

CS-E4830 Kernel Methods in Machine Learning

Assignment 1

Kunal Ghosh, 546247

September 20, 2016

1 Solution to Question 1

1.1 (a) True

Its given that K_1 and K_2 are Positive semi-definite matrices. So, for any vector v we have

$$\begin{aligned} v^T K_1 v &\geq 0 \text{ and } v^T K_2 v \geq 0 \\ \text{from question } K &= aK_1 + bK_2 \text{ where } a, b \in \mathcal{R}^+ \\ \text{then, } v^T K v &= v^T (aK_1 + bK_2) v \\ &= v^T (aK_1) v + v^T (bK_2) v \\ &= a(v^T K_1 v) + b(v^T K_2 v) \end{aligned} \tag{1}$$

Since, $a, b > 0$ and $v^T K_1 v, v^T K_2 v \geq 0$

$$\begin{aligned} \implies a(v^T K_1 v) + b(v^T K_2 v) &\geq 0 \\ \implies v^T (aK_1 + bK_2) v &\geq 0 \\ \implies v^T K v &\geq 0 \end{aligned}$$

So, K is Positive semi-definite. Hence a Kernel Matrix.

1.2 (b) False

If matrix K is defined as $K = K_1 - K_2$, where K_1, K_2 are positive semi-definite matrices, then

$$\begin{aligned} v^T K v &= v^T (K_1 - K_2) v \\ &= v^T (K_1) v - v^T (K_2) v \end{aligned} \tag{2}$$

We know that, $v^T K_1 v, v^T K_2 v \geq 0$. However, the difference of two non-negative numbers is not always non-negative.

$$\begin{aligned} \implies v^T (K_1) v - v^T (K_2) v &\not\geq 0 \\ \implies v^T K v &\not\geq 0 \end{aligned} \tag{3}$$

Hence $K = K_1 - K_2$ is not always positive semi-definite. Hence K need not be a Kernel Matrix.

1.3 (c) False

If matrix K is defined as $K = K_1 K_2$ where K_1, K_2 are positive semi-definite matrices, then
[Counter Example] : Product of two symmetric matrices is not always symmetric.

$$\text{for, } K_1 = \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix} \text{ and } K_2 = \begin{bmatrix} -2 & 1 \\ 1 & -3 \end{bmatrix}$$

where $\det(K_1) > 0$ and $\det(K_2) > 0$ so all the eigen values of K_1 and K_2 are positive. Hence K_1 and K_2 are positive semi-definite matrices. But,

$$K = K_1 K_2 = \begin{bmatrix} -3 & 4 \\ 3 & -4 \end{bmatrix}$$

is not a symmetric matrix. Therefore $K = K_1 K_2$ need not be a Kernel Matrix.

1.4 (d) False

[Counter Example] : Examples of K_1 and K_2 given above (previous section) are valid Kernel Matrices but don't have all positive entries.

1.5 (e) True

e). $K_1(x_i, x_j) = \Phi_1(x_i)^T \Phi_1(x_j) = \sum_{l=1}^d \phi_{1l}(x_i) \cdot \phi_{1l}(x_j)$
 $K_2(x_i, x_j) = \Phi_2(x_i)^T \Phi_2(x_j) = \sum_{m=1}^d \phi_{2m}(x_i) \cdot \phi_{2m}(x_j)$
 so $K(x_i, x_j) = K_1(x_i, x_j) \cdot K_2(x_i, x_j) = \left[\sum_{l=1}^d \phi_{1l}(x_i) \cdot \phi_{1l}(x_j) \right] \left[\sum_{m=1}^d \phi_{2m}(x_i) \cdot \phi_{2m}(x_j) \right]$

Expanding out a few terms.

$$= \left[\phi_{11}(x_i) \cdot \phi_{11}(x_j) + \phi_{12}(x_i) \cdot \phi_{12}(x_j) + \dots + \phi_{1d}(x_i) \cdot \phi_{1d}(x_j) \right] \text{ has } d \text{ terms} \\
\times \left[\phi_{21}(x_i) \cdot \phi_{21}(x_j) + \phi_{22}(x_i) \cdot \phi_{22}(x_j) + \dots + \phi_{2d}(x_i) \cdot \phi_{2d}(x_j) \right] \text{ has } d \text{ terms}$$

