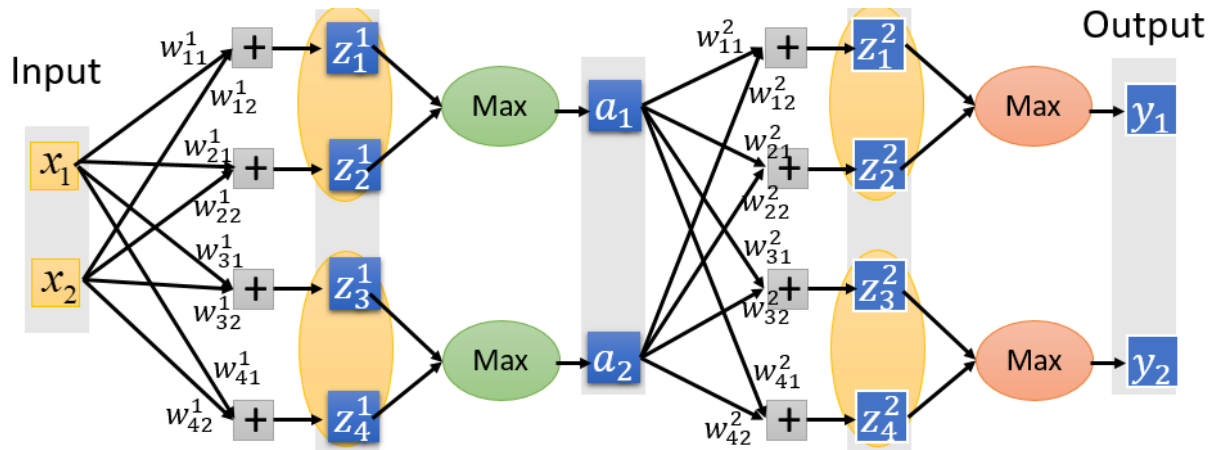




**4. (20 pts) Maxout Network**

Consider the following maxout network

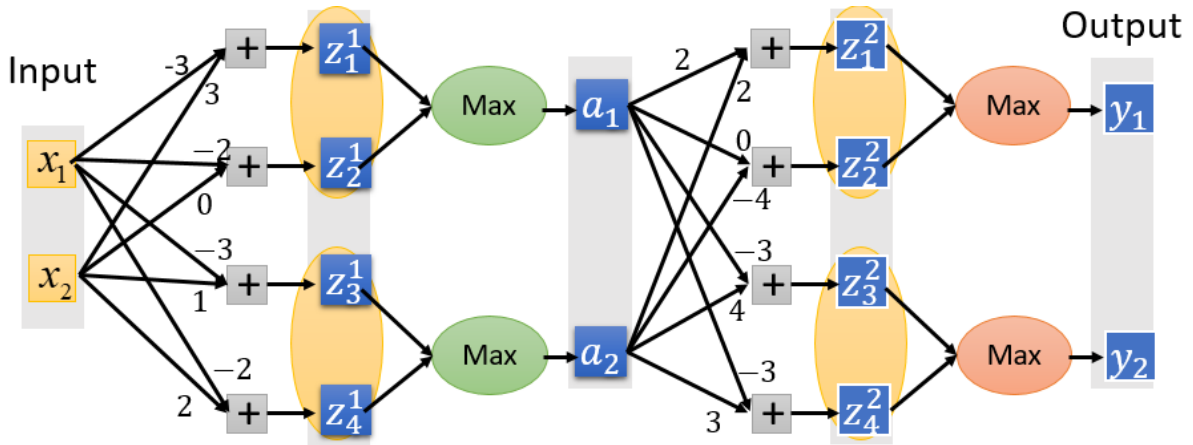


The above neural network can be represented as a function  $f_\theta$ , namely

$$\begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = f_\theta \left( \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \right),$$

where parameter  $\theta$  records all the weights  $w_{ij}^k$ .

(1) (6 pts) Suppose the weights are initialized as follows:



If  $(x_1, x_2) = (1, 10)$ , please fill out the following table. No derivation required.

Variable	$z_1^1$	$z_2^1$	$z_3^1$	$z_4^1$	$a_1$	$a_2$	$z_1^2$	$z_2^2$	$z_3^2$	$z_4^2$	$y_1$	$y_2$
Value												

(2) (6 pts) Continuing (1), if the ground truth is  $(\hat{y}_1, \hat{y}_2) = (-8, -1)$ , and the L1-loss is adopted, namely

$$L(\theta) = \left\| \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \end{bmatrix} - f_\theta \left( \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \right) \right\|_1 = |y_1 - \hat{y}_1| + |y_2 - \hat{y}_2|$$

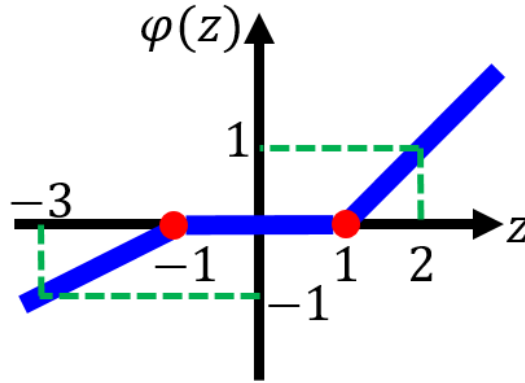
Please compute  $\frac{\partial L}{\partial w_{12}^1}$  and  $\frac{\partial L}{\partial w_{41}^2}$ .

