

Problem 1: Kernels

$$1.1 \quad k_1(x, x') = \phi_1(x)^T \phi_1(x'); \quad k_2(x, x') = \phi_2(x)^T \phi_2(x')$$

either ① design a feature map ϕ using ϕ_1, ϕ_2 such that $k(x, x') = \phi(x)^T \phi(x')$
 or ② explain why it cannot exists.

$$\text{Let } k_1: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R} \quad ; \quad k_2: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R} \\ (\phi_1(x)) \quad (\phi_2(x)) \\ (\phi_1(x')) \quad (\phi_2(x'))$$

k_1 & k_2 are positive semi-definite

$$a) \quad k(x, x') = c k_1(x, x') \text{ for any } c \geq 0.$$

Since k_1 is kernel.

$\Rightarrow K_1 = \phi_1(x) \cdot \phi_1(x)^T$ would be corresponding kernel matrix for $k_1(x, x')$

$$\text{where } (K_1)_{ij} = \phi_1(x_i)^T \phi_1(x_j)$$

Similarly,

$K_2 = \phi_2(x) \cdot \phi_2(x)^T$ would be corresponding to kernel matrix for
 $k_2(x, x')$, where $(K_2)_{ij} = \phi_2(x_i)^T \phi_2(x_j)$.

For any vector
 $v \in \mathbb{R}^n$

$$\Rightarrow K = CK_1$$

$$\Rightarrow V^T K V = V^T (C K_1) V = C V^T K_1 V \geq 0$$

Hence we showed that K is positive semi-definite matrix

Hence K is a kernel matrix $\Rightarrow k(x, x')$ is kernel function.



扫描全能王 创建

$$b) K(x, x') = k_1(x, x') + k_2(x, x')$$

As the set up above,

$$K_1 = \phi_1(x) \phi_1^T(x)$$

$$K_2 = \phi_2(x) \phi_2^T(x)$$

$$K = K_1 + K_2$$

$$V^T K V = V^T (K_1 + K_2) V = (V^T K_1 V + V^T K_2 V)$$

Since K_1 & K_2 are PSD

for any vector $V \in \mathbb{R}^n$

$$\Rightarrow V^T K_1 V \geq 0 \text{ & } V^T K_2 V \geq 0$$

$$\Rightarrow V^T K V = V^T K_1 V + V^T K_2 V \geq 0$$

$\Rightarrow K'$ is kernel matrix

$\Rightarrow K(x, x')$ is kernel function.

$$c) K(x, x') = k_1(x, x') - k_2(x, x')$$

for any vector $V \in \mathbb{R}^n$

$$K' = K_1 - K_2 \quad V^T K_1 V \geq 0 \text{ & } V^T K_2 V \geq 0$$

$$V^T K' V = \underbrace{V^T K_1 V}_{\geq 0} - \underbrace{V^T K_2 V}_{\geq 0}$$

this is not guaranteed to be ≥ 0

Hence, K' might not be kernel but if $V^T K_1 V \geq V^T K_2 V$, then K' could be kernel.



扫描全能王 创建

$$(d) \cdot k(x, x') = k_1(x, x') \cdot k_2(x, x') = \phi_1^T(x) \phi_1(x) \phi_2^T(x) \phi_2(x)$$

we can construct $\phi_k(x) = \phi_1(x) \otimes \phi_2(x)$.

$$\Rightarrow k(x, x') = \phi_k^T(x) \phi_k(x') = (\phi_1(x) \otimes \phi_2(x))^T (\phi_1(x') \otimes \phi_2(x'))$$

where \otimes denotes the tensor product.

Hence, it is a valid kernel.

$$(e) k(x, x') = k(f(x), f(x')) = \phi_f(f(x))^T \phi_f(f(x'))$$

$$\text{let } \phi_k(x) = \phi_f(f(x))$$

$$\text{Hence } k(x, x') = \phi_k^T(x) \phi_k(x') = \phi_f(f(x))^T \phi_f(f(x'))$$

$\Rightarrow k(x, x')$ is a valid kernel.

1.2 Show that $k'(x, x') = \sum_{i=1}^d 2_i k(x, x')$ is a valid kernel for any valid kernel k if $2_i \geq 0$ for all $i \in \{0, 1\}$

$$k'(x, x') = \sum_{i=1}^d 2_i k(x, x')^i = 2_1 k(x, x')^1 + 2_2 k(x, x')^2 + \dots$$

according to a) $k' = C \circ k$ is still valid kernel

d) $k' = k_1 \cdot k_2$ is still valid, hence

$$k(x, x')^2 = k(x, x') \cdot k(x, x') \xrightarrow{\text{kernel}} \boxed{\text{all are valid kernel}}$$

$$k(x, x')^3 = k(x, x')^2 \cdot k(x, x') \xrightarrow{\text{kernel}} \boxed{\text{all are valid kernel}}$$

$$k(x, x')^d = k(x, x)^{d-1} \cdot k(x, x) \xrightarrow{\text{kernel}}$$



扫描全能王 创建

Then according to b), the sum of valid kernels are still kernel

$\Rightarrow K(x, x') = \sum_{i=1}^n d_i k(x, x')$ is a valid kernel function for all $x, x' \in \mathbb{R}^2$.