= When we find the product of the above two terms it is going to be a term with $d \times d$ terms.

$$= \left[\begin{aligned} &\phi_{11}(x_i) \phi_{11}(x_j) \phi_{21}(x_i) \phi_{21}(x_j) + \phi_{11}(x_i) \phi_{11}(x_j) \phi_{22}(x_i) \phi_{22}(x_j) + \dots + \phi_{11}(x_i) \phi_{11}(x_j) \phi_{2d}(x_i) \phi_{2d}(x_j) \\ &+ \phi_{12}(x_i) \phi_{12}(x_j) \phi_{21}(x_i) \phi_{21}(x_j) + \dots \\ &+ \phi_{1d}(x_i) \phi_{1d}(x_j) \phi_{21}(x_i) \phi_{21}(x_j) + \dots + \phi_{1d}(x_i) \phi_{1d}(x_j) \phi_{2d}(x_i) \phi_{2d}(x_j) \end{aligned} \right] \text{ } d \text{ terms}$$

Individual terms among these
 These $d \times d$ terms can be written as.

$$\left[\begin{aligned} &(\phi_{11}(x_i) \phi_{21}(x_i)) (\phi_{11}(x_j) \phi_{21}(x_j)) + (\phi_{11}(x_i) \phi_{22}(x_i)) (\phi_{11}(x_j) \phi_{22}(x_j)) + \dots + (\phi_{11}(x_i) \phi_{2d}(x_i)) (\phi_{11}(x_j) \phi_{2d}(x_j)) \\ &+ (\phi_{12}(x_i) \phi_{21}(x_i)) (\phi_{12}(x_j) \phi_{21}(x_j)) + (\phi_{12}(x_i) \phi_{22}(x_i)) (\phi_{12}(x_j) \phi_{22}(x_j)) + \dots + (\phi_{12}(x_i) \phi_{2d}(x_i)) (\phi_{12}(x_j) \phi_{2d}(x_j)) \\ &+ \dots \\ &+ (\phi_{1d}(x_i) \phi_{21}(x_i)) (\phi_{1d}(x_j) \phi_{21}(x_j)) + (\phi_{1d}(x_i) \phi_{22}(x_i)) (\phi_{1d}(x_j) \phi_{22}(x_j)) + \dots + (\phi_{1d}(x_i) \phi_{2d}(x_i)) (\phi_{1d}(x_j) \phi_{2d}(x_j)) \end{aligned} \right]$$

This sum can be re-written as a dot-product of a vector $\Phi(x_i)$ with itself. i.e. $\Phi(x_i)^T \Phi(x_j)$ such that where $\Phi(x)$ is defined as follows.

$$\left[\begin{aligned} &\phi_{11}(x) \phi_{21}(x), \phi_{11}(x) \phi_{22}(x), \dots, \phi_{11}(x) \phi_{2d}(x), \\ &\phi_{12}(x) \phi_{21}(x), \phi_{12}(x) \phi_{22}(x), \dots, \phi_{12}(x) \phi_{2d}(x), \\ &\vdots \\ &\phi_{1d}(x) \phi_{21}(x), \dots, \phi_{1d}(x) \phi_{2d}(x) \end{aligned} \right]^T$$

This is exactly what $\Phi(x)$ would be, if it were defined as.

$\Phi(x) = \Phi_1(x) \otimes \Phi_2(x)$. Hence if $K(x_i, x_j) = K_1(x_i, x_j) \times K_2(x_i, x_j)$, then $K(x_i, x_j)$ can be written as $\Phi(x_i)^T \Phi(x_j)$ where $\Phi(x) = \Phi_1(x) \otimes \Phi_2(x)$.

