

1

a) $x = \{0; 1.5; 3; 4.5; 6\}$ $\phi_j(x) = x^j$
 $t = \{1, 0; -1, 0, 1\}$ $y(x, w) = \sum_{j=0}^3 w_j \cdot \phi_j(x) = w_0 + w_1 \cdot x + w_2 \cdot x^2 + w_3 \cdot x^3$

Design Matrix: $m \times (d+1)$
 \rightarrow # input features $\rightarrow 5 \times 4$
 \rightarrow # examples \rightarrow rows \rightarrow cols
 \rightarrow our regression uses 3 features!

As such, we have the matrix.

$$\Phi = \begin{bmatrix} 1 & \phi_1(x_1) & \phi_2(x_1) & \phi_3(x_1) \\ 1 & \phi_1(x_2) & \phi_2(x_2) & \phi_3(x_2) \\ 1 & \phi_1(x_3) & \phi_2(x_3) & \phi_3(x_3) \\ 1 & \phi_1(x_4) & \phi_2(x_4) & \phi_3(x_4) \\ 1 & \phi_1(x_5) & \phi_2(x_5) & \phi_3(x_5) \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 1 & 1.5 & 2.25 & 3.375 \\ 1 & 3 & 9 & 27 \\ 1 & 4.5 & 20.25 & 91.125 \\ 1 & 6 & 36 & 216 \end{bmatrix}$$

$\rightarrow \phi_1(x) = x^1 \rightarrow \phi_3(x) = x^3$
 $\rightarrow \phi_2(x) = x^2$

b)

With our Design matrix & our Target vector

$$X = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 1 & 1.5 & 2.25 & 3.375 \\ 1 & 3 & 9 & 27 \\ 1 & 4.5 & 20.25 & 91.125 \\ 1 & 6 & 36 & 216 \end{bmatrix} \quad Y = \begin{bmatrix} 1 \\ 0 \\ -1 \\ 0 \\ 1 \end{bmatrix}$$

We want the weight vector $w = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ w_3 \end{bmatrix}$, which minimizes the sum of squared errors.

Knowing that $w = (X^T X)^{-1} X^T Y$, we have:

$$w = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ 0 & 1.5 & 3 & 4.5 & 6 \\ 0 & 2.25 & 9 & 20.25 & 36 \\ 0 & 3.375 & 27 & 91.125 & 216 \end{bmatrix}^{-1} \cdot \begin{bmatrix} 1 & 0 & 0 & 0 \\ 1 & 1.5 & 2.25 & 3.375 \\ 1 & 3 & 9 & 27 \\ 1 & 4.5 & 20.25 & 91.125 \\ 1 & 6 & 36 & 216 \end{bmatrix} \cdot \begin{bmatrix} 1 \\ 0 \\ -1 \\ 0 \\ 1 \end{bmatrix}$$

$$= \begin{bmatrix} 5 & 15 & 67.5 & 337.5 \\ 15 & 67.5 & 337.5 & 1792.125 \\ 67.5 & 337.5 & 1792.125 & 9871.835 \\ 337.5 & 1792.125 & 9871.835 & 55700.156 \end{bmatrix}^{-1} \cdot \begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ 0 & 1.5 & 3 & 4.5 & 6 \\ 0 & 2.25 & 9 & 20.25 & 36 \\ 0 & 3.375 & 27 & 91.125 & 216 \end{bmatrix} \cdot \begin{bmatrix} 1 \\ 0 \\ -1 \\ 0 \\ 1 \end{bmatrix}$$

1

b (cont.)

$$= \begin{bmatrix} 0.99 & 0.06 & -0.09 & 0.06 & -0.01 \\ -1.00 & 1.08 & 0.38 & 0.70 & 0.23 \\ 0.29 & -0.49 & -0.06 & 0.41 & -0.16 \\ -0.02 & 0.05 & 0 & -0.05 & 0.02 \end{bmatrix} \cdot \begin{bmatrix} 1 \\ 0 \\ -1 \\ 0 \\ 1 \end{bmatrix}$$

$$\Rightarrow w = \begin{bmatrix} 1.057 \\ -1.143 \\ 0.190 \\ 0 \end{bmatrix} \approx \begin{bmatrix} 1.06 \\ -1.14 \\ 0.19 \\ 0 \end{bmatrix} \quad \checkmark$$

c) To compute the gradient of $E(w) = \frac{1}{2} \sum_{k=1}^M (t_k - w^T \phi_k)^2 + \frac{\lambda}{2} \|w\|_2^2$

we have.

