=1}^n (1 - \delta(y_t, f(\mathbf{x}; \theta))) = \frac{1}{n} \sum_{t=1}^n \text{Loss}(y_t, f(\mathbf{x}; \theta))$$

where $\delta(y, y') = 1$ if $y = y'$ and 0 if otherwise.

Perceptron update rule: Let k donates the number of parameter updates we have performed and $\theta^{(k)}$ is the parameter vector after k updates. Initially $k = 0$, and $\theta^{(k)} = \mathbf{0}$. We the loop through all the training instances (\mathbf{x}_t, y_t) , and updates the parameters only in response to mistakes,

$$\begin{cases} \theta^{(k+1)} \leftarrow \theta^{(k)} + y_t \mathbf{x}_t & \text{if } y_t (\theta^{(k+1)})^\top \mathbf{x}_t < 0 \\ \text{The parameters unchanged} \end{cases}$$

Geometry intuition of Perceptron

8

```
1 begin
2   n = 1000 # sample size
3   d = 2; # dimensionality of data
4   μ = 5 # mean
5   Σ = 8 # variance
6 end
```

```
points1_train =
2x500 Matrix{Float64}:
 9.63346  5.39397  6.67285  12.4209  ...  1.0458  8.02516  12.5985  4.6004
 8.10227  3.11485  -2.07827  4.62534  ... -3.02241  9.0917  8.29807  2.27181
```

```
1 points1_train = rand(MvNormal([Σ, μ], 5 .* [2 (μ - d)/Σ; (μ - d)/Σ d]), n ÷ 2)
```

```
points2_train =
2x500 Matrix{Float64}:
 0.0865891 -0.172354 -7.39168 -1.825 ... 0.597771 -4.65396 -2.17438 4.42893
 9.37609 10.483 6.2005 11.6646 12.0652 6.94485 4.4954 10.4031
```

```
1 points2_train = rand(MvNormal([-μ+d, Σ+d], 5 .* [3 (μ - d)/μ; (μ - d)/μ d]), n ÷ 2)
```

```
points1_test =
2x500 Matrix{Float64}:
 10.4694  7.27176  10.6064  10.3368  ...  10.399  6.32523  3.67006  4.26559
 -1.53045  3.94667  4.20854  4.10918  ...  3.6226  10.3159  3.40714  2.8007
```

```
1 points1_test = rand(MvNormal([Σ, μ], 5 .* [2 (μ - d)/Σ; (μ - d)/Σ d]), n ÷ 2)
```

```
points2_test =
2x500 Matrix{Float64}:
 -5.15394 -9.12328 -4.35624 -2.15082 ... 1.09356 -7.3331 -6.76925 1.35783
 9.48395 4.7303 6.69915 7.73966 6.90306 7.76199 7.93746 6.75461
```

```
1 points2_test = rand(MvNormal([-μ+d, Σ+d], 5 .* [3 (μ - d)/μ; (μ - d)/μ d]), n ÷ 2)
```

Todo

Your task here is implement the PLA (1 point). You can modify your own code in the area bounded by START YOUR CODE and END YOUR CODE.

```
1 md"""
2 !!! todo
3 Your task here is implement the PLA (1 point). You can modify your own code in the
  area bounded by START YOUR CODE and END YOUR CODE.
4 """
```

pla

Perceptron learning algorithm (PLA) implement function.

Fields

- pos_data::Matrix{Float64}: Input features for postive class (+1)
- neg_data::Matrix{Float64}: Input features for negative class (-1)
- n_epochs::Int64=10000: Maximum training epochs. Default is 10000
- η ::Float64=0.03: Learning rate. Default is 0.03

```
1 """
2     Perceptron learning algorithm (PLA) implement function.
3
4     ### Fields
5     - pos_data::Matrix{Float64}: Input features for postive class (+1)
6     - neg_data::Matrix{Float64}: Input features for negative class (-1)
7     - n_epochs::Int64=10000: Maximum training epochs. Default is 10000
8     -  $\eta$ ::Float64=0.03: Learning rate. Default is 0.03
9 """
10 function pla(pos_data::Matrix{Float64}, neg_data::Matrix{Float64},
11             n_epochs::Int64=10000,  $\eta$ ::Float64=0.03)
12     # START YOUR CODE
13
14     # END YOUR CODE
15     # return  $\theta$ 
16 end
```

$\theta_{ml} =$

```
1  $\theta_{ml} =$  pla(points1_train, points2_train)
```

draw_pla

Decision boundary visualization function for PLA

Fields

- θ : PLA parameters
- pos_data::Matrix{Float64}: Input features for postive class (+1)
- neg_data::Matrix{Float64}: Input features for negative class (-1)

```
1  """
2      Decision boundary visualization function for PLA
3
4  ### Fields
5  -  $\theta$ : PLA paramters
6  - pos_data::Matrix{Float64}: Input features for postive class (+1)
7  - neg_data::Matrix{Float64}: Input features for negative class (-1)
8  """
9  function draw_pla( $\theta$ , pos_data::Matrix{Float64}, neg_data::Matrix{Float64})
10     plt = scatter(pos_data[1, :], pos_data[2, :], label="y = 1")
11     scatter!(plt, neg_data[1, :], neg_data[2, :], label="y = -1")
12
13     b =  $\theta$ [3]
14      $\theta_{ml}$  =  $\theta$ [1:2]
15
16     decision(x) =  $\theta_{ml}' * x + b$ 
17
18     D = ([
19         tuple.(eachcol(pos_data), 1)
20         tuple.(eachcol(neg_data), -1)
21     ])
22
23     x_min = minimum(map((p) -> p[1][1], D))
24     y_min = minimum(map((p) -> p[1][2], D))
25     x_max = maximum(map((p) -> p[1][1], D))
26     y_max = maximum(map((p) -> p[1][2], D))
27
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
```

```
1 # Uncomment this line below when you finish your implementation
2 # draw_pla( $\theta_{ml}$ , points1_train, points2_train)
```

tpfptnfn_cal

Calculating values for True Positives (TP), False Positives (FP), True Negatives (TN), and False Negatives (FN)