2 Solution to Question 2

Question 2: To show, $h(x) = \text{sign}(\|\phi(x) - c_-\|^2 - \|\phi(x) - c_+\|^2)$ can be re-written as.

$$\text{sign}\left(\sum_{i=1}^m \alpha_i K(x, x_i) + b\right).$$

we know that $\|\phi(x) - c_-\|^2 = \langle \phi(x) - c_-, \phi(x) - c_- \rangle$

$$\text{so } h(x) = \text{sign}(\langle \phi(x) - c_-, \phi(x) - c_- \rangle - \langle \phi(x) - c_+, \phi(x) - c_+ \rangle)$$

Using the property of linearity of inner product. according to which.
 $\langle x+y, z \rangle = \langle x, z \rangle + \langle y, z \rangle$ and $\langle x, y+z \rangle = \langle x, y \rangle + \langle x, z \rangle$.

$$\begin{aligned} h(x) &= \text{sign}(\langle \phi(x), \phi(x) - c_- \rangle - \langle c_-, \phi(x) - c_- \rangle - \{\langle \phi(x), \phi(x) - c_+ \rangle - \langle c_+, \phi(x) - c_+ \rangle\}) \\ &= \text{sign}(\langle \phi(x), \phi(x) \rangle - \langle \phi(x), c_- \rangle - \{c_-, \phi(x) - c_-\} - \{\langle \phi(x), \phi(x) \rangle - \langle \phi(x), c_+ \rangle\} \\ &\quad + \{\langle c_+, \phi(x) \rangle - \langle c_+, c_+ \rangle\}) \\ &= \text{sign}(\langle \phi(x), \phi(x) \rangle - \langle \phi(x), c_- \rangle - \langle c_-, \phi(x) \rangle + \langle c_-, c_- \rangle - \langle \phi(x), \phi(x) \rangle + \langle \phi(x), c_+ \rangle \\ &\quad + \langle c_+, \phi(x) \rangle - \langle c_+, c_+ \rangle) \end{aligned}$$

Continues in next page...

also $\langle a, b \rangle \equiv \langle b, a \rangle$.

$$\begin{aligned} \therefore h(x) &= \text{sign} (2 \langle \phi(x), c_+ \rangle - 2 \langle \phi(x), c_- \rangle + \langle c_-, c_- \rangle - \langle c_+, c_+ \rangle) \\ &= \text{sign} (2 \langle \phi(x), c_+ - c_- \rangle + \langle c_-, c_- \rangle - \langle c_+, c_+ \rangle) \quad \text{--- (1)} \\ &= \text{sign} (2 \langle \phi(x), c_+ - c_- \rangle + \end{aligned}$$

c_- and c_+ are vectors. $\therefore \langle c_-, c_- \rangle = \left[\frac{1}{m_-} \sum_{i \in I^-} \phi(x_i) \right]^T \left[\frac{1}{m_-} \sum_{i \in I^-} \phi(x_i) \right]$

$$\begin{aligned} &= \frac{1}{m_-^2} \sum_{i, j \in I^-} \phi(x_i)^T \phi(x_j) = \frac{1}{m_-^2} \sum_{i, j \in I^-} \phi(x_i)^T \phi(x_j) \quad \text{where } x_i = x_j \text{ if } i \in I^- \\ &= \frac{1}{m_-^2} \sum_{i, j \in I^-} \langle \phi(x_i), \phi(x_j) \rangle \\ &= \frac{1}{m_-^2} \sum_{i, j \in I^-} k(x_i, x_j) \quad \text{using kernel trick.} \quad \text{--- (2)} \end{aligned}$$

Similarly, $\langle c_+, c_+ \rangle = \frac{1}{m_+^2} \sum_{i, j \in I^+} k(x_i, x_j)$. --- (3).

nd.

$$\begin{aligned} \langle \phi(x), c_+ - c_- \rangle &= \phi(x)^T \left[\frac{1}{m_+} \sum_{i \in I^+} \phi(x_i) - \frac{1}{m_-} \sum_{i \in I^-} \phi(x_i) \right] \\ &= \frac{1}{m_+} \sum_{i \in I^+} \phi(x)^T \phi(x_i) - \frac{1}{m_-} \sum_{i \in I^-} \phi(x)^T \phi(x_i). \quad \text{--- (4)} \end{aligned}$$