$$\begin{aligned} \frac{\partial E(w)}{\partial w} &= \frac{\partial \left[\frac{1}{2} \sum_{k=1}^M (t_k - w^T \phi_k)^2 + \frac{\lambda}{2} \|w\|_2^2 \right]}{\partial w} = \frac{\partial \left[\frac{1}{2} \sum_{k=1}^M (t_k - w^T \phi_k)^2 \right]}{\partial w} + \frac{\partial \left[\frac{\lambda}{2} \|w\|_2^2 \right]}{\partial w} \\ &= \underbrace{\frac{1}{2} \frac{\partial \left[\sum_{k=1}^M (t_k - w^T \phi_k)^2 \right]}{\partial w}}_A + \underbrace{\frac{\lambda}{2} \frac{\partial [\|w\|_2^2]}{\partial w}}_B \end{aligned}$$

Dealing with A & B individually we have.

$$\begin{aligned} \textcircled{A} \frac{\partial \left[\sum_{k=1}^M (t_k - w^T \phi_k)^2 \right]}{\partial w} &\rightarrow \sum_{k=1}^M (t_k - w^T \phi_k)^2 \\ &= (T - Xw)^T (T - Xw) \\ &= \frac{\partial (T - Xw)^T (T - Xw)}{\partial w} \\ &= \left(\frac{\partial}{\partial w} (T - Xw)^T \right) (T - Xw) + (T - Xw)^T \left(\frac{\partial}{\partial w} (T - Xw) \right) \\ &= (-X^T) (T - Xw) + (T - Xw)^T (-X) \\ &= -2X^T (T - Xw) \end{aligned}$$

$$\frac{\delta(w^T A w)}{\delta w} = (A + A^T) w \quad \hookrightarrow I$$

1

C (cont.)

$$\begin{aligned} \textcircled{\beta} \quad \frac{\delta \|w\|_2^2}{\delta w} &= \frac{\delta \left[\sum_i w_i^2 \right]}{\delta w} = \frac{\delta [w^T w]}{\delta w} \quad \nearrow \sum w_i^2 = w^T w \quad \uparrow \\ &= \frac{\delta [w^T \cdot \overset{\nearrow = w}{I} w]}{\delta w} = (I + I^T) w = 2I w = 2w \quad \checkmark \end{aligned}$$

Joining things together we have

$$\frac{\delta E(w)}{\delta w} = \frac{1}{2} \cdot (2 X^T (T - Xw)) + \frac{\lambda}{2} (2w)$$

$$= -X^T (T - Xw) + \lambda w \quad \checkmark$$

$$d) \quad E(w) = \frac{1}{2} \sum_{n=1}^N (t_n - w^T x_n)^2 + \frac{\lambda}{2} \|w\|_2^2$$

$$\left(\begin{array}{l} \text{Ridge Regression:} \\ E(w) = \frac{1}{2} \sum_{n=1}^N (t_n - w^T x_n)^2 + \|w\|_2^2 \\ \quad \cdot \frac{1}{2} \sum_{n=1}^N (t_n - w^T x_n)^2 + I \frac{\lambda}{2} \|w\|_2^2 \\ \quad \cdot \frac{1}{2} \sum_{n=1}^N (t_n - w^T x_n)^2 + I \frac{\lambda}{2} \|w\|_2^2 \end{array} \right) \quad \times$$

$$\text{we have } X = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 1 & 1.5 & 2.25 & 3.375 \\ 1 & 3 & 9 & 27 \\ 1 & 4.5 & 20.25 & 91.125 \\ 1 & 6 & 36 & 216 \end{bmatrix} \quad Y = \begin{bmatrix} 1 \\ 0 \\ -1 \\ 0 \\ 1 \end{bmatrix}$$

$$\log\left(\frac{\lambda}{2}\right) = 0 \Leftrightarrow \frac{\lambda}{2} = 10^0 \Leftrightarrow \frac{\lambda}{2} = 1 \Leftrightarrow \lambda = 2$$

from c)

$$\frac{\delta E(w)}{\delta w} = 0 \Leftrightarrow -X^T (T - Xw) + \lambda w = 0 \quad \Leftrightarrow -X^T T + X^T X w + \lambda w = 0$$

$$\Leftrightarrow X^T X w + \lambda w = X^T T \quad \Leftrightarrow (X^T X + I \lambda) w = X^T T$$

$$\Leftrightarrow w = (X^T X + I \lambda)^{-1} X^T T$$

HW 2 - Page 4

14 April 2021

1

d (cont.)

