

1.

1) kernels

$$a) \quad k(\vec{u}, \vec{v}) = (\langle u, v \rangle + 1)^3 = (\langle u, v \rangle)^3 + 3(\langle u, v \rangle)^2 + 3\langle u, v \rangle + 1.$$

$$= \sum_{(j_1, \dots, j_d) : \sum j_i = 3} \binom{3}{j_1, \dots, j_d} (u^{(1)})^{j_1} \dots (u^{(d)})^{j_d} (v^{(1)})^{j_1} \dots (v^{(d)})^{j_d} \\ + 3 \cdot \sum_{(k_1, \dots, k_d) : \sum k_i = 2} \binom{2}{k_1, \dots, k_d} (u^{(1)})^{k_1} \dots (u^{(d)})^{k_d} (v^{(1)})^{k_1} \dots (v^{(d)})^{k_d} \\ + 3 \sum_{i=1}^d u^{(i)} v^{(i)} + 1.$$

$$= \langle \Phi(u), \Phi(v) \rangle.$$

$$\text{Thus, } \Phi(u) = \left[\dots, \sqrt{\binom{3}{j_1, \dots, j_d}} (u^{(1)})^{j_1} \dots (u^{(d)})^{j_d}, \dots, \right. \\ \left. \dots, \sqrt{\binom{2}{k_1, \dots, k_d}} (u^{(1)})^{k_1} \dots (u^{(d)})^{k_d}, \dots, \right. \\ \left. \dots, \sqrt{3} (u^{(1)}), \dots, \sqrt{3} (u^{(d)}), 1 \right]$$

$$b) \quad i) \quad k(\vec{x}, \vec{z}) = k_1(\vec{x}, \vec{z}) + k_2(\vec{x}, \vec{z}) \\ = k_1(\vec{z}, \vec{x}) + k_2(\vec{z}, \vec{x}) = k(\vec{z}, \vec{x}).$$

$\Rightarrow k$ is symmetric

$$\text{Suppose } A = B + C, \text{ where } B = \begin{bmatrix} k_1(x_1, x_1) & \dots & k_1(x_1, x_n) \\ \vdots & & \vdots \\ k_1(x_n, x_1) & \dots & k_1(x_n, x_n) \end{bmatrix}.$$

B is PSD matrix for k_1 .

and C is PSD matrix for k_2 .

$$C = \begin{bmatrix} k_2(x_1, x_1) & \dots & k_2(x_1, x_n) \\ \vdots & & \vdots \\ k_2(x_n, x_1) & \dots & k_2(x_n, x_n) \end{bmatrix}$$

Thus, for $\forall \vec{z} \in \mathbb{R}^n$,

$$\vec{z}^T A \vec{z} = \vec{z}^T B \vec{z} + \vec{z}^T C \vec{z} \geq 0. \Rightarrow A \text{ is PSD matrix}$$

Thus $k(\vec{x}, \vec{z})$ is a ^{symmetric} positive-definite kernel.



(ii) $k(\vec{x}, \vec{z})$ may not be a positive-definite kernel.

$$\text{let } k_2(\vec{x}, \vec{z}) = -k_1(\vec{x}, \vec{z})$$

$$\text{Then } k(\vec{x}, \vec{z}) = -k_1(\vec{x}, \vec{z}).$$

$$\text{Let } A = \begin{bmatrix} k_1(x_1, x_1) & \dots & k_1(x_1, x_n) \\ \vdots & & \vdots \\ k_1(x_n, x_1) & \dots & k_1(x_n, x_n) \end{bmatrix} \text{ is the PSD matrix for } k_1.$$

Then $-A$ can't be PSD and thus $k(\vec{x}, \vec{z})$ isn't a positive-definite kernel.

$$(iii) \quad k(\vec{x}, \vec{z}) = a k_1(\vec{x}, \vec{z}) = a k_1(\vec{z}, \vec{x}) = k(\vec{z}, \vec{x}).$$

$\Rightarrow k(\vec{x}, \vec{z})$ is symmetric.

$$\text{let } B = \begin{bmatrix} k(x_1, x_1) & \dots & k(x_1, x_n) \\ \vdots & & \vdots \\ k(x_n, x_1) & \dots & k(x_n, x_n) \end{bmatrix} = a \cdot \begin{bmatrix} k_1(x_1, x_1) & \dots & k_1(x_1, x_n) \\ \vdots & & \vdots \\ k_1(x_n, x_1) & \dots & k_1(x_n, x_n) \end{bmatrix}$$

$$= a \cdot A.$$

Thus A is the PSD matrix for k_1

$$\forall \vec{z} \in \mathbb{R}^n. \quad \vec{z}^T A \vec{z} = a (\vec{z}^T B \vec{z}) \geq 0. \quad (\because a \in \mathbb{R}^+)$$

$\Rightarrow A$ is PSD matrix.

