

50.021 Artificial Intelligence

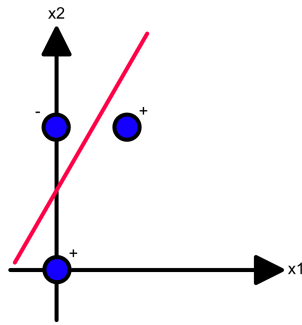
Theory Homework 1

Due: every Monday, 4PM before class starts

[Q1]. Find *any* separating hyperplane equation for these three sample points: $\mathbf{x}_1 = (0, 0), y_1 = +1$, $\mathbf{x}_2 = (1, 3), y_2 = +1$, and $\mathbf{x}_3 = (0, 3), y_3 = -1$. Draw (by hand) or plot (using Python, see matplotlib) the result.

Solution:

Any line that separates the points into two classes is ok. Below is one example.



[Q2]. Find by hand the value of optimum weights $\hat{\mathbf{w}}$ and bias \hat{b} using linear regression for the four following sample points: $\mathbf{x}_1 = (1, 1), y_1 = +1$, $\mathbf{x}_2 = (2, 2), y_2 = +1$, $\mathbf{x}_3 = (1, 3), y_3 = -1$, $\mathbf{x}_4 = (2, 3), y_4 = -1$. Show your working.

Solution:

We can begin by making the matrix X ,

$$X = \begin{bmatrix} 1 & 1 & 1 \\ 2 & 2 & 1 \\ 1 & 3 & 1 \\ 2 & 3 & 1 \end{bmatrix}$$

and the vector $\mathbf{y} = [1 \ 1 \ -1 \ -1]^T$. We know that the optimum weight is $\hat{\mathbf{w}} = [w_1 \ w_2 \ \hat{b}]^T = (X^T \cdot X)^{-1} X^T \cdot Y$. Substituting the values above,

$$\hat{\mathbf{w}} = \left(\begin{bmatrix} 1 & 2 & 1 & 2 \\ 1 & 2 & 3 & 3 \\ 1 & 1 & 1 & 1 \end{bmatrix} \cdot \begin{bmatrix} 1 & 1 & 1 \\ 2 & 2 & 1 \\ 1 & 3 & 1 \\ 2 & 3 & 1 \end{bmatrix} \right)^{-1} \begin{bmatrix} 1 & 2 & 1 & 2 \\ 1 & 2 & 3 & 3 \\ 1 & 1 & 1 & 1 \end{bmatrix} \cdot \begin{bmatrix} 1 \\ 1 \\ -1 \\ -1 \end{bmatrix}$$

$$X^T \cdot X = \begin{bmatrix} 10 & 14 & 6 \\ 14 & 23 & 9 \\ 6 & 9 & 4 \end{bmatrix}, \quad (X^T \cdot X)^{-1} = \begin{bmatrix} 1.1 & -0.2 & -1.2 \\ -0.2 & 0.4 & -0.6 \\ -1.2 & -0.6 & 3.4 \end{bmatrix}$$

Hence,

$$\hat{\mathbf{w}} = \begin{bmatrix} 1.1 & -0.2 & -1.2 \\ -0.2 & 0.4 & -0.6 \\ -1.2 & -0.6 & 3.4 \end{bmatrix} \cdot \begin{bmatrix} 1 & 2 & 1 & 2 \\ 1 & 2 & 3 & 3 \\ 1 & 1 & 1 & 1 \end{bmatrix} \cdot \begin{bmatrix} 1 \\ 1 \\ -1 \\ -1 \end{bmatrix} = \begin{bmatrix} 0.6 \\ -1.2 \\ 1.8 \end{bmatrix}$$

[Q3]. In the lecture notes, we solve the objective function:

$$\hat{\mathbf{w}} = \operatorname{argmin}_{\mathbf{w}} \sum_{i=1}^n (y_i - \mathbf{w} \cdot \mathbf{x}_i)^2 \quad (1)$$

by hand, and we have the analytical solution for optimum weights \mathbf{w}^* in linear regression,

$$\hat{\mathbf{w}} = (X^T \cdot X)^{-1} X^T \cdot Y$$

$$\mathbf{x}, \hat{\mathbf{w}} \in \mathbb{R}^d$$

Now suppose instead of using \mathbf{x}_i directly, we want to use some basis function $\phi(\mathbf{x}_i)$ on each dataset i , and suppose that we use a slightly different squared error loss function than the lecture notes,

$$L(y, f(\mathbf{x})) = \frac{1}{2N} (y - f(\mathbf{x}))^2,$$

$$f(\mathbf{x}) = \phi(\mathbf{x}) \cdot \mathbf{w}$$

1. Rewrite the objective function in equation (1) using the new loss function $L(y, f(x))$ above.
2. Show that, even when the loss function is changed from what's shown in the lecture notes to the above, and although we apply basis function $\phi(\mathbf{x}_i)$, the solution for optimum weight $\hat{\mathbf{w}}$ still takes the similar form,

$$\hat{\mathbf{w}} = (\Phi^T \cdot \Phi)^{-1} \Phi^T \cdot Y,$$

where,

$$\Phi = \begin{bmatrix} \phi_1(\mathbf{x}_1) & \dots & \phi_d(\mathbf{x}_1) \\ \vdots & \ddots & \vdots \\ \phi_1(\mathbf{x}_N) & \dots & \phi_d(\mathbf{x}_N) \end{bmatrix},$$

and d is number of dimensions of each sample \mathbf{x} , and N is the number of samples.

