

Solution of Assignment 3

November 6, 2014

1. Show that in regularized linear regression, $\mathbf{w}^* = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_d)^{-1} \mathbf{X}^\top \mathbf{r} = \mathbf{X}^\top (\mathbf{X} \mathbf{X}^\top + \lambda \mathbf{I}_N)^{-1} \mathbf{r}$, a linear combination of the examples.

(Hint: use the matrix inversion in block form: <http://www.cs.nthu.edu.tw/~jang/book/addenda/matinv/matinv/>)

Answer :

$$(\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_d)^{-1} \mathbf{X}^\top$$

Since the matrix inversion lemma, let $A = \lambda \mathbf{I}_d$, $B = \mathbf{X}^\top$, $C = \mathbf{I}_N$, $D = \mathbf{X}$

$$\begin{aligned} &= \frac{1}{\lambda} \mathbf{X}^\top - \frac{1}{\lambda^2} \mathbf{X}^\top (\mathbf{I}_N + \frac{1}{\lambda} \mathbf{X} \mathbf{X}^\top)^{-1} \mathbf{X} \mathbf{X}^\top \\ &= \frac{1}{\lambda} (\mathbf{X}^\top - \mathbf{X}^\top (\lambda)^{-1} (\mathbf{I}_N + \frac{1}{\lambda} \mathbf{X} \mathbf{X}^\top)^{-1} \mathbf{X} \mathbf{X}^\top) \\ &= \frac{1}{\lambda} (\mathbf{X}^\top - \mathbf{X}^\top (\lambda \mathbf{I}_N + \mathbf{X} \mathbf{X}^\top)^{-1} \mathbf{X} \mathbf{X}^\top) \\ &= \frac{1}{\lambda} (\mathbf{X}^\top (\lambda \mathbf{I}_N + \mathbf{X} \mathbf{X}^\top)^{-1} (\lambda \mathbf{I}_N + \mathbf{X} \mathbf{X}^\top) - \mathbf{X}^\top (\lambda \mathbf{I}_N + \mathbf{X} \mathbf{X}^\top)^{-1} \mathbf{X} \mathbf{X}^\top) \\ &= \frac{1}{\lambda} ((\mathbf{X}^\top (\lambda \mathbf{I}_N + \mathbf{X} \mathbf{X}^\top)^{-1}) (\lambda \mathbf{I}_N + \mathbf{X} \mathbf{X}^\top - \mathbf{X} \mathbf{X}^\top)) \\ &= \mathbf{X}^\top (\mathbf{X} \mathbf{X}^\top + \lambda \mathbf{I}_N)^{-1} \end{aligned}$$

2. Show that in an RKHS, the inner product $\langle f, g \rangle := \sum_{i=1}^n \sum_{j=1}^m \alpha^{(i)} \beta^{(j)} k(\mathbf{x}^{(i)}, \mathbf{y}^{(j)})$ for any $f = \sum_{i=1}^n \alpha^{(i)} k(\mathbf{x}^{(i)}, \cdot)$ and $g = \sum_{j=1}^m \beta^{(j)} k(\mathbf{y}^{(j)}, \cdot)$ is well-defined; i.e., it satisfies

- (a) *symmetry*: $\langle f, g \rangle = \langle g, f \rangle$;
- (b) *linearity*: $\langle af + bg, h \rangle = a \langle f, h \rangle + b \langle g, h \rangle$ for any $a, b \in \mathbb{R}$; and
- (c) *positive definiteness*: $\langle f, f \rangle \geq 0$ with equality holds iff $f(\cdot) = 0(\cdot)$.

Answer :

- (a) $\langle f, g \rangle = \sum_{i=1}^n \sum_{j=1}^m \alpha^{(i)} \beta^{(j)} k(\mathbf{x}^{(i)}, \mathbf{y}^{(j)})$
 $= \sum_{j=1}^m \sum_{i=1}^n \beta^{(j)} \alpha^{(i)} k(\mathbf{y}^{(j)}, \mathbf{x}^{(i)}) = \langle g, f \rangle$
- (b) $\langle af + bg, h \rangle = \sum_{i=1}^n \sum_{j=1}^m \alpha^{(i)} \beta^{(j)} k(ax^{(i)} + by^{(i)}, z^{(j)}) = \sum_{i=1}^n \sum_{j=1}^m \alpha^{(i)} \beta^{(j)} (k(ax^{(i)}, z^{(j)}) + k(by^{(i)}, z^{(j)}))$
 $= \sum_{i=1}^n \sum_{j=1}^m \alpha^{(i)} \beta^{(j)} (a \times k(x^{(i)}, z^{(j)}) + b \times k(y^{(i)}, z^{(j)})) = a \langle f, h \rangle + b \langle g, h \rangle.$