$\Rightarrow k$ is SPD kernel.

$$(iv). \quad k(\vec{x}, \vec{z}) = k_1(\vec{x}, \vec{z}) k_2(\vec{x}, \vec{z}) = k_1(\vec{z}, \vec{x}) \cdot k_2(\vec{z}, \vec{x}) = k(\vec{z}, \vec{x}).$$

$\Rightarrow k$ is symmetric

Suppose $k_1(\vec{x}, \vec{z}) = \Phi_1(x)^T \Phi_1(z)$ (since SPD kernel \Leftrightarrow inner product kernel.)

$$k_2(\vec{x}, \vec{z}) = \Phi_2(x)^T \Phi_2(z)$$

Where $\Phi_1(x) \in \mathbb{R}^n$ $\Phi_2(x) \in \mathbb{R}^m$, $\vec{x}, \vec{z} \in \mathbb{R}^d$.

$$\begin{aligned} \Rightarrow k(\vec{x}, \vec{z}) &= \langle \Phi_1(x), \Phi_1(z) \rangle * \langle \Phi_2(x), \Phi_2(z) \rangle \\ &= \left(\sum_{i=1}^n \Phi_1^{(i)}(x) \Phi_1^{(i)}(z) \right) * \left(\sum_{j=1}^m \Phi_2^{(j)}(x) \Phi_2^{(j)}(z) \right) \end{aligned}$$



$$k(\vec{x}, \vec{z}) = \sum_{i=1}^n \sum_{j=1}^m (\Phi_i^{(1)}(x) \Phi_j^{(1)}(x)) \cdot (\Phi_i^{(1)}(z) \cdot \Phi_j^{(1)}(z))$$

$\Rightarrow k(\vec{x}, \vec{z})$ is a Inner Product Kernel

Thus $k(\vec{x}, \vec{z})$ is a SPD kernel.

$$(v) \quad k(\vec{x}, \vec{z}) = f(\vec{x}) f(\vec{z}) = f(\vec{z}) f(\vec{x}) = k(\vec{z}, \vec{x})$$

$\Rightarrow k$ is symmetric.

$$K = \begin{bmatrix} k(\vec{x}_1, \vec{x}_1) & \dots & k(\vec{x}_1, \vec{x}_n) \\ \vdots & & \vdots \\ k(\vec{x}_n, \vec{x}_1) & \dots & k(\vec{x}_n, \vec{x}_n) \end{bmatrix}_{n \times n} = [f(x_1), \dots, f(x_n)]^T \cdot [f(x_1), \dots, f(x_n)]$$

$$\Rightarrow \forall \vec{z} \in \mathbb{R}^n, \quad \vec{z}^T K \vec{z} = \vec{z}^T [f(x_1), \dots, f(x_n)]^T [f(x_1), \dots, f(x_n)] \vec{z} \\ = \| \vec{z} \cdot [f(x_1), \dots, f(x_n)] \|^2 \geq 0.$$

$\Rightarrow K$ is PSD matrix
symmetric

Thus k is positive-definite kernel.

$$(vi) \quad k(\vec{x}, \vec{z}) = P(\phi_1(\vec{x}, \vec{z})) = P(\phi_1(\vec{z}, \vec{x})) = k(\vec{z}, \vec{x}).$$

$\Rightarrow k$ is symmetric.

According conclusion of (iv). Let $k_2(x, z) = k_1(x, z)$.

Then $k(x, z) = k_1^2(x, z)$ is SPD kernel.

By induction, $k(x, z) = k_1^n(x, z)$ is still SPD kernel.

According conclusion of (iii). $\forall a_n \in \mathbb{R}^+$.

$k(x, z) = a_n \cdot k_1^n(x, z)$ is still SPD kernel

According conclusion of (i). $\forall a_n, a_m \in \mathbb{R}^+$

$k(x, z) = a_n k_1^n(x, z) + a_m k_1^m(x, z)$ is still SPD kernel.

Thus, $k(\vec{x}, \vec{z}) = P(\phi_1(\vec{x}, \vec{z})) = \sum_{i=1}^p a_i k_1^i(x, z)$ ($a_i \in \mathbb{R}^+$) is SPD kernel.



