

**Machine Learning Course**  
**Second intermediate assessment – May 31, 2018**

**Students should do all the exercises to get the maximum score.**  
**If you solve all the three exercises correctly, you get 33 points.**  
**Please, justify carefully each answer.**

**Name:** ..... **Surname:** ..... **Student ID:** .....

**Exercise 1 (13 points)**

Given the  $n=6$  two-dimensional data points  $\mathbf{x}$ , and their labels  $y$

$$\mathbf{x} = \begin{bmatrix} -2 & 0 \\ 0 & -2 \\ -1 & -1 \\ 1 & 2 \\ 1 & 1 \\ 2 & 0 \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} 1 \\ 1 \\ 1 \\ -1 \\ -1 \\ -1 \end{bmatrix},$$

find the linear discriminant function using a **batch gradient-descent algorithm** to minimize the following objective function (using the so-called *ramp loss*):

$$L(\mathbf{w}, b) = \frac{1}{n} \sum_{i=1}^n [\max(0, 1 - y_i f(\mathbf{x}_i)) - \max(0, -y_i f(\mathbf{x}_i))] + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w}, \text{ where } f(\mathbf{x}_i) = \mathbf{w}^T \mathbf{x}_i + b.$$

The gradient of  $l(yf(\mathbf{x})) = [\max(0, 1 - yf(\mathbf{x})) - \max(0, -yf(\mathbf{x}))]$  with respect to the classifier parameters is not null only when  $0 < yf(\mathbf{x}) < 1$ , and in that case it equals  $-y \mathbf{x}$  and  $-y$ , respectively when computed with respect to  $\mathbf{w}$  and  $b$ .

Initialize  $\mathbf{w} = [0.2, -1]^T$ ,  $b = 0$ , the gradient step size  $\eta = 0.5$ , the parameter  $\lambda = 0.5$ , and the threshold on the termination condition  $\theta = 0.7$ .

Use the l1 norm to compute  $|\nabla_{\mathbf{w}} L(\mathbf{w}, b)| + |\nabla_b L(\mathbf{w}, b)|$  in the termination condition.

- State the gradient-descent learning algorithm.
- Compute  $\mathbf{w}, b$  for the first two iterations of the algorithm, and check if it converges.
- Plot the initial decision boundary along with the training points, and how it changes during the first two iterations of the algorithm.
- Plot the ramp loss function  $l(z) = [\max(0, 1 - z) - \max(0, -z)]$ , with respect to  $z$ . Explain why the gradient is not null only when  $0 < z < 1$ .
- The ramp loss is more robust to outliers in the training data. Can you explain why?

**Exercise 2 (12 points)**

Given the two-dimensional data points  $\mathbf{x}$ , and the initial  $k=3$  centroids  $\mathbf{v}$ ,

$$\mathbf{x} = \begin{bmatrix} -1 & 0 \\ 0 & -1 \\ 2 & 2 \\ 1 & 0 \\ 3 & 3 \end{bmatrix}, \quad \mathbf{v} = \begin{bmatrix} 1 & 1 \\ 2 & 2 \\ 1 & -1 \end{bmatrix},$$

cluster the data points  $\mathbf{x}$  using the *k-means* algorithm, reporting the clustering labels, the updated centroids and the objective function at each iteration of the algorithm.

For simplicity, use the  $l_1$  (Manhattan) distance instead of the  $l_2$  (Euclidean) distance, both for computing the objective function  $\sum_{i=1}^n \|\mathbf{x}_i - \mathbf{v}_k\|_1$  (being  $\mathbf{v}_k$  the closest centroid to  $\mathbf{x}_i$ ) and for computing the distances between the data points  $\mathbf{x}$  and the centroids  $\mathbf{v}$ . If a point has the same distance with respect to a number of centroids, assign it to the centroid with the *lowest* class index in this set (e.g., if the point has the same distance w.r.t. centroid 0 and 2, assign it to centroid 0).

Make a two-dimensional plot displaying the data points (with a clear indication to explain to which cluster each point belongs to, after the last iteration) and the final centroids.

Plot the decision boundaries of the nearest mean centroid classifier that uses the final centroids of the *k-means* algorithm as the estimated centroids of each class.