(3) (5 pts) Let  $g_1(x), g_2(x), \dots, g_k(x)$  be convex functions, prove that

$$h(x) = \max(g_1(x), g_2(x), \dots, g_k(x))$$

is also a convex function.

- (4) **(3 pts)** Is it possible to implement the following piecewise-linear activation function  $\varphi$  with maxout network? If possible, please show how to implement; If impossible, please prove why it is impossible.



**5. (20 pts) Gradient Boosting**

Consider the binary classification problem, where we are given training data set  $\{(\mathbf{x}_i, y_i)\}_{i=1}^N$  with  $\mathbf{x}_i \in \mathbb{R}^d$  and  $y_i \in \{\pm 1\}$ . Let  $\mathcal{F}$  be a collection of binary classifiers, each mapping from  $\mathbb{R}^d$  to  $\{\pm 1\}$ . Given number of epochs  $T \in \mathbb{N}$ , suppose we want to find the function

$$g(\mathbf{x}) = \sum_{t=1}^T \alpha_t f_t(\mathbf{x})$$

Where  $f_t \in \mathcal{F}$  and  $\alpha_t \in \mathbb{R}$  for all  $t = 1, \dots, T$ , by which the aggregated classifier is given by

$$h(\mathbf{x}) = \begin{cases} 1 & \text{if } g(\mathbf{x}) > 0 \\ -1 & \text{if } g(\mathbf{x}) \leq 0 \end{cases}$$

Please apply gradient boosting to show how the functions  $f_t$  and coefficients  $\alpha_t$  are computed with an aim to minimize the following loss function

$$L(g) = \sum_{i=1}^N \log(1 + e^{-y_i g(\mathbf{x}_i)})$$

**6. (20 pts) Expectation Maximization Interpretation behind Semi-Supervised Learning**

Given  $N$  samples  $\mathbf{x}_1, \dots, \mathbf{x}_N \in \mathbb{R}^m$  as well as their labels  $y_1, \dots, y_N \in \{0, 1, \dots, K\}$ . Consider the generative model where each sample  $\mathbf{x}_i$  is generated independently according to Gaussian mixture model that depends on the label  $y_i$ , as represented by random variable

$$X_i \sim \begin{cases} \sum_{j=1}^K \pi_j \mathcal{N}(\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j) & \text{if } y_i = 0 \\ \mathcal{N}(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) & \text{if } y_i = k \neq 0 \end{cases}$$

where  $\pi_1 + \dots + \pi_K = 1$ , and  $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  denotes the Gaussian distribution with mean  $\boldsymbol{\mu}$  and covariance matrix  $\boldsymbol{\Sigma}$ , with probability density function

$$\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{\sqrt{(2\pi)^m |\boldsymbol{\Sigma}|}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right)$$

We would like to apply Expectation Maximization algorithm to find the maximum likelihood

estimation of parameters  $\theta = \{(\pi_k, \mu_k, \Sigma_k)\}_{k=1}^K$ .

(1) (16 pts) Please write down the E-step and M-step and show that the parameters are updated from

$\theta^{(t)} = \{(\pi_k^{(t)}, \mu_k^{(t)}, \Sigma_k^{(t)})\}_{k=1}^K$  to  $\theta^{(t+1)} = \{(\pi_k^{(t+1)}, \mu_k^{(t+1)}, \Sigma_k^{(t+1)})\}_{k=1}^K$  in the following form:

$$\pi_k^{(t+1)} = \frac{N_k + \sum_{i:y_i=0} \delta_{ik}^{(t)}}{N}$$

$$\mu_k^{(t+1)} = \frac{\sum_{i:y_i=k} \mathbf{x}_i + \sum_{i:y_i=0} \delta_{ik}^{(t)} \mathbf{x}_i}{N_k + \sum_{i:y_i=0} \delta_{ik}^{(t)}}$$

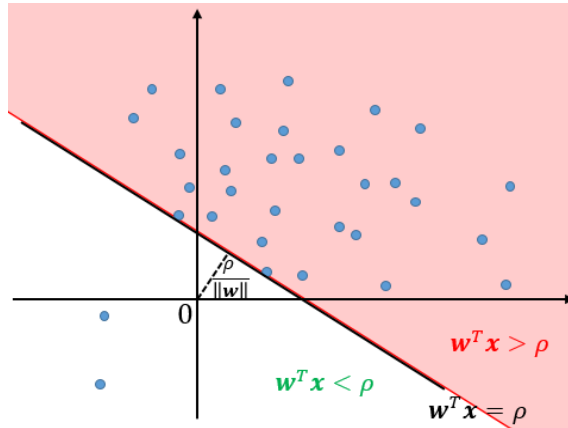
$$\Sigma_k^{(t+1)} = \frac{\sum_{i:y_i=k} (\mathbf{x}_i - \mu_k^{(t+1)}) (\mathbf{x}_i - \mu_k^{(t+1)})^T + \sum_{i:y_i=0} \delta_{ik}^{(t)} (\mathbf{x}_i - \mu_k^{(t+1)}) (\mathbf{x}_i - \mu_k^{(t+1)})^T}{N_k + \sum_{i:y_i=0} \delta_{ik}^{(t)}}$$

where  $N_k = \sum_{i:y_i=k} 1$  is the number of samples in class k. Please show your derivations.