Fields

- y_test: Actual labels
- y_pred: Predicted labels

```
1  """
2      Calculating values for True Positives (TP), False Positives (FP), True Negatives
3      (TN), and False Negatives (FN)
4  """
5  ### Fields
6  - y_test: Actual labels
7  - y_pred: Predicted labels
8  """
9  function tpfptnfn_cal(y_test, y_pred, positive_class=1)
10     true_positives = 0
11     false_positives = 0
12     true_negatives = 0
13     false_negatives = 0
14
15     # Calculate true positives, false positives, false negatives, and true negatives
16     for (true_label, predicted_label) in zip(y_test, y_pred)
17         if true_label == positive_class && predicted_label == positive_class
18             true_positives += 1
19         elseif true_label != positive_class && predicted_label == positive_class
20             false_positives += 1
21         elseif true_label == positive_class && predicted_label != positive_class
22             false_negatives += 1
23         elseif true_label != positive_class && predicted_label != positive_class
24             true_negatives += 1
25         end
26     end
27
28     return true_positives, false_positives, true_negatives, false_negatives
29 end
```


eval_pla

Evaluation function for PLA to calculate accuracy

Fields

- θ : PLA parameters
- pos_data::Matrix{Float64}: Input features for postive class (+1)
- neg_data::Matrix{Float64}: Input features for negative class (-1)

```
1  """
2      Evaluation function for PLA to calculate accuracy
3
4  ### Fields
5  -  $\theta$ : PLA paramters
6  - pos_data::Matrix{Float64}: Input features for postive class (+1)
7  - neg_data::Matrix{Float64}: Input features for negative class (-1)
8  """
9  function eval_pla( $\theta$ , pos_data, neg_data)
10     n = size(pos_data, 2)
11     X = vcat(hcat(pos_data, neg_data), ones(n * 2)')
12
13     y_test = vcat(ones(n), -ones(n))'
14     y_pred = [sign(x) for x in  $\theta'$  * X]
15
16     # START YOUR CODE
17     # TODO: acc, p, r, f1???
18
19     # END YOUR CODE
20
21     print(" acc: $acc\n precision: $p\n recall: $r\n f1_score: $f1\n")
22     return acc, p, r, f1
23 end
```

```
1 # Uncomment this line below when you finish your implementation
2 # eval_pla( $\theta_{ml}$ , points1_test, points2_test)
```

Convergence Proof

Assume that all the training instances have bounded Euclidean norms, i.e $\|\mathbf{x}\| \leq R$. Assume that exists a linear classifier in class of functions with finite parameter values that correctly classifies all the training instances. For precisely, we assume that there is some $\gamma > 0$ such that $\mathbf{y}_t(\boldsymbol{\theta}^*)^\top \mathbf{x}_t \geq \gamma$ for all $t = 1 \dots n$.

The convergence proof is based on combining two results:

- **Result 1:** we will show that the inner product $(\boldsymbol{\theta}^*)^\top \boldsymbol{\theta}^{(k)}$ increases at least linearly with each update.

Todo

Your task here is show the proof of result 1. (0.25 point)

```
1 md"""
2 !!! todo
3 Your task here is show the proof of result 1. (0.25 point)
4 """
```

START YOUR PROOF

<content>

END YOUR PROOF

- **Result 2:** The squared norm $\|\theta^{(k)}\|^2$ increases at most linearly in the number of updates k .

Todo

Your task here is show the proof of result 2. (0.25 point)

```
1 md"""
2 !!! todo
3 Your task here is show the proof of result 2. (0.25 point)
4 """
```

START YOUR PROOF

<content>

END YOUR PROOF

We can now combine parts 1) and 2) to bound the cosine of the angle between $\theta^{(k)}$ and θ^* . Since cosine is bounded by one, thus

$$1 \geq \frac{k\gamma}{\sqrt{kR^2} \|\theta^{(*)}\|} \leftrightarrow k \leq \frac{R^2 \|\theta^{(*)}\|^2}{\gamma^2}$$

By combining the two we can show that the cosine of the angle between $\theta^{(k)}$ and θ^* has to increase by a finite increment due to each update. Since cosine is bounded by one, it follows that we can only make a finite number of updates.