$$\begin{aligned}
 W &= (X^T X + I \lambda)^{-1} X^T T = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ 0 & 1.5 & 3 & 4.5 & 6 \\ 0 & 2.25 & 9 & 20.25 & 36 \\ 0 & 3.375 & 27 & 91.125 & 216 \end{bmatrix} \cdot \begin{bmatrix} 1 & 0 & 0 & 0 \\ 1 & 1.5 & 2.25 & 3.375 \\ 1 & 3 & 9 & 27 \\ 1 & 4.5 & 20.25 & 91.125 \\ 1 & 6 & 36 & 216 \end{bmatrix} + \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \cdot 2 \cdot \begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ 0 & 1.5 & 3 & 4.5 & 6 \\ 0 & 2.25 & 9 & 20.25 & 36 \\ 0 & 3.375 & 27 & 91.125 & 216 \end{bmatrix} \cdot \begin{bmatrix} 1 \\ 0 \\ -1 \\ 0 \\ 1 \end{bmatrix} \\
 &= \begin{bmatrix} 5 & 15 & 67.5 & 337.5 \\ 15 & 67.5 & 337.5 & 1742.125 \\ 67.5 & 337.5 & 1742.125 & 9871.835 \\ 337.5 & 1742.125 & 9871.835 & 55700.156 \end{bmatrix} + \begin{bmatrix} 2 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 \\ 0 & 0 & 2 & 0 \\ 0 & 0 & 0 & 2 \end{bmatrix} \cdot \begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ 0 & 1.5 & 3 & 4.5 & 6 \\ 0 & 2.25 & 9 & 20.25 & 36 \\ 0 & 3.375 & 27 & 91.125 & 216 \end{bmatrix} \cdot \begin{bmatrix} 1 \\ 0 \\ -1 \\ 0 \\ 1 \end{bmatrix} \\
 &= \begin{bmatrix} 0.29 & -0.07 & -0.02 & 0.00 \\ -0.07 & 0.33 & -0.15 & 0.01 \\ -0.02 & -0.15 & 0.09 & -0.01 \\ 0.00 & 0.01 & -0.01 & 0.00 \end{bmatrix} \cdot \begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ 0 & 1.5 & 3 & 4.5 & 6 \\ 0 & 2.25 & 9 & 20.25 & 36 \\ 0 & 3.375 & 27 & 91.125 & 216 \end{bmatrix} \cdot \begin{bmatrix} 1 \\ 0 \\ -1 \\ 0 \\ 1 \end{bmatrix} \Rightarrow W = \begin{bmatrix} 0.295 \\ -0.132 \\ -0.138 \\ 0.030 \end{bmatrix}
 \end{aligned}$$

e) To compute the gradient of $E(w) = \frac{1}{2} \sum_{k=1}^M (t_k - w^T \phi_k)^2 + \lambda \|w\|_1$

we have.

$$\begin{aligned}
 \frac{\partial E(w)}{\partial w} &= \frac{\partial}{\partial w} \left[\frac{1}{2} \sum_{k=1}^M (t_k - w^T \phi_k)^2 + \lambda \|w\|_1 \right] = \frac{\partial}{\partial w} \left[\frac{1}{2} \sum_{k=1}^M (t_k - w^T \phi_k)^2 \right] + \frac{\partial}{\partial w} [\lambda \|w\|_1] \\
 &= \underbrace{\frac{1}{2} \frac{\partial}{\partial w} \left[\sum_{k=1}^M (t_k - w^T \phi_k)^2 \right]}_A + \underbrace{\lambda \frac{\partial}{\partial w} [\|w\|_1]}_B
 \end{aligned}$$

Dealing with A & B individually we have.

$$\begin{aligned}
 \textcircled{A} \quad \frac{\partial}{\partial w} \left[\sum_{k=1}^M (t_k - w^T \phi_k)^2 \right] &\rightarrow \sum_{k=1}^M (t_k - w^T \phi_k)^2 \\
 &= (T - Xw)^T (T - Xw) \\
 &= \frac{\partial}{\partial w} (T - Xw)^T (T - Xw) \\
 &= \left(\frac{\partial}{\partial w} (T - Xw)^T \right) (T - Xw) + (T - Xw)^T \left(\frac{\partial}{\partial w} (T - Xw) \right) \\
 &= (-X^T) (T - Xw) + (T - Xw)^T (-X) \\
 &= -2X^T (T - Xw)
 \end{aligned}$$

1

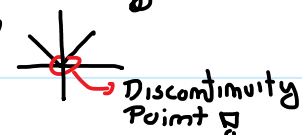
e (cont.)

$$\textcircled{\beta} \frac{\partial \|w\|_1}{\partial w} = \frac{\partial \left[\sum_i |w_i| \right]}{\partial w} \rightarrow \frac{\partial [1^T |w|]}{\partial w}$$

↑ now vector of 1's

$$= 1^T \cdot \frac{\partial [|w|]}{\partial w} \quad \text{which is not differentiable!}$$