(c). k is IP kernel $\Rightarrow k$ is SPD kernel

$$k(\vec{u}, \vec{v}) = \langle \Phi(\vec{u}), \Phi(\vec{v}) \rangle = \langle \Phi(\vec{v}), \Phi(\vec{u}) \rangle = k(\vec{v}, \vec{u})$$

$\Rightarrow k$ is symmetric.

$$\text{Suppose } A = \begin{pmatrix} k(x_1, x_1) & \dots & k(x_1, x_n) \\ \vdots & & \vdots \\ k(x_n, x_1) & \dots & k(x_n, x_n) \end{pmatrix} = \begin{pmatrix} \Phi(x_1)^T \Phi(x_1) & \dots & \Phi(x_1)^T \Phi(x_n) \\ \vdots & & \vdots \\ \Phi(x_n)^T \Phi(x_1) & \dots & \Phi(x_n)^T \Phi(x_n) \end{pmatrix}$$

$$\forall \vec{z} \in \mathbb{R}^n, \quad \vec{z}^T A \vec{z} = \sum_{i,j=1}^n z_i z_j \langle \Phi(x_i), \Phi(x_j) \rangle$$

$$= \left\langle \sum_{i=1}^n z_i \Phi(x_i), \sum_{j=1}^n z_j \Phi(x_j) \right\rangle$$

$$= \left\| \sum_{i=1}^n z_i \Phi(x_i) \right\|^2 \geq 0.$$

$\Rightarrow A$ is PSD matrix.

Thus, k is SPD kernel.

2). Kernel Ridge Regression.

$$(a). \text{ obj. fun} = \min_{\vec{w}, b} \frac{1}{n} \sum_{i=1}^n (y_i - \vec{w}^T \vec{x}_i - b)^2 + \lambda \|\vec{w}\|^2$$

$$\text{solution: } \hat{\vec{w}} = (\hat{X}^T \hat{X} + n\lambda I)^{-1} \hat{X}^T \hat{\vec{y}}$$

$$\hat{b} = \bar{y} - \hat{\vec{w}}^T \bar{\vec{x}}$$

where

$$\hat{\vec{y}} = \begin{bmatrix} \hat{y}_1 \\ \vdots \\ \hat{y}_n \end{bmatrix}, \quad \hat{X} = \begin{bmatrix} \hat{\vec{x}}_1^T \\ \vdots \\ \hat{\vec{x}}_n^T \end{bmatrix}$$

$$\hat{\vec{x}}_i = \vec{x}_i - \bar{\vec{x}} \in \mathbb{R}^d.$$

$$\hat{y}_i = y_i - \bar{y} \in \mathbb{R}.$$

$$\bar{\vec{x}} = \frac{1}{n} \sum_{i=1}^n \vec{x}_i$$

$$\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$$

$$\text{Thus the obj. fun} = \|\hat{\vec{y}} - \hat{X} \vec{w}\|^2 + n\lambda \|\vec{w}\|^2$$

According to the lemma in notes.

$$(\hat{X}^T \hat{X} + n\lambda I)^{-1} = \frac{1}{n\lambda} [I - \hat{X}^T (\frac{1}{n\lambda} I + \hat{X} \hat{X}^T)^{-1} \hat{X}]$$

$$\Rightarrow \hat{\vec{w}} = (\hat{X}^T \hat{X} + n\lambda I)^{-1} \hat{X}^T \hat{\vec{y}} = \frac{1}{n\lambda} [\hat{X}^T - \hat{X}^T (\frac{1}{n\lambda} I + \hat{X} \hat{X}^T)^{-1} \hat{X} \hat{X}^T] \hat{\vec{y}}$$



$$\hat{W} = \frac{1}{n\lambda} [\hat{X}^T - \tilde{X}^T (\tilde{G} + n\lambda I)^{-1} \tilde{G}] \tilde{y}$$

where $\tilde{G} = \tilde{X} \tilde{X}^T$.

$$\begin{aligned} \hat{b} &= \bar{y} - \hat{W}^T \bar{x} = \bar{y} - \frac{1}{n\lambda} \tilde{y}^T [\tilde{X} - \tilde{G} (\tilde{G} + n\lambda I)^{-1} \tilde{X}] \bar{x} \\ &= \bar{y} - \tilde{y}^T (\tilde{G} + n\lambda I)^{-1} \tilde{X} \bar{x} \\ &= \bar{y} - \tilde{y}^T (\tilde{G} + n\lambda I)^{-1} \tilde{g}(\bar{x}) \end{aligned}$$

where $\tilde{G} = \begin{bmatrix} \langle \tilde{x}_1, \tilde{x}_1 \rangle & \dots & \langle \tilde{x}_1, \tilde{x}_n \rangle \\ \vdots & & \vdots \\ \langle \tilde{x}_n, \tilde{x}_1 \rangle & \dots & \langle \tilde{x}_n, \tilde{x}_n \rangle \end{bmatrix}$ $\tilde{g}(\bar{x}) = \begin{bmatrix} \langle \tilde{x}_1, \bar{x} \rangle \\ \vdots \\ \langle \tilde{x}_n, \bar{x} \rangle \end{bmatrix}$

substitute $\langle \tilde{x}_i, \tilde{x}_j \rangle$ with $\langle \tilde{\Phi}(x_i), \tilde{\Phi}(x_j) \rangle$.

