

# CS4487: Home Assignment №3

Department of Computer Science  
City University of Hong Kong

**Due on** October 30, 2019, 7pm

## Exercise 1

[5 points]. In Section 5.4 of the lecture note, we have given a detailed derivation of the dual form of SVM with soft margin. With simpler arguments, derive the dual form of SVM with hard margin

$$\begin{aligned} \min_{w,b} \quad & \frac{1}{2} w^T w \\ \text{s.t.} \quad & y^{(i)}(w^T x^{(i)} + b) \geq 1, \quad i = 1, \dots, m. \end{aligned}$$

## Exercise 2

[5 points]. Considering the two-way partitioning problem

$$\begin{aligned} \min_x \quad & \frac{1}{2} x^T W x \\ \text{s.t.} \quad & x_i^2 = 1, \quad i = 1, \dots, n, \end{aligned}$$

where  $x \in \mathbb{R}^n$  is the parameter vector to be optimized and  $W \in \mathbb{R}^{n \times n}$  is a fixed weight matrix. One interpretation of this problem is to partition  $\{1, 2, \dots, n\}$  in two sets;  $W_{ij}$  is cost of assigning  $i, j$  to the same set;  $-W_{ij}$  is cost of assigning to different sets.

- [1 point]. Is it a convex optimization problem? Give your reasons.
- [3 points]. Derive the Lagrange dual function of the above problem. Hint: recall the definition of a positive definite matrix<sup>1</sup>.

---

<sup>1</sup>[https://en.wikipedia.org/wiki/Definiteness\\_of\\_a\\_matrix](https://en.wikipedia.org/wiki/Definiteness_of_a_matrix)

## Exercise 1

[5 points]. In Section 5.4 of the lecture note, we have given a detailed derivation of the dual form of SVM with soft margin. With simpler arguments, derive the dual form of SVM with hard margin.

$$\begin{aligned} \min_{w,b} \quad & \frac{1}{2} w^T w \\ \text{s.t.} \quad & y^{(i)}(w^T x^{(i)} + b) \geq 1, \quad i = 1, \dots, m. \end{aligned}$$

Lagrangian:

$$\begin{aligned} L(\vec{w}, \vec{b}, \vec{\alpha}) &= \frac{1}{2} \|\vec{w}\|^2 - \sum_{i=1}^m \alpha_i [y^{(i)} (\vec{w}^T \vec{x}^{(i)}) + b] - 1] \\ &= \frac{1}{2} \vec{w}^T \vec{w} - \sum_{i=1}^m \alpha_i y^{(i)} \vec{w}^T \vec{x}^{(i)} - \sum_{i=1}^m \alpha_i y^{(i)} b \\ &\quad + \sum_{i=1}^m \alpha_i \end{aligned}$$

$$\nabla_{\vec{w}} L = \vec{w} - \sum_{i=1}^m \alpha_i y^{(i)} \vec{x}^{(i)} = 0$$

$$\nabla_b L = - \sum_{i=1}^m \alpha_i y^{(i)} = 0$$

Lagrange dual function:

$$g(\alpha) = \inf_{\vec{w}, \vec{b}} (L(\vec{w}, \vec{b}, \vec{\alpha}))$$

$$\begin{aligned} &= \frac{1}{2} \left( \sum_{i=1}^m \alpha_i y^{(i)} \vec{x}^{(i)} \right)^T \left( \sum_{i=1}^m \alpha_i y^{(i)} \vec{x}^{(i)} \right) \\ &\quad - \sum_{i=1}^m \alpha_i y^{(i)} \left( \sum_{j=1}^m \alpha_j y^{(j)} \vec{x}^{(j)} \right)^T \vec{x}^{(i)} + \sum_{i=1}^m \alpha_i \\ &= \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m y^{(i)} y^{(j)} \alpha_i \alpha_j (\vec{x}^{(i)})^T \vec{x}^{(j)}. \end{aligned}$$

Dual form:

$$\max_{\vec{\alpha}} \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m y^{(i)} y^{(j)} \alpha_i \alpha_j (\vec{x}^{(i)})^T (\vec{x}^{(j)})$$

$$\text{subject to } \nabla_{\vec{\alpha}} L = - \sum_{i=1}^m \alpha_i y^{(i)} = 0.$$

$$\alpha_i \geq 0$$

## Exercise 2

[5 points]. Considering the two-way partitioning problem

$$\min_{\vec{x}} \frac{1}{2} \vec{x}^T W \vec{x} = \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n W_{ij} x_i x_j$$

$$\text{s.t. } x_i^2 = 1, \quad i = 1, \dots, n,$$

where  $\vec{x} \in \mathbb{R}^n$  is the parameter vector to be optimized and  $W \in \mathbb{R}^{n \times n}$  is a ~~seeded~~ weight matrix. One interpretation of this problem is to partition  $\{1, 2, \dots, n\}$  in two sets;  $W_{ij}$  is cost of assigning  $i, j$  to the same set;  $-W_{ij}$  is cost of assigning to different sets.

a) [1 point]. Is it a convex optimization problem? Give your reasons.