Bonus: Show that $k'(x, x') = \exp(\min(x, x'))$ is a valid kernel over input space $\mathcal{X} = [0, 1] \subset \mathbb{R}$.

$$K_{ij} = k'(x_i, x_j) = \exp(\min(x_i, x_j))$$

For any vector $v \in \mathbb{R}^n$,

$$\begin{aligned} v^T K v &= \begin{bmatrix} v_1 \\ \vdots \\ v_n \end{bmatrix}^T K \begin{bmatrix} v_1 \\ \vdots \\ v_n \end{bmatrix} = [v_1 \dots v_n] \begin{bmatrix} K_{11} & K_{12} & \cdots & K_{1n} \\ \vdots & \ddots & & \vdots \\ K_{n1} & \cdots & \ddots & K_{nn} \end{bmatrix} \begin{bmatrix} v_1 \\ \vdots \\ v_n \end{bmatrix} \\ &= \sum_{i=1}^n \sum_{j=1}^n v_i v_j \exp(\min(x_i, x_j)) \end{aligned}$$

Since the exponential function is always positive

$$\Rightarrow f(x) = e^x \geq 0 \text{ for all } x \in \mathbb{R}$$

and since $x \in [0, 1] \Rightarrow \min\{x_i, x_j\} \in [0, 1]$ for all $x_i, x_j \in [0, 1]$.

$\Rightarrow \exp(\min(x_i, x_j)) \geq 0 \Rightarrow$ the summation of positive terms is still positive.

$$\Rightarrow v^T K v \geq 0 \text{ for any } v \in \mathbb{R}^n$$

$\Rightarrow K'$ is kernel.



扫描全能王 创建

Bohm's continuity.
 $K'(x, x') = \exp(-\min(x, x'))$. K' is the kernel matrix for $K(x, x')$.

We just need to show $\exp(K_1)$ is kernel and
 $\exp(\min(x, x'))$ is kernel.

Let $c_1, \dots, c_n \in C$ and $x_1, \dots, x_n \in R$. Let $k(x, x') = \min(x, x')$.

$$\begin{aligned} \sum_{i,j=1}^n c_i c_j K(x_i, x_j) &= \sum_{i,j=1}^n c_i c_j \mathbb{1}_{[S \leq x_i]} \mathbb{1}_{[S \leq x_j]} ds \\ &= \int_0^\infty \left| \sum_{i=1}^n c_i \mathbb{1}_{[S \leq x_i]} \right|^2 ds \geq 0. \end{aligned}$$

which means $C^T K_C \geq 0$ is indeed kernel.
And K_1 .

Then we just need to show $K = \exp(K_1)$ is kernel.

$$\begin{aligned} \exp(K_1) &= \exp(0) + \exp(0)K_1 + \frac{\exp(0)}{2!} K_1^2 + \dots + \\ &= 1 + K_1 + \frac{1}{2} K_1^2 + \dots \end{aligned}$$

According to 1.4, $K = \exp(K_1)$ is a kernel.

Hence, we showed that $K(x, x') = \exp(-\min(x, x'))$

is kernel.



扫描全能王 创建

Problem 2. Gradient Descent

$$\min_w \frac{1}{m} \sum_{i=1}^m \log(1 + \exp(-y_i w^T x_i)) + R(w)$$

Show properties about when this objective is convex, L-smooth, or possibly M-strongly convex

$$2.1 R(w) = 0$$

$$\Rightarrow F(w) = \frac{1}{m} \sum_{i=1}^m \log(1 + \exp(-y_i w^T x_i))$$

$$WTS: \nabla F(w) \geq M I$$

$$\text{Prof: } \nabla F(w) = \frac{1}{m} \sum_{i=1}^m \frac{\partial (\log(1 + \exp(-y_i w^T x_i)))}{\partial w}$$

$$\nabla^2 F(w) = \frac{1}{m} \sum_{i=1}^m \frac{-y_i x_i e^{-y_i w^T x_i}}{1 + e^{-y_i w^T x_i}}$$

$$\nabla^2 F(w) = \frac{1}{m} \sum_{i=1}^m \frac{y_i^2 x_i^T x_i e^{-y_i w^T x_i}}{(1 + e^{-y_i w^T x_i})^2}$$

