

HW 3

Sherin Karuvallil Saji

q1 + q3 answers are
inside a ML-HW3-q1andq3 doc
code for q3 is inside
q3.ipynb. (just click
on run all)

Answers for q2, q4, q5
are inside this
handwritten pdf

2)

kernels over $\mathbb{R}^n \times \mathbb{R}^n$

$$1. K(x, z) = k_1(x, z) k_2(x, z)$$

To prove that $K(x, z)$ is also a kernel, we need to show that the kernel matrix is positive semi-definite and that $K(x, z)$ is symmetric.

Let's denote kernel matrix corresponding to k_1 as M_1 and kernel matrix corresponding to k_2 as M_2 .

kernel matrix M , corresponding to $K(x, z)$ is given by the element-wise multiplication of M_1 and M_2

$$M = M_1 \odot M_2$$

Also for any vector v :

$$v^T M v = v^T (M_1 \odot M_2) v = (v^T M_1 v) * (v^T M_2 v)$$

Since M_1 and M_2 are positive semi-definite,

$$v^T M_1 v \geq 0 \text{ and}$$

$$v^T M_2 v \geq 0 \text{ for any vector } v$$

\therefore their product

$$(v^T M_1 v) * (v^T M_2 v)$$

is also non-negative.

$\therefore M$ is positive semi-definite

To show that $K(x, z)$ is symmetric:

We need to show that

$$K(x, z) = K(z, x)$$

$$K(x, z) = k_1(x, z) k_2(x, z)$$

$$K(z, x) = k_1(z, x) k_2(z, x)$$

$$= k_1(x, z) k_2(x, z)$$

$$= K(x, z)$$

because k_1 and k_2 are kernels and definitely symmetric meaning

$$\therefore K(x, z) = K(z, x)$$

$$k_1(x, z) = k_1(z, x)$$

$$k_2(x, z) = k_2(z, x)$$

$\therefore K$ is valid kernel

q2)
 2. $k(x, z) = a k_1(x, z) + b k_2(x, z)$ where $a, b > 0$

for kernel $k(x, z)$ let the corresponding kernel matrix be

$k(x, z)$	M
$k_1(x, z)$	M_1
$k_2(x, z)$	M_2

To prove that $k(x, z)$ is also a kernel, we need to show that the kernel matrix is positive semi-definite and that $k(x, z)$ is symmetric.

For any vector v :

$$v^T M v = v^T (a * M_1 + b * M_2) v$$

$$= a * (v^T M_1 v) + b * (v^T M_2 v)$$

Since M_1 & M_2 are positive semi-definite,

$$v^T M_1 v \geq 0 \text{ \& } v^T M_2 v \geq 0 \text{ for any vector } v$$

As a & b are positive real numbers,

the combination of $a * (v^T M_1 v) + b * (v^T M_2 v)$ will also be non-negative.

Proving $k(x, z) = k(z, x)$:

$$k(z, x) = a k_1(z, x) + b k_2(z, x)$$

$$= a k_1(x, z) + b k_2(x, z)$$

$$= k(x, z)$$

k_1 and k_2 are kernels and definitely symmetric meaning

$$k_1(x, z) = k_1(z, x)$$

$$k_2(x, z) = k_2(z, x)$$

Thus, $k(x, z) = a k_1(x, z) + b k_2(x, z)$ satisfies Mercer's theorem and is a valid kernel.

3. $k(x, z) = aK_1(x, z) - bK_2(x, z)$ where $a, b > 0$

To prove that $K(x, z)$ is also a kernel, we need to show that the kernel matrix is positive semi-definite and that $K(x, z)$ is symmetric.

for kernel $K(x, z)$ let the corresponding kernel matrix be

$K(x, z)$	M
$K_1(x, z)$	M_1
$K_2(x, z)$	M_2

For any vector v :

$$v^T M v = v^T (a * M_1 - b * M_2) v$$

$$= a * (v^T M_1 v) - b * (v^T M_2 v)$$

Since M_1 & M_2 are positive semi-definite,

$$v^T M_1 v \geq 0 \text{ \& } v^T M_2 v \geq 0 \text{ for any vector } v$$

As a & b are positive real numbers, the combination of

$$a * (v^T M_1 v) - b * (v^T M_2 v)$$

may not always be non-negative.

\therefore this case does not satisfy the conditions of Mercer's theorem & is not a valid kernel.

An example of where $K(x, z)$ is not a kernel is:

$$a \underset{1 \times 2}{[0 \ 1]} \underset{2 \times 2}{\begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}} \begin{bmatrix} 0 \\ 1 \end{bmatrix} - b \underset{1 \times 2}{[0, 1]} \begin{bmatrix} 4 & 3 \\ 2 & 1 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

$$= a \underset{1 \times 2}{[3 \ 4]} \underset{2 \times 1}{\begin{bmatrix} 0 \\ 1 \end{bmatrix}} - b \underset{1 \times 2}{[2 \ 1]} \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

$$= a \underset{1 \times 1}{(4)} - b \underset{1 \times 1}{(1)}$$

$$= 4a - b$$

So for the example that

$$M_1 = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}, M_2 = \begin{bmatrix} 4 & 3 \\ 2 & 1 \end{bmatrix}$$

$$v^T = [0, 1], v = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

$$a = 1, b = 8$$

These values will give a $K(x, z)$ that is not a valid kernel.

q2)

4. To show that $k(x, z) = f(x)f(z)$ is a valid kernel, we need to demonstrate that for any dataset $X = \{x_1, x_2, \dots, x_n\}$ and any coefficients c_1, c_2, \dots, c_n , the Gram matrix G formed using the kernel is positive semi-definite, $c^T G c \geq 0$ for any vector c .