$f(w) = |w| \rightarrow$



Joining things together we have

$$\frac{\partial E(w)}{\partial w} = \frac{1}{2} \cdot (-2 X^T (T - Xw)) + \lambda \cdot 1^T \frac{\partial [|w|]}{\partial w}$$

f)

The LASSO regression lacks a closed form solution. By definition, a closed form solution is one in which there are no limit, **differentiation** or integration. To perform this regularization we would have to perform $\frac{\partial E(w)}{\partial w} = 0$ where $E(w) = \frac{1}{2} \sum_{k=1}^n (t_k - w_k \cdot \phi_k)^2 + \lambda \|w\|_1$. From exercise 1.e), we've

already seen that when computing the gradient of $E(w)$, we get a formula that still maintains a differentiation (due to $\|w\|_1$ not being differentiable). As such, LASSO regression, does not possess a closed form solution.

2

$$\text{a) } w^{[0]} = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 & 1 \end{bmatrix} \quad b^{[1]} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} \quad // \quad w^{[2]} = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix} \quad b^{[2]} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad // \quad w^{[3]} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} \quad b^{[3]} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

$$f(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}} = \tanh(x) \quad m = 0.1 \quad x = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}, \quad t = \begin{bmatrix} 1 \\ -1 \end{bmatrix}$$

1 Forward Propagation

$$z^{[1]} = w^{[1]} x^{[0]} + b^{[1]} \Rightarrow \begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \end{bmatrix} = \begin{bmatrix} 5 \\ 5 \end{bmatrix}$$

↑ element wise

$$x^{[1]} = f(z^{[1]}) = \begin{bmatrix} \tanh(5) \\ \tanh(5) \\ \tanh(5) \\ \tanh(5) \\ \tanh(5) \end{bmatrix} \approx \begin{bmatrix} 0.9999 \\ 0.9999 \\ 0.9999 \\ 0.9999 \\ 0.9999 \end{bmatrix}$$

2

a (cont.)

$$\bullet z^{[2]} = w^{[2]} x^{[1]} + b^{[2]} \Rightarrow \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix} \cdot \begin{bmatrix} 0.9999 \\ 0.9999 \\ 0.9999 \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \end{bmatrix} = \begin{bmatrix} 2.9997 \\ 2.9997 \end{bmatrix}$$

$$x^{[2]} = f(z^{[2]}) = \begin{bmatrix} \tanh(2.9997) \\ \tanh(2.9997) \end{bmatrix} = \begin{bmatrix} 0.9951 \\ 0.9951 \end{bmatrix}$$

$$\bullet z^{[3]} = w^{[3]} x^{[2]} + b^{[3]} \Rightarrow \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} \cdot \begin{bmatrix} 0.9951 \\ 0.9951 \end{bmatrix} + \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 2.9902 \\ 2.9902 \end{bmatrix}$$

$$x^{[3]} = f(z^{[3]}) = \begin{bmatrix} \tanh(2.9902) \\ \tanh(2.9902) \end{bmatrix} = \begin{bmatrix} 0.995 \\ 0.995 \end{bmatrix}$$

2 Back Propagation

Squared Error loss $E(t, x) = \frac{1}{2}(x - t)^2$

$$\rightarrow \frac{\partial E}{\partial x^{[L]}}(t, x^{[L]}) = \frac{\partial E}{\partial (x^{[L]} - t)^2} \cdot \frac{\partial (x^{[L]} - t)^2}{\partial (x^{[L]} - t)} \cdot \frac{\partial (x^{[L]} - t)}{\partial x^{[L]}} = \frac{1}{2} [2(x^{[L]} - t)] = x^{[L]} - t$$

$$\rightarrow \frac{\partial x^{[L]}}{\partial z^{[L]}}(z^{[L]}) = 1 - \tanh(z^{[L]})^2$$

$$\rightarrow \frac{\partial z^{[L]}}{\partial w^{[L]}}(w^{[L]}, b^{[L]}, x^{[L-1]}) = x^{[L-1]}$$

$$\rightarrow \frac{\partial z^{[L]}}{\partial b^{[L]}}(w^{[L]}, b^{[L]}, x^{[L-1]}) = 1$$

$$\rightarrow \frac{\partial z^{[L]}}{\partial x^{[L-1]}}(w^{[L]}, b^{[L]}, x^{[L-1]}) = w^{[L]}$$

Let us begin the recursion.