(c) $\langle f, f \rangle = \left\langle \sum_{i=1}^n \alpha^{(i)} k(x^{(i)}, \cdot), \sum_{j=1}^n \alpha^{(j)} k(x^{(j)}, \cdot) \right\rangle = \sum_{i=1}^n \sum_{j=1}^n \alpha^{(i)} \alpha^{(j)} k(x^{(i)}, x^{(j)}) \geq 0$ (by definition of kernel).

Next, we show that $\langle f, f \rangle = 0$ iff $\forall x, \sum_{i=1}^n \alpha^{(i)} k(x^{(i)}, \cdot) = 0(\cdot)$

$\Rightarrow 0 = \sum_{i=1}^n \sum_{j=1}^n \alpha^{(i)} \alpha^{(j)} k(x^{(i)}, x^{(j)}) = \sum_{j=1}^n \alpha^{(j)} (\sum_{i=1}^n \alpha^{(i)} k(x^{(i)}, x^{(j)})) \Rightarrow f = 0(\cdot).$

$\Leftarrow Trivial.$

3. Show that in a large-margin linear classifier, the margin between the hyperplanes $\{\mathbf{x} : \mathbf{w}^\top \mathbf{x} - b = 1\}$ and $\{\mathbf{x} : \mathbf{w}^\top \mathbf{x} - b = -1\}$ is $2/\|\mathbf{w}\|$.

Answer :

Let p be a point on h_1 , i.e., $\mathbf{w}^\top p - b = 1$, and q be the point on h_2 , $\mathbf{w}^\top q - b = -1$, which has minimum distance to p . Since $\frac{\mathbf{w}}{\|\mathbf{w}\|}$ is the unit normal vector of h_1 , we have $q = p + t \frac{\mathbf{w}}{\|\mathbf{w}\|}$. We can solve t by $\mathbf{w}^\top q - b = -1 \Rightarrow \mathbf{w}^\top (q = p + t \frac{\mathbf{w}}{\|\mathbf{w}\|}) - b = -1 \Rightarrow t = \frac{(-1 - (\mathbf{w}^\top p - b))}{\|\mathbf{w}\|} = \frac{-2}{\|\mathbf{w}\|}$. So, the distance between p and q is $\|p - q\| = \|p - (p + t \frac{\mathbf{w}}{\|\mathbf{w}\|})\| = \frac{-2}{\|\mathbf{w}\|}$.

4. Prove the Semiparametric Representer theorem below:

Theorem 1. Let $\tilde{g} := g + b\psi$, where $g \in \mathcal{H}$, $b \in \mathbb{R}$, and $\psi : \mathcal{I} \rightarrow \mathbb{R}$. Then each minimizer \tilde{h} of the regularized risk functional:

$$\arg \min_{\tilde{g}} L((\mathbf{x}^{(1)}, r^{(1)}, \tilde{g}(\mathbf{x}^{(1)})), \dots, (\mathbf{x}^{(N)}, r^{(N)}, \tilde{g}(\mathbf{x}^{(N)}))) + \Omega(\|g\|_{\mathcal{RKHS}})$$

subject to $C_k((\mathbf{x}^{(1)}, r^{(1)}, \tilde{g}(\mathbf{x}^{(1)})), \dots, (\mathbf{x}^{(N)}, r^{(N)}, \tilde{g}(\mathbf{x}^{(N)}))) \leq 0, \forall k$

admits the form $\tilde{h}(\mathbf{x}) = \sum_{t=1}^N c_t k(\mathbf{x}^{(t)}, \mathbf{x}) + b\psi(\mathbf{x})$.

Answer :

(a) Consider $\tilde{\Omega}(\|g\|_{\mathcal{RKHS}}^2)$ for convenience.

(b) By the fundamental theorem of linear algebra, decompose any $g \in H$ into two vectors parallel and orthogonal to $V = \text{span}(k(\cdot, x_i))_{i=1, \dots, m}$ respectively. We have $g(x) = f_p(x) + f_o(x) \Rightarrow \tilde{g}(x) = f_p(x) + f_o(x) + b\psi(x)$, where $f_p \in V$, and $f_o \in V^\perp$.

