

# Introduction to Machine Learning, Fall 2023

## Homework 4

(Due Tuesday Dec.19 at 11:59pm (CST))

1. [15 points] [Maximum Margin Classifier] Consider a data set of  $nd$ -dimensional sample points,  $\{X_1, \dots, X_n\}$ . Each sample point,  $X_i \in \mathbb{R}^d$ , has a corresponding label,  $y_i$ , indicating to which class that point belongs. For now, we will assume that there are only two classes and that every point is either in the given class ( $y_i = 1$ ) or not in the class ( $y_i = -1$ ). Consider the linear decision boundary defined by the hyperplane

$$\mathcal{H} = \{x \in \mathbb{R}^d : x \cdot w + \alpha = 0\}.$$

The maximum margin classifier maximizes the distance from the linear decision boundary to the closest training point on either side of the boundary, while correctly classifying all training points.

- (a) An in-class sample point is correctly classified if it is on the positive side of the decision boundary, and an out-of-class sample is correctly classified if it is on the negative side. Write a set of  $n$  constraints to ensure that all  $n$  points are correctly classified. [3 points]
- (b) The maximum margin classifier aims to maximize the distance from the training points to the decision boundary. Derive the distance from a point  $X_i$  to the hyperplane  $\mathcal{H}$ . [3 points]
- (c) Assuming all the points are correctly classified, write an inequality that relates the distance of sample point  $X_i$  to the hyperplane  $\mathcal{H}$  in terms of only the normal vector  $w$ . [3 points]
- (d) For the maximum margin classifier, the training points closest to the decision boundary on either side of the boundary are referred to as support vectors. What is the distance from any support vector to the decision boundary? [3 points]
- (e) Using the previous parts, write an optimization problem for the maximum margin classifier. [3 points]

**Solution:**

$$(a). \quad y_i \cdot (x_i \cdot w + \alpha) > 0. \quad i=1 \dots n.$$

$$(b) \quad r_i = \frac{y_i (x_i \cdot w + \alpha)}{\|w\|} \quad i=1 \dots n. \quad (X_i \text{ to the } H: \frac{|Ax+By+C|}{\sqrt{A^2+B^2}})$$

(c). Assuming that all the points are correctly classified, we can suppose that  $|x_i \cdot w + \alpha| \geq 1$   $i=1 \dots n$  with normalized  $w$  and  $\alpha$ . So that margin is 1, then  $|x_i \cdot w + \alpha| \geq 1, i=1 \dots n$

$$\therefore d_i = \frac{y_i (x_i \cdot w + \alpha)}{\|w\|} \geq \frac{1}{\|w\|}$$

d) The distance from any support vector to the decision boundary is the minimum distance of all the points to the boundary. So that.

$$d_{\text{support}} = \min d_i = \min \frac{y_i(x_i \cdot w + \alpha)}{\|w\|} = \frac{1}{\|w\|}$$

(If margin is 1 with normalized  $w$  and  $\alpha$ .)

e)  $\max \frac{1}{\|w\|}$

$$\text{s.t. } y_i(x_i \cdot w + \alpha) \geq 1 \quad \forall i = 1 \dots n.$$

2. [15 points] 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.

Solution:

Suppose  $\mathbf{X}$  to be centralized. The  $\mathbf{V} \in \mathbb{R}^{d \times p}$  is the projected matrix and  $\mathbf{V}^T \mathbf{V} = \mathbf{I}$ .

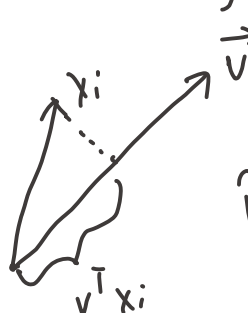
PCA based on projected variance maximization is.

$$\max_{\mathbf{V}} \frac{1}{n} \sum_{i=1}^n (\mathbf{V}^T \mathbf{x}_i)^2.$$

PCA based on projected error (Euclidean error) minimization is

$$\min_{\mathbf{V}} \frac{1}{n} \sum_{i=1}^n \|\mathbf{x}_i - (\mathbf{V}^T \mathbf{x}_i) \mathbf{V}\|^2.$$

According to the 勾股定理 (Pythagorean Theorem).