(a) The objective (cost function):

$$\min_{\vec{x}} \frac{1}{2} \vec{x}^T W \vec{x} = \min_{\vec{x}} \frac{1}{2} \sum_{i,j=1}^n x_i x_j W_{ij}.$$

$$\text{where } x_i^2 \geq 1, \quad i = 1, 2, \dots, n.$$

$\frac{1}{2} \vec{x}^T W \vec{x}$  is nonconvex,  $x_i^2 \geq 1$  is convex function for  $\vec{x}$ , since  $\frac{1}{2} (\theta \vec{x} + (1-\theta) \vec{a})^T W (\theta \vec{x} + (1-\theta) \vec{a})$

$$= \frac{1}{2} (\theta \vec{x}^T W \vec{x} + \theta (1-\theta) \vec{x}^T W \vec{y} + \theta (1-\theta) \vec{y}^T W \vec{x} + (1-\theta)^2 \vec{y}^T W \vec{y})$$

$$= \frac{1}{2} [\theta^2 \vec{x}^T W \vec{x} + 2\theta(1-\theta) \vec{x}^T W \vec{y} + (1-\theta)^2 \vec{y}^T W \vec{y}]$$

$$\leq \frac{1}{2} \theta \vec{x}^T W \vec{x} + \frac{1}{2} (1-\theta) \vec{y}^T W \vec{y}$$

proof:

$$\theta \vec{x}^T W \vec{x} + (1-\theta) \vec{y}^T W \vec{y} \geq \theta^2 \vec{x}^T W \vec{x} + 2\theta(1-\theta) \vec{x}^T W \vec{y} + (1-\theta)^2 \vec{y}^T W \vec{y}$$

$$\theta(1-\theta) \vec{x}^T W \vec{x} + \theta(1-\theta) \vec{y}^T W \vec{y} - 2\theta(1-\theta) \vec{x}^T W \vec{y} \geq 0.$$

since  $0 < \theta < 1, \quad \theta(1-\theta) > 0$ 

$$\vec{x}^T W \vec{x} + \vec{y}^T W \vec{y} - 2 \vec{x}^T W \vec{y} \geq 0.$$

$$(\vec{x} - \vec{y})^T W (\vec{x} - \vec{y}) \geq 0.$$

W might be a negative matrix.

for e.g.  $\begin{bmatrix} -1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & -1 & 0 \\ 0 & 0 & 0 & -1 \end{bmatrix}$ , therefore

we cannot prove the inequality.

(b).

b) [3 points]. Derive the Lagrange dual function of the above problem. Hint: recall the definition of a positive definite matrix<sup>1</sup>.<sup>1</sup>[https://en.wikipedia.org/wiki/Definiteness\\_of\\_a\\_matrix](https://en.wikipedia.org/wiki/Definiteness_of_a_matrix)Lagrangian:  $\vec{x}^T \vec{\alpha} \vec{x}$ 

$$L(\vec{x}, \vec{\alpha}) = \frac{1}{2} \vec{x}^T W \vec{x} + \sum_{i=1}^n \alpha_i (x_i^2 - 1)$$

$$= \frac{1}{2} \vec{x}^T W \vec{x} + \sum_{i=1}^n \alpha_i x_i^2 - \sum_{i=1}^n \alpha_i$$

$$= \frac{1}{2} \vec{x}^T W \vec{x} + \vec{x}^T A \vec{x} - \sum_{i=1}^n \text{tr}(A) \quad \text{where } A \text{ is a diagonal matrix}$$

$$\nabla_{\vec{x}} L = 2 \vec{x} \left( \frac{1}{2} W + A \right)$$

Lagrange dual function.