Geometric margin & SVM Motivation

There is a question? Does $\frac{\|\theta^{(*)}\|^2}{\gamma^2}$ relate to how difficult the classification problem is? Its inverse, i.e., $\frac{\gamma^2}{\|\theta^{(*)}\|^2}$ is the smallest distance in the vector space from any samples to the decision boundary specified by $\theta^{(*)}$. In other words, it serves as a measure of how well the two classes of data are separated (by a linear boundary). We call this is geometric margin, denoted by γ_{geom} . As a result, the bound on the number of perceptron updates can be written more succinctly in terms of the geometric margin γ_{geom} (You know that man, Vapnik–Chervonenkis Dimension)

$$k \leq \left(\frac{R}{\gamma_{geom}} \right)^2$$

. We note some interesting thing about the result:

- Does not depend (directly) on the dimension of the data, nor
- number of training instances

Can't we find such a large margin classifier directly? YES, in this homework, you will do it with Support Vector Machine :)

Linear Support Vector Machine (Maximum 6 points)

From the problem statement section, we are given

$$\{(x_i, y_i) | x_i \in \mathbb{R}^d, y_i \in \{-1, 1\}\}_{i=1}^n$$

And based on previous section, we want to find the "maximum-geometric margin" that divides the space into two parts so that the distance between the hyperplane and the nearest point from either class is maximized. Any hyperplane can be written as the set of data points \mathbf{x} satisfying

$$\theta^T \mathbf{x} + b = 0$$

```
1 md"""
2 ## Linear Support Vector Machine (Maximum 6 points)
3
4 From the problem statement section, we are given
5
6  $\{(x_i, y_i) \mid x_i \in \mathbb{R}^d, y_i \in \{-1, 1\}\}_{i=1}^n$ 
7
8 And based on previous section, we want to find the "maximum-geometric margin" that
9 divides the space into two parts so that the distance between the hyperplane and the
10 nearest point from either class is maximized. Any hyperplane can be written as the
11 set of data points  $\mathbf{x}$  satisfying
12
13  $\theta^T \mathbf{x} + b = 0$ 
14 """
```

Hard-margin

The goal of SVM is to choose two parallel hyperplanes that separate the two classes of data in order to maximize the distance between them. The region defined by these two hyperplanes is known as the "margin," and the maximum-margin hyperplane is the one located halfway between them. And these hyperplane can be described as

$$\theta^T \mathbf{x} + b = 1 (\text{anything on or above this boundary is of one class, with label } 1)$$

and

$$\theta^T \mathbf{x} + b = -1 (\text{anything on or below this boundary is of the other class, with label } -1)$$

Geometrically, the distance between these two hyperplanes is $\frac{2}{\|\theta\|}$

Todo

Your task here is show that the distance between these two hyperplanes is $\frac{2}{\|\theta\|}$ (1 point). You can modify your own code in the area bounded by START YOUR PROOF and END YOUR PROOF.

START YOUR PROOF

<content>

END YOUR PROOF

So we want to maximize the distance between these two hyperplanes? Right? Equivalently, we minimize $\|\theta\|$. We also have to prevent data points from falling into the margin, we add the following constraint: for each i either

$$\theta^T \mathbf{x}_i + b \geq 1 \text{ if } y_i = 1$$

and

$$\theta^T \mathbf{x} + b \leq -1 \text{ if } y_i = -1$$

And, we can rewrite this as

$$y_i(\theta^T \mathbf{x}_i + b) \geq 1, \forall i \in \{1 \dots n\}$$

Finally, the optimization problem is

$$\begin{aligned} \min_{\theta, b} \quad & \frac{1}{2} \|\theta\|^2 \\ \text{s.t.} \quad & y_i(\theta^T \mathbf{x}_i + b) - 1 \geq 0, \forall i = 1 \dots n \end{aligned}$$

The parameters θ and b that solve this problem determine the classifier

$$\mathbf{x} \rightarrow \text{sign}(\theta^T \mathbf{x}_i + b)$$

Todo

Your task here is implement the hard-margin SVM solving the primal formulation using gradient descent (3 points). You can modify your own code in the area bounded by START YOUR CODE and END YOUR CODE.

```
1 md"""
2 !!! todo
3 Your task here is implement the hard-margin SVM solving the primal formulation using
  gradient descent (3 points). You can modify your own code in the area bounded by
  START YOUR CODE and END YOUR CODE.
4 """
```

hardmargin_svm

SVM solving the primal formulation using gradient descent (hard-margin)

Fields

- `pos_data::Matrix{Float64}`: Input features for postive class (+1)
- `neg_data::Matrix{Float64}`: Input features for negative class (-1)
- `η::Float64=0.03`: Learning rate. Default is 0.03
- `n_epochs::Int64=10000`: Maximum training epochs. Default is 10000

```
1 """
2 SVM solving the primal formulation using gradient descent (hard-margin)
3 ### Fields
4 - pos_data::Matrix{Float64}: Input features for postive class (+1)
5 - neg_data::Matrix{Float64}: Input features for negative class (-1)
6 - η::Float64=0.03: Learning rate. Default is 0.03
7 - n_epochs::Int64=10000: Maximum training epochs. Default is 10000
8 """
9 function hardmargin_svm(pos_data, neg_data, η=0.04, n_epochs=10000)
10     # START YOUR CODE
11
12     ## Create variables for the separating hyperplane  $w'x = b$ .
13
14     ## Loss function
15
16     # Train using gradient descent
17     ## For each epoch
18     ### For each training instance  $\in D$ 
19
20     ## Update weight
21
22     # END YOUR CODE
23     ## Return hyperplane parameters
24     #return w, b
25 end
```

```
1 # Uncomment this line below when you finish your implementation
2 # w, b = hardmargin_svm(points1_train, points2_train)
```

draw

Visualization function for SVM solving the primal formulation using gradient descent (hard-margin)