$$\|\mathbf{x}_i\|^2 = \|(\mathbf{V}^T \mathbf{x}_i) \mathbf{V}\|^2 + \|\mathbf{x}_i - (\mathbf{V}^T \mathbf{x}_i) \mathbf{V}\|^2 \quad i=1 \dots n.$$

Proof:  $\|\mathbf{x}_i - (\mathbf{V}^T \mathbf{x}_i) \mathbf{V}\|^2 = \|\mathbf{x}_i\|^2 + \|(\mathbf{V}^T \mathbf{x}_i) \mathbf{V}\|^2 - 2 \mathbf{V}^T \mathbf{x}_i \mathbf{V}^T \mathbf{x}_i$

$$= \|\mathbf{x}_i\|^2 + \|(\mathbf{V}^T \mathbf{x}_i) \mathbf{V}\|^2 - \mathbf{x}_i^T \mathbf{V}^T \mathbf{x}_i \mathbf{V} - \mathbf{V}^T \mathbf{x}_i \mathbf{V}^T \mathbf{x}_i$$

$$= \|\mathbf{x}_i\|^2 + \|(\mathbf{V}^T \mathbf{x}_i) \mathbf{V}\|^2 - 2(\mathbf{V}^T \mathbf{x}_i)^2$$

$$= \|\mathbf{x}_i\|^2 + \|\mathbf{V}^T \mathbf{x}_i\|^2 - 2\|\mathbf{V}^T \mathbf{x}_i\|^2.$$

$$= \|\mathbf{x}_i\|^2 - \|\mathbf{V}^T \mathbf{x}_i\|^2.$$

$$\therefore \|\mathbf{x}_i\|^2 = \|\mathbf{V}^T \mathbf{x}_i\|^2 \|\mathbf{V}\|^2 + \|\mathbf{x}_i\|^2 - \|\mathbf{V}^T \mathbf{x}_i\|^2 = \|\mathbf{x}_i\|^2 \neq$$

$$\therefore \max_{\mathbf{V}} \frac{1}{n} \sum_{i=1}^n (\mathbf{V}^T \mathbf{x}_i)^2 = \max_{\mathbf{V}} \frac{1}{n} \sum_{i=1}^n [\|\mathbf{x}_i\|^2 - \|\mathbf{x}_i - (\mathbf{V}^T \mathbf{x}_i) \mathbf{V}\|^2]$$

As  $\frac{1}{n} \sum_{i=1}^n \|\mathbf{x}_i\|^2$  is a constant.

$$\Leftrightarrow \min_{\mathbf{V}} \frac{1}{n} \sum_{i=1}^n \|\mathbf{x}_i - (\mathbf{V}^T \mathbf{x}_i) \mathbf{V}\|^2$$

Above all, PCA based on projected variance maximization is equal to PCA based on projected error (Euclidean error) minimization.

3. [15 points] [Performing PCA by Hand] Let's do principal components analysis (PCA)! Consider this sample of six points  $X_i \in \mathbb{R}^2$ .

$$\left\{ \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \end{bmatrix}, \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 \\ 2 \end{bmatrix}, \begin{bmatrix} 2 \\ 1 \end{bmatrix}, \begin{bmatrix} 2 \\ 2 \end{bmatrix} \right\}.$$

(a) [4 pts] Compute the mean of the sample points and write the centered design matrix  $\dot{X}$ .

Hint: The sample mean is

Hint: By subtracting the mean from each sample, we form the centered design matrix

$$\dot{X} =$$

(b) [5 pts] Find all the principal components of this sample. Write them as unit vectors.

Hint: The principal components of our dataset are the eigenvectors of the matrix

$$\dot{X}^\top \dot{X} =$$

The characteristic polynomial of this symmetric matrix is

$$\det(sI - \dot{X}^\top \dot{X})$$

(c) [6 pts]

Which of those two principal components would be preferred if you use only one? [2 pts]

What information does the PCA algorithm use to decide that one principal components is better than another? [2 pts]

From an optimization point of view, why do we prefer that one? [2 pts]

**Solution:**

a) Hint: The sample mean is  $\begin{bmatrix} 1 \\ 1 \end{bmatrix}$ .

