

# Neural Networks for Classification

---

*EE698V - Machine Learning for Signal Processing*

Vipul Arora



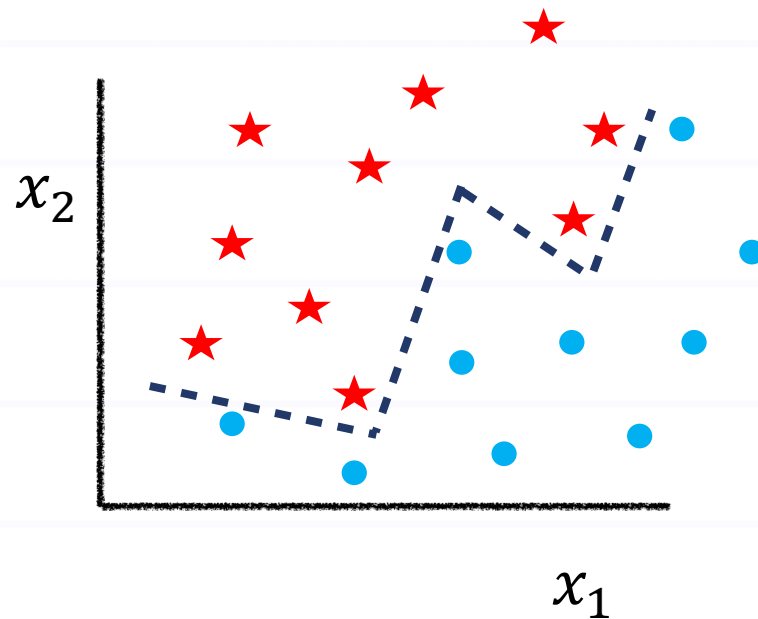
# Non-linear Classification

**Reference:** PRML Sec 5.1, 5.2

# Non-linear Classifiers

---

- Remember, non-linear models can be more powerful than linear ones



# Multi-class Classification

$$\begin{array}{c}
 \mathbb{R} \\
 \begin{matrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{matrix} \\
 \begin{matrix} i_0 & i_1 & i_2 \end{matrix}
 \end{array}
 \quad
 \mathbf{h} = \sum_{i=1}^C \omega_{i,c} v_{i,c}
 \quad
 \mathbf{y} = \sigma(\mathbf{h})$$

- The output of model is  $\mathbf{h}$  or  $h_c \in \mathbb{R}, c = 0, 1, 2, \dots, C-1$
- $y_c = \operatorname{argmax}_c h_c$ ,  $\mathbf{y}$  is a one-hot vector
- $\operatorname{argmax}$  is not differentiable, so we need a soft version

- $\mathbf{y} = \operatorname{softmax}(\mathbf{h})$ ; i.e.,  $y_c = \frac{e^{h_c}}{\sum_{c'} e^{h_{c'}}}$



$$\mathbf{h} = \begin{bmatrix} 2.1 \\ 0.3 \\ 0.8 \end{bmatrix}
 \quad
 \mathbf{y} = \begin{bmatrix} e^{2.1}/\text{den} \\ e^{0.3}/\text{den} \\ e^{0.8}/\text{den} \end{bmatrix} = \begin{bmatrix} .7 \\ .11 \\ .18 \end{bmatrix}
 \quad
 \text{den} = e^{2.1} + e^{0.3} + e^{0.8}$$

$$\mathbf{y} = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$$

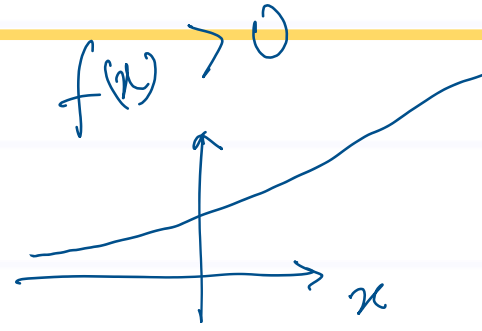
# Why exp in softmax?

$$y_c = \frac{h_c}{\sum_{c'} h_{c'}} \leftarrow \text{could be -ve or 0}$$

$$= \frac{2.1}{2.1 + 0.8 + 0.3}$$

$$= \frac{\sigma(2.1) \rightarrow 1}{\sigma(2.1) + \sigma(0.8) + \sigma(0.3)}$$

$\downarrow \quad \downarrow \quad \downarrow$   
 $1 \quad 0.9 \quad 0.5$

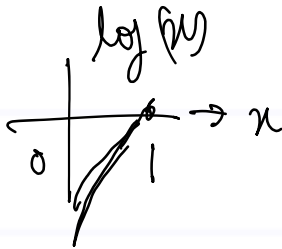
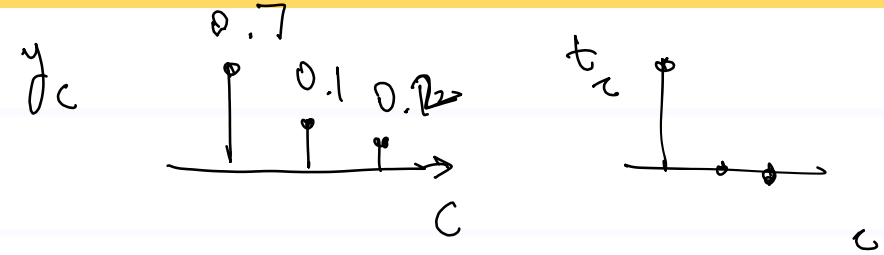


exp.