Fields

- `w` & `b`: SVM parameters
- `pos_data::Matrix{Float64}`: Input features for positive class (+1)
- `neg_data::Matrix{Float64}`: Input features for negative class (-1)

```
1  """
2      Visualization function for SVM solving the primal formulation using gradient
3      descent (hard-margin)
4
5  """
6  ### Fields
7  - w & b: SVM parameters
8  - pos_data::Matrix{Float64}: Input features for positive class (+1)
9  - neg_data::Matrix{Float64}: Input features for negative class (-1)
10 """
11 function draw(w, b, pos_data, neg_data)
12     plt = scatter(pos_data[1, :], pos_data[2, :], label="y = 1")
13     scatter!(plt, neg_data[1, :], neg_data[2, :], label="y = -1")
14
15     hyperplane(x) = w' * x + b
16
17     D = ([
18         tuple.(eachcol(pos_data), 1)
19         tuple.(eachcol(neg_data), -1)
20     ])
21
22     x_min = minimum(map((p) -> p[1][1], D))
23     y_min = minimum(map((p) -> p[1][2], D))
24     x_max = maximum(map((p) -> p[1][1], D))
25     y_max = maximum(map((p) -> p[1][2], D))
26
27     contour!(plt, x_min:0.1:x_max, y_min:0.1:y_max,
28         (x, y) -> hyperplane([x, y]),
29         levels=[-1],
30         linestyle=:dash,
31         colorbar_entry=false, color=:red, label = "Negative points")
32     contour!(plt, x_min:0.1:x_max, y_min:0.1:y_max,
33         (x, y) -> hyperplane([x, y]),
34         levels=[0], linestyle=:solid, label="SVM prediction",
35         colorbar_entry=false, color=:green)
36     contour!(plt, x_min:0.1:x_max, y_min:0.1:y_max,
37         (x, y) -> hyperplane([x, y]), levels=[1], linestyle=:dash,
38         colorbar_entry=false, color=:blue, label = "Positive points")
39 end
```

```
1 # Uncomment this line below when you finish your implementation
2 # draw(w, b, points1_train, points2_train)
```

eval_svm

Evaluation function for hard-margin & soft-margin SVM to calculate accuracy

Fields

- θ : PLA parameters
- pos_data::Matrix{Float64}: Input features for positive class (+1)
- neg_data::Matrix{Float64}: Input features for negative class (-1)

```
1  """
2      Evaluation function for hard-margin & soft-margin SVM to calculate accuracy
3
4  ### Fields
5  -  $\theta$ : PLA parameters
6  - pos_data::Matrix{Float64}: Input features for positive class (+1)
7  - neg_data::Matrix{Float64}: Input features for negative class (-1)
8  """
9  function eval_svm(w, b, pos_data, neg_data)
10     n = size(pos_data, 2)
11     X = hcat(pos_data, neg_data)
12
13     # Actual labels, and predicted labels
14     y_test = vcat(ones(n), -ones(n))'
15     y_pred = [sign(x) for x in w' * X .+ b]
16
17     # START YOUR CODE
18     # TODO: acc, p, r, f1???
19
20     # END YOUR CODE
21
22     print(" acc: $acc\n precision: $p\n recall: $r\n f1_score: $f1\n")
23
24     return acc, p, r, f1
25 end
```

```
1 # Uncomment this line below when you finish your implementation
2 # eval_svm(w, b, points1_test, points2_test)
```

Soft-margin

The limitation of Hard Margin SVM is that it only works for data that can be separated linearly. In reality, however, this would not be the case. In practice, the data will almost certainly contain noise and may not be linearly separable. In this section, we will talk about soft-margin SVM (an relaxation of the optimization problem).

Basically, the trick here is very simple, we add slack variables ς_i to the constraint of the optimization problem.

$$y_i(\theta^\top \mathbf{x}_i + b) \geq 1 - \varsigma_i, \forall i = 1 \dots n$$

The regularized optimization problem become as

$$\begin{aligned} \min_{\theta, b, \varsigma} \quad & \frac{1}{2} \|\theta\|^2 + \sum_{i=1}^n \varsigma_i \\ \text{s.t.} \quad & y_i(\theta^\top \mathbf{x}_i + b) \geq 1 - \varsigma_i, \forall i = 1 \dots n \end{aligned}$$

Furthermore, we add a regularization parameter C to determine how important ς should be. And, we got it :)

$$\begin{aligned} \min_{\theta, b, \varsigma} \quad & \frac{1}{2} \|\theta\|^2 + C \sum_{i=1}^n \varsigma_i \\ \text{s.t.} \quad & y_i(\theta^\top \mathbf{x}_i + b) \geq 1 - \varsigma_i, \varsigma_i \geq 0, \forall i = 1 \dots n \end{aligned}$$

Todo

Your task here is implement the soft-margin SVM solving the primal formulation using gradient descent (3 points). You can modify your own code in the area bounded by START YOUR CODE and END YOUR CODE.

softmargin_svm

SVM solving the primal formulation using gradient descent (soft-margin)

