

CS/ECE/ME532 Activity 21

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

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

$$\hat{y} = x_1^2 w_1 + x_2^2 w_2 + \sqrt{2} x_1 x_2 w_3 + \sqrt{2} x_1 w_4 + \sqrt{2} x_2 w_5 + w_6$$

- a) Show that $\hat{y} = \boldsymbol{\phi}^T(\mathbf{x})\mathbf{w}$ and find $\boldsymbol{\phi}, \mathbf{w}$.
 - b) Show that the “kernel” $\boldsymbol{\phi}^T(\mathbf{x}_i)\boldsymbol{\phi}(\mathbf{x}_j)$ is identical to $(\mathbf{x}_i^T \mathbf{x}_j + 1)^2$.
 - c) The number of multiplications may be used as a crude measure of computational complexity. Compare the number of multiplications required to compute $\boldsymbol{\phi}^T(\mathbf{x}_i)\boldsymbol{\phi}(\mathbf{x}_j)$ (ignoring the $\sqrt{2}$ terms) to that required to compute $(\mathbf{x}_i^T \mathbf{x}_j + 1)^2$.
2. You are given N observations $y_i, \mathbf{x}_i, i = 1, 2, \dots, N$ and solve the ridge-regression problem

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

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

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

- 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

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

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

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

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

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

- d) Assume $\phi(\mathbf{x})$ is defined as in Problem 1 and find $[\mathbf{K}]_{i,j}$ as a function of $\mathbf{x}_i^T \mathbf{x}_j$.
- e) Recall from Problem 1 that $\hat{y}(\mathbf{x}) = \phi^T(\mathbf{x})\hat{\mathbf{w}}$. Thus, $\hat{y}(\mathbf{x}) = \phi^T(\mathbf{x})\Phi^T(\Phi\Phi^T + \lambda\mathbf{I})^{-1}\mathbf{y}$. Show that

$$\hat{y}(\mathbf{x}) = \sum_{j=1}^N K(\mathbf{x}, \mathbf{x}_j)\alpha_j$$

where $K(\mathbf{x}, \mathbf{x}_j) = (\mathbf{x}^T \mathbf{x}_j + 1)^2$.

3. Suppose $\phi(\mathbf{x}) = \mathbf{x}$. Use the results of the previous problem.
- a) Find the expression for the corresponding kernel $K(\mathbf{x}, \mathbf{x}_j)$.
- b) Express $\hat{y}(\mathbf{x})$ in terms of α_j and your expression for $K(\mathbf{x}, \mathbf{x}_j)$. How does each training sample influence the prediction $\hat{y}(\mathbf{x})$ at some new value \mathbf{x} ?
4. The results we developed in this exercise so far show that regression can be expressed entirely in terms of the kernel function $K(\mathbf{x}, \mathbf{x}_j)$:

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

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

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

For simplicity this problem assumes \mathbf{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 _____ sum of _____ functions centered on the _____.