

Optimization and Machine Learning, Spring 2020

Homework 5

(Due Tuesday, June 2 at 11:59pm (CST))

1. Show that weighted Euclidean distance in \mathbb{R}^p ,

$$d_e^{(w)}(x_i, x_{i'}) = \frac{\sum_{l=1}^p w_l (x_{il} - x_{i'l})^2}{\sum_{l=1}^p w_l},$$

satisfies

$$d_e^{(w)}(x_i, x_{i'}) = d_e(z_i, z_{i'}) = \sum_{l=1}^p (z_{il} - z_{i'l})^2,$$

where

$$z_{il} = x_{il} \left(\frac{w_l}{\sum_{l=1}^p w_l} \right)^{1/2}.$$

Thus weighted Euclidean distance based on x is equivalent to unweighted Euclidean distance based on z . (15 points)

Solution:

$$\begin{aligned} d_e^{(w)}(x_i, x_{i'}) &= \frac{\sum_{l=1}^p w_l (x_{il} - x_{i'l})^2}{\sum_{l=1}^p w_l} \\ &= \sum_{l=1}^p \frac{w_l}{\sum_{l=1}^p w_l} (x_{il} - x_{i'l})^2 \\ &= \sum_{l=1}^p \left(\frac{w_l}{\sum_{l=1}^p w_l} (x_{il} - x_{i'l}) \right)^2 \\ &= \sum_{l=1}^p \left(\frac{w_l}{\sum_{l=1}^p w_l}^{1/2} x_{il} - \frac{w_l}{\sum_{l=1}^p w_l}^{1/2} x_{i'l} \right)^2 \\ &= \sum_{l=1}^p (z_{il} - z_{i'l})^2 \\ &= d_e(z_i, z_{i'}) \end{aligned}$$

2. Consider a dataset of n observations $\mathbf{X} \in \mathbb{R}^{n \times d}$, and our goal is to project the data onto a subspace having dimensionality p , $p < d$. Prove that PCA based on projected variance maximization is equivalent to PCA based on projected error (Euclidean error) minimization. (20 points)

Solution: Suppose \mathbf{X} has been centralized. Let $\mathbf{V} \in \mathbb{R}^{d \times p}$ represent the projected matrix and $\mathbf{V}^T \mathbf{V} = \mathbf{I}$. Then we have two different optimization goals. The one based on projected variance maximization is

$$\text{maximize } \|\mathbf{XV}\|_F^2 = \text{Tr}(\mathbf{VV}^T \mathbf{X}^T \mathbf{XV}) = \text{Tr}(\mathbf{X}^T \mathbf{XV})$$

, the other one based on projected error minimization is

$$\text{minimize } \|\mathbf{X} - \mathbf{XV}\|_F^2 = \text{Tr}((\mathbf{X} - \mathbf{XV})^T (\mathbf{X} - \mathbf{XV}))$$

And we have

$$\begin{aligned} \text{Tr}((\mathbf{X} - \mathbf{XV})^T (\mathbf{X} - \mathbf{XV})) &= \text{Tr}((\mathbf{I} - \mathbf{VV}^T)^T \mathbf{X}^T \mathbf{X} (\mathbf{I} - \mathbf{VV}^T)) \\ &= \text{Tr}(\mathbf{X}^T \mathbf{X} (\mathbf{I} - \mathbf{VV}^T)) \\ &= \text{Tr}(\mathbf{X}^T \mathbf{X}) - \text{Tr}(\mathbf{X}^T \mathbf{XV}) \end{aligned}$$

Since $\text{Tr}(\mathbf{X}^T \mathbf{X})$ is a constant value, so projected variance maximization is equivalent to projected error minimization.

3. Show that the conventional linear PCA algorithm is recovered as a special case of kernel PCA if we choose the linear kernel function given by $k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T \mathbf{x}'$. (15 points)

Solution: Suppose \mathbf{X} has been centralized. The kernel function is given by $k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T \mathbf{x}'$, so we have the matrix form $\mathbf{K} = \mathbf{X}^T \mathbf{X}$. We can represent principal components $\mathbf{v} = \sum_{i=1}^n a_i \mathbf{x}_i = \mathbf{X} \alpha$. And to solve the conventional linear PCA, we have

$$\begin{aligned} & \frac{1}{n} \mathbf{X} \mathbf{X}^T \mathbf{v} = \lambda \mathbf{v} \\ \Rightarrow & \mathbf{X}^T \frac{1}{n} \mathbf{X} \mathbf{X}^T \mathbf{v} = \mathbf{X}^T \lambda \mathbf{v} \\ \Rightarrow & \frac{1}{n} \mathbf{X}^T \mathbf{X} \mathbf{X}^T \mathbf{X} \alpha = \lambda \mathbf{X}^T \mathbf{X} \alpha \\ \Rightarrow & \frac{1}{n} \mathbf{K} \mathbf{K} \alpha = \lambda \mathbf{K} \alpha \\ \Rightarrow & \frac{1}{n} \mathbf{K} \alpha = \lambda \alpha \end{aligned}$$

Then we show the conventional linear PCA is a special case of kernel PCA with the linear kernel function given by $k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T \mathbf{x}'$.