$$y = \begin{bmatrix} 1/2.8 \\ 0.9/2.4 \\ 0.5/2.4 \end{bmatrix}$$

# Multi-class Classification: Loss function

- Mean Squared Error?
- We need a stronger pull to 0 or 1
- $y_c$  can be seen as  $P(c)$  as  $\sum_c y_c = 1$ , so also  $t_c$
- Categorical cross-entropy

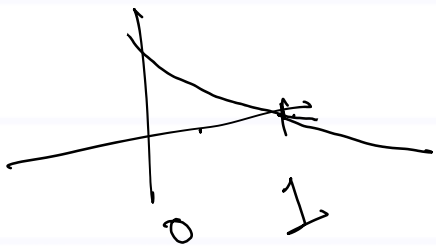
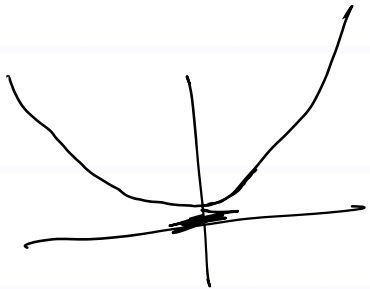


$$E_{Xent} = - \sum_c t_c \log(y_c) = - (1 \log 0.7 + 0 + 0)$$

# Why not MSE?

$$E_{MSE} = \frac{1}{2} \sum_c (t_c - y_c)^2$$

$$E_{xent} = \sum_c -t_c \log y_c$$



$$\frac{\partial E_{MSE}}{\partial y_{c'}} = (t_c - y_{c'}) \left( -\frac{\partial y_{c'}}{\partial y_{c'}} \right) = -0.3 \quad y_c \approx 0.7$$

$$\frac{\partial E_{xent}}{\partial y_{c'}} = -\frac{\delta_{cc'}}{y_c} \frac{\partial y_{c'}}{\partial y_{c'}} = -\frac{1}{0.7} = -1.4$$

$$\frac{\partial E_{xent}}{\partial y_{c'}} = -2$$

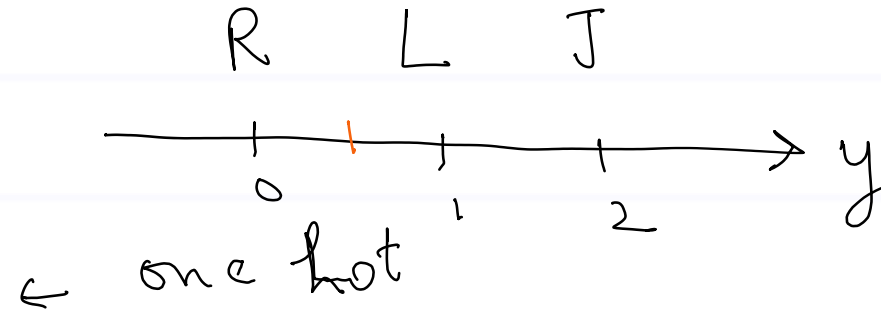
$$\frac{\partial E_{MSE}}{\partial y_{c'}} = -0.5$$

$$y_c \approx 0.5$$

# Why not represent target as $\{0,1,2\}$ ?

$$t = \{0, 1, 2\}$$

$$t = \begin{Bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{Bmatrix}$$



flower  $\sim R, L$

0.5  
0.5  
0

$\sim R, J$

0.5  
0  
0.5



# Multi-label Classification

attributes

|         |   |   |
|---------|---|---|
| beard   | 0 | 1 |
| glasses | 0 | 1 |
| hair    | 0 | 1 |

- The output of model is  $\mathbf{h}$  or  $h_l \in \mathbb{R}, l = 0, 1, 2, \dots, L - 1$

$$h_{i_2} = \sum_{i_1} w_{i_2 i_1} v_{i_1}$$

- $y_l = \text{sign}(h_l)$  or  $y_l = (\text{sign}(h_l) + 1) / 2$

- sign is not differentiable, so we need a soft version

- $y_l = \text{sigmoid}(h_l)$ ,  $\mathbf{y}$  is a multi-hot vector

|   |   |   |   |   |   |
|---|---|---|---|---|---|
| 1 | 0 | 0 | 1 | 1 | 0 |
| 1 | 1 | 0 | 0 | 1 | 0 |
| 1 | 1 | 1 | 1 | 0 | 0 |

$$\mathbf{h} = \begin{bmatrix} 2.1 \\ 0.8 \\ -0.3 \end{bmatrix} \quad \mathbf{y} = \begin{bmatrix} 0.8 \\ 0.25 \\ 0.4 \end{bmatrix}$$

# Multi-label Classification: Loss function

- Probabilistic interpretation:

- $P(l = 1) = y_l$

- $P(l = 0) = 1 - y_l$

$E = 0$

when  $y = 1$  for  $t = 1$

and  $y = 0$  for  $t = 0$

- Binary Cross-entropy

- $E_{binXent} = -\sum_l (t_l \log y_l + (1 - t_l) \log (1 - y_l))$

beard,  $t = 1$   
 $y = 0.7$

$$E_{binXent} = -1 \log 0.7$$

$t = 0$   
 $y = 0.7$   
 $E = -(1 - 0) \log 0.3$

# Optimization

---

- You can find optimal model parameters (or weights) using gradient descent
- All you need is  $\partial E / \partial \mathbf{w}$
- Can you find  $\partial E_{xent} / \partial h_c$  or  $\partial E_{binxent} / \partial h_l$ ?

# Perceptron Algorithm

---

- DIY
- Reference: PRML Section 4.1.7