$$\text{Since } y_i \in \{+1, -1\} \Rightarrow y_i^2 = 1$$

$$\Rightarrow \nabla^2 F(w) = \frac{1}{m} \sum_{i=1}^m \frac{x_i^T x_i e^{-(y_i w^T x_i)}}{(1 + e^{-y_i w^T x_i})^2}$$

$$x_i^T x_i = \|x_i\|_2^2 \leq 1$$

$$0 \leq \nabla^2 F(w) = \frac{1}{m} \sum_{i=1}^m \frac{x_i^T x_i e^{-(y_i w^T x_i)}}{(1 + e^{-y_i w^T x_i})^2}$$

since sigmoid function
 $\sigma(x) = \frac{1}{1 + e^{-x}}$ has
 codomain (0, 1)

$$\leq \frac{1}{m} \sum_{i=1}^m \frac{e^{-(y_i w^T x_i)}}{(1 + e^{-y_i w^T x_i})^2}$$

$$\leq \frac{1}{m} \sum_{i=1}^m \frac{1}{(1 + e^{-y_i w^T x_i})^2}$$

This is sigmoid.

$$S = \frac{1}{m} \cdot m \cdot 0 \leq \left(-\frac{1}{m} \sum_{i=1}^m \boxed{\frac{1}{1 + e^{-(y_i w^T x_i)}}} \right) \leq \frac{1}{m} \cdot m I = I$$



However, we can't find its lower bound but only $\nabla^2 F(w) \succ 0$.
it is

\Rightarrow We cannot say it is strong convex but only it is convex.
i.e. there doesn't exist a M s.t. $\nabla^2 F(w) \geq M I$.

2.2. In order to show it's 1-smooth.

We have gotten the $\nabla^2 F(w)$ previous part:

$$\nabla^2 F(w) = \frac{1}{m} \sum_{i=1}^m \frac{1}{1+e^{-y_i w^T x_i}} I_m$$

As we see, $\frac{1}{1+e^{-y_i w^T x_i}}$ is a squashing (Sigmoid) Function.

$$\sigma(w) = \frac{1}{1+e^{-y_i w^T x_i}} \in (0, 1)$$

Hence $\nabla F(w) = \frac{1}{m} \sum_{i=1}^m \frac{1}{1+e^{-y_i w^T x_i}} \in (\frac{1}{m} \cdot m \cdot 0, \frac{1}{m} \cdot m \cdot I)$
 $\Rightarrow (0, I)$.

$$\Rightarrow \nabla^2 F(w) \leq \begin{bmatrix} 1 & \\ & 1 \end{bmatrix} = L I = 1 \cdot I.$$

which means $L \leq 1$.

$\Rightarrow F(w)$ is 1-smooth function.



扫描全能王 创建

2.3 Convergence rate of GD when $R(w) = 0$

$$\|W_{T+1} - w^*\|_2 \leq \varepsilon$$

$$F(W_{T+1}) - F(W^*) \leq \frac{L}{2T}$$

$$\frac{L}{2T} = \varepsilon \Rightarrow \frac{1}{2T} = \frac{\varepsilon}{L}$$

$$T = \frac{1}{2\varepsilon} \sim O\left(\frac{1}{2\varepsilon}\right)$$

$$\Rightarrow \text{convergence rate } \boxed{O\left(\frac{1}{2T}\right)}$$

$$2.4. R(w) = \sum_{j=1}^d \lambda_j w_j^2 \quad (\text{regularizer})$$

Show that the objective with $R(w)$ is μ -strongly convex and L -smooth

$$\text{for } \mu = 2 \min_{j \in [d]} \lambda_j \text{ and } L = 1 + 2 \max_{j \in [d]} \lambda_j$$

$$\Rightarrow F(w) = \underbrace{\frac{1}{m} \sum_{i=1}^m \log(1 + e^{-y_i w^T x_i})}_{F_0(w)} + \sum_{j=1}^d \lambda_j w_j^2$$

$$\textcircled{1} \quad \nabla^2 F(w) = \nabla^2 F_0(w) + \nabla^2 \left(\sum_{j=1}^d \lambda_j w_j^2 \right)$$

$$\nabla^2 \left(\sum_{j=1}^d \lambda_j w_j^2 \right) = \nabla^2 \left(\sum_{j=1}^d 2\lambda_j w_j \right) = \sum_{j=1}^d 2\lambda_j \quad \text{sigmoid function}$$

$$\Rightarrow \nabla^2 F(w) = \frac{1}{m} \sum_{i=1}^m \frac{1}{1 + e^{-y_i w^T x_i}} + \sum_{j=1}^d 2\lambda_j$$

$$= \frac{1}{m} \sum_{i=1}^m \underbrace{\frac{1}{1 + e^{-y_i w^T x_i}}}_{\text{Sigmoid}} + 2 \sum_{j=1}^d \lambda_j \in (0 + 2 \sum_{j=1}^d \lambda_j, 1 + 2 \sum_{j=1}^d \lambda_j)$$