The Gram matrix G is defined as follows:

$$G_{ij} = k(x_i, x_j) = f(x_i)f(x_j)$$

Now, let's compute $c^T G c$:

$$\begin{aligned} c^T G c &= \sum_{i=1}^n \sum_{j=1}^n G_{ij} c_i c_j \\ &= \sum_{i=1}^n \sum_{j=1}^n f(x_i)f(x_j) c_i c_j \\ &= \sum_{i=1}^n \left(f(x_i) c_i \sum_{j=1}^n f(x_j) c_j \right) \\ &= \sum_{i=1}^n \left(f(x_i) c_i \right)^2 \end{aligned}$$

Since the square of any real number is non-negative,

$$c^T G c \geq 0.$$

$\therefore k(x, z) = f(x)f(z)$ is positive semi-definite.

Also, $k(x, z)$ is symmetric because:

$$\begin{aligned} k(x, z) &= f(x)f(z) \\ k(z, x) &= f(z)f(x) \\ &= f(x)f(z) \\ &= k(x, z) \end{aligned}$$

\therefore By Mercer's theorem,

$k(x, z)$ is a kernel

$$3. e^{(t)}(\theta) = \log(1 + e^{-y^{(t)}(\theta \cdot x^{(t)} + \theta_0)})$$

$$u = -y^{(t)}(\theta \cdot x^{(t)} + \theta_0)$$

$$e^{(t)}(\theta) = \log(1 + e^u)$$

$$\frac{\nabla e^{(t)}(\theta, \theta_0)}{\nabla \theta} = \frac{\frac{du}{d\theta} e^u}{1 + e^u}$$

$$= \frac{(\frac{du}{d\theta} e^u)(e^{-u})}{(1 + e^u)(e^{-u})}$$

$$= \frac{\frac{du}{d\theta}}{e^{-u} + 1}$$

$$= \frac{-y^{(t)} x^{(t)}}{1 + e^{y^{(t)}(\theta \cdot x^{(t)})}}$$

$$= \frac{-y^{(t)} x^{(t)}}{1 + \exp(y^{(t)}(\theta \cdot x^{(t)} + \theta_0))}$$

$$\frac{\nabla e^{(t)}(\theta, \theta_0)}{\nabla \theta_0} = \frac{\frac{du}{d\theta_0} e^u}{1 + e^u}$$

$$= \frac{\frac{du}{d\theta_0}}{e^{-u} + 1}$$

$$= \frac{-y^{(t)}}{1 + \exp(y^{(t)}(\theta \cdot x^{(t)} + \theta_0))}$$

$$u = -y^{(t)}(\theta \cdot x^{(t)} + \theta_0)$$

$$\frac{du}{d\theta} = -y^{(t)} x^{(t)}$$

$$u = -y^{(t)}(\theta \cdot x^{(t)} + \theta_0)$$

$$\frac{du}{d\theta_0} = -y^{(t)}$$

4.

$$y = \underline{\omega_0} + \sum_{i=1}^n \omega_i x_i + \omega_i x_i^2$$

$$E = \frac{1}{2} \sum_j (y_j - y_j^*)^2$$

constant

$$\frac{\partial E}{\partial \omega_0} = \sum_j (y_j - y_j^*) \left(\frac{\partial y_j}{\partial \omega_0} \right)$$

$$= \sum_j (y_j - y_j^*)$$

$$\frac{\partial E}{\partial \omega_i} = \sum_j (y_j - y_j^*) \left(\frac{\partial y_j}{\partial \omega_i} \right)$$

differentiated ω_0 w.r.t. ω_i

$$\frac{\partial y_j}{\partial \omega_i} = 0 + \sum_{i=1}^n x_i + x_i^2$$

$$\frac{\partial E}{\partial \omega_i} = \sum_j (y_j - y_j^*) \left(\sum_{i=1}^n x_i + x_i^2 \right)$$

∴ The update rule:

$$\omega_0_{\text{new}} = \omega_0_{\text{old}} - \text{learning rate} \times \sum_j (y_j - y_j^*)$$

$$\omega_i_{\text{new}} = \omega_i_{\text{old}} - \text{learning rate} \times \sum_j (y_j - y_j^*) \left(\sum_{i=1}^n x_i + x_i^2 \right)$$

5) For 2 features:

Naive Bayes needs to estimate the following probabilities:

1. $P(X_1 | y=0)$
2. $P(X_1 | y=1)$
3. $P(X_2 | y=0)$
4. $P(X_2 | y=1)$
5. $P(Y=0)$
6. $P(Y=1) = 1 - P(Y=0)$

↳ so not needed

∴ Total parameters = 5

The goal is to predict the class label (0 or 1) based on the values of X_1 and X_2 .

For 3 features:

Naive Bayes needs to estimate the following probabilities:

1. $P(X_1 | y=0)$
2. $P(X_1 | y=1)$
3. $P(X_2 | y=0)$
4. $P(X_2 | y=1)$
5. $P(X_3 | y=0)$
6. $P(X_3 | y=1)$
7. $P(Y=0)$
8. $P(Y=1) = 1 - P(Y=0)$

↳ so not needed, just find from $P(Y=0)$

∴ Total parameters = 7

For n features:

For every feature i in $1, 2, \dots, n$:

Need to find:

$$P(X_i | y=0)$$

$$P(X_i | y=1)$$

on top of this,
need to find

$$P(Y=0)$$

$$\& P(Y=1) = 1 - P(Y=0)$$

↳ not needed, just find from $P(Y=0)$

∴ Total parameters = $2n + 1$