To kernelize, observe: $\langle \tilde{x}_i, \tilde{x}_j \rangle = \langle x_i - \bar{x}, x_j - \bar{x} \rangle$

$$\begin{aligned} &= \langle x_i, x_j \rangle - \frac{1}{n} \sum_{r=1}^n \langle x_i, x_r \rangle - \frac{1}{n} \sum_{s=1}^n \langle x_s, x_j \rangle \\ &\quad + \frac{1}{n^2} \sum_{r=1}^n \sum_{s=1}^n \langle x_r, x_s \rangle \end{aligned}$$

and $\langle \tilde{x}_i, \bar{x} \rangle = \langle x_i - \bar{x}, \bar{x} \rangle$

$$= \frac{1}{n} \sum_{r=1}^n \langle x_i, x_r \rangle - \frac{1}{n^2} \sum_{r=1}^n \sum_{s=1}^n \langle x_r, x_s \rangle$$

Now substitute $\langle x, x' \rangle$ with $k(x, x')$, Define the mean-centered feature map

$$\tilde{\Phi}(x) := \Phi(x) - \frac{1}{n} \sum_{l=1}^n \Phi(x_l)$$

$$\Rightarrow \hat{b} = \bar{y} - \tilde{y}^T (\tilde{K} + n\lambda I)^{-1} \tilde{k}(\bar{x}) \quad \text{where } \tilde{K} = [\langle \tilde{\Phi}(x_i), \tilde{\Phi}(x_j) \rangle]_{1 \leq i, j \leq n}$$

$$\tilde{k}(\bar{x}) = \begin{bmatrix} \langle \tilde{\Phi}(x_1), \frac{1}{n} \sum_{l=1}^n \Phi(x_l) \rangle \\ \vdots \\ \langle \tilde{\Phi}(x_n), \frac{1}{n} \sum_{l=1}^n \Phi(x_l) \rangle \end{bmatrix}$$



(b). $\tilde{K} = K - K O - O K + O K O$

Determine that: $\tilde{K}'_{n \times m} = K' - K O_1 - O_2 K' + O_2 K O_1$

c) where: O_1 is a $n \times m$ matrix where all entries equal to $\frac{1}{n}$
 O_2 is a $n \times n$ matrix where all entries equal to $\frac{1}{n}$

To check this: $\tilde{K}'_{ij} = \langle \tilde{\Phi}(x_i), \tilde{\Phi}(x'_j) \rangle = (\Phi(x_i) - \frac{1}{n} \sum_{p=1}^n \Phi(x_p)) * (\Phi(x'_j) - \frac{1}{n} \sum_{q=1}^n \Phi(x_q))$

$$= \Phi(x_i) \cdot \Phi(x'_j) - \frac{1}{n} \Phi(x_i) \cdot \sum_{q=1}^n \Phi(x_q) - \frac{1}{n} \Phi(x'_j) \cdot \sum_{p=1}^n \Phi(x_p) + \frac{1}{n^2} \sum_{p,q=1}^n \Phi(x_p) \Phi(x_q)$$

$$= K'_{ij} - (K O_1)_{ij} - (O_2 K')_{ij} + (O_2 K O_1)_{ij} \quad \checkmark$$

prediction formula =

$$\hat{f}(x'_i) = \hat{w}^T x'_i + \hat{b} = \bar{y} + \tilde{y}^T (\tilde{K} + n\lambda I)^{-1} \tilde{K}(x'_i) \quad \text{where } \tilde{K}(x'_i) =$$

$$\frac{1}{n\lambda} \tilde{y}^T [\tilde{X} - \tilde{K} (\tilde{K} + n\lambda I)^{-1} \tilde{X}] x'_i$$

$$\begin{bmatrix} \langle \tilde{\Phi}(x_1), \tilde{\Phi}(x'_i) \rangle \\ \vdots \\ \langle \tilde{\Phi}(x_n), \tilde{\Phi}(x'_i) \rangle \end{bmatrix}$$

$$\hat{f}(X') = \begin{bmatrix} 1 \\ \vdots \\ \vdots \\ 1 \end{bmatrix}_{m \times 1} * \bar{y} + \tilde{K}'^T (\tilde{K} + n\lambda I)^{-1} \tilde{y}$$

(where x' is $m \times d$ test data matrix)

3). Support Vector Regression

a. The obj. fun: $\frac{1}{2} \|w\|^2 + \frac{c}{n} \sum_{i=1}^n (\xi_i^+ + \xi_i^-)$

let $t_i = w^T x_i - b$. Observe that if $\xi_i^+ > 0$, $y_i - t_i = \epsilon + \xi_i^+$

That is, the inequality constraint is an equality because increasing ξ_i^+ beyond equality unnecessarily increases the obj. fun. Since

this equality can only hold if $y_i - t_i \geq \epsilon$, $\xi_i^+ = \max\{0, y_i - t_i - \epsilon\}$

Similarly we may deduce that $\xi_i^- = \max\{0, t_i - y_i - \epsilon\}$



4.