$$\bullet \delta^{[3]} = \frac{\partial E}{\partial x^{[3]}} \circ \frac{\partial x^{[3]}}{\partial z^{[3]}} = (x^{[3]} - t) \circ (1 - \tanh(z^{[3]})^2) = \begin{bmatrix} 0.995 \\ 0.995 \end{bmatrix} \cdot \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 0.995 \\ 0.995 \end{bmatrix} \approx \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

$$\bullet \delta^{[2]} = \frac{\partial z^{[3]}}{\partial x^{[2]}}^T \cdot \delta^{[3]} \circ \frac{\partial x^{[2]}}{\partial z^{[2]}} = (w^{[3]})^T \cdot \delta^{[3]} \circ (1 - \tanh(z^{[2]})^2) = \begin{bmatrix} 1 & 1 \end{bmatrix}^T \cdot \begin{bmatrix} 1 \\ 1 \end{bmatrix} \circ \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 1 & 1 \end{bmatrix}$$

$$\bullet \delta^{[1]} = \frac{\partial z^{[2]}}{\partial x^{[1]}}^T \cdot \delta^{[2]} \circ \frac{\partial x^{[1]}}{\partial z^{[1]}} = (w^{[2]})^T \cdot \delta^{[2]} \circ (1 - \tanh(z^{[1]})^2) = \begin{bmatrix} 1 & 1 & 1 \end{bmatrix}^T \cdot \begin{bmatrix} 1 & 1 \end{bmatrix} \circ \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 2 \\ 2 \\ 2 \end{bmatrix}$$

2

a (cont.) \rightarrow Update

$$\begin{aligned}
 &= \delta^{[1]} \cdot \frac{\partial z^{[1]T}}{\partial w^{[1]}} \cdot \delta^{[1]} \cdot (x^{[0]})^T = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \cdot \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}^T = \begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} \\
 \bullet \quad w^{[1]} &= w^{[1]} - \eta \frac{\partial E}{\partial w^{[1]}} = w^{[1]} - 0.1 \begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} \Rightarrow w^{[1]} = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \end{bmatrix} \\
 \bullet \quad b^{[1]} &= b^{[1]} - \eta \frac{\partial E}{\partial b^{[1]}} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} - 0.1 \cdot \begin{bmatrix} 0 \\ 0 \end{bmatrix} = \boxed{b^{[1]} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}} \\
 &= \delta^{[1]} \cdot \frac{\partial z^{[1]T}}{\partial b^{[1]}} = \delta^{[1]}
 \end{aligned}$$

$$\begin{aligned}
 \bullet \quad w^{[2]} &= w^{[2]} - \eta \frac{\partial E}{\partial w^{[2]}} = w^{[2]} - \eta (\delta^{[2]} x^{[1]T}) = w^{[2]} - \eta \left(\begin{bmatrix} 0.0001 \\ 0.0001 \end{bmatrix} \cdot \begin{bmatrix} 0.9999 & 0.9999 \end{bmatrix}^T \right) \\
 &= w^{[2]} - 0.1 \begin{bmatrix} 0.0001 & 0.0001 & 0.0001 \\ 0.0001 & 0.0001 & 0.0001 \end{bmatrix} \Rightarrow \boxed{w^{[2]} \approx \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}}
 \end{aligned}$$

$$\bullet \quad b^{[2]} = b^{[2]} - \eta \frac{\partial E}{\partial b^{[2]}} = b^{[2]} - \eta \delta^{[2]} = b^{[2]} - 0.1 \begin{bmatrix} 0.0001 \\ 0.0001 \end{bmatrix} \Rightarrow \boxed{b^{[2]} \approx \begin{bmatrix} 0 \\ 0 \end{bmatrix}}$$

$$\begin{aligned}
 \bullet \quad w^{[3]} &= w^{[3]} - \eta \frac{\partial E}{\partial w^{[3]}} = w^{[3]} - \eta (\delta^{[3]} x^{[2]T}) = w^{[3]} - \eta \left(\begin{bmatrix} 0 \\ 0.02 \end{bmatrix} \cdot \begin{bmatrix} 0.9951 & 0.9951 \end{bmatrix}^T \right) \\
 &= w^{[3]} - 0.1 \begin{bmatrix} 0 & 0 & 0.02 \\ 0 & 0.02 & 0.02 \end{bmatrix} = \boxed{w^{[3]} = \begin{bmatrix} 1 & 1 & 1 \\ 0.998 & 0.998 & 1 \end{bmatrix}}
 \end{aligned}$$

$$\bullet \quad b^{[3]} = b^{[3]} - \eta \frac{\partial E}{\partial b^{[3]}} = b^{[3]} - \eta \delta^{[3]} = b^{[3]} - 0.1 \begin{bmatrix} 0 \\ 0.02 \end{bmatrix} \Rightarrow \boxed{b^{[3]} = \begin{bmatrix} 1 \\ 0.998 \end{bmatrix}}$$

b) Cross-Entropy loss measures the performance of a classification model whose output is a probability value between 0 and 1. It aims to minimize the distance between two **probability distributions**, making it a good loss function for classification problems. Looking at our target vector given in this exercise, we can see that it constitutes of $[1, -1]^T$, and as such, it clearly does not fall under this category. For Cross-Entropy Loss to be favourable, our outcome vector would have to have all its values between 0 and 1, and they would need to sum up to a total of 1 (making them probabilistic values).

