CS/ECE/ME532 Activity 21

*Estimated Time: 20 minutes for P1, 20 minutes for P2, 10 minutes for P3, 15 minutes for P4.*

Handwritten (top right):

$$[x_{i_1}^2 \quad x_{i_2}^2 \quad \sqrt{2}x_{i_1}x_{i_2} \quad \sqrt{2}x_{i_1} \quad \sqrt{2}x_{i_2} 1] \begin{bmatrix} x_{j_1}^2 \\ x_{j_2}^2 \\ \sqrt{2}x_{j_1}x_{j_2} \\ \sqrt{2}x_{j_1} \\ \sqrt{2}x_{j_2} \\ 1 \end{bmatrix}$$

$= x_{i_1}^2 x_{j_1}^2 + x_{i_2}^2 x_{j_2}^2 + 2x_{i_1}x_{i_2}x_{j_1}x_{j_2}$
$+ 2x_{i_1}x_{j_1} + 2x_{i_2}x_{j_2} + 1$

$(x_{i_1}x_{j_1} + x_{i_2}x_{j_2}+1)(x_{i_1}x_{j_1}+x_{i_2}x_{j_2}+1)$
$= x_{i_1}^2x_{j_1}^2 + x_{i_1}x_{i_2}x_{j_1}x_{j_2} + x_{i_1}x_{j_1} +$
$x_{i_2}x_{j_2}x_{i_1}x_{j_1} + x_{i_2}^2x_{j_2}^2 + x_{i_2}x_{j_2} + x_{i_1}x_{j_1}$
$+ x_{i_2}x_{j_2} + 1$
$= x_{i_1}^2x_{j_1}^2 + x_{i_2}^2x_{j_2}^2$
$+ 2x_{i_1}x_{i_2}x_{j_1}x_{j_2}$
$+ 2x_{i_1}x_{j_1}$
$+ 2x_{i_2}x_{j_2} + 1$

1. Consider performing regression using all quadratic and lower order functions of a 2-dimensional observation $\boldsymbol{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$

$$\hat{y} = x_1^2 w_1 + x_2^2 w_2 + \sqrt{2}x_1x_2w_3 + \sqrt{2}x_1w_4 + \sqrt{2}x_2w_5 + w_6$$

a) Show that $\hat{y} = \boldsymbol{\phi}^T(\boldsymbol{x})\boldsymbol{w}$ and find $\boldsymbol{\phi}, \boldsymbol{w}$.

$\phi = [x_1^2 \quad x_2^2 \quad \sqrt{2}x_1x_2 \quad \sqrt{2}x_1 \quad \sqrt{2}x_2 1]^T$
$w = [w_1 \quad w_2 \quad w_3 \quad w_4 \quad w_5 \quad w_6]^T$

b) Show that the "kernel" $\boldsymbol{\phi}^T(\boldsymbol{x}_i)\boldsymbol{\phi}(\boldsymbol{x}_j)$ is identical to $(\boldsymbol{x}_i^T\boldsymbol{x}_j + 1)^2$.

$\hookrightarrow [x_{i_1}x_{i_2}]\begin{bmatrix} x_{j_1} \\ x_{j_2} \end{bmatrix} = x_{i_1}x_{j_1}+x_{i_2}x_{j_2}+1$

c) The number of multiplications may be used as a crude measure of computational complexity. Compare the number of multiplications required to compute $\boldsymbol{\phi}^T(\boldsymbol{x}_i)\boldsymbol{\phi}(\boldsymbol{x}_j)$ (ignoring the $\sqrt{2}$ terms) to that required to compute $(\boldsymbol{x}_i^T\boldsymbol{x}_j + 1)^2$.

$\underbrace{\phi^T(x_i)\phi(x_j)}_{\hookrightarrow \text{6 multiplications}}$ $\underbrace{(x_i^Tx_j+1)^2}_{\hookrightarrow \text{3 multiplications,}}$
2 for $x_i^Tx_j$,
1 for square

2. You are given $N$ observations $y_i, \boldsymbol{x}_i, i = 1, 2, \ldots, N$ and solve the ridge-regression problem

$$\arg\min_{\boldsymbol{w}} \|\boldsymbol{y} - \boldsymbol{\Phi}\boldsymbol{w}\|_2^2 + \lambda\|\boldsymbol{w}\|_2^2$$

where $\boldsymbol{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}$ and $\boldsymbol{\Phi} = \begin{bmatrix} \boldsymbol{\phi}^T(\boldsymbol{x}_1) \\ \boldsymbol{\phi}^T(\boldsymbol{x}_2) \\ \vdots \\ \boldsymbol{\phi}^T(\boldsymbol{x}_N) \end{bmatrix}$. You know the solution may be expressed in standard form as

$\Phi^T\Phi\Phi^T + \lambda\Phi^T$

$$\hat{\boldsymbol{w}} = (\boldsymbol{\Phi}^T\boldsymbol{\Phi} + \lambda\boldsymbol{I})^{-1}\boldsymbol{\Phi}^T\boldsymbol{y}$$