4. Let $S = \{x^{(1)}, \dots, x^{(n)}\}$ be a dataset of n samples with 2 features, i.e. $x^{(i)} \in \mathbb{R}^2$. The samples are classified into 2 categories with labels $y^{(i)} \in \{0, 1\}$. A scatter plot of the dataset is shown in Figure 1.

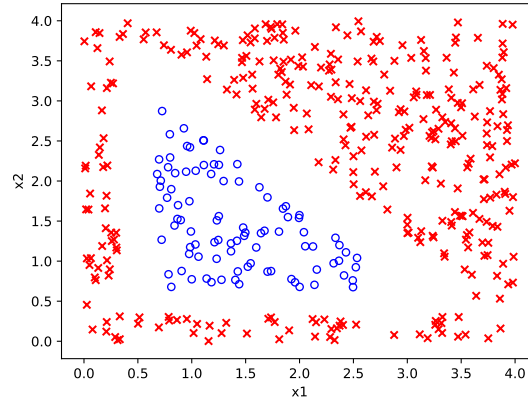


Figure 1: Plot of dataset S .

The examples in class 1 are marked as “ \times ” and examples in class 0 are marked as “ \circ ”. We want to perform binary classification using a simple neural network with the architecture shown in Figure 2.

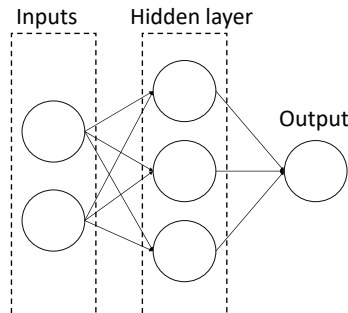


Figure 2: Architecture for our simple neural network.

Denote the two features x_1 and x_2 , the three neurons in the hidden layer h_1, h_2 , and h_3 , and the output neuron as o . Let the weight from x_i to h_j be $w_{i,j}^{[1]}$ for $i \in \{1, 2\}, j \in \{1, 2, 3\}$, and the weight from h_j to o be

$w_j^{[2]}$. Finally, denote the intercept weight for h_j as $w_{0,j}^{[1]}$, and the intercept weight for o as $w_0^{[2]}$. For the loss function, we'll use average squared loss instead of the usual negative log-likelihood:

$$l = \frac{1}{n} \sum_{i=1}^n (o^{(i)} - y^{(i)})^2,$$

where $o^{(i)}$ is the result of the output neuron for example i .

- (a) Suppose we use the sigmoid function as the activation function for h_1, h_2, h_3 and o . What is the gradient descent update to $w_{2,1}^{[1]}$, assuming we use a learning rate of α ? Your answer should be written in terms of $x^{(i)}$, $o^{(i)}$, $y^{(i)}$, and the weights. (10 points)

Solution: Let z^{h_i} for $i \in \{1, 2\}$ and z^o be the input to the sigmoid function at the hidden and output layers.

$$\begin{aligned} \frac{\partial l}{\partial w_{2,1}^{[1]}} &= \frac{\partial l}{\partial o} \frac{\partial o}{\partial z^o} \frac{\partial z^o}{\partial h_1} \frac{\partial h_1}{\partial z^{h_1}} \frac{\partial z^{h_1}}{\partial w_{2,1}^{[1]}} \\ &= \frac{1}{n} \sum_{i=1}^n 2o^{(i)}(o^{(i)} - y^{(i)})(1 - o^{(i)})w_1^{[2]}h_1(1 - h_1)x_2^{(i)}, \end{aligned}$$

where we exploit the fact that the derivative (w.s.t. z) of the sigmoid function is $\frac{\partial \sigma}{\partial z} = \sigma(z)(1 - \sigma(z))$.

- (b) Now, suppose instead of using the sigmoid function for the activation function for h_1, h_2, h_3 and o , we instead used the step function $f(x)$, defined as

$$f(x) = \begin{cases} 1, & x \geq 0, \\ 0, & x < 0. \end{cases}$$

Is it possible to have a set of weights that allow the neural network to classify this dataset with 100% accuracy? If it is possible, please provide a set of weights that enable 100% accuracy and explain your reasoning for those weights in your PDF. If it is not possible, please explain your reasoning in your PDF. (10 points) (Hint: There are three sides to a triangle, and there are three neurons in the hidden layer.)

Solution: Based on the hinted triangle, we construct the following matrix multiplication as

$$\begin{bmatrix} -1 & 5 & 0 \\ -1 & 0 & 5 \\ 4.2 & -1 & -1 \end{bmatrix} \begin{bmatrix} 1 \\ x_1 \\ x_2 \end{bmatrix}. \quad (1)$$

When the data point is within the triangle, $f(z) = [1, 1, 1]^T$. Otherwise, it's one of the seven binary vectors (e.g., $[0, 1, 0]^T$).

