

# Mid-term exam solution key

Prayag  
Neural networks and learning systems-I

May 21, 2019

## Problem 1.

*Solution.* 1. Consider the update equation

$$\bar{W}(n+1) = \bar{W}(n) - \eta \Delta \bar{W} \quad (1)$$

where  $\bar{W}(n)$  is the weight vector at time  $n$ ,  $\eta$  is the learning rate, and  $\Delta \bar{W}$  is the gradient of the cost function with respect to the weight vector  $\bar{W}(n)$ .

- (a) The learning rate must be increased when the derivative of the cost function with respect to a weight vector has the same algebraic sign to accelerate the convergence of the algorithm.
  - (b) The learning rate must be decreased when the algebraic sign of the derivative of cost function alternates with consecutive iterations. This reduces the oscillations during the convergence of the algorithm.
2. This follows from the composition of linear maps. Consider an affine linear map  $g(x) = ax + b$  where  $a$  and  $b$  are non-zero constants. Similarly  $f(x) = a'x + b'$  with  $a'$  and  $b'$  being non-zero constants. Consider the composition  $f(g(x)) = cx + d$  where  $c = aa'$  and  $d = a'b + b'$  which is again an affine linear map.
3. Linear regression:  $Y = \sigma_0 + \sigma_1 X + \epsilon$ , Quartic regression:  $Y = \sigma_0 + \sigma_1 X + \sigma_2 X^2 + \sigma_3 X^3 + \sigma_4 X^4 + \epsilon$ , LRSS: training residual sum of squares for linear regression, and QRSS: training residual sum of squares for quartic regression.

**Case 1:** Let  $\epsilon \neq 0$  and  $\mathbb{E} \epsilon = 0$ . In the case of linear regression,  $Y = f(x) + \epsilon$  we have

$$\begin{aligned} \mathbb{E}((y - \hat{y})^2) &= \mathbb{E}\left(f(x) + \epsilon - \hat{f}(x)\right)^2 \\ &= \mathbb{E}\left(f(x) - \hat{f}(x)\right)^2 + \mathbb{E}(\epsilon^2) + 2\mathbb{E}(\epsilon)\mathbb{E}\left(f(x) - \hat{f}(x)\right) \\ &= \underbrace{\mathbb{E}\left(f(x) - \hat{f}(x)\right)^2}_{\text{can be minimized}} + \underbrace{\mathbb{E}(\epsilon^2)}_{\text{cannot be minimized}} \end{aligned} \quad (2)$$

Since  $\text{Var}(x)$  cannot be minimized using linear regression and on the other hand the quartic regression is more flexible and can even fit the points with noise, therefore

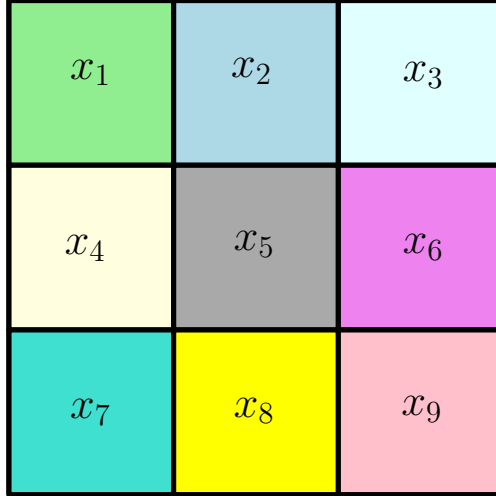


Figure 1: A  $3 \times 3$  image.

QRSS  $\leq$  LRSS.

**Case 2:** With  $\epsilon = 0$ , the relationship is truly linear i.e.,  $Y = \sigma_0 + \sigma_1 X$ . In such cases, LRSS and QRSS will be equal since both of the models can fit the data exactly.