$$c) \quad w^{[0]} = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 & 1 \end{bmatrix} \quad b^{[0]} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} \quad // \quad w^{[1]} = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix} \quad b^{[1]} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad // \quad w^{[2]} = \begin{bmatrix} 1 & 1 & 1 \\ 0.998 & 0.998 & 1 \end{bmatrix} \quad b^{[2]} = \begin{bmatrix} 1 \\ 0.998 \end{bmatrix}$$

$$\begin{aligned}
 f(x) &= \frac{e^x - e^{-x}}{e^x + e^{-x}} = \tanh(x) & \eta &= 0.1 & x &= \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}, \quad t = \begin{bmatrix} 1 \\ 0 \end{bmatrix} \\
 \text{Softmax}(z) &= x \rightarrow x, \quad = \frac{e^{z_i}}{\sum_{j=1}^n e^{z_j}}
 \end{aligned}$$

2

C (cont.)

1 Find derivative of a given x_i with respect to z_i (softMax)

$$\frac{\partial x_i}{\partial z_j} = \frac{\partial}{\partial z_j} \frac{e^{z_i}}{\sum_{k=1}^d e^{z_k}}$$

$$\begin{aligned} \text{For } i=j: \quad \frac{\partial x_i}{\partial z_j} &= \frac{\partial}{\partial z_j} \frac{e^{z_i}}{\sum_{k=1}^d e^{z_k}} = \left(\frac{\partial}{\partial z_j} e^{z_i} \right) \left(\sum_{k=1}^d e^{z_k} \right) - (e^{z_i}) \left(\frac{\partial}{\partial z_j} \sum_{k=1}^d e^{z_k} \right) \\ &= \frac{e^{z_i} \sum_{k=1}^d e^{z_k} - e^{z_i} e^{z_j}}{\left(\sum_{k=1}^d e^{z_k} \right)^2} = \frac{e^{z_i} \left(\sum_{k=1}^d e^{z_k} - e^{z_j} \right)}{\left(\sum_{k=1}^d e^{z_k} \right)^2} = \frac{e^{z_i}}{\sum_{k=1}^d e^{z_k}} \left(\frac{\sum_{k=1}^d e^{z_k} - e^{z_j}}{\sum_{k=1}^d e^{z_k}} \right) \\ &= x_i (1 - x_j) = x_i (1 - x_i) \end{aligned}$$

$$\begin{aligned} \text{For } i \neq j: \quad \frac{\partial x_i}{\partial z_j} &= \frac{\partial}{\partial z_j} \frac{e^{z_i}}{\sum_{k=1}^d e^{z_k}} = e^{z_i} \frac{\partial}{\partial z_j} \frac{1}{\sum_{k=1}^d e^{z_k}} = e^{z_i} \left(- \frac{1}{\left(\sum_{k=1}^d e^{z_k} \right)^2} \right) e^{z_j} \\ &= -x_i x_j \end{aligned}$$

2 Forward Propagation

$$z^{[1]} = W^{[1]} x^{[0]} + b^{[1]} \Rightarrow \begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} = \begin{bmatrix} 5 \\ 5 \\ 5 \end{bmatrix}$$

$$x^{[1]} = \text{element wise } f(z^{[1]}) = \begin{bmatrix} \tanh(5) \\ \tanh(5) \\ \tanh(5) \end{bmatrix} \approx \begin{bmatrix} 0.9999 \\ 0.9999 \\ 0.9999 \end{bmatrix}$$

$$z^{[2]} = W^{[2]} x^{[1]} + b^{[2]} \Rightarrow \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} 0.9999 \\ 0.9999 \\ 0.9999 \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \end{bmatrix} = \begin{bmatrix} 2.9997 \\ 2.9997 \end{bmatrix}$$

$$x^{[2]} = f(z^{[2]}) = \begin{bmatrix} \tanh(2.9997) \\ \tanh(2.9997) \end{bmatrix} = \begin{bmatrix} 0.9950 \\ 0.9950 \end{bmatrix}$$

$$z^{[3]} = W^{[3]} x^{[2]} + b^{[3]} \Rightarrow \begin{bmatrix} 1 & 1 \\ 0.998 & 0.998 \end{bmatrix} \begin{bmatrix} 0.995 \\ 0.995 \end{bmatrix} + \begin{bmatrix} 1 \\ 0.998 \end{bmatrix} = \begin{bmatrix} 2.990 \\ 2.984 \end{bmatrix}$$

$$x^{[3]} = \text{softmax}(z) = \begin{bmatrix} \frac{e^{2.990}}{e^{2.990} + e^{2.984}} \\ \frac{e^{2.984}}{e^{2.990} + e^{2.984}} \end{bmatrix} = \begin{bmatrix} 0.5015 \\ 0.4985 \end{bmatrix}$$