Note that if $\xi_i^+ > 0$ then $\xi_i^- = 0$, and vice versa. Thus,

$$\xi_i^- + \xi_i^+ = \max\{0, |y_i - t_i| - \varepsilon\},$$

and the optimization problem reduce to

$$\min_{\vec{w}, b} \frac{1}{2} \|\vec{w}\|^2 + \frac{C}{n} \sum_{i=1}^n \max\{0, |y_i - t_i| - \varepsilon\}$$

or setting $\lambda = \frac{1}{2C}$

$$\min_{\vec{w}, b} \frac{1}{n} \sum_{i=1}^n \ell_\varepsilon(y_i, \vec{w}^T \mathbf{x}_i + b) + \lambda \|\vec{w}\|^2$$

b. The optimization problem is:

$$\min_{\vec{w}, b, \xi^+, \xi^-} \frac{1}{2} \|\vec{w}\|^2 + \frac{C}{n} \sum_{i=1}^n (\xi_i^+ + \xi_i^-)$$

$$\text{s.t. } y_i - \vec{w}^T \mathbf{x}_i - b \leq \varepsilon + \xi_i^+ \quad \forall i$$

$$\vec{w}^T \mathbf{x}_i + b - y_i \leq \varepsilon + \xi_i^- \quad \forall i$$

$$\xi_i^+ \geq 0 \quad \forall i$$

$$\xi_i^- \geq 0 \quad \forall i$$

Lagrangian: $L(\vec{w}, b, \xi^+, \xi^-, \alpha, \beta, \lambda, \gamma) = \frac{1}{2} \|\vec{w}\|^2 + \frac{C}{n} \sum_{i=1}^n (\xi_i^+ + \xi_i^-)$
 $+ \sum_{i=1}^n \alpha_i (y_i - \vec{w}^T \mathbf{x}_i - b - \varepsilon - \xi_i^+)$
 $+ \sum_{i=1}^n \beta_i (\vec{w}^T \mathbf{x}_i + b - y_i - \varepsilon - \xi_i^-)$
 $- \sum_{i=1}^n \lambda_i \xi_i^+ - \sum_{i=1}^n \gamma_i \xi_i^-$

Dual Function: $L_D(\alpha, \beta, \lambda, \gamma) = \min_{\vec{w}, b, \xi^+, \xi^-} L(\vec{w}, b, \xi^+, \xi^-, \alpha, \beta, \lambda, \gamma)$

original problem is quadratic w.r.t. \vec{w} , and linear w.r.t. ξ^+, ξ^- , it's convex, we can get optimization of Dual Function above by:

$$\frac{\partial L}{\partial \vec{w}} = \vec{w} - \sum_{i=1}^n \alpha_i \mathbf{x}_i + \sum_{i=1}^n \beta_i \mathbf{x}_i \stackrel{!}{=} 0 \Rightarrow \vec{w}^* = \sum_{i=1}^n \alpha_i^* \mathbf{x}_i - \sum_{i=1}^n \beta_i^* \mathbf{x}_i$$

$$\frac{\partial L}{\partial b} = -\sum_{i=1}^n \alpha_i + \sum_{i=1}^n \beta_i \stackrel{!}{=} 0 \Rightarrow \sum_{i=1}^n \alpha_i^* = \sum_{i=1}^n \beta_i^*$$



$$\frac{\partial L}{\partial \xi_i^+} = \frac{c}{n} - \alpha_i - \lambda_i \stackrel{!}{=} 0 \Rightarrow \frac{c}{n} = \lambda_i^* + \alpha_i^*$$

$$\frac{\partial L}{\partial \xi_i^-} = \frac{c}{n} - \beta_i - \gamma_i \stackrel{!}{=} 0 \Rightarrow \frac{c}{n} = \beta_i^* + \gamma_i^*$$

$$\Rightarrow L_D(\alpha, \beta, \lambda, \gamma) = \frac{1}{2} \|w^*\|_2^2 + \frac{c}{n} \sum_{i=1}^n (\xi_i^+ + \xi_i^-) + \sum_{i=1}^n \alpha_i^* (y_i - w^{*T} x_i - b - \varepsilon - \xi_i^+) \\ + \sum_{i=1}^n \beta_i^* (w^{*T} x_i + b - y_i - \varepsilon - \xi_i^-) - \sum_{i=1}^n \lambda_i^* \xi_i^+ - \sum_{i=1}^n \gamma_i^* \xi_i^-$$