(2) (4 pts) What is the closed form expression of  $\delta_{ik}^{(t)}$ ? Please show your derivations.

## 7. (25 pts) One-Class Support Vector Machine

Given unlabeled training data  $\mathbf{x}_1, \dots, \mathbf{x}_N \in \mathbb{R}^d$ , the goal of *one-class SVM* is to find the “smallest” half space  $S = \{\mathbf{x} \in \mathbb{R}^d: \mathbf{w}^T \mathbf{x} \geq \rho\}$  that contains the most data points, as illustrated below.



Formally speaking, with a given hyper-parameter  $0 < \nu < 1$ , the *one-class SVM* aims to solve the following optimization problem

$$\begin{aligned} & \text{minimize} && \frac{1}{2} \|\mathbf{w}\|^2 + \frac{1}{\nu N} \sum_{i=1}^N \xi_i - \rho \\ & \text{subject to} && \mathbf{w}^T \mathbf{x}_i \geq \rho - \xi_i, \quad \forall i = 1, \dots, N \\ & && \xi_i \geq 0 \\ & \text{variables} && \mathbf{w} \in \mathbb{R}^d, \rho \in \mathbb{R}, \xi_1, \dots, \xi_N \in \mathbb{R} \end{aligned}$$

Which can be rewritten in the form of **primal problem**

$$\begin{aligned} & \text{minimize} && f(\mathbf{w}, \rho, \xi) = \frac{1}{2} \|\mathbf{w}\|^2 + \frac{1}{\nu N} \sum_{i=1}^N \xi_i - \rho \\ & \text{subject to} && g_{1,i}(\mathbf{w}, \rho, \xi) = \rho - \xi_i - \mathbf{w}^T \mathbf{x}_i \leq 0, \quad \forall i = 1, \dots, N \\ & && g_{2,i}(\mathbf{w}, \rho, \xi) = -\xi_i \leq 0 \\ & \text{variables} && \mathbf{w} \in \mathbb{R}^d, \rho \in \mathbb{R}, \xi_1, \dots, \xi_N \in \mathbb{R} \end{aligned}$$

(1) **(7 pts)** Show that the **dual problem** of *one-class SVM* can be written as

$$\begin{aligned}
 &\text{maximize} && \theta(\alpha) = -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j \mathbf{x}_i^T \mathbf{x}_j \\
 &\text{subject to} && \sum_{i=1}^N \alpha_i = 1 \\
 &\text{variables} && 0 \leq \alpha_i \leq \frac{1}{vN}, \quad i = 1, \dots, N
 \end{aligned}$$

Let  $(\bar{\mathbf{w}}, \bar{b}, \bar{\xi})$  be a primal optimal solution,  $\bar{\alpha}$  be a dual optimal solution.

(2) **(8 pts)** Which of the following statements are true? Please justify your answers.

- i.  $f(\bar{\mathbf{w}}, \bar{\rho}, \bar{\xi}) = \theta(\bar{\alpha})$ .
- ii. If  $\bar{\mathbf{w}}^T \mathbf{x}_i < \bar{\rho}$ , then  $\bar{\alpha}_i = 1/(vN)$ .
- iii. If  $\bar{\mathbf{w}}^T \mathbf{x}_i > \bar{\rho}$ , then  $\bar{\alpha}_i = 0$ .
- iv.  $\bar{\xi}_i = \max(\bar{\rho} - \bar{\mathbf{w}}^T \mathbf{x}_i, 0)$ .

(3) **(5 pts)** Show that the optimal bias  $\bar{\rho}$  can be expressed in terms of the optimal weighting vector  $\bar{\mathbf{w}}$  as

$$\bar{\rho} = \operatorname{argmin}_{\rho \in \mathbb{R}} \left( \frac{1}{vN} \sum_{i=1}^N \max(\rho - \bar{\mathbf{w}}^T \mathbf{x}_i, 0) - \rho \right)$$

(Hint: Rewrite the primal problem as an unconstrained optimization problem over variables  $\mathbf{w}, \rho$ )

(4) **(5 pts)** Suppose  $k - 1 < vN < k$ , for which  $k$  is a positive integer, then how many among the  $N$  data points  $\mathbf{x}_1, \dots, \mathbf{x}_N$  satisfy  $\bar{\mathbf{w}}^T \mathbf{x} \geq \bar{\rho}$ ? In other words, please compute  $\sum_i: \bar{\mathbf{w}}^T \mathbf{x}_i \geq \bar{\rho} 1$ . Please show your derivations.

(Hint: Solve the minimization problem in (3) by taking upper-derivatives)