

ML HW3

SNIR HORDAN 205689581, ALEXANDER SHENDER 328626114

1. Let $\hat{w} = V(\Sigma^+)^2 \Sigma^T U^T y$ where $\Sigma \in \mathbb{R}^{m \times d}$ with the singular values derived from the SVD, as the entries on the diagonal, and U, V square orthogonal matrices of order $m \times m$ and $d \times d$, respectively.

Notice that

$$\begin{aligned} X\hat{w} &= U\Sigma \underbrace{V^T V}_{=I_d} (\Sigma^+)^2 \Sigma^T U^T y = \\ &= U \underbrace{\Sigma(\Sigma^+)^2 \Sigma^T}_{\text{see (1)}} U^T y = UU^T y = y \end{aligned}$$

$$(1) \Sigma = \begin{pmatrix} \sigma_1 & 0 & 0 & 0 & 0 & 0 \\ 0 & \sigma_2 & 0 & 0 & 0 & 0 \\ 0 & 0 & . & 0 & 0 & 0 \\ 0 & 0 & 0 & \sigma_r & 0 & 0 \end{pmatrix} \text{ and } \Sigma^+ = \begin{pmatrix} \frac{1}{\sigma_1} & 0 & 0 & 0 \\ 0 & . & 0 & 0 \\ 0 & 0 & . & 0 \\ 0 & 0 & 0 & \frac{1}{\sigma_r} \end{pmatrix} \text{ then } \Sigma \Sigma^+ = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & . & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \end{pmatrix} \text{ Then } \Sigma(\Sigma^+)^2 \Sigma^T = I_m$$

Then $\|X\hat{w} - y\|_2^2 = 0$

The metric $\|\cdot\|_2 : \mathbb{R}^d \rightarrow \mathbb{R}^+$ is non-negative then $\hat{w} \in \argmin_w \|Xw - y\|_2^2$

2. 1.1.

$$\text{Define } \phi_3 : \mathbb{R}^d \rightarrow \mathbb{R}^{2d}, \phi_3(\vec{x}) = \begin{pmatrix} \phi_1(\vec{x}) \\ \vdots \\ \phi_2(\vec{x}) \end{pmatrix}$$

Define $K_3 : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R} \quad K_3(x, x') = \langle \phi_3(x), \phi_3(x') \rangle$

Then

$$\begin{aligned} K_3(x, x') &= \langle \phi_3(x), \phi_3(x') \rangle = \left\langle \begin{pmatrix} \phi_1(x) \\ \vdots \\ \phi_2(x) \end{pmatrix}, \begin{pmatrix} \phi_1(x') \\ \vdots \\ \phi_2(x') \end{pmatrix} \right\rangle = \left\langle \begin{pmatrix} \phi_1(x) \\ \vdots \end{pmatrix}, \begin{pmatrix} \phi_1(x') \\ \vdots \end{pmatrix} \right\rangle + \left\langle \begin{pmatrix} \phi_2(x) \\ \vdots \end{pmatrix}, \begin{pmatrix} \phi_2(x') \\ \vdots \end{pmatrix} \right\rangle = \\ &= K_1(x, x') + K_2(x, x') \end{aligned}$$

K_3 is a valid kernel function iff its Gram matrix is Positive Semi-Definite.

Let $0_m \neq z \in \mathbb{R}^m$ and $x_1, \dots, x_m \in \mathbb{R}^d$

Then $z^T G_{K_3} z = z^T (G_{K_1} + G_{K_2}) z = z^T G_{K_1} z + z^T G_{K_2} z \geq 0$, because K_1, K_2 are valid kernel functions.

Therefore K_3 is a valid kernel function.

2.1.2

Define $\phi_4 : \mathbb{R} \rightarrow \mathbb{R}$, $\phi_4(\vec{x}) = f(\vec{x}) \phi_1(\vec{x})$ and let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be a function.

Define $K_4 : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R} \quad K_4(x, x') = \langle \phi_4(x), \phi_4(x') \rangle$

Then $K_4(x, x') = \langle \phi_4(x), \phi_4(x') \rangle = \langle f(x) \phi_1(x), f(x') \phi_1(x') \rangle = f(x) f(x') \langle \phi_1(x), \phi_1(x') \rangle = f(x) f(x') K_1(x, x')$

K_4 is a valid kernel function iff its Gram matrix is Positive Semi-Definite.

Let $0_m \neq z \in \mathbb{R}^m$ and $x_1, \dots, x_m \in \mathbb{R}^d$

Then $z^T G_{K_4} z = z^T \text{diag}(f(x_1), \dots, f(x_m))^T G_{K_1} \text{diag}(f(x_1), \dots, f(x_m)) z \geq 0$ because $\text{diag}(f(x_1), \dots, f(x_m)) z \in \mathbb{R}$ and G_{K_1} is PSD because K_1 is a valid kernel function.