let's consider the Hessian matrix for $\sum_{j=1}^d \lambda_j w_j^2$

$$\text{upperbd of } \nabla^2 F_0(w) = 2I + H = \begin{bmatrix} 2I & H \\ H & 2I \end{bmatrix} \leq \begin{bmatrix} 2\lambda_{\max} I & \\ & 2\lambda_{\max} I \end{bmatrix} = L \cdot I \Rightarrow L = 1 + 2 \max_{j \in [d]} \lambda_j$$

$$2 \min_{j \in [d]} \lambda_j \leq \nabla^2 F(w) \leq 1 + \sum_{j=1}^d 2\lambda_j \leq 1 + 2 \max_{j \in [d]} \lambda_j \Rightarrow \text{It's } L = (2 \max_{j \in [d]} \lambda_j + 1) \text{-smooth function.}$$



扫描全能王 创建

and

$$\nabla^2 F(w) \geq (2 \min_{j \in [d]} \lambda_j) I - M I$$

Lower bound
 $\nabla^2 F(w) \geq \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_d \end{bmatrix} \geq \begin{bmatrix} \alpha_{\min} & & \\ & \ddots & \\ & & \alpha_{\min} \end{bmatrix}$

where $M = 2 \max_{j \in [d]} \lambda_j \Rightarrow$ it's $(2 \min_{j \in [d]} \lambda_j)$ -strongly convex function.

Hence we proved $F(w)$ with $\|w\|_2$ is both L -smooth &

With $M = 2 \min_{j \in [d]} \lambda_j$ & $L = 2 \max_{j \in [d]} \lambda_j + 1$ M -strongly convex.

Bonus.

$$\Rightarrow \|w_T - w^*\|_2^2 \leq (1 - \frac{M}{L})^T \|w_1 - w^*\|_2^2$$

$$F(w) = \frac{1}{m} \sum_{i=1}^m \log(1 + \exp(-y_i w^T x_i)) + \frac{\alpha}{2} \sum_{j=1}^d \lambda_j w_j^2$$

$$(1 - \frac{M}{L})^T \|w_1 - w^*\|_2^2$$

Assume $\|w^*\|_2 = 1$
& $w_1 = 0$.

$$= (1 - \frac{M}{L})^T \|w_1 - w^*\|_2^2$$

$$1 - \frac{M}{L} = \log \varepsilon$$

$$\Rightarrow T = \frac{L}{m} \log(\frac{1}{\varepsilon})$$

$$M = 2 \min_{j \in [d]} \lambda_j$$

$$L = 2 \max_{j \in [d]} \lambda_j + 1$$

$$= \frac{2 \max_{j \in [d]} \lambda_j + 1}{2 \min_{j \in [d]} \lambda_j} \log(\frac{1}{\varepsilon})$$

$$= \frac{\max_{j \in [d]} \lambda_j + \frac{1}{2}}{\min_{j \in [d]} \lambda_j} \log(\frac{1}{\varepsilon}).$$

$$\frac{M}{L} = \frac{\min_{j \in [d]} \lambda_j}{\max_{j \in [d]} \lambda_j + 1}$$

$$1 - \frac{M}{L} = \frac{\frac{1}{2} \max_{j \in [d]} \lambda_j - \min_{j \in [d]} \lambda_j}{\max_{j \in [d]} \lambda_j + 1}$$

$$\Rightarrow T \sim O\left(\frac{\max_{j \in [d]} \lambda_j + \frac{1}{2}}{\min_{j \in [d]} \lambda_j} \log(\frac{1}{\varepsilon})\right)$$

$$\sim O\left(\log \frac{1}{\varepsilon}\right)$$

$\log \frac{1}{\varepsilon} \rightarrow$ convergence rate

$O(\log \frac{1}{\varepsilon})$



扫描全能王 创建

contingy Bows for λ

$$T = \frac{\frac{1}{z} + \max \lambda_j}{\min \lambda_j} \log(\frac{1}{z})$$

$$\log z = -T \frac{\min \lambda_j}{\frac{1}{z} + \max \lambda_j}$$

$$z = e^{-T \left(\frac{\min \lambda_j}{\frac{1}{z} + \max \lambda_j} \right)} = e^{-\left(\frac{\min \lambda_j}{\max \lambda_j} \right) T} \sim O(e^{-\left(\frac{\min \lambda_j}{\max \lambda_j} \right) T})$$

$$\Rightarrow z \sim O(e^{-\left[\left(\frac{\min \lambda_j}{\frac{1}{z} + \max \lambda_j} \right) T \right]})$$

$$T \sim O\left(\frac{\max \lambda_j + 0.5}{\min \lambda_j} \cdot \log(\frac{1}{z})\right)$$



扫描全能王 创建