(c) Since $f_o \in V^\perp \Rightarrow f_o(x) = 0$, $\tilde{g}(x) = f_p(x) + b\psi(x)$.

(d) Suppose the minimizer \tilde{h} has the form $\tilde{h}(x) = f_p(x) + f_o(x) + b\psi(x)$; however, $\tilde{h} - f_o$ is always a better solution.

i. First, $\tilde{h} - f_o$ satisfies all constraints C'_k s as $(\tilde{h} - f_o)(x^{(i)}) = \tilde{h}(x^{(i)})$ for all $1 \leq i \leq N$.

ii. Due to the same reason, $\tilde{h} - f_o$ has the same loss score from L as h .

iii. $\tilde{\Omega}(\|h - f_o\|^2) = \tilde{\Omega}(\|f_p\|^2) \leq \tilde{\Omega}(\|f_p\|^2 + \|f_o\|^2)$
 $= \tilde{\Omega}(\|f_p + f_o\|^2) = \tilde{\Omega}(\|h\|^2)$

(e) By contradiction, $\tilde{h}(x) = f_p(x) + b\psi(x) = \sum_{t=1}^N c_t k(x^{(t)}, x) + b\psi(x)$.

5. Consider the necessary optimality condition for the SVM dual (that is, if the following holds, then $\alpha = \alpha^*$):

$$\exists \rho \in \mathbb{R} \text{ such that } \max_{i \in I_{up}} r^{(i)} g_i \leq \rho \leq \min_{j \in I_{down}} r^{(j)} g_j \quad (1)$$

where $I_{up} := \{i | r^{(i)} \alpha_i < B^{(i)}\}$ and $I_{down} := \{j | A^{(j)} < r^{(j)} \alpha_j\}$.

I_{up} :

正Instance : $r^{(i)} \alpha_i < C \cdot \alpha_i < C \cdot r^{(i)}$ 在線上或線外

負Instance : $r^{(i)} \alpha_i < 0 \cdot \alpha_i > 0 \cdot r^{(i)}$ 在線上或線內

I_{down} :

正Instance : $r^{(i)} \alpha_i > 0 \cdot \alpha_i > 0 \cdot r^{(i)}$ 在線上或線內

負Instance : $r^{(i)} \alpha_i > -C \cdot \alpha_i < C \cdot r^{(i)}$ 在線上或線外

正Instance : $rg > \rho \cdot g > \rho \cdot r\alpha = C \cdot \alpha = C \cdot SV$
 負Instance : $rg > \rho \cdot g < \rho \cdot r\alpha = 0 \cdot \alpha = 0 \cdot \text{Non SV}$
 正Instance又在 I_{down} 才是SV

正Instance : $rg < \rho \cdot g < \rho \cdot r\alpha = 0 \cdot \alpha = 0 \cdot \text{Non SV}$
 負Instance : $rg < \rho \cdot g > \rho \cdot r\alpha = -C \cdot \alpha = C \cdot SV$
 負Instance又在 I_{up} 才是SV

(a) Show that the above condition can be rewritten as

$$\exists \rho \in \mathbb{R} \text{ such that } \forall t, \begin{cases} r^{(t)}\alpha_t^* = B^{(t)}, & \text{if } r^{(t)}g_t^* > \rho \\ r^{(t)}\alpha_t^* = A^{(t)}, & \text{if } r^{(t)}g_t^* < \rho \end{cases}, \begin{matrix} I_{down} \text{的正Ins或} I_{up} \text{的非SV} \\ I_{up} \text{的SV或} I_{down} \text{的非SV} \end{matrix} \quad (2)$$

which is equivalent to

$$\text{結論} \cdot g > \rho \text{ 才是SV} \cdot \text{此時} \alpha = C \quad \exists \rho \in \mathbb{R} \text{ such that } \forall t, \begin{cases} \alpha_t^* = C, & \text{if } g_t^* > r^{(t)}\rho \text{ 是SV} \\ \alpha_t^* = 0, & \text{if } g_t^* < r^{(t)}\rho \end{cases}. \quad (3)$$