In order to satisfy KKT condition, we still need:

$$\alpha_i^* (y_i - w^{*T} x_i - b - \varepsilon - \xi_i^+) = 0 \quad \forall i$$

$$\beta_i^* (w^{*T} x_i + b - y_i - \varepsilon - \xi_i^-) = 0 \quad \forall i \quad \text{where } x_i^* \text{ is primal optimal}$$

$$L_D(\alpha, \beta, \lambda, \gamma) = \frac{1}{2} \left\| \sum_{i=1}^n (\alpha_i^* - \beta_i^*) x_i \right\|^2 + \frac{c}{n} \sum_{i=1}^n (\xi_i^+ + \xi_i^-) \\ + \sum_{i=1}^n \alpha_i^* (y_i - w^{*T} x_i - b - \varepsilon) - \sum_{i=1}^n \alpha_i^* \xi_i^+ - \sum_{i=1}^n \lambda_i^* \xi_i^+ \\ + \sum_{i=1}^n \beta_i^* (w^{*T} x_i + b - y_i - \varepsilon) - \sum_{i=1}^n \beta_i^* \xi_i^- - \sum_{i=1}^n \gamma_i^* \xi_i^- \\ = \frac{1}{2} \sum_{i,j=1}^n (\alpha_i^* - \beta_i^*) (\alpha_j^* - \beta_j^*) x_i^T x_j + \sum_{i=1}^n \alpha_i^* (y_i - w^{*T} x_i - \varepsilon) \\ + \sum_{i=1}^n \beta_i^* (w^{*T} x_i - y_i - \varepsilon) \\ = \frac{1}{2} \sum_{i,j=1}^n (\alpha_i^* - \beta_i^*) (\alpha_j^* - \beta_j^*) x_i^T x_j + \sum_{i=1}^n \alpha_i^* (y_i - \sum_{j=1}^n (\alpha_j^* - \beta_j^*) x_j^T x_i - \varepsilon) \\ + \sum_{i=1}^n \beta_i^* (\sum_{j=1}^n (\alpha_j^* - \beta_j^*) x_j^T x_i - y_i - \varepsilon) \\ = -\frac{1}{2} \sum_{i,j=1}^n (\alpha_i^* - \beta_i^*) (\alpha_j^* - \beta_j^*) x_i^T x_j - \sum_{i=1}^n (\alpha_i^* + \beta_i^*) \varepsilon$$

s.t. $\alpha_i^* \geq 0, \beta_i^* \geq 0, \forall i \quad \lambda_i^* \geq 0, \gamma_i^* \geq 0, \forall i.$

$$\sum_{i=1}^n \alpha_i^* = \sum_{i=1}^n \beta_i^*, \quad \frac{c}{n} = \lambda_i^* + \alpha_i^* = \beta_i^* + \gamma_i^*.$$



Thus, the Dual Optimization Problem is:

$$\max_{\alpha, \beta} -\frac{1}{2} \sum_{i,j=1}^n (\alpha_i - \beta_j) (\alpha_j - \beta_i) x_i^T x_j - \sum_{i=1}^n (\alpha_i + \beta_i) \varepsilon$$

$$\text{st. } \alpha_i \geq 0, \beta_i \geq 0; \lambda_i \geq 0, \gamma_i \geq 0, \forall i \left\{ \begin{array}{l} \sum_{i=1}^n \alpha_i = \sum_{i=1}^n \beta_i, \quad \frac{C}{n} = \lambda_i + \alpha_i = \beta_i + \gamma_i \end{array} \right.$$

above constraints equal to: $0 \leq \alpha_i \leq \frac{C}{n}, 0 \leq \beta_i \leq \frac{C}{n}, \sum_{i=1}^n \alpha_i = \sum_{i=1}^n \beta_i$

c. Let α^*, β^* be dual opt. $\Rightarrow w^* = \sum_{i=1}^n (\alpha_i^* - \beta_i^*) x_i$

from complementary slackness,

$$\alpha_i^* (y_i - w^{*T} x_i - b - \varepsilon - \zeta_i^+) = 0, \quad \forall i$$

$$\beta_i^* (w^{*T} x_i + b - y_i - \varepsilon - \zeta_i^-) = 0$$

consider any i such that $0 < \alpha_i^* < \frac{C}{n}$, since $\alpha_i^* > 0$, we know

$$y_i - w^{*T} x_i - b = \varepsilon + \zeta_i^{+*}$$

Since $\alpha_i^* < \frac{C}{n}$, we know $\lambda_i > 0 \Rightarrow \zeta_i^+ = 0$.