$$g(\vec{\alpha}) = \inf_{\vec{x}} L(\vec{x}, \vec{\alpha})$$

$$= \inf_{\vec{x}} \vec{x}^T \left( \frac{1}{2} W + A \right) \vec{x} - \text{tr}(A) = \begin{cases} -\text{tr}(A) & \text{when } \frac{1}{2} W + A \succeq 0 \\ -\infty & \text{otherwise.} \end{cases}$$

- c) [1 point]. Try to give a nontrivial lower bound (i.e., other than  $-\infty$ ) of the optimal value  $p^*$  of the above problem. Hint: use the lower bound property.

$$P^* = \frac{1}{2} (x^*)^T w x^*$$

$$\geq L(x^*, a) \geq g$$

$$\geq -\text{tr}(A) \quad (\text{if } \nabla_w A \geq 0)$$

Therefore, the lower bound of  $\frac{1}{2} x^T w x$   
is  $-\text{tr}(A)$ , where  $A$  is the diagonal  
matrix  $\begin{bmatrix} a_1 & & & \\ & a_2 & & \\ & & \ddots & \\ & & & a_n \end{bmatrix}$  where  $a_i$  is a Lagrange  
multiplier.

### Exercise 3

[5 points]. Prove the following matrix identity

$$(P^{-1} + B^T R^{-1} B)^{-1} B^T R^{-1} = P B^T (B P B^T + R)^{-1}, \quad (1)$$

where  $P \in \mathbb{R}^{n \times n}$ ,  $B \in \mathbb{R}^{m \times n}$ , and  $R \in \mathbb{R}^{m \times m}$ .  $P$  and  $R$  are invertible. Note that if  $m \ll n$ , it will be much cheaper to evaluate the right-hand side than the left-hand side. Hint: right multiplying both sides by  $(B P B^T + R)$ . Using Eq. (1) to prove a special case

$$(I + AB)^{-1} A = A(I + BA)^{-1},$$

where  $A \in \mathbb{R}^{n \times m}$  and  $B \in \mathbb{R}^{m \times n}$ .

$$\begin{aligned} \text{Proof: L.H.S.} &= (P^{-1} + B^T R^{-1} B)^{-1} B^T R^{-1} \\ &= [P + B^T R(B^T)^{-1}] B^T R^{-1} \\ &= P B^T R^{-1} + B^T R(B^T)^{-1} B^T R^{-1} \\ &= P B^T R^{-1} + B^T \\ &= P B^T R^{-1} + P B^T (P B^T)^{-1} B^{-1} \\ &= P B^T (R^{-1} + (P B^T)^{-1} B^{-1}) \\ &= P B^T (B P B^T + R)^{-1} = \text{R.H.S.} \quad * \end{aligned}$$

### Exercise 4 - Optional

For a pair of random variables  $(X, Y)$ , sampled from the joint distribution  $P(X, Y) = P(X)P(Y|X)$ , the mutual information between  $X$  and  $Y$  is defined as

$$\begin{aligned} I(X; Y) &= \sum_{x \in X} \sum_{y \in Y} p(X=x, Y=y) \log \frac{p(X=x, Y=y)}{p(X=x)p(Y=y)} \\ &= \sum_{x \in X} \sum_{y \in Y} p(x, y) \log \frac{p(x, y)}{p(x)p(y)}. \end{aligned}$$

Prove that for fixed  $P(Y|X)$  the mutual information  $I(X; Y)$  is a concave function of  $P(X)$ . Hint: recall (In Home Assignment 1) the relationship between mutual information and entropy which is defined as

$$H(X) = - \sum_{x \in X} p(X=x) \log p(X=x).$$

$$\begin{aligned} I(X; Y) &= \sum_{x \in X} \sum_{y \in Y} p(x, y) \log \frac{p(x, y)}{p(x)p(y)} \\ &= \sum_{x \in X} \sum_{y \in Y} p(x, y) \log p(Y|x) \\ &\quad - \sum_{x \in X} \sum_{y \in Y} p(x, y) \log p(y) \\ &= \sum_{x \in X} \sum_{y \in Y} p(x, y) \log p(Y|x) - \sum_{y \in Y} p(y) \log p(y) \\ &= \sum_{x \in X} p(x) \sum_{y \in Y} p(y|x) \log p(y|x) - \sum_{y \in Y} p(y) \log p(y). \\ &= H(Y) - H(Y|X) \end{aligned}$$

