

# **Softm** **ression**

## **Classification by minimising cross-entropy loss**

Assignment Project Exam Help

<https://eduassistpro.github.io/>

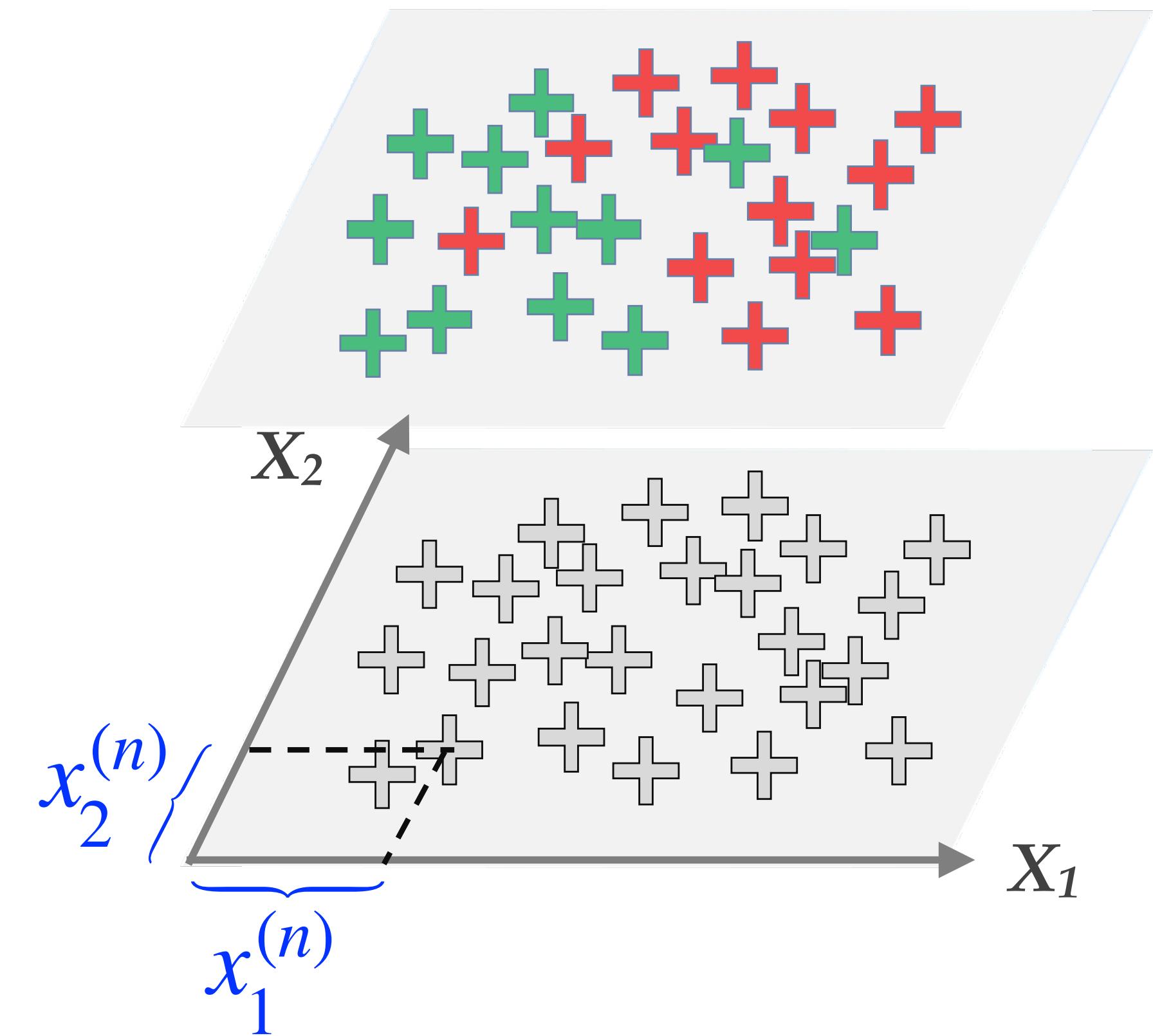
Add WeChat edu\_assist\_pro

**Srinandan Dasmahapatra**

# Classification: discrete output

Minimise deviation of prediction from annotation

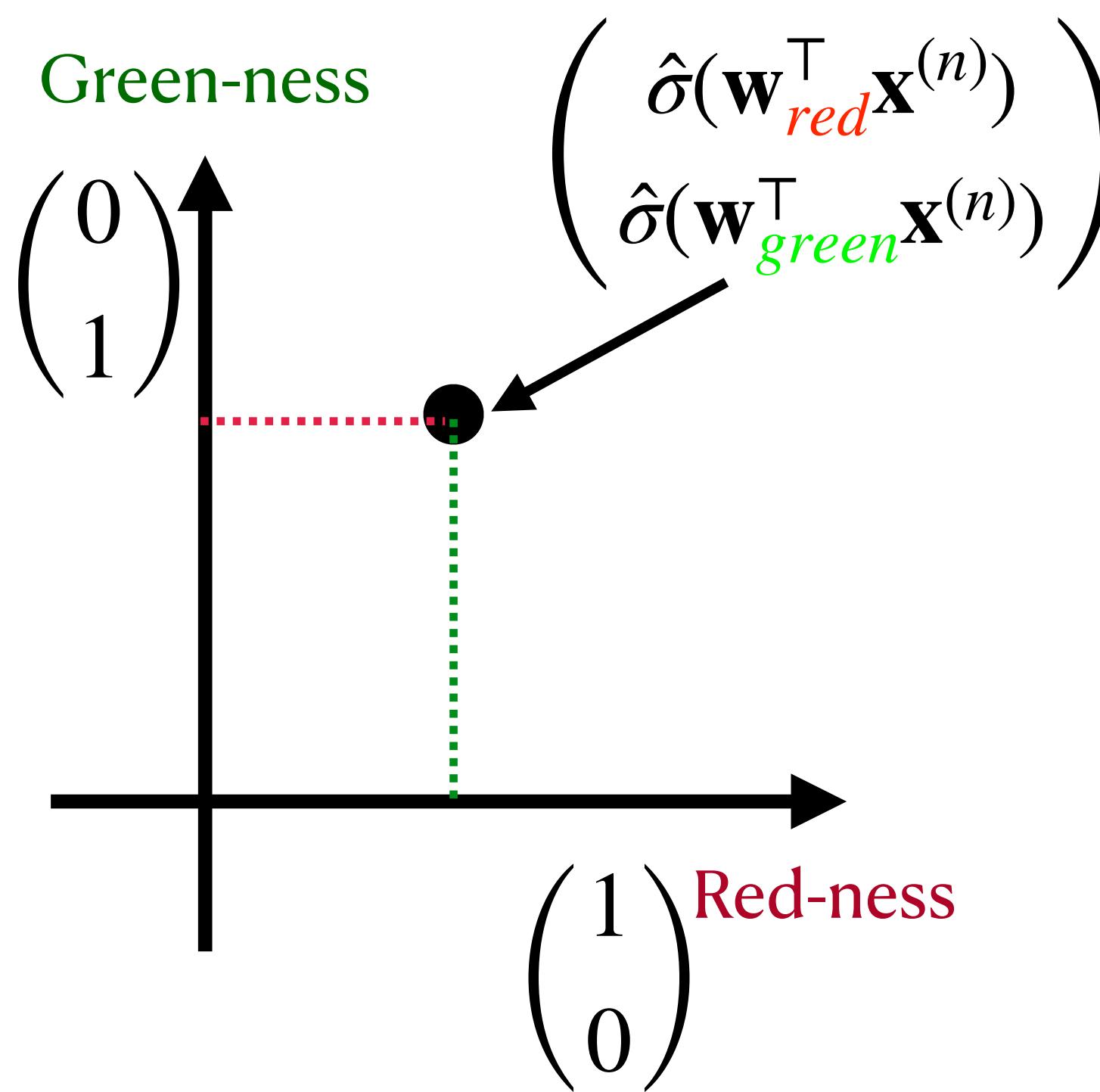
- Given training set represented by points labelled green and red, variable  $Y$  ...
- ... where each point has two features
- $\mathcal{D} := \{((x_1^{(n)}, x_2^{(n)}), y^{(n)}) \mid n = 1, \dots, N\}$
- Task: find function  $f(x_1^{(n)}, x_2^{(n)}) = \hat{y}^{(n)}$  that reproduces given labels



Assignment Project Exam Help  
<https://eduassistpro.github.io/>  
Add WeChat edu\_assist\_pro