4. Consider a  $3 \times 3$  image as shown in Figure 1 and a  $2 \times 2$  kernel as shown in Figure 2. Moving the kernel over the  $3 \times 3$  image we get 4 outputs which is given by

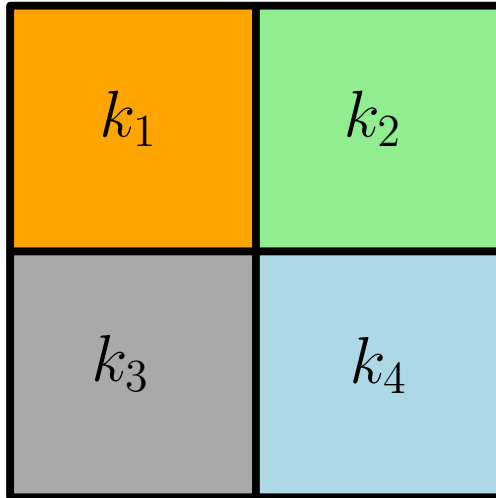


Figure 2: A  $2 \times 2$  kernel.

$$\begin{aligned}
 a_1 &= x_1 k_1 + x_2 k_2 + x_4 k_3 + x_5 k_4 \\
 a_2 &= x_2 k_1 + x_3 k_2 + x_5 k_3 + x_6 k_4 \\
 a_3 &= x_4 k_1 + x_5 k_2 + x_7 k_3 + x_8 k_4 \\
 a_4 &= x_5 k_1 + x_6 k_2 + x_8 k_3 + x_9 k_4.
 \end{aligned}$$