**Exercise 3 (8 points)**

Given the two-dimensional training points  $\mathbf{x}_{\text{tr}}$ , along with their labels  $\mathbf{y}_{\text{tr}}$ , and a set of test examples  $\mathbf{x}_{\text{ts}}$ , with their labels  $\mathbf{y}_{\text{ts}}$

$$\mathbf{x}_{\text{tr}} = \begin{bmatrix} -1 & 0 \\ 0 & -1 \\ 2 & 2 \\ 1 & 0 \\ 3 & 3 \end{bmatrix}, \quad \mathbf{y}_{\text{tr}} = \begin{bmatrix} 2 \\ 2 \\ 0 \\ 1 \\ 0 \end{bmatrix}, \quad \mathbf{x}_{\text{ts}} = \begin{bmatrix} 1 & 1 \\ 2 & 2 \\ 0 & 0 \\ 1 & 2 \end{bmatrix}, \quad \mathbf{y}_{\text{ts}} = \begin{bmatrix} 0 \\ 0 \\ 2 \\ 1 \end{bmatrix},$$

classify the points in  $\mathbf{x}_{\text{ts}}$  with a *k-NN* algorithm with  $k=1$ , using the  $l_2$  distance as the distance metric. The distance matrix computed by comparing  $\mathbf{x}_{\text{ts}}$  against  $\mathbf{x}_{\text{tr}}$  is given below:

$$\begin{bmatrix} 2.24 & 2.24 & 1.41 & 1.00 & 2.83 \\ 3.61 & 3.61 & 0.00 & 2.24 & 1.41 \\ 1.00 & 1.00 & 2.83 & 1.00 & 4.24 \\ 2.83 & 3.16 & 1.00 & 2.00 & 2.24 \end{bmatrix}$$

- Compute the classification error. In case of equal (minimum) distances between a given test sample and a subset of the training points, assign the test sample to the class of the first point of the training set (from left to right in the distance matrix).
- Plot the decision function of the given *k-NN* classifier.

## EXERCISE 1 - SOLUTION

The algorithm is:

```

begin initialize  $\mathbf{w}, \theta, \eta, k=0$ 
  repeat
     $\mathbf{w} = \mathbf{w} - \eta \nabla_{\mathbf{w}} L(\mathbf{w}, b)$ 
     $b = b - \eta \nabla_b L(\mathbf{w}, b)$ 
  until  $\eta (|\nabla_{\mathbf{w}} L(\mathbf{w}, b)| + |\nabla_b L(\mathbf{w}, b)|) < \theta$ 
  
```

We need to compute the derivatives of the objective function w.r.t.  $\mathbf{w}$  and  $b$ :

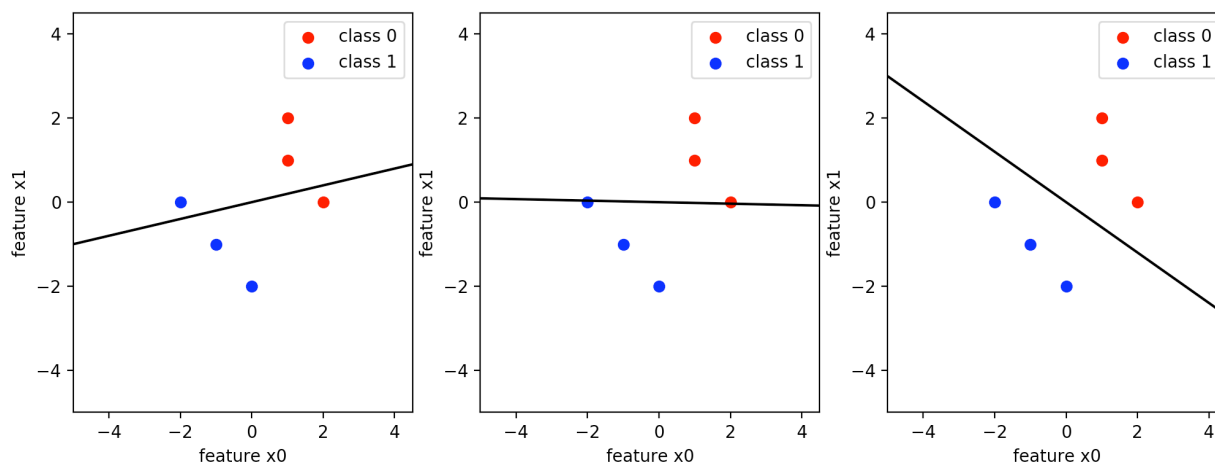
$$\nabla_{\mathbf{w}} L = \frac{1}{n} \sum_{i: 0 < y_i f(\mathbf{x}_i) < 1} -y_i \mathbf{x}_i + \lambda \mathbf{w}$$

$$\nabla_b L = \frac{1}{n} \sum_{0 < y_i f(\mathbf{x}_i) < 1} -y_i$$

iter	$L(\mathbf{w}, b)^*$	$\nabla_{\mathbf{w}} L$	$\nabla_b L$	$\mathbf{w}$	$b$	term. cond.	$\theta$
0	0.89	[ 0.43 -0.17]	0.0	[-0.02 -0.92]	0.0	0.3	0.7
1	1.25	[ 0.99 -0.12]	0.0	[-0.51 -0.85]	0.0	0.56	0.7

(\*) The objective function here is computed for  $\mathbf{w}, b$  after update

The algorithm has converged already after the first iteration.