2

C (cont.) \rightarrow Back Propagation

Cross-Entropy loss $E(t, x) = \sum_{i=1}^d t_i \log x_i$

$$\begin{aligned} \rightarrow \frac{\partial E}{\partial z_i}(t, x^{[3]}) &= \delta_i^{[3]} = \frac{\partial}{\partial z_i} \left(- \sum_{k=1}^d t_k \log x_k^{[3]} \right) = - \sum_{k=1}^d t_k \frac{\partial}{\partial z_i} \log x_k^{[3]} \\ &= - \sum_{k=1}^d t_k \frac{1}{x_k^{[3]}} \frac{\partial x_k^{[3]}}{\partial z_i} = - \sum_{k=1}^d t_k \frac{1}{x_k^{[3]}} \frac{\partial x_k^{[3]}}{\partial z_i} - \sum_{k \neq i} t_k \frac{1}{x_k^{[3]}} \frac{\partial x_k^{[3]}}{\partial z_i} \\ &= -t_i \frac{1}{x_i^{[3]}} (x_i^{[3]} (1 - x_i^{[3]})) - \sum_{k \neq i} t_k \frac{1}{x_k^{[3]}} (-x_k^{[3]} x_i^{[3]}) \\ &= -t_i + t_i x_i^{[3]} + \sum_{k \neq i} t_k x_i^{[3]} = -t_i + x_i^{[3]} (t_i + \sum_{k \neq i} t_k) = -t_i + x_i^{[3]} \left(\sum_{k=1}^d t_k \right) \\ &= x_i^{[3]} - t_i \end{aligned}$$

$$\rightarrow \frac{\partial x^{[1]}}{\partial z^{[0]}}(z^{[0]}) = 1 - \tanh(z^{[0]})^2$$

$$\rightarrow \frac{\partial z^{[1]}}{\partial w^{[0]}}(w^{[0]}, b^{[0]}, x^{[0-1]}) = x^{[0-1]}$$

$$\rightarrow \frac{\partial z^{[0]}}{\partial b^{[-1]}}(w^{[-1]}, b^{[-1]}, x^{[-1-0]}) = 1$$

$$\rightarrow \frac{\partial z^{[1]}}{\partial x^{[0-1]}}(w^{[0]}, b^{[0]}, x^{[0-1]}) = w^{[0]}$$

To start the recursion we need the $\delta^{[3]}$

$$\delta^{[3]} = \begin{bmatrix} \delta_1^{[3]} \\ \delta_2^{[3]} \end{bmatrix} = \begin{bmatrix} x_1^{[3]} - t_1 \\ x_2^{[3]} - t_2 \end{bmatrix} = \begin{bmatrix} 0.5015 \\ 0.4985 \end{bmatrix} - \begin{bmatrix} 1 \\ 0 \end{bmatrix} = \begin{bmatrix} -0.4985 \\ 0.4985 \end{bmatrix}$$

Now we can compute the remaining δ

$$\begin{aligned} \bullet \delta^{[2]} &= \frac{\partial z^{[3]}}{\partial x^{[2]}} \delta^{[3]} \circ \frac{\partial x^{[3]}}{\partial z^{[2]}} = (w^{[2]})^T \cdot \delta^{[3]} \circ (1 - \tanh(z^{[2]})^2) \\ &= \begin{bmatrix} 1 & 0.998 \\ 1 & 0.998 \end{bmatrix} \begin{bmatrix} -0.4985 \\ 0.4985 \end{bmatrix} \circ \left(\begin{bmatrix} 1 \\ 1 \end{bmatrix} - \begin{bmatrix} \tanh(0.9950)^2 \\ \tanh(0.9950)^2 \end{bmatrix} \right) = \begin{bmatrix} -0.0004 \\ -0.0004 \end{bmatrix} \end{aligned}$$

$$\begin{aligned} \bullet \delta^{[1]} &= \frac{\partial z^{[2]}}{\partial x^{[1]}} \delta^{[2]} \circ \frac{\partial x^{[2]}}{\partial z^{[1]}} = (w^{[1]})^T \cdot \delta^{[2]} \circ (1 - \tanh(z^{[1]})^2) \\ &= \begin{bmatrix} 1 & 1 & 1 \end{bmatrix} \cdot \begin{bmatrix} -0.0004 \\ -0.0004 \\ -0.0004 \end{bmatrix} \circ \left(\begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} - \begin{bmatrix} \tanh(5)^2 \\ \tanh(5)^2 \\ \tanh(5)^2 \end{bmatrix} \right) \approx \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} \end{aligned}$$