Hint: By subtracting the mean from each sample, we form the centered design matrix

$$\dot{X} = \begin{bmatrix} -1 & -1 \\ -1 & 0 \\ 0 & -1 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \end{bmatrix}$$

b)

$$\dot{X}^\top \dot{X} = \begin{bmatrix} -1 & -1 & 0 & 0 & 1 & 1 \\ -1 & 0 & -1 & 1 & 0 & 1 \end{bmatrix} \begin{bmatrix} -1 & -1 \\ -1 & 0 \\ 0 & -1 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \end{bmatrix}$$

$$= \begin{bmatrix} 4 & 2 \\ 2 & 4 \end{bmatrix}$$

$$\det (sI - \dot{X}^T \dot{X}) = \begin{vmatrix} s-4 & -2 \\ -2 & s-4 \end{vmatrix} = (s-4)^2 - 4 = (s-2)(s-6).$$

$$s_1 = 2, \quad s_2 = 6.$$

$$1^\circ \quad s_1 = 2.$$

$$2^\circ \quad s_2 = 6$$

$$\begin{bmatrix} 4-2 & 2 \\ 2 & 4-2 \end{bmatrix} = \begin{bmatrix} 2 & 2 \\ 2 & 2 \end{bmatrix} \quad \begin{bmatrix} -2 & 2 \\ 2 & -2 \end{bmatrix} v = 0.$$

$$\begin{bmatrix} 2 & 2 \\ 2 & 2 \end{bmatrix} v = 0$$

$$v_2 = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

$$\exists -1c: \quad v_2 = \begin{bmatrix} \frac{\sqrt{2}}{2} \\ \frac{\sqrt{2}}{2} \end{bmatrix}$$

$$v_1 = \begin{bmatrix} 1 \\ -1 \end{bmatrix}$$

$$\exists -1c: \quad v_1 = \begin{bmatrix} \frac{\sqrt{2}}{2} \\ -\frac{\sqrt{2}}{2} \end{bmatrix}$$

$$c). \quad 1^\circ \text{ use } v_2 = \begin{bmatrix} \frac{\sqrt{2}}{2} \\ \frac{\sqrt{2}}{2} \end{bmatrix}$$

2<sup>o</sup> If  $\lambda$  is bigger, use  $\lambda$ 's principal components.

3<sup>o</sup>. From an optimization point of view, we prefer of which variance is bigger, so that we can keep more original information when reducing the same dimensions.

$$\text{that is } \max_{1 \leq i \leq n} \frac{1}{n} (v^T x_i)^2 = v^T X X^T v$$

$$\Rightarrow \max_v v^T X X^T v \quad \text{s.t.} \quad v^T v = 1.$$

$$\Rightarrow \max_v v^T X X^T v - (\lambda v^T v - 1) \lambda \geq 0$$

$$\Rightarrow \frac{\partial (v^T X X^T v - \lambda v^T v)}{\partial v} = 0. \quad (X X^T - \lambda I) v = 0$$

$$\Rightarrow X X^T v = \lambda v.$$

$$v^T X X^T v = v^T \lambda v = \lambda v^T v = \lambda$$

Now we only need to maximize  $\lambda$ , of which here is  $s_2$ .

4. [15 points] [Backpropagation on an Arithmetic Expression] Consider an arithmetic network with the inputs  $a, b$ , and  $c$ , which computes the following sequence of operations, where  $s(\gamma) = \frac{1}{1+e^{-\gamma}}$  is the logistic (sigmoid) function and  $r(\gamma) = \max\{0, \gamma\}$  is the hinge function used by ReLUs.

$$d = ab \quad e = s(d) \quad f = r(a) \quad g = 3a \quad h = 2e + f + g \quad i = ch \quad j = f + i^2$$