$$\begin{aligned} A &= \theta I(X; Y) + (1-\theta) [I(X; Y) - H(Y)] \\ &= \sum_{x \in X} [\theta p(x) + (1-\theta) q(x)] \sum_{y \in Y} p(y|x) \log p(y|x) \\ &\quad - \theta \sum_{y \in Y} (p(y|x) p(x)) \log (p(y|x) p(x)) \\ &\quad - (1-\theta) \sum_{y \in Y} (p(y|x) q(x)) \log (p(y|x) q(x)) \end{aligned}$$

$$\begin{aligned} B &= I(X; Y) - (\theta p + (1-\theta) q) \\ &= \sum_{x \in X} [\theta p(x) + (1-\theta) q(x)] \sum_{y \in Y} p(y|x) \log p(y|x) \\ &\quad - \sum_{y \in Y} p(y|x) (\theta p(x) + (1-\theta) q(x)) \log (\theta p(x) + (1-\theta) q(x)) \end{aligned}$$

$$\begin{aligned} B-A &= \sum_{y \in Y} [\theta p(y|x) p(x) \log (p(y|x) p(x)) + (1-\theta) p(y|x) q(x) \log (p(y|x) q(x))] \\ &\quad - p(y|x) (\theta p(x) + (1-\theta) q(x)) \log (\theta p(x) + (1-\theta) q(x)) \end{aligned}$$

$I$  is concave wrt.  $p(x)$ . If  $B-A \geq 0$ .

For each  $y \in Y$ , if the Big Agg 30 holds, we can prove it.

For each  $y \in Y$ , denotes  $p(y|x) = k$ ,  $p(x) = p$ ,  $q(x) = q$ .

$$\begin{aligned} &\theta k p \log(kp) + (1-\theta) k q \log(kq) \\ &- k(\theta p + (1-\theta) q) \log(k(\theta p + (1-\theta) q)) \geq 0 \\ &k > 0 \end{aligned}$$

$$\begin{aligned} &\Leftrightarrow \theta p \log p + p \log k + (1-\theta) [q \log q + q \log k] \\ &- (\theta p + (1-\theta) q) \log (\theta p + (1-\theta) q) \\ &- \theta p \log k - (1-\theta) q \log k \geq 0 \end{aligned}$$

$$\Leftrightarrow \theta(\log p) + (1-\theta)(q \log q) \geq \theta p + (1-\theta) q \log(\theta p + (1-\theta) q).$$

$$\begin{aligned} &\Leftrightarrow f'(x) = x \log x \text{ is convex. } (x > 0) \\ &f''(x) = \frac{d}{dx} (1 + \log x) = \frac{1}{x} > 0. \quad ** \end{aligned}$$

- c) [1 point]. Try to give a nontrivial lower bound (*i.e.*, other than  $-\infty$ ) of the optimal value  $p^*$  of the above problem. Hint: use the lower bound property.

### Exercise 3

[5 points]. Prove the following matrix identity

$$(P^{-1} + B^T R^{-1} B)^{-1} B^T R^{-1} = P B^T (B P B^T + R)^{-1}, \quad (1)$$

where  $P \in \mathbb{R}^{n \times n}$ ,  $B \in \mathbb{R}^{m \times n}$ , and  $R \in \mathbb{R}^{m \times m}$ .  $P$  and  $R$  are invertible. Note that if  $m \ll n$ , it will be much cheaper to evaluate the right-hand side than the left-hand side. Hint: right multiplying both sides by  $(B P B^T + R)$ . Using Eq. (1) to prove a special case

$$(I + AB)^{-1} A = A(I + BA)^{-1},$$

where  $A \in \mathbb{R}^{n \times m}$  and  $B \in \mathbb{R}^{m \times n}$ .

### Exercise 4 - Optional

For a pair of random variables  $(X, Y)$ , sampled from the joint distribution  $P(X, Y) = P(X)P(Y|X)$ , the mutual information between  $X$  and  $Y$  is defined as

$$\begin{aligned} I(X; Y) &= \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p(X = x, Y = y) \log \frac{p(X = x, Y = y)}{p(X = x)p(Y = y)} \\ &= \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p(x, y) \log \frac{p(x, y)}{p(x)p(y)}. \end{aligned}$$

Prove that for fixed  $P(Y|X)$ , the mutual information  $I(X; Y)$  is a concave function of  $P(X)$ . Hint: recall (in Home Assignment I) the relationship between mutual information and entropy which is defined as

$$H(X) = - \sum_{x \in \mathcal{X}} p(X = x) \log p(X = x).$$