2

c (cont.) 4 Update

$$\begin{aligned}
 \bullet \quad w^{[1]} &= w^{[1]} - \eta \frac{\partial E}{\partial w^{[1]}} = w^{[1]} - \eta \left(\delta^{[1]} \frac{\partial z^{[1]}}{\partial w^{[1]}} \right) \\
 &= w^{[1]} - \eta \left(\delta^{[1]} x^{[0]T} \right) = w^{[1]} - \eta \left(\begin{bmatrix} 0 \\ 0 \end{bmatrix} [1 \ 1 \ 1 \ 1 \ 1] \right) \\
 &= w^{[1]} - \eta \left(\begin{bmatrix} 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix} \right) \Rightarrow \boxed{w^{[1]} = w^{[1]} = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 & 1 \end{bmatrix}}
 \end{aligned}$$

$$\bullet \quad b^{[1]} = b^{[1]} - \eta \frac{\partial E}{\partial b^{[1]}} = b^{[1]} - \eta \delta^{[1]} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} - 0.1 \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Rightarrow \boxed{b^{[1]} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}}$$

$$\begin{aligned}
 \bullet \quad w^{[2]} &= w^{[2]} - \eta \frac{\partial E}{\partial w^{[2]}} = w^{[2]} - \eta \left(\delta^{[2]} \frac{\partial z^{[2]}}{\partial w^{[2]}} \right) \\
 &= w^{[2]} - \eta \left(\delta^{[2]} x^{[1]T} \right) = w^{[2]} - \eta \left(\begin{bmatrix} -0.0004 \\ -0.0004 \end{bmatrix} \begin{bmatrix} 0.9999 & 0.9999 & 0.9999 & 0.9999 \end{bmatrix}^T \right) \\
 &= w^{[2]} - 0.1 \begin{bmatrix} -0.0004 & -0.0004 & -0.0004 & -0.0004 \\ -0.0004 & -0.0004 & -0.0004 & -0.0004 \end{bmatrix} \Rightarrow \boxed{w^{[2]} = \begin{bmatrix} 1.0004 & 1.0004 & 1.0004 & 1.0004 \\ 1.0004 & 1.0004 & 1.0004 & 1.0004 \end{bmatrix}}
 \end{aligned}$$

$$\bullet \quad b^{[2]} = b^{[2]} - \eta \frac{\partial E}{\partial b^{[2]}} = b^{[2]} - \eta \delta^{[2]} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} - 0.1 \begin{bmatrix} -0.0004 \\ -0.0004 \end{bmatrix} \Rightarrow \boxed{b^{[2]} = \begin{bmatrix} 0.0004 \\ 0.0004 \end{bmatrix}}$$

$$\begin{aligned}
 \bullet \quad w^{[3]} &= w^{[3]} - \eta \frac{\partial E}{\partial w^{[3]}} = w^{[3]} - \eta \left(\delta^{[3]} \frac{\partial z^{[3]}}{\partial w^{[3]}} \right) = w^{[3]} - \eta \left(\delta^{[3]} x^{[2]T} \right) \\
 &= w^{[3]} - \eta \left(\begin{bmatrix} -0.4965 \\ 0.4965 \end{bmatrix} \cdot \begin{bmatrix} 0.9950 & 0.9950 \end{bmatrix}^T \right) = \begin{bmatrix} 1.050 & 1.050 \\ 0.9484 & 0.9484 \end{bmatrix} - 0.1 \begin{bmatrix} -0.4960 & -0.4960 \\ 0.4960 & 0.4960 \end{bmatrix} \\
 \Rightarrow \boxed{w^{[3]} = \begin{bmatrix} 1.050 & 1.050 \\ 0.9484 & 0.9484 \end{bmatrix}}
 \end{aligned}$$

2

C (cont.)

$$\bullet \quad b^{[3]} = b^{[2]} - \eta \frac{\partial E}{\partial b^{[2]}} = b^{[2]} - \eta \cdot \delta^{[2]} = \begin{bmatrix} 1 \\ 0.998 \end{bmatrix} - 0.1 \begin{bmatrix} 0.4985 \\ 0.4985 \end{bmatrix} \Rightarrow \boxed{b^{[3]} = \begin{bmatrix} 1.04985 \\ 0.94815 \end{bmatrix}}$$

♡