- (b) Show that Condition (1) is also sufficient by letting $\mathbf{w}^* = \sum_{t=1}^N \alpha_t^* r^{(t)} \Phi(\mathbf{x}^{(t)})$, $b^* = \rho$, and $\xi_t^* = \max\{0, g_t^* - r^{(t)}\rho\}$.
 (Hint: show that there is no duality gap, i.e., $\text{primal}(\mathbf{w}^*, b^*, \xi^*) - \text{dual}(\alpha^*) = 0$, by Condition (3))

Answer :

不等於B只會<B

- (a) If $r^{(t)}g_t^* > \rho$ and $r^{(t)}\alpha_t^* < B^{(t)}$, then the instance t must be in I_{up} , meaning $\max_{i \in I_{up}} r^{(i)}g_i > \rho$ in Condition 1, a contradiction. Similarly, if $r^{(t)}g_t^* < \rho$ and $r^{(t)}\alpha_t^* > A^{(t)}$, then the instance t must be in I_{down} , leading to another contradiction $\rho > \min_{j \in I_{down}} r^{(j)}g_j$ to Condition 1. Condition (3) is just a simple rewrite of Condition (2).

- (b) Recall that

$$\text{primal}(w, b, \xi) = \frac{1}{2} \|w\|^2 + C \sum_{t=1}^N \xi_t,$$

$$\text{dual}(\alpha) = \alpha 1 - \frac{1}{2} \alpha^\top \tilde{K} \alpha = \sum_{t=1}^N \alpha_t - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N r^{(i)} r^{(j)} \alpha_i \alpha_j K_{i,j}.$$

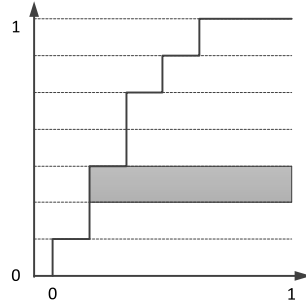
We show that there is no duality gap, i.e., $\text{primal}(\mathbf{w}^*, b^*, \xi^*) - \text{dual}(\alpha^*) = 0$, by letting $\mathbf{w}^* = \sum_{t=1}^N \alpha_t^* r^{(t)} \Phi(\mathbf{x}^{(t)})$, $b^* = \rho$, and $\xi_t^* = \max\{0, g_t^* - r^{(t)}\rho\}$.

$$\begin{aligned} & \text{primal}(\mathbf{w}^*, b^*, \xi^*) - \text{dual}(\alpha^*) \\ &= \left(\frac{1}{2} \|\mathbf{w}^*\|^2 + C \sum_{t=1}^N \xi_t^* \right) - \left(\sum_{t=1}^N \alpha_t^* - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N r^{(i)} r^{(j)} \alpha_i^* \alpha_j^* K_{i,j} \right) \\ &= C \sum_{t=1}^N \xi_t^* - \sum_{t=1}^N \alpha_t^* + \sum_{i=1}^N \sum_{j=1}^N r^{(i)} r^{(j)} \alpha_i^* \alpha_j^* K_{i,j} \quad (\text{due to } \frac{1}{2} \|\mathbf{w}^*\|^2 = \frac{1}{2} (\alpha^*)^\top \tilde{K} \alpha^*) \\ &= C \sum_{t=1}^N \xi_t^* - \sum_{i=1}^N \alpha_i^* g_i^* \quad (\text{due to } g_t^* = 1 - r^{(t)} \sum_{i=1}^N r^{(i)} \alpha_i^* K_{i,t}) \\ &= \sum_{t=1}^N (C \xi_t^* - \alpha_t^* g_t^*) \\ &= \sum_{t=1}^N (-\rho r^{(t)} \alpha_t^*) \quad (\text{due to Condition (3)}) \\ &= 0 \quad (\text{due to } r^\top \alpha = 0). \end{aligned}$$

The optimality holds.

6. Show that the *Area Under the ROC Curve* (AUC) is equal to the probability that a classifier ranks a randomly chosen positive instance higher than a randomly chosen negative one.
 (Hint: by partitioning the AUC horizontally)

Answer :



- (a) Consider the horizontal partition of AUC based on each positive instance. The height of each partition (shown as the shaded area in the above figure) is the probability that a positive instance is chosen. On the other hand, the width of the partition is the conditional probability that given a positive instance is chosen, a randomly chosen negative instance is ranked after that positive instance (note that each negative instance ranked before that positive instance contributes to the portion of the horizontal bar before the shaded area). Therefore, by summing up all the partitions, we have the joint probability that a randomly chosen positive instance is ranked before a randomly chosen negative one.