Graphical illustration of all the connections is shown in Figure 3

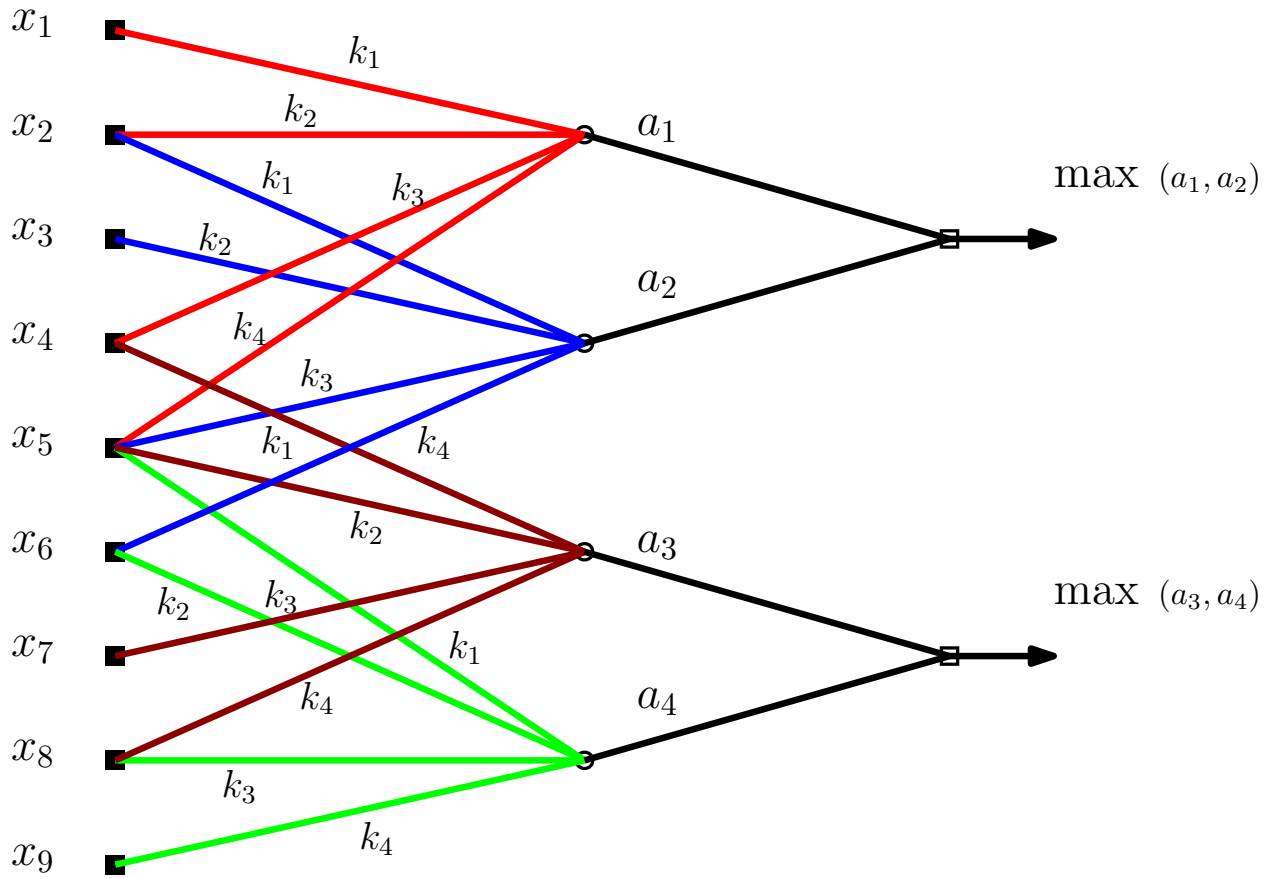


Figure 3

■

**Problem 2.**

*Solution.* The intermediate variables in the network shown in Figure 4 is given by

$$\begin{aligned}
 v_1 &= x_1 + x_2 + x_3 - 0.5 \\
 v_2 &= x_1 + x_2 + x_3 - 1.5 \\
 v_3 &= x_1 + x_2 + x_3 - 2.5 \\
 y_1 &= \phi(v_1) \\
 y_2 &= \phi(v_2) \\
 y_3 &= \phi(v_3)
 \end{aligned} \tag{3}$$

where

$$\begin{cases} 1 & x \geq 0 \\ 0 & \text{Otherwise} \end{cases} \tag{4}$$

■

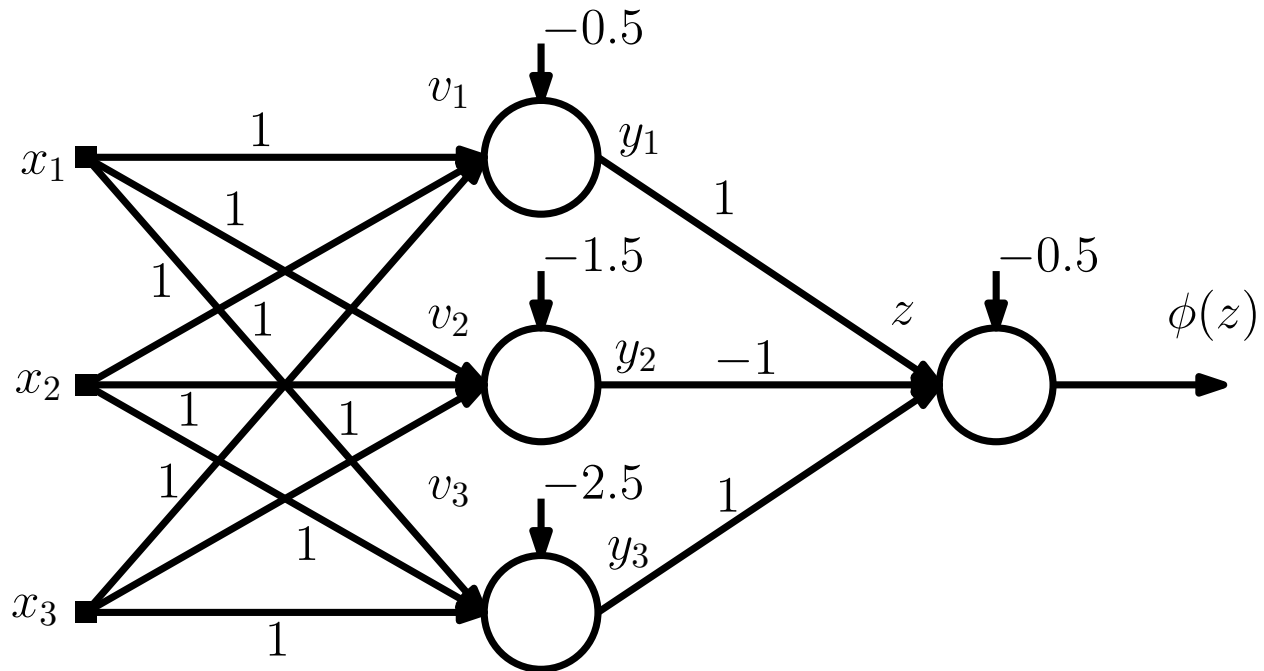


Figure 4: 3-bit XOR network.