# Analogy with seeing in colour

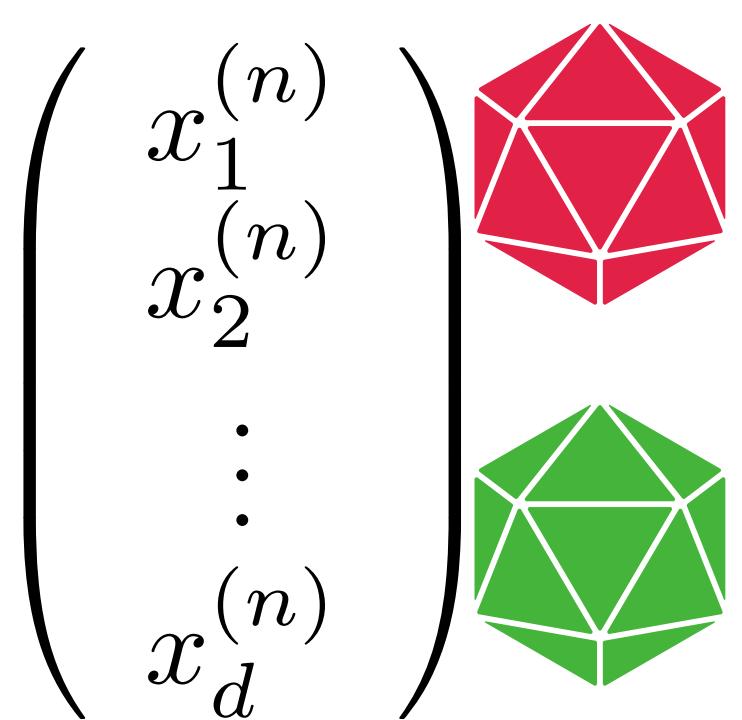
**Opsins (photopigments) in cones respond to colour preferentially**



Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

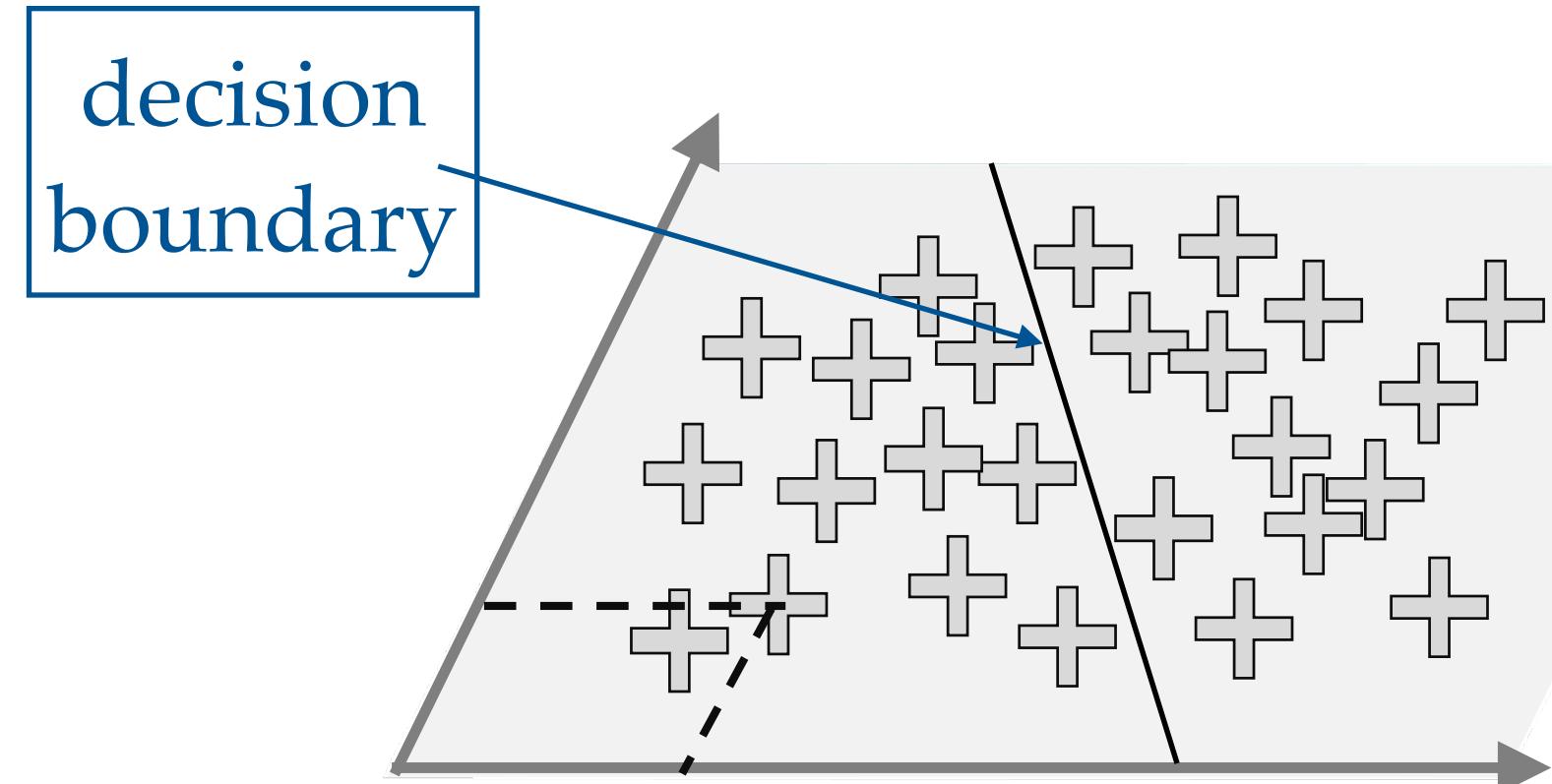


$$\mathbf{w}_{red} := (\mathbf{w}_{red,1}, \mathbf{w}_{red,2}, \dots, \mathbf{w}_{red,d})$$

$$\mathbf{w}_{green} := (\mathbf{w}_{green,1}, \mathbf{w}_{green,2}, \dots, \mathbf{w}_{green,d})$$

# Find equation for decision boundary

Assign probability for each point being green/red



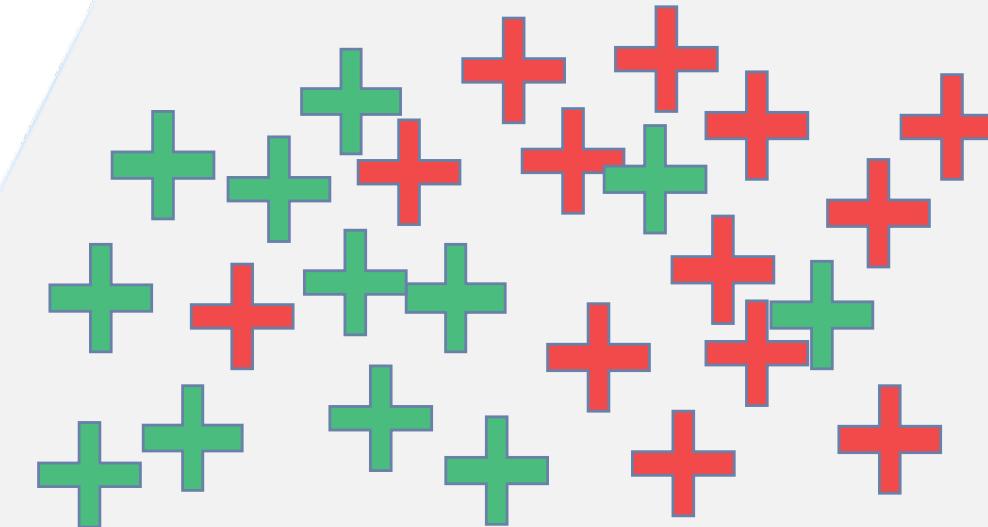
$$f(x_1, x_2; \mathbf{w}) = w_0 + w_1 x_1 + w_2 x_2$$

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

$$\sigma(f) = \frac{1}{1 + \exp(-f)}$$



Learning = adjusting weights until agreement with data



# Constructing scale for comparing predictions with training labels

- Output  $f(\mathbf{x}^{(n)}; \mathbf{w}^i) \triangleq f_i^{(n)}$ ,  $i = \text{red/green}$ ,  $C = 2$ ,  $\hat{y}^{(n)} = \sigma(f^{(n)})$
  - $0 \leq \sigma(f_i^{(n)}) \leq 1$  probability, with  $\sum_{i=1}^C \sigma(f_i^{(n)}) = 1$
  - Let  $\hat{y}_i^{(n)} = \sigma(f_i^{(n)})$
  - $y^{(n)} = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$  or  $\begin{pmatrix} 0 \\ 1 \end{pmatrix}$  for red/green
  - Evaluation of classification:  $\text{cost}(y^{(n)}, \hat{y}^{(n)})$
  - Compare two different costs – quadratic and logarithmic
  - Logarithm penalises mistakes more, also has a sharper drop (large gradient to guide weights to lower loss)
- Assignment Project Exam Help  
<https://eduassistpro.github.io/>  
Add WeChat edu\_assist\_pro
- 0.5
- $-y^{(n)} \ln \hat{y}^{(n)}$
- $y^{(n)} - \hat{y}^{(n)})^2$
- $f^{(n)}$
- 
- For a one component (scalar) output

# Multiclass classification

**Input: images 32 x 32 x 3 dimensions, Output: one-hot encodings: 10 dimensions**

CIFAR-10: Example dataset for  
multi class classification

Assignment Project Exam Help

<https://eduassistpro.github.io/>  
Add WeChat edu\_assist\_pro