$\Phi^T(\Phi\Phi^T+\lambda I) = (\Phi^T\Phi + \lambda I)\Phi^T$

a) Factor $\boldsymbol{\Phi}^T$ from the left and the right of $\boldsymbol{\Phi}^T\boldsymbol{\Phi}\boldsymbol{\Phi}^T + \lambda\boldsymbol{\Phi}^T$ to show that

$$(\boldsymbol{\Phi}^T\boldsymbol{\Phi} + \lambda\boldsymbol{I})^{-1}\boldsymbol{\Phi}^T = \boldsymbol{\Phi}^T(\boldsymbol{\Phi}\boldsymbol{\Phi}^T + \lambda\boldsymbol{I})^{-1}$$

$\Phi^T(\Phi\Phi^T+\lambda I)^{-1} = (\Phi^T\Phi+\lambda I)^{-1}\Phi^T$
$\Phi^T = \Phi^T$

*Hint:* we did this a previous activity and you used the result in the breast cancer classification assignment.

$\hat{w} = (\Phi^T\Phi + \lambda I)^{-1}\Phi^Ty$

b) Use the result of the previous part to show that

$(\Phi^T\Phi+\lambda I)^{-1}\Phi^T = \Phi^T(\Phi\Phi^T+\lambda I)^{-1}$

$$\hat{\boldsymbol{w}} = \boldsymbol{\Phi}^T(\boldsymbol{\Phi}\boldsymbol{\Phi}^T + \lambda\boldsymbol{I})^{-1}\boldsymbol{y}$$

$\hat{w} = \Phi^T(\Phi\Phi^T+\lambda I)^{-1}y$

c) Let the kernel matrix $\boldsymbol{K} = \boldsymbol{\Phi}\boldsymbol{\Phi}^T$. Express the $i, j$ element of $\boldsymbol{K}$, $[\boldsymbol{K}]_{i,j}$ using $\boldsymbol{\phi}(\boldsymbol{x})$.

$[K]_{i,j} = \phi^T(x_i)\phi(x_j)$

$$[K]_{i,j} = (x_i^T x_j + 1)^2$$

**d)** Assume $\phi(x)$ is defined as in Problem 1 and find $[K]_{i,j}$ as a function of $x_i^T x_j$.

**e)** Recall from Problem 1 that $\hat{y}(x) = \phi^T(x)\hat{w}$. Thus, $\hat{y}(x) = \phi^T(x)\Phi^T(\Phi\Phi^T + \lambda I)^{-1}y$. Show that

$$\hat{y}(x) = \sum_{j=1}^{N} K(x, x_j)\alpha_j \qquad \alpha = (\Phi\Phi^T + \lambda I)^{-1}y$$

where $K(x, x_j) = (x^T x_j + 1)^2$. $\longrightarrow \phi^T(x)\phi(x_j)$

$$\hat{y}(x) = \sum_{j=1}^{N} (x^T x_j + 1)^2 \alpha_j = \sum_{j=1}^{N} (x^T x_j + 1)^2 (\Phi\Phi^T + \lambda I)^{-1}y = \phi^T(x)\Phi^T(\Phi\Phi^T + \lambda I)^{-1}y$$

**3.** Suppose $\phi(x) = x$. Use the results of the previous problem.

**a)** Find the expression for the corresponding kernel $K(x, x_j)$. $= X^T X_j$

**b)** Express $\hat{y}(x)$ in terms of $\alpha_j$ and your expression for $K(x, x_j)$. How does each training sample influence the prediction $\hat{y}(x)$ at some new value $x$?

$$\hat{y}(x) = \sum_{j=1}^{N} \alpha_j (x^T x_j)$$

$x_j$ is used as part of each multiplication that sums to $\hat{y}(x)$, so each training sample is used in the prediction.

**4.** The results we developed in this exercise so far show that regression can be expressed entirely in terms of the kernel function $K(x, x_j)$:

$$\hat{y}(x) = \sum_{j=1}^{n} K(x, x_j)\alpha_j$$

where $\alpha_j$ is a function of the kernel matrix $K$, regularization parameter $\lambda$, and data $y$. This form allows us to perform regression when the high dimensional feature vector $\phi(x)$ is not easily defined, but $K(x, x_j) = \phi^T(x)\phi(x_j)$ is easily defined. One such case is the Gaussian kernel,

$$K(x, x_j) = \exp\left\{-\frac{||x - x_j||_2^2}{2\sigma}\right\}$$

For simplicity this problem assumes $x$ is one dimensional, that is $\hat{y}(x)$ describes a graph of a function of one variable.

**a)** Suppose $x_1 = -2, x_2 = 0$, and $x_3 = 2$. Sketch $K(x, x_j)$ as a function of $x$ for $j = 1, 2, 3$ assuming $\sigma = 1$.

**b)** Now sketch $\hat{y}(x)$ assuming $\alpha_1 = -1, \alpha_2 = 2$, and $\alpha_3 = 1$.

**c)** Fill in the blanks. The expression $\hat{y}(x) = \sum_{j=1}^{n} K(x, x_j)\alpha_j$ interpolates a value $y$ corresponding to $x$ as a **weighted** sum of __n__ functions centered on the **y-axis**.

a).