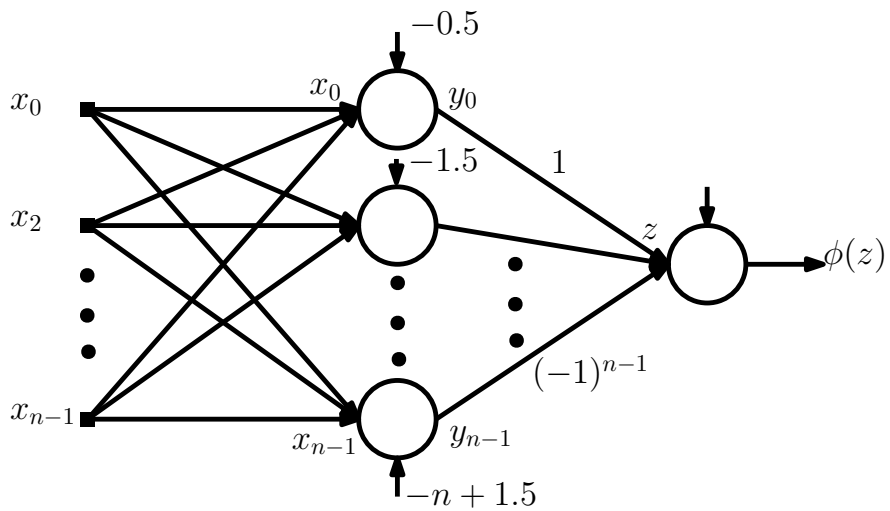


Figure 5:  $n$ -bit XOR network.

Table 1: Intermediate variables in the network

$x_1$	$x_2$	$x_3$	$v_1$	$v_2$	$v_3$	$y_1$	$y_2$	$y_3$	$z$	$\phi(z)$
0	0	0	-0.5	-1.5	-2.5	0	0	0	-0.5	0
0	0	1	0.5	-0.5	-1.5	1	0	0	0.5	1
0	1	0	0.5	-0.5	-1.5	1	0	0	0.5	1
0	1	1	1.5	0.5	-0.5	1	1	0	-0.5	0
1	0	0	0.5	-1.5	-1.5	1	0	0	0.5	1
1	0	1	1.5	0.5	-0.5	1	1	0	-0.5	0
1	1	0	1.5	0.5	-0.5	1	1	0	-0.5	0
1	1	1	2.5	1.5	0.5	1	1	1	0.5	1

**Problem 3.**

*Solution.* The cost function is given by

$$J(\bar{W}) = \sum_{x \in \mathcal{H}} \left( -\bar{W}^T x \right) \quad (5)$$

where  $\mathcal{H}$  is the of misclassified inputs. Differentiating  $J(\bar{W})$  with respect to  $\bar{W}(n)$  we get

$$\nabla_{\bar{W}} J(\bar{W}) = \sum_{x \in \mathcal{H}} -x \quad (6)$$

The weight vector is updated as follows:

$$\begin{aligned} \bar{W}(n+1) &= \bar{W}(n) - \eta(n) \nabla_{\bar{W}} J(\bar{W}) \\ &= \bar{W}(n) + \eta \sum_{x \in \mathcal{H}} x \end{aligned} \quad (7)$$

where  $\eta$  is assumed to remain same for all  $n$ , say  $\eta = 1$ . Let the initial weight vector be  $\bar{W}(0) = \bar{0}$ . Consider  $\bar{W}^T \bar{x} \leq 0$  and  $\bar{x}(n) \in \mathcal{H}$  is the set of misclassified samples. We know that

$$\begin{aligned} \bar{W}(n+1) &= \bar{W}(n) + \sum_{\bar{x} \in \mathcal{H}} \bar{x} \\ &= \sum_{\bar{x} \in \mathcal{H}_0} \bar{x} + \dots + \sum_{\bar{x} \in \mathcal{H}_n} \bar{x} \end{aligned} \quad (8)$$

let us consider a  $\bar{W}_0$  such that  $\bar{W}_0^T \bar{x} > 0$  for all  $\bar{x}(n)$  belongs to class 1. Pre-multiplying the above equation by  $\bar{W}_0^T$  we get

$$\bar{W}_0^T \bar{W}(n+1) = \sum_{\bar{x} \in \mathcal{H}_0} \bar{W}_0^T \bar{x} + \dots + \sum_{\bar{x} \in \mathcal{H}_n} \bar{W}_0^T \bar{x}. \quad (9)$$