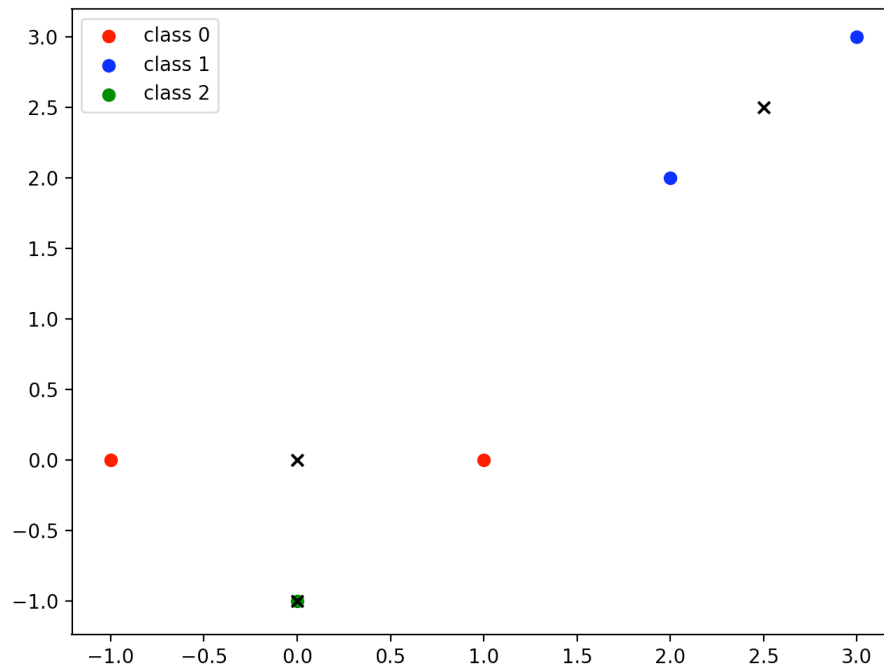


The ramp loss is flat when  $yf(\mathbf{x})$  is greater than 1 or lower than 0. Therefore, the gradient is null in those regions, while the loss is equal to  $1 - yf(\mathbf{x})$  when  $0 < yf(\mathbf{x}) < 1$ . Being equal to 1 when  $yf(\mathbf{x}) < 0$ , this loss *bounds* the influence of outliers in training data (it is a non-increasing function for decreasing  $yf(\mathbf{x})$ ), and it is thus less sensitive to misclassified points that are arbitrarily far from the rest of the training data.

## EXERCISE 2 – SOLUTION

<i>iter</i>	$\sum_{i=1}^n   \mathbf{x}_i - \mathbf{v}_k  _1$	<i>distance matrix</i>	<i>cluster assignments</i>	<i>current <math>\mathbf{v}</math></i>
0	7.0	$\begin{bmatrix} 3. & 5. & 3. \\ 3. & 5. & 1. \\ 2. & 0. & 4. \\ 1. & 3. & 1. \\ 4. & 2. & 6. \end{bmatrix}$	0 2 1 0 1	$\begin{bmatrix} 1. & 1. \\ 2. & 2. \\ 1. & -1. \end{bmatrix}$
1	4.0	$\begin{bmatrix} 1. & 6. & 2. \\ 1. & 6. & 0. \\ 4. & 1. & 5. \\ 1. & 4. & 2. \\ 6. & 1. & 7. \end{bmatrix}$	0 2 1 0 1	$\begin{bmatrix} 0. & 0. \\ 2.5 & 2.5 \\ 0. & -1. \end{bmatrix}$

After iteration 1, the cluster assignments do not change anymore. Therefore, the algorithm stops. The final clustering (along with the centroids) is shown below.



### EXERCISE 3 – SOLUTION

It is not difficult to see that the minimum distances per row are:

```
[[ 2.24  2.24  1.41  1.  2.83]
 [ 3.61  3.61  0.  2.24  1.41]
 [ 1.  1.  2.83  1.  4.24]
 [ 2.83  3.16  1.  2.  2.24]]
```

This corresponds to classify the test samples as  $y_c = [1 \ 0 \ 2 \ 0]$ .

The true labels are:  $[0 \ 0 \ 2 \ 1]$ , and thus, the classification error is 50%.