Then, we construct the second matrix multiplication as

$$\begin{bmatrix} -1 & -1 & -1 & -2.5 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}. \quad (2)$$

Note that, the 1 in the first row of the all-one vector is the bias. Therefore, only when a data point is inside the triangle, it produces the result less than 0, and hence categorized as class 0. Otherwise, it will be classified as class 1. Hence, the accuracy is 100%.

- (c) Let the activation functions for h_1, h_2, h_3 be the linear function $f(x) = x$ and the activation function for o be the same step function as before. Is it possible to have a set of weights that allow the neural network to classify this dataset with 100% accuracy? If it is possible, please provide a set of weights that enable 100% accuracy and explain your reasoning for those weights in your PDF. If it is not possible, please explain your reasoning in your PDF. (10 points)

Solution: No, it is impossible. Because the current classes cannot be linearly separated. If the activation is linear, then it's an affine transformation from 2D to 3D spaces, and the classification problem will still be non-linearly in the 3D space, which cannot be solved by a step function.

5. Convolutional neural networks targets the processing of 2-D features instead of the 1-D ones in multi-layer perceptron (MLP), the structure of which is depicted in Fig. 3.

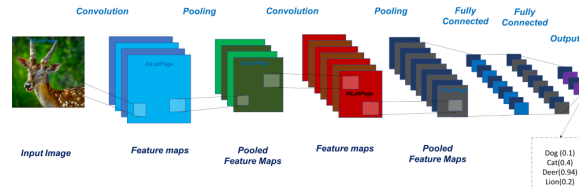


Figure 3: <https://mc.ai/how-does-convolutional-neural-network-work/>

- (a) Kernel convolution is process where we take a kernel (or filter), we pass it over our image and transform it based on the values from filter. Now, you are given the following formula

$$G(m, n) = (f * h)(m, n) = \sum_j \sum_k h(j, k) f(m - j, n - k),$$

where the input image is denoted by f and the kernel by h . The indexes of rows and columns of the result matrix are marked with m and n respectively. Please calculate the feature maps, if you are given the following 3×3 image matrix and 2×2 kernel matrix. (5 points)

1	2	3
4	5	6
7	8	9

Table 1: 3×3 -Image Matrix.

1	0
0	1

Table 2: 2×2 -Kernel Matrix.

Solution: By the formula, we have the feature maps, as shown in Tab. 3.

6	8
12	14

Table 3: 2×2 -feature maps.

- (b) Assume the input with the size of (width = 28, height = 28) and a filter with the size of (width = 5, height = 5) and the convolutional layer parameters are $S = 1$ (the stride), $P = 0$ (the amount of zero padding). What is the exact size of the convolution output? (5 points)

Solution: The Conv layer:

- Accepts the input with the size 28×28
- Requires the hyperparameters:
 - * A filter with the size 5×5 ,
 - * the stride $S = 1$,
 - * the amount of zero padding $P = 0$.
- Produces a output of size $W_{\text{out}} \times H_{\text{out}}$ where
 - * $W_{\text{out}} = (28 - 5 + 0)/1 + 1 = 24$
 - * $H_{\text{out}} = (28 - 5 + 0)/1 + 1 = 24$.

6. Consider the following grid environment. Starting from any unshaded square, you can move up, down, left, or right. Actions are deterministic and always succeed (e.g. going left from state 1 goes to state 0) unless they will cause the agent to run into a wall. The thicker edges indicate walls, and attempting to move in the direction of a wall results in staying in the same square. Taking any action from the green target square (no. 5) earns a reward of +5 and ends the episode. Taking any action from the red square of death (no. 11) earns a reward of -5 and ends the episode. Otherwise, each move is associated with some reward $r \in \{-1, 0, +1\}$. Assume the discount factor $\gamma = 1$ unless otherwise specified.

0	1	2	3	4
5	6	7	8	9
10	11	12	13	14
15	16	17	18	19
20	21	22	23	24

- (a) Define the reward r for all states (except state 5 and state 11 whose rewards are specified above) that would cause the optimal policy to return the shortest path to the green target square (no. 5). (3 points)

Solution: Let all rewards be -1 . This is because it penalizes those schemes moving many times to get to the green target square. In other words, the fewer times you move, the less negative rewards you receive.

- (b) Using r from part (a), find the optimal value function for each square. (5 points)

Solution: With converse thinking, one can start from the green square to obtain the optimal value functions for other squares.

-4	-3	-2	-1	0
5	4	3	2	1
4	-5	2	1	0
-5	-4	-3	-2	-1
-6	-5	-4	-3	-2

- (c) Does setting $\gamma = 0.8$ change the optimal policy? Why or why not? (2 points)

Solution: No. Changing γ changes the value function but not the relative order.