Let  $\alpha = \min_i \sum_{\bar{x} \in \mathcal{H}_i} \bar{W}_0^T \bar{x}$ . Using Cauchy-Schwartz inequality we get

$$\begin{aligned} \|\bar{W}_0\|^2 \|\bar{W}_{n+1}\|^2 &\geq (n+1)^2 \alpha^2 \\ \|\bar{W}_{n+1}\|^2 &\geq \frac{(n+1)^2 \alpha^2}{\|\bar{W}_0\|^2} \end{aligned} \quad (10)$$

We know that,

$$\begin{aligned} \bar{W}(n+1) &= \bar{W}(n) + \sum_{\bar{x} \in \mathcal{H}_n} \bar{x} \\ \|\bar{W}(n+1)\|^2 &= \|\bar{W}(n)\|^2 + \left\| \sum_{\bar{x} \in \mathcal{H}_n} \bar{x} \right\|^2 + 2 \underbrace{\bar{W}^T(n) \sum_{\bar{x} \in \mathcal{H}_n} \bar{x}}_{\bar{W}^T(n) \bar{x} \leq 0} \\ \|\bar{W}(n+1)\|^2 &\leq \|\bar{W}(n)\|^2 + \left\| \sum_{\bar{x} \in \mathcal{H}_n} \bar{x} \right\|^2 \\ &\leq \sum_1^n \left\| \sum_{\bar{x} \in \mathcal{H}_n} \bar{x} \right\|^2 \\ &\leq (n+1)\beta \end{aligned} \quad (11)$$

where  $\beta = \max_i \left\| \sum_{\bar{x} \in \mathcal{H}_i} \bar{x} \right\|^2$ . Using equations (10) and (11) we get

$$n_{\max} = \left( \frac{\beta}{\alpha^2} \|\bar{W}_0\|^2 \right) - 1. \quad (12)$$

Therefore the batch perceptron algorithm converges after  $n_{\max}$  epochs. ■

#### Problem 4.

*Solution.* Consider through the Taylor series expansion

$$E_{\text{avg}}(\bar{W}(n) + \Delta \bar{W}(n)) = E_{\text{avg}}(\bar{W}(n)) + \bar{g}^T(n) \Delta \bar{W}(n) + \frac{1}{2} \Delta \bar{W}(n)^T \mathbf{H}(n) \Delta \bar{W}(n) + \text{h.o.t} \quad (13)$$

and neglecting the h.o.t we get

$$E_{\text{avg}}(\bar{W}(n) + \Delta \bar{W}(n)) = E_{\text{avg}}(\bar{W}(n)) + \bar{g}^T(n) \Delta \bar{W}(n) + \frac{1}{2} \Delta \bar{W}(n)^T \mathbf{H}(n) \Delta \bar{W}(n). \quad (14)$$

Differentiating the above equation with respect to  $\Delta \bar{W}(n)^T$  we get

$$\begin{aligned} \frac{\partial E_{\text{avg}}(\bar{W}(n) + \Delta \bar{W}(n))}{\partial \Delta \bar{W}(n)^T} &= 0 \\ \mathbf{H}(n) \Delta \bar{W}(n) &= -\bar{g}(n) \\ \Delta \bar{W}(n)^T &= -\mathbf{H}^{-1}(n) \bar{g}(n) \end{aligned} \quad (15)$$

provided  $\mathbf{H}^{-1}(n)$  exists. One can also get a pseudo-inverse of  $\mathbf{H}$  in case of singularity. The advantages of Hessian are as follows:

1. Accelerated convergence.
2. Possibly low rank approximations over  $\mathbf{H}$  to obtain low complexity algorithms (i.e., there is control on complexity).