Fields

- `pos_data::Matrix{Float64}`: Input features for positive class (+1)
- `neg_data::Matrix{Float64}`: Input features for negative class (-1)
- `C`: relaxation variable control slack variables ς
- `η::Float64=0.03`: Learning rate. Default is 0.03
- `n_epochs::Int64=10000`: Maximum training epochs. Default is 10000

```
1  """
2      SVM solving the primal formulation using gradient descent (soft-margin)
3  """
4  - pos_data::Matrix{Float64}: Input features for positive class (+1)
5  - neg_data::Matrix{Float64}: Input features for negative class (-1)
6  - C: relaxation variable control slack variables  $\varsigma$ 
7  - η::Float64=0.03: Learning rate. Default is 0.03
8  - n_epochs::Int64=10000: Maximum training epochs. Default is 10000
9  """
10 function softmargin_svm(pos_data, neg_data, n_epochs=10000, C=0.12, η=0.01)
11     # START YOUR CODE
12
13     ## Create variables for the separating hyperplane  $w'x = b$ .
14
15     ## Loss function
16
17     # Train using gradient descent
18     ## For each epoch
19     ### For each training instance  $\epsilon \in D$ 
20
21     #### Calculate slack variables  $\varsigma$ 
22
23     ## Update weight
24
25     # END YOUR CODE
26     ## Return hyperplane parameters
27     # return w, b
28 end
```

```
1 # Uncomment this line below when you finish your implementation
2 # sw, sb = softmargin_svm(points1_train, points2_train)
```

```
1 # Uncomment this line below when you finish your implementation
2 # draw(sw, sb, points1_train, points2_train)
```

```
1 # Uncomment this line below when you finish your implementation
2 # eval_svm(sw, sb, points1_test, points2_test)
```

Computing the SVM classifier (To get beyond 8.5 points)

We should know about some popular kernel types we could use to classify the data such as linear kernel, polynomial kernel, Gaussian, sigmoid and RBF (radial basis function) kernel.

- Linear Kernel: $K(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i^\top \mathbf{x}_j$
- Polynomial kernel: $K(\mathbf{x}_i, \mathbf{x}_j) = (1 + \mathbf{x}_i^\top \mathbf{x}_j)^p$
- Gaussian: $K(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{2\sigma^2}\right)$
- Sigmoid: $K(\mathbf{x}_i, \mathbf{x}_j) = \tanh(\beta_0 \mathbf{x}_i^\top \mathbf{x}_j + \beta_1)^p$
- RBF kernel: $K(\mathbf{x}_i, \mathbf{x}_j) = \exp(-\gamma \|\mathbf{x}_i - \mathbf{x}_j\|^2)$

```
1 md"""
2 ## Computing the SVM classifier (To get beyond 8.5 points)
3
4 We should know about some popular kernel types we could use to classify the data such
  as linear kernel, polynomial kernel, Gaussian, sigmoid and RBF (radial basis
  function) kernel.
5 - Linear Kernel:  $K(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i^\top \mathbf{x}_j$ 
6 - Polynomial kernel:  $K(\mathbf{x}_i, \mathbf{x}_j) = (1 + \mathbf{x}_i^\top \mathbf{x}_j)^p$ 
7 - Gaussian:  $K(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{2\sigma^2}\right)$ 
8 - Sigmoid:  $K(\mathbf{x}_i, \mathbf{x}_j) = \tanh(\beta_0 \mathbf{x}_i^\top \mathbf{x}_j + \beta_1)^p$ 
9 - RBF kernel:  $K(\mathbf{x}_i, \mathbf{x}_j) = \exp(-\gamma \|\mathbf{x}_i - \mathbf{x}_j\|^2)$ 
10 """
```

two_spirals

Function for creating two spirals dataset.

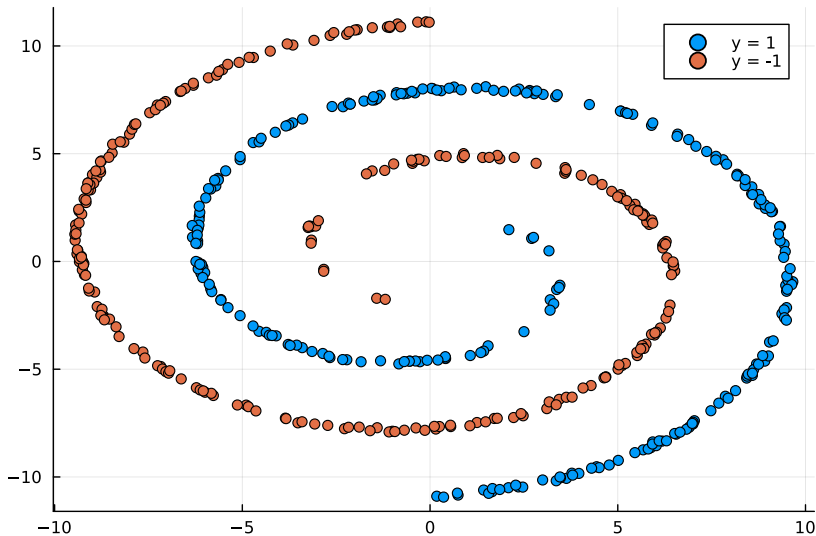
You can check the MATLAB implement here: 6 functions for generating artificial datasets, <https://www.mathworks.com/matlabcentral/fileexchange/41459-6-functions-for-generating-artificial-datasets>