With the same reason, we have that for any j s.t. $0 < \beta_j^* < \frac{C}{n}$.

$$w^{*T} x_j + b - y_j = \varepsilon + \zeta_j^{-*} \quad \& \quad \zeta_j^- = 0$$

Solve for $b^* \Rightarrow b^* = y_i - w^{*T} x_i - \varepsilon$ or $b^* = y_j - w^{*T} x_j + \varepsilon$

From the solution of (b), we can easily find that the dual optimization problem is a inner product of $\langle x_i, x_j \rangle$. We can substitute it with $k(x_i, x_j) = \langle \Phi(x_i), \Phi(x_j) \rangle$, Then SVR solves

$$\max_{\alpha, \beta} -\frac{1}{2} \sum_{i,j=1}^n (\alpha_i - \beta_j) (\alpha_j - \beta_i) k(x_i, x_j) - \sum_{i=1}^n (\alpha_i + \beta_i) \varepsilon$$

And SVR is expressed:

$$f(\vec{x}) = \sum_{i=1}^n (\alpha_i^* - \beta_i^*) k(x_i, x_j) + b^*$$

where b^* is computed as above.



d. If x_i satisfies: $y_i - w^{*T} x_i - b = \varepsilon + \xi_i^+$

or $w^{*T} x_i + b - y_i = \varepsilon + \xi_i^-$,

we call x_i a support vector.

Therefore, if x_i is not a support vector, then $\alpha_i^* = 0$ or $\beta_i^* = 0$.

Primal:
$$\min_{\tilde{w}, b, \xi_i^+, \xi_i^-} \frac{1}{2} \|w\|_2^2 + \frac{C}{n} \sum_{i=1}^n (\xi_i^+ + \xi_i^-)$$

s.t. $y_i - w^T x_i - b \leq \varepsilon + \xi_i^+ \quad \forall i$

$w^T x_i + b - y_i \leq \varepsilon + \xi_i^- \quad \forall i$

$\xi_i^+ \geq 0 \quad \forall i$

$\xi_i^- \geq 0 \quad \forall i$

There are 5 cases:

① if $\left. \begin{aligned} y_i - w^{*T} x_i - b^* < \varepsilon \\ \text{and } w^{*T} x_i + b^* - y_i < \varepsilon \end{aligned} \right\}$, then $\left. \begin{aligned} \xi_i^+ &= 0 \\ \text{and } \xi_i^- &= 0 \end{aligned} \right\}$ and x_i is not a support vector
 $\xi_i^{+*} = 0, \xi_i^{-*} = 0$.

② if $y_i - w^{*T} x_i - b^* = \varepsilon$, then $w^{*T} x_i + b^* - y_i \leq 0 \leq \varepsilon$, $\xi_i^+ = 0$, and x_i is a support vector

③ if $w^{*T} x_i + b^* - y_i = \varepsilon$, then $y_i - w^{*T} x_i - b^* \leq 0 \leq \varepsilon$, $\xi_i^- = 0$, $\xi_i^+ = 0$ and x_i is a support vector

④ if $y_i - w^{*T} x_i - b^* > \varepsilon$, then $w^{*T} x_i + b^* - y_i \leq 0 \leq \varepsilon$, $\xi_i^+ > 0$, $\xi_i^- = 0$, and x_i is a support vector

⑤ if $w^{*T} x_i + b^* - y_i > \varepsilon$, then $y_i - w^{*T} x_i - b^* \leq 0 \leq \varepsilon$, $\xi_i^- > 0$, $\xi_i^+ = 0$ and x_i is a support vector

Since $w^* = \sum_{i=1}^n (\alpha_i^* - \beta_i^*) x_i$, if x_i is not a support vector, $\alpha_i^* = \beta_i^* = 0$,

there is no contribution to compute w^* for non-support vector. Also b^* is decided by a support vector,

Thus the final predictor only depend on x_i if it's a support vector, they are a subset of training examples.



EECS 545 Homework 4

Yuan Yin

November 1, 2018

Problem 1

(c)

The result is:

```
MSE for training data is: 19.016860651808518
MSE for test data is: 23.535860253927392
offset b is: 19.429475366627464
```

(d)

The result is:

```
MSE for training data is: 19.249564400341065
MSE for test data is: 26.396905605851654
```

Finding that for kernel ridge regression without offset, the errors both for training and test data are a little larger than kernel ridge regression with offset.