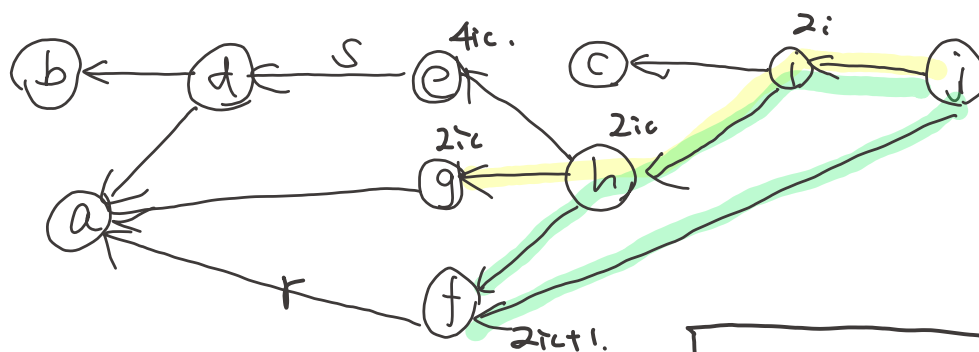
We want to find the partial derivatives of  $j$  with respect to every other variable  $a$  through  $i$ , in backpropagation style. This means that for each variable  $z$ , we want you to write  $\partial j / \partial z$  in two forms: (1) in terms of derivatives involving each variable that directly uses the value of  $z$ , and (2) in terms of the inputs and intermediate values  $a \dots i$ , as simply as possible but with no derivative symbols. For example, we write

$$\frac{\partial j}{\partial i} = \frac{dj}{di} = 2i \quad (\text{no chain rule needed for this one only})$$

$$\frac{\partial j}{\partial h} = \frac{\partial j}{\partial i} \frac{\partial i}{\partial h} = 2ic \quad (\text{chain rule, then backprop the derivative expressions})$$

(a) Now, please write expressions for  $\partial j / \partial g$ ,  $\partial j / \partial f$ ,  $\partial j / \partial e$ ,  $\partial j / \partial d$ ,  $\partial j / \partial c$ ,  $\partial j / \partial b$ , and  $\partial j / \partial a$  as we have written  $\partial j / \partial h$  above. If they are needed, express the derivative  $s'(\gamma)$  in terms of  $s(\gamma)$  and express the derivative  $r'(\gamma)$  as the indicator function  $1(\gamma \geq 0)$ . (Hint:  $f$  is used in two places and  $a$  is used in three, so they will need a multivariate chain rule. It might help you to draw the network as a directed graph, but it's not required.)

**Solution:**



a).  $\frac{\partial j}{\partial g} = \frac{\partial j}{\partial h} \frac{\partial h}{\partial g} = 2ic \cdot 1 = 2ic.$

$$\frac{\partial j}{\partial f} = \frac{\partial j}{\partial i} \frac{\partial i}{\partial h} \frac{\partial h}{\partial f} + \frac{\partial j}{\partial f} = 2ic + 1$$

$$\frac{\partial j}{\partial e} = \frac{\partial j}{\partial i} \frac{\partial i}{\partial h} \frac{\partial h}{\partial e} = 2ic \cdot 2 = 4ic$$

$$\frac{\partial j}{\partial d} = \frac{\partial j}{\partial e} \cdot \frac{\partial e}{\partial d} = 4ic \cdot s'(d) = 4ic \cdot s(d) [1 - s(d)].$$

$$\frac{\partial j}{\partial c} = \frac{\partial j}{\partial i} \cdot \frac{\partial i}{\partial c} = 2ih.$$

$$\frac{\partial j}{\partial b} = \frac{\partial j}{\partial d} \cdot \frac{\partial d}{\partial b} = 4ic \cdot s'(d) \cdot a = 4ica s(d) [1 - s(d)]$$

$$\frac{\partial j}{\partial a} = \frac{\partial j}{\partial g} \frac{\partial g}{\partial a} + \frac{\partial j}{\partial f} \frac{\partial f}{\partial a} + \frac{\partial j}{\partial d} \frac{\partial d}{\partial a} = 6ic + (2ic+1)r'(a)$$

$$\begin{aligned} s'(r) &= \left( \frac{1}{1+e^{-x}} \right)' \\ &= \frac{-e^{-x}}{-(1+e^{-x})^2} \\ &= \frac{e^{-x}}{(1+e^{-x})^2} \\ &= s(r) [1 - s(r)]. \end{aligned}$$

$$+ 4i_c s(d)[1-s(d)]b \quad \left| \quad r'(r)=1 \quad (r \geq 0) \right.$$

$$= b_i^c + 4i_c s(d)[1-s(d)]b + (2i_c + 1) \mathbb{1}_{(a \geq 0)}$$