3. Show that if we define a function,

$$\mathcal{K}(\mathbf{x}_i, \mathbf{x}_r) = \phi(\mathbf{x}_i) \cdot \phi(\mathbf{x}_r),$$

then $f(\mathbf{x})$ can be written only in terms of the function \mathcal{K} above without the need to specify ϕ explicitly.

Hint: Let $\mathbf{w} = \Phi^T \mathbf{v}$, where \mathbf{v} is some new parameter vector that you now maximise.

4. Write down the analytical form of the optimum parameter for $L(y, f(\mathbf{x}))$ using your new expression of $f(\mathbf{x})$ (that now contains \mathcal{K} in part (3)). This can be done by reusing some of your answer in part (2).

Solution:

1. The function we are going to minimise is the sum of the loss functions over the entire training samples,

$$\hat{\mathbf{w}} = \operatorname{argmin}_{\mathbf{w}} \sum_{i=1}^N \frac{1}{2N} (y_i - f(\mathbf{x}_i))^2$$

$$= \operatorname{argmin}_{\mathbf{w}} \sum_{i=1}^N \frac{1}{2N} (y_i - \phi(\mathbf{x}_i) \cdot \mathbf{w})^2$$

In vectorised form, it is,

$$\hat{\mathbf{w}} = \operatorname{argmin}_{\mathbf{w}} \frac{1}{2N} \|\mathbf{y} - \Phi \mathbf{w}\|^2,$$

where $\mathbf{y} = (y_1, \dots, y_N)^T$, and $\mathbf{w} = (w_1, \dots, w_d)^T$.

2. We can differentiate the equation in part 1 with respect to \mathbf{w} and set it to zero. In short, we can use the vector derivative,

$$\nabla_{\mathbf{w}} \left[\frac{1}{2N} \|\mathbf{y} - \Phi \mathbf{w}\|^2 \right] = 0$$

Hence,

$$\begin{aligned} 0 &= \frac{1}{2N} \nabla_{\mathbf{w}} [(\mathbf{y} - \Phi \mathbf{w})^T (\mathbf{y} - \Phi \mathbf{w})] \\ &= \frac{1}{2N} \nabla_{\mathbf{w}} [\mathbf{y}^T \mathbf{y} - 2\mathbf{y}^T \Phi \mathbf{w} + \mathbf{w}^T \Phi^T \Phi \mathbf{w}] \\ &= \frac{1}{2N} (-2\Phi^T \mathbf{y} + 2\Phi^T \Phi \mathbf{w}) \\ &\quad \text{(Google "Matrix Calculus" if you're lost on how to get the above expression)} \\ &= -\Phi^T \mathbf{y} + \Phi^T \Phi \mathbf{w} \\ \therefore \hat{\mathbf{w}} &= (\Phi^T \Phi)^{-1} \Phi^T \mathbf{y} \end{aligned}$$

3. The key is to first realise that the possible values of \mathbf{w} lies in the span of data points $\mathbf{x}_i \forall i = 1, \dots, N$. In other words, the optimal $\hat{\mathbf{w}}$ depends on *all* of the training data,

$$\mathbf{w} = \Phi^T \mathbf{v},$$

for some vector $\mathbf{v} = (v_1, \dots, v_N)$, $\mathbf{v} \in \mathbb{R}^N$.

We know that $f(\mathbf{x}) = \phi(\mathbf{x}) \cdot \mathbf{w}$. Plugging in the new expression for \mathbf{w} above,

$$\begin{aligned} f(\mathbf{x}) &= \phi(\mathbf{x}) \cdot (\Phi^T \mathbf{v}) \\ &= \phi(\mathbf{x}) \cdot \left(\sum_{i=1}^N \phi(\mathbf{x}_i) v_i \right) \\ &= \sum_{i=1}^N (\phi(\mathbf{x}) \cdot \phi(\mathbf{x}_i)) v_i \\ &= \sum_{i=1}^N \mathcal{K}(\mathbf{x}, \mathbf{x}_i) v_i \end{aligned}$$

This is called the *dual* form of linear regression. It is more efficient to compute \mathcal{K} than to compute ϕ , because ϕ . In the literature, \mathcal{K} is known as a *kernel*.

4. Using results from part (2),

$$\begin{aligned} \Phi^T \mathbf{y} &= \Phi^T \Phi \hat{\mathbf{w}}, \\ \Phi^T \mathbf{y} &= \Phi^T \Phi (\Phi^T \hat{\mathbf{v}}), \\ \therefore \mathbf{y} &= \Phi \Phi^T \hat{\mathbf{v}}, \\ \hat{\mathbf{v}} &= (\Phi \Phi^T)^{-1} \mathbf{y}, \\ \hat{\mathbf{v}} &= (\mathbf{K})^{-1} \mathbf{y}, \end{aligned}$$

where,

$$\mathbf{K} = \begin{bmatrix} \mathcal{K}(\mathbf{x}_1, \mathbf{x}_1) & \dots & \mathcal{K}(\mathbf{x}_1, \mathbf{x}_N) \\ \vdots & \ddots & \vdots \\ \mathcal{K}(\mathbf{x}_N, \mathbf{x}_1) & \dots & \mathcal{K}(\mathbf{x}_N, \mathbf{x}_N) \end{bmatrix}$$