(e)

The code for (c) is:

```
import numpy as np
import scipy.io as sio

# Load the data and preprocessed the data
bodyfat_data = sio.loadmat('bodyfat_data.mat')
x = bodyfat_data['X']
y = bodyfat_data['y']
n,d = x.shape
train_num = 150; m = n - train_num
x_train = x[:train_num,:]; x_test = x[train_num:,:]
```



```

y_train = y[:train_num]; y_test = y[train_num:]
y_train_bar = np.mean(y_train)

# Helper function
def dist2(x,c):
    ndata,dimx = x.shape
    ncenters, dimc = c.shape

    xsum = np.sum(x**2,axis = 1)
    xsum = xsum[:,np.newaxis]
    csum = np.sum(c**2,axis = 1)
    csum = csum[:,np.newaxis]

    n2 = xsum.dot(np.ones([1,ncenters]))+ np.ones([ndata,1]).dot(csum.T)- 2*x.dot(c.T)
    return n2

sigma = 15; lamda = 0.003
kernel = np.exp(-1/2/sigma**2 * dist2(x_train, x_train))
kernel_prime = np.exp(-1/2/sigma**2 * dist2(x_train, x_test))
O1 = 1/train_num * np.ones([train_num, m]); O2 = 1/train_num * np.ones([train_num,
    train_num])
kernel_tilda = kernel - kernel.dot(O2) - O2.dot(kernel) + (O2.dot(kernel)).dot(O2)
kernel_prime_tilda = kernel_prime - kernel.dot(O1) - O2.dot(kernel_prime) +
    (O2.dot(kernel)).dot(O1)

# Kernel ridge regression with offset
y_pre_train = []
for i in range(train_num):
    y_pre_train.append(y_train_bar + (np.mat(y_train - y_train_bar).T
        * np.mat(kernel_tilda + train_num * lamda *
            np.eye(train_num)).I
        * np.mat(kernel_tilda[:,i]).T)[0,0])

y_pre_test = []
for i in range(m):
    y_pre_test.append(y_train_bar + (np.mat(y_train - y_train_bar).T
        * np.mat(kernel_tilda + train_num * lamda *
            np.eye(train_num)).I
        * np.mat(kernel_prime_tilda[:,i]).T)[0,0])

# Compute MSE
MSE_train = sum([(y_train[i] - y_pre_train[i])**2 for i in range(train_num)]) / train_num
MSE_test = sum([(y_test[i] - y_pre_test[i])**2 for i in range(n-train_num)]) / (n-train_num)
print("MSE for training data is: ", MSE_train[0])
print("MSE for test data is: ", MSE_test[0])

```

```

# Compute offset b
b = y_train_bar - (np.mat(y_train - y_train_bar).T
                    * np.mat(kernel_tilda + train_num * lamda * np.eye(train_num)).I
                    * np.mat(kernel.dot(O2)[: ,0] - ((O2.dot(kernel)).dot(O2))[0,0]).T)[0,0]
print("offset b is: ", b)

```

And the code for (d) is:

```

import numpy as np
import scipy.io as sio

# Load the data and preprocessed the data
bodyfat_data = sio.loadmat('bodyfat_data.mat')
x = bodyfat_data['X']
y = bodyfat_data['y']
n,d = x.shape
train_num = 150; m = n - train_num
x_train = x[:train_num,:]; x_test = x[train_num:,:]
y_train = y[:train_num]; y_test = y[train_num:]

# Helper function
def dist2(x,c):
    ndata,dimx = x.shape
    ncenters, dimc = c.shape

    xsum = np.sum(x**2,axis = 1)
    xsum = xsum[:,np.newaxis]
    csum = np.sum(c**2,axis = 1)
    csum = csum[:,np.newaxis]

    n2 = xsum.dot(np.ones([1,ncenters]))+ np.ones([ndata,1]).dot(csum.T)- 2*x.dot(c.T)
    return n2

sigma = 15; lamda = 0.003
kernel = np.exp(-1/2/sigma**2 * dist2(x_train, x_train))
kernel_prime = np.exp(-1/2/sigma**2 * dist2(x_train, x_test))

# Kernel ridge regression without offset
y_pre_train = []
for i in range(train_num):
    y_pre_train.append(np.array(np.mat(kernel + train_num * lamda * np.eye(train_num)).I *
                                np.mat(y_train)).T.dot(kernel[i,:].T))

```

```

y_pre_test = []
for i in range(m):
    y_pre_test.append(np.array(np.mat(kernel + train_num * lamda * np.eye(train_num)).I *
                                np.mat(y_train)).T.dot(kernel_prime[:,i]))

# Compute MSE
MSE_train = sum([(y_train[i] - y_pre_train[i])**2 for i in range(train_num)])/train_num
MSE_test = sum([(y_test[i] - y_pre_test[i])**2 for i in range(n-train_num)])/(n-train_num)
print("MSE for training data is: ", MSE_train[0])
print("MSE for test data is: ", MSE_test[0])

```