Let $\alpha_i = \pm 1 = \begin{cases} \frac{1}{m_+} & \text{if } y_i = +1 \\ \frac{1}{m_-} & \text{if } y_i = -1 \end{cases}$

$$\alpha_i = \begin{cases} \frac{1}{m_+} & \text{if } y_i = +1 \\ \frac{1}{m_-} & \text{if } y_i = -1 \end{cases}$$

$$\therefore h(x) = \text{sign} \left(2 \sum_{i=1}^m \alpha_i \phi(x)^T \phi(x_i) + \frac{1}{m_+^2} \sum_{i, j \in I^+} k(x_i, x_j) + \frac{1}{m_-^2} \sum_{i, j \in I^-} k(x_i, x_j) \right)$$

b.

$$h(x) = \text{sign} \left(\sum_{i=1}^m \alpha_i \phi(x)^T \phi(x_i) + b \right)$$

taking the 2. out, common doesn't affect the sign.

$$b = \frac{1}{2m_+^2} \sum_{i, j \in I^+} k(x_i, x_j) + \frac{1}{2m_-^2} \sum_{i, j \in I^-} k(x_i, x_j)$$

$$\alpha_i = \begin{cases} \frac{1}{m_+} & \text{if } y_i = +1. \\ -\frac{1}{m_-} & \text{if } y_i = -1. \end{cases}$$

2 Solution to Question 3

```
function kernel_mat = gaussian_kernel(X, Z, sigma)
    % X is (n,d)
    % Z is (m,d)
    % sigma is scalar
    % gaussian kernel = exp( -0.5 * (1/sigma^2)* || phi(x1) - phi(x2) ||^2)
    % Assuming, here X and Z are already projected in phi (feature) space.
    % So in this case, gaussian kernel = exp( -0.5 * (1/sigma^2)* || x1 - x2 ||^2)
    % Dimension of kernel_mat is (n,m)
    [n,d] = size(X);
    [m,d] = size(Z);
    kernel_mat = zeros(n,m);
    for row = 1:n
        kernel_mat(row,:) = sum(bsxfun(@minus,Z,X(row,:)).^2,2)';
    end
    kernel_mat = exp(kernel_mat * (-0.5 / (sigma^2)));
end
```

3 Solution to Question 4

```
function y_pred = parzen_classify(Kx_train, Kx_train_test, y_train)
    m_pos = sum(y_train == 1);
    m_neg = length(y_train) - m_pos;
    % size(Kx_train) (800,800)
    % size(Kx_train_test) (800,200)
    b = zeros(2,1);
    negidxs = find(y_train == -1);
    for i = negidxs'
        b(1) = b(1) + sum(Kx_train(i,negidxs));
    end
    b(1) = b(1)/(2*m_neg^2);

    posidxs = find(y_train == 1);
    for i = posidxs'
        b(2) = b(2) + sum(Kx_train(i,posidxs));
    end
    b(2) = b(2)/(2*m_pos^2);

    const = b(1) - b(2);

    % here alpha is (800,1)
    alpha = (1/m_pos)*ones(size(y_train));
    alpha(negidxs) = (-1/m_neg);

    y_pred = sign((alpha' * Kx_train_test) + const);
end
```

4 Solution to Question 5

In the caption in the figures below σ is a parameter of the gaussian kernel which is used for the task of classification in this problem. Among the figures, fig-1 shows the learning curves and fig-2 shows the decision boundary for various values of σ . More description in the captions of the figures.