FIELDS

- `n_samples`: number of samples you want :)
- `noise`: noise rate for creating process you want :)

```
1  """
2      Function for creating two spirals dataset.
3
4      You can check the MATLAB implement here: 6 functions for generating artificial
5      datasets, https://www.mathworks.com/matlabcentral/fileexchange/41459-6-functions-
6      for-generating-artificial-datasets
7
8  """
9  function two_spirals(n_samples, noise::Float64=0.2)
10     start_angle =  $\pi$  / 2
11     total_angle =  $3\pi$ 
12
13     N1 = floor(Int, n_samples / 2)
14     N2 = n_samples - N1
15
16     n = start_angle .+ sqrt.(rand(N1, 1)) .* total_angle
17     d1 = [-cos.(n) .* n + rand(N1, 1) .* noise sin.(n) .* n + rand(N1, 1) .* noise]
18
19     n = start_angle .+ sqrt.(rand(N2, 1)) .* total_angle
20     d2 = [cos.(n) .* n + rand(N2, 1) * noise -sin.(n) .* n + rand(N2, 1) .* noise]
21
22     return d1', d2'
23 end
```

```
► (2×250 adjoint(::Matrix{Float64}) with eltype Float64:
  -0.400044  9.12106  -2.66661  -5.44424  ...  6.29711  6.9828  6.65789  8.85921  -6.
1  # create two spirals which are not linearly seperable
2  sp_points1, sp_points2 = two_spirals(500)
```



```
1 scatter!(scatter(sp_points1[1, :], sp_points1[2, :], label="y = 1"), sp_points2[1, :], sp_points2[2, :], label="y = -1")
```

```
γ = 0.2
```

```
1 # Kernel function: in this lab, we use RBF kernel function, you want to do more experiment, please try again at home
```

```
2 γ = 1 / 5
```

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

```
1 K(x, y) = exp(-γ * (x - y)' * (x - y))
```

SMO algorithm

For more detail, you should read: Platt, J. (1998). Sequential minimal optimization: A fast algorithm for training support vector machines.

Wikipedia just quite good for describes this algorithm: MO is an iterative algorithm for solving the optimization problem. MO breaks this problem into a series of smallest possible sub-problems, which are then solved analytically. Because of the linear equality constraint involving the Lagrange multipliers λ_i , the smallest possible problem involves two such multipliers.

The SMO algorithm proceeds as follows:

- Step 1: Find a Lagrange multiplier α_1 that violates the Karush–Kuhn–Tucker (KKT) conditions for the optimization problem.
- Step 2: Pick a second multiplier α_2 and optimize the pair (α_1, α_2)
- Step 3: Repeat steps 1 and 2 until convergence.

Dual SVM - Hard-margin

If you want to find minimum of a function f under the equality constraint g , we can use Lagrangian function

$$f(x) - \lambda g(x) = 0$$

where λ is Lagrange multiplier.

In terms of SVM optimization problem

$$\begin{aligned} \min_{\theta, b} \quad & \frac{1}{2} \|\theta\|^2 \\ \text{s.t.} \quad & y_i(\theta^\top \mathbf{x}_i + b) - 1 \geq 0, \forall i = 1 \dots n \end{aligned}$$

The equality constraint is $g(\theta, b) = y_i(\theta^\top \mathbf{x}_i + b) - 1, \forall i = 1 \dots n$

Then the Lagrangian function is

$$\mathcal{L}(\theta, b, \lambda) = \frac{1}{2} \|\theta\|^2 + \sum_i^n \lambda_i (y_i(\theta^\top \mathbf{x}_i + b) - 1)$$

Equivalently, Lagrangian primal problem is formulated as

$$\begin{aligned} \min_{\theta, b} \quad & \max \mathcal{L}(\theta, b, \lambda) \\ \text{s.t.} \quad & \lambda_i \geq 0, \forall i = 1 \dots n \end{aligned}$$

Note

We need to MINIMIZE the MAXIMIZATION of $\mathcal{L}(\theta, b, \lambda)$? What we are doing???

Danger

More precisely, λ here should be KKT (Karush-Kuhn-Tucker) multipliers

$$\lambda[-y_i(\theta^\top \mathbf{x}_i + b) + 1] = 0, \forall i = 1 \dots n$$

```

1 md"""
2 ### Dual SVM - Hard-margin
3
4 If you want to find minimum of a function  $f$  under the equality constraint  $g$ , we
  can use Lagrangian function
5
6  $f(x) - \lambda g(x) = 0$ 
7 where  $\lambda$  is Lagrange multiplier.
8
9 In terms of SVM optimization problem
10
11 
$$\begin{gather*} \underset{\{\theta, b\}}{\text{min}} \frac{1}{2} \|\theta\|^2 \\ \text{s.t.} \quad y_i (\mathbf{\theta}^{\text{top}} \mathbf{x}_i + b) - 1 \geq 0, \text{ for all } i = 1 \dots n \end{gather*}$$

12
13 The equality constraint is  $g(\theta, b) = y_i (\mathbf{\theta}^{\text{top}} \mathbf{x}_i + b) - 1, \text{ for all } i = 1 \dots n$ 
14
15 Then the Lagrangian function is
16
17 
$$\mathcal{L}(\theta, b, \lambda) = \frac{1}{2} \|\theta\|^2 + \sum_{i=1}^n \lambda_i (y_i (\mathbf{\theta}^{\text{top}} \mathbf{x}_i + b) - 1)$$

18
19 Equivalently, Lagrangian primal problem is formulated as
20
21 
$$\begin{gather*} \underset{\{\theta, b\}}{\text{min}} \quad \underset{\{\lambda\}}{\text{max}} \quad \mathcal{L}(\theta, b, \lambda) \\ \text{s.t.} \quad \lambda_i \geq 0, \text{ for all } i = 1 \dots n \end{gather*}$$

22
23 !!! note
24
25 We need to MINIMIZE the MAXIMIZATION of  $\mathcal{L}(\theta, b, \lambda)$ ? What
  we are doing???
26
27 !!! danger
28
29 More precisely,  $\lambda$  here should be KKT (Karush-Kuhn-Tucker) multipliers
30
31 
$$\lambda [-y_i (\mathbf{\theta}^{\text{top}} \mathbf{x}_i + b) + 1] = 0, \text{ for all } i = 1 \dots n$$

32
33 """