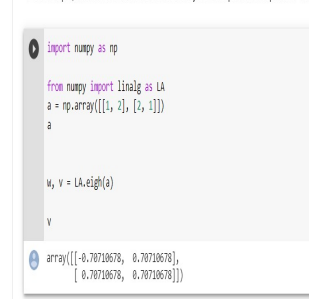
Therefore K_4 is a valid kernel function.

2.2.1

$$\text{Observe } \begin{pmatrix} K(x_1, x_1) & K(x_1, x_2) \\ K(x_2, x_1) & K(x_2, x_2) \end{pmatrix} = \begin{pmatrix} 1 & 2 \\ 2 & 1 \end{pmatrix}$$

$$\left| \begin{pmatrix} 1-\lambda & 2 \\ 2 & 1-\lambda \end{pmatrix} \right| = (1-\lambda)^2 - 4 = -3 - 2\lambda + \lambda^2 = (\lambda-3)(\lambda+1)$$

Then eigenvalues are 3, -1.



```
import numpy as np
from numpy import linalg as LA
a = np.array([[1, 2], [2, 1]])
a

w, v = LA.eigh(a)
v
array([[0.70710678, 0.70710678],
       [0.70710678, 0.70710678]])
```

Above are the eigenvectors of norm 1.

2.2.2

Define $\hat{v} = \begin{pmatrix} -1 \\ 1 \end{pmatrix}$. v is an eigenvector of the Gram matrix. Note, $y_1 = 1, y_2 = -1$

Let $\lambda \geq 0$.

$$L_{SVM}(c\hat{v}) = \lambda \begin{pmatrix} -c & c \end{pmatrix} \begin{pmatrix} 1 & 2 \\ 2 & 1 \end{pmatrix} \begin{pmatrix} -c \\ c \end{pmatrix} + \frac{1}{2} \sum_{i=1}^2 \max\{0, 1 - y_i (cGv)_i\}$$

$$= -\lambda 2c^2 + \frac{1}{2} [\max\{0, 1 - 1(+c)\} + \max\{0, 1 - (-1)(-c)\}] = \dots$$

$$\dots = -\lambda 2c^2 + \frac{1}{2} [2 + 2c] = -\lambda c^2 + 1 - c \xrightarrow{c \rightarrow \infty} -\infty$$

Non-PSD kernels are “problematic” because there exists an eigenvector with a negative eigenvalue. Then, exists a contradiction to the existence of a minimal value over the set $\{L_{SVM}(\alpha) \mid \alpha \in \mathbb{R}^m\}$.

2.2.3. Assume $G \geq 0$ and let $\lambda \geq 0$.

2.2.3.1. Note $G \geq 0$ if $\forall z \neq 0_m, z^T G z \geq 0$

$$\frac{1}{m} \sum_{i=1}^m \max\{0, 1 - y_i (G\alpha)_i\} \geq 0 \text{ by definition of the hinge-loss function.}$$

$\lambda \alpha^T G \alpha \geq 0$ because G is PSD.

$$\text{Then } L_{SVM}(\alpha) = \min_{\alpha} \lambda \alpha^T G \alpha + \frac{1}{m} \sum_{i=1}^m \max\{0, 1 - y_i (G\alpha)_i\}$$

2.2.3.2.

By theorem optimal solution $\hat{w} = \sum_{i=1}^m \alpha_i x_i$.

Let $w = \sum_{i=1}^m \alpha_i x_i$ be a separating hyperplane.

By the Representer Theorem, if w is an optimal solution to the Soft-SVM optimization problem then $y_i \langle w, x_i \rangle = \sum_{j=1}^m \alpha_j \langle x_i, x_j \rangle \geq 1$ Then $\alpha_i = 0$.

Then we can assume above claim holds for w , because we are attempting to minimize the loss function over $\alpha \in \mathbb{R}^m$.

Therefore,

$$L_{SVM}(\alpha) = \lambda \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j \langle x_i, x_j \rangle + \frac{1}{m} \sum_{i=1}^m \max\left\{0, 1 - \sum_{j=1}^m \alpha_j \langle x_i, x_j \rangle\right\}$$

$$\stackrel{\text{Representor Thm.}}{=} \lambda \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j \langle x_i, x_j \rangle + 1 - \frac{1}{m} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j \langle x_i, x_j \rangle = 1 + \left(\lambda - \frac{1}{m} \right) \alpha^T G \alpha$$

Note $G \geq 0$ if $\forall z \neq 0_m, z^T G z \geq 0$

Then $L_{SVM}(\alpha) = 1 + \left(\lambda - \frac{1}{m} \right) \alpha^T G \alpha \implies \min_{\alpha} L_{SVM}(\alpha) \leq 1$ because $\|\alpha\|_2$ can be arbitrarily small.