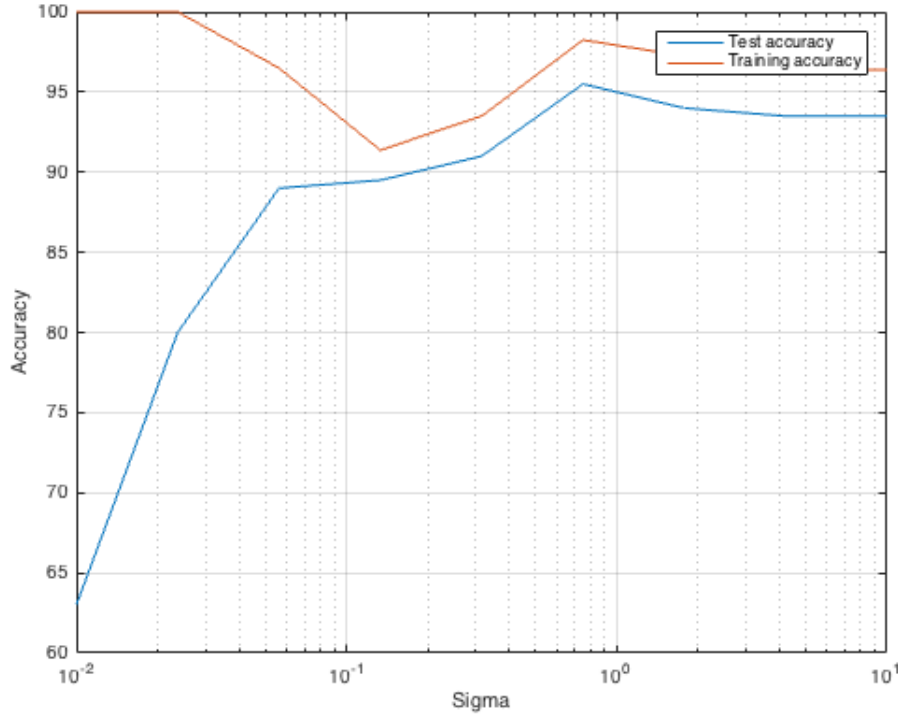


Figure 1: Learning curves: Initially, with low values of sigma there is overfitting (high training accuracy, low test accuracy) as the sigma values increase the model generalizes better until sigma reaches ≈ 0.7 then the model starts to underfit.

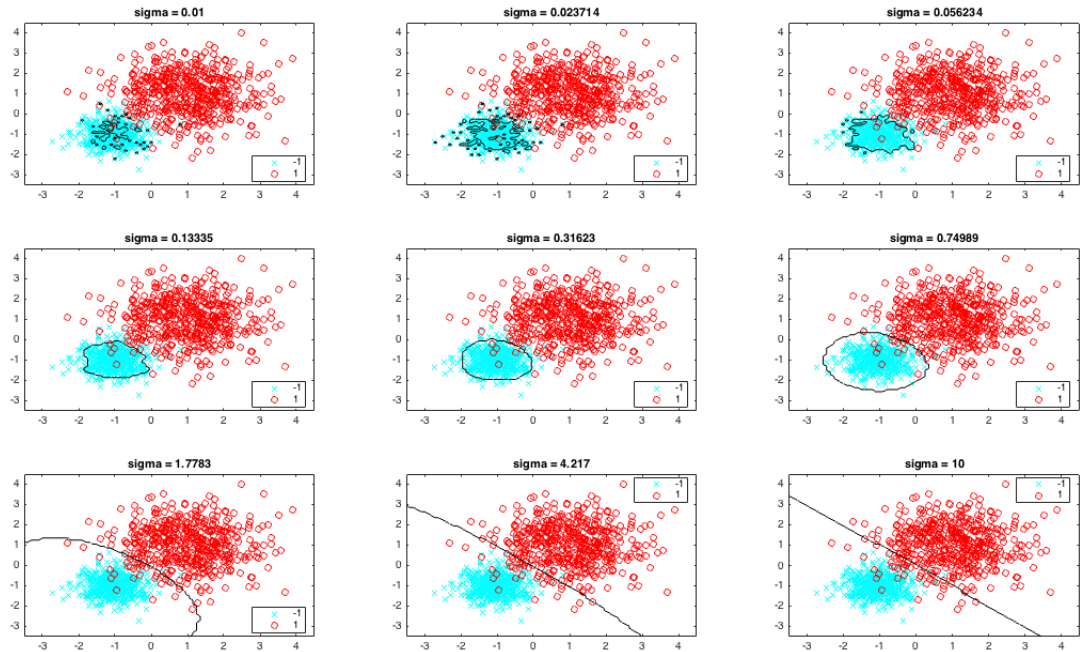


Figure 2: Decision Boundaries: For low values of sigma overfitting is clearly visible (small "circles" tightly encircling certain datapoints). As the sigma increases the model generalizes better, for $\sigma \approx 0.7$ the test accuracy is the best (least number of red-circles inside the decision boundary). Finally, as the sigma keeps on increasing the decision boundary tends to a linear decision boundary which isn't the most ideal for this dataset as can be seen in the learning curves Fig-1 above where both the training and test accuracy drops as the sigma is increased beyond ≈ 0.7 .