```

With the Lagrangian function

$$\begin{aligned} \min_{\theta, b} \max \mathcal{L}(\theta, b, \lambda) &= \frac{1}{2} \|\theta\|^2 + \sum_{i=1}^n \lambda_i (y_i(\theta^\top \mathbf{x}_i + b - 1)) \\ \text{s.t. } \lambda_i &\geq 0, \forall i = 1 \dots n \end{aligned}$$

Setting derivatives to 0 yield:

$$\begin{aligned} \nabla_{\theta} \mathcal{L}(\theta, b, \lambda) &= \theta - \sum_{i=1}^n \lambda_i y_i \mathbf{x}_i = 0 \Leftrightarrow \theta^* = \sum_{i=1}^n \lambda_i y_i \mathbf{x}_i \\ \nabla_b \mathcal{L}(\theta, b, \lambda) &= - \sum_{i=1}^n \lambda_i y_i = 0 \end{aligned}$$

We substitute them into the Lagrangian function, and get

$$W(\lambda, b) = \sum_{i=1}^n \lambda_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \lambda_i \lambda_j y_i y_j \mathbf{x}_i \mathbf{x}_j$$

So, dual problem is stated as

$$\begin{aligned} \max_{\lambda} \sum_{i=1}^n \lambda_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \lambda_i \lambda_j y_i y_j \mathbf{x}_i \mathbf{x}_j \\ \text{s.t. } \lambda_i \geq 0, \forall i = 1 \dots n, \sum_{i=1}^n \lambda_i y_i = 0 \end{aligned}$$

To solve this one has to use quadratic optimization or **sequential minimal optimization**

```

draw_nl (generic function with 1 method)
1 function draw_nl( $\lambda$ , b, pos_data, neg_data)
2   plt = scatter(pos_data[1, :], pos_data[2, :], label="y = 1")
3   scatter!(plt, neg_data[1, :], neg_data[2, :], label="y = -1")
4
5   D = ([
6     tuple.(eachcol(pos_data), 1)
7     tuple.(eachcol(neg_data), -1)
8   ])
9
10  X = [x for (x, y) in D]
11  Y = [y for (x, y) in D]
12
13  k(x, y) = exp(-1 / 5 * (x - y)' * (x - y))
14
15  hyperplane(x) = ( $\lambda$  .* Y) . k.(X, Ref(x)) + b
16
17  x_min = minimum(map((p) -> p[1][1], D))
18  y_min = minimum(map((p) -> p[1][2], D))
19  x_max = maximum(map((p) -> p[1][1], D))
20  y_max = maximum(map((p) -> p[1][2], D))
21
22  contour!(plt, x_min:0.1:x_max, y_min:0.1:y_max,
23    (x, y) -> hyperplane([x, y]),
24    levels=[-1],
25    linestyle=:dash,
26    colorbar_entry=false, color=:red, label = "Negative points")
27  contour!(plt, x_min:0.1:x_max, y_min:0.1:y_max,
28    (x, y) -> hyperplane([x, y]),
29    levels=[0], linestyle=:solid, label="SVM prediction",
30    colorbar_entry=false, color=:green)
31  contour!(plt, x_min:0.1:x_max, y_min:0.1:y_max,
32    (x, y) -> hyperplane([x, y]), levels=[1], linestyle=:dash,
33    colorbar_entry=false, color=:blue, label = "Positive points")
34 end

```

Todo

Your task here is implement the hard-margin SVM solving the dual formulation using sequential minimal optimization (2 points). You can modify your own code in the area bounded by START YOUR CODE and END YOUR CODE.


```

dualsvm_smo_hard (generic function with 4 methods)
1 function dualsvm_smo_hard(pos_data, neg_data, n_epochs=100,  $\lambda_{tol}=0.0001$ ,
  err_tol=0.0001)
2   # You do not need implement kernel, please use the K(.) kernel function in
  previous cell code.
3
4   # START YOUR CODE
5   # Step 1: Data preparation
6   # First you construct and shuffle to obtain dataset D in a stochastically manner
7
8   # For more easily access to data point
9   X = [x for (x, y) ∈ D]
10  Y = [y for (x, y) ∈ D]
11
12  # Step 2: Initialization
13  # Larangian multipliers, and bias
14  λ = zeros(length(D))
15  b = 0
16  n = length(λ)
17
18  # Step 3: Training loop
19
20
21  # END YOUR CODE
22  ## Return hyperplane parameters
23  # return λ, b
24 end

```

```

1 # Uncomment this line below when you finish your implementation
2 # λh, bh = dualsvm_smo_hard(points1train, points2train)

```

```

1 # Uncomment this line below when you finish your implementation
2 # draw_n1(λh, bh, points1train, points2train)

```

Dual SVM - Soft-margin

As we know that, the regularized optimization problem in the case of soft-margin as

$$\begin{aligned}
 & \min_{\theta, b, \varsigma} \frac{1}{2} \|\theta\|^2 + C \sum_{i=1}^n \varsigma_i \\
 \text{s.t. } & y_i(\theta^T \mathbf{x}_i + b) \geq 1 - \varsigma_i, \varsigma_i \geq 0, \forall i = 1 \dots n
 \end{aligned}$$

We use Larangian multipliers, and transform to a dual problem as

$$\begin{aligned}
 & \max_{\lambda} \sum_{i=1}^n \lambda_i - \frac{1}{2} \sum_i^n \sum_j^n \lambda_i \lambda_j y_i y_j \mathbf{x}_i \mathbf{x}_j \\
 \text{s.t. } & 0 \leq \lambda_i \leq C, \forall i = 1 \dots n, \sum_{i=1}^n \lambda_i y_i = 0
 \end{aligned}$$

Todo

Your task here is implement the soft-margin SVM solving the dual formulation using sequential minimal optimization (2 points). You can modify your own code in the area bounded by START YOUR CODE and END YOUR CODE.

dualsvm_smo_soft (generic function with 5 methods)

```
1 function dualsvm_smo_soft(pos_data, neg_data, n_epochs=100, C=1000,  $\lambda_{tol}$ =0.0001,
2   err_tol=0.0001)
3   # START YOUR CODE
4
5   # Step 1: Data preparation
6   # First you construct and shuffle to obtain dataset D in a stochastically manner
7
8   # For more easily access to data point
9   X = [x for (x, y) ∈ D]
10  Y = [y for (x, y) ∈ D]
11
12  # Step 2: Initialization
13  # Lagrangian multipliers, and bias
14  λ = zeros(length(D))
15  b = 0
16  n = length(λ)
17
18  # Step 3: Training loop
19
20  # END YOUR CODE
21  ## Return hyperplane parameters
22  # return λ, b
23 end
```

```
1 # Uncomment this line below when you finish your implementation
2 # λs, bs = dualsvm_smo_soft(sp_points1, sp_points2)
```

```
1 # Uncomment this line below when you finish your implementation
2 # draw_n(λs, bs, sp_points1, sp_points2)
```

Multi-classes classification problem with SVMs (To get beyond 10.0 points)

```
1 md"""
2 ## Multi-classes classification problem with SVMs (To get beyond 10.0 points)
3 """
```

Load MNIST dataset

```
1 md"""
2 ### Load MNIST dataset
3 """
```

10000

```
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
```

```

784x60000 Matrix{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
0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 ... 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0

```

```

1 begin
2   train_x = Array{Float64}(undef, 28^2, NUMBER_TRAIN_SAMPLES)
3   train_y = Array{Int64}(undef, NUMBER_TRAIN_SAMPLES)
4
5   io_images = open(train_x_dir)
6   io_labels = open(train_y_dir)
7
8   for i ∈ 1:NUMBER_TRAIN_SAMPLES
9     seek(io_images, (i-1)*28^2 + 16) # offset 16 to skip header
10    seek(io_labels, (i-1)*1 + 8) # offset 8 to skip header
11    train_x[:,i] = convert(Array{Float64}, read(io_images, 28^2))
12    train_y[i] = convert(Int, read(io_labels, UInt8))
13  end
14  close(io_images)
15  close(io_labels)
16
17  train_x = train_x
18 end

```

```

784x10000 Matrix{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
0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 ... 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0

```

```

1 begin
2   test_x = Array{Float64}(undef, 28^2, NUMBER_TEST_SAMPLES)
3   test_y = Array{Int64}(undef, NUMBER_TEST_SAMPLES)
4
5   io_images_test = open(test_x_dir)
6   io_labels_test = open(test_y_dir)
7
8   for i ∈ 1:NUMBER_TEST_SAMPLES
9     seek(io_images_test, (i-1)*28^2 + 16) # offset 16 to skip header
10    seek(io_labels_test, (i-1)*1 + 8) # offset 8 to skip header
11    test_x[:,i] = convert(Array{Float64}, read(io_images_test, 28^2))
12    test_y[i] = convert(Int, read(io_labels_test, UInt8))
13  end
14  close(io_images_test)
15  close(io_labels_test)
16
17  test_x = test_x
18 end

```

► ((784, 60000), (60000), (784, 10000), (10000))

```
1 size(train_x), size(train_y), size(test_x), size(test_y)
```

Training SVMs

```
1 md"""
2 ### Training SVMs
3 """
```

```
1 # START YOUR CODE
2
3 # END YOUR CODE
```

Evaluation

```
1 md"""
2 ### Evaluation
3 """
```

```
1 # START YOUR CODE
2
3 # END YOUR CODE
```

This is the end of Lab 05. However, there still a lot of things that you can learn about SVM. There are many open tasks to do in your spare time such as how to deal with multi-class, or Bayesian SVM. :) Hope all you will enjoy SVM. Good luck!

References

-
- [1] Boyd, S. P., & Vandenberghe, L. (2004). Convex optimization. Cambridge university press.
 - [2] Griva, I., Nash, S. G., & Sofer, A. (2008). Linear and Nonlinear Optimization 2nd Edition. Society for Industrial and Applied Mathematics.
 - [3] Schölkopf, B., & Smola, A. J. (2002). Learning with kernels: support vector machines, regularization, optimization, and beyond. MIT press.
 - [4] Lab 3, Logistic Regresion, Introduction to Machine Learning course, Department of Computer Science, Faculty of Information Technology, Ho Chi Minh University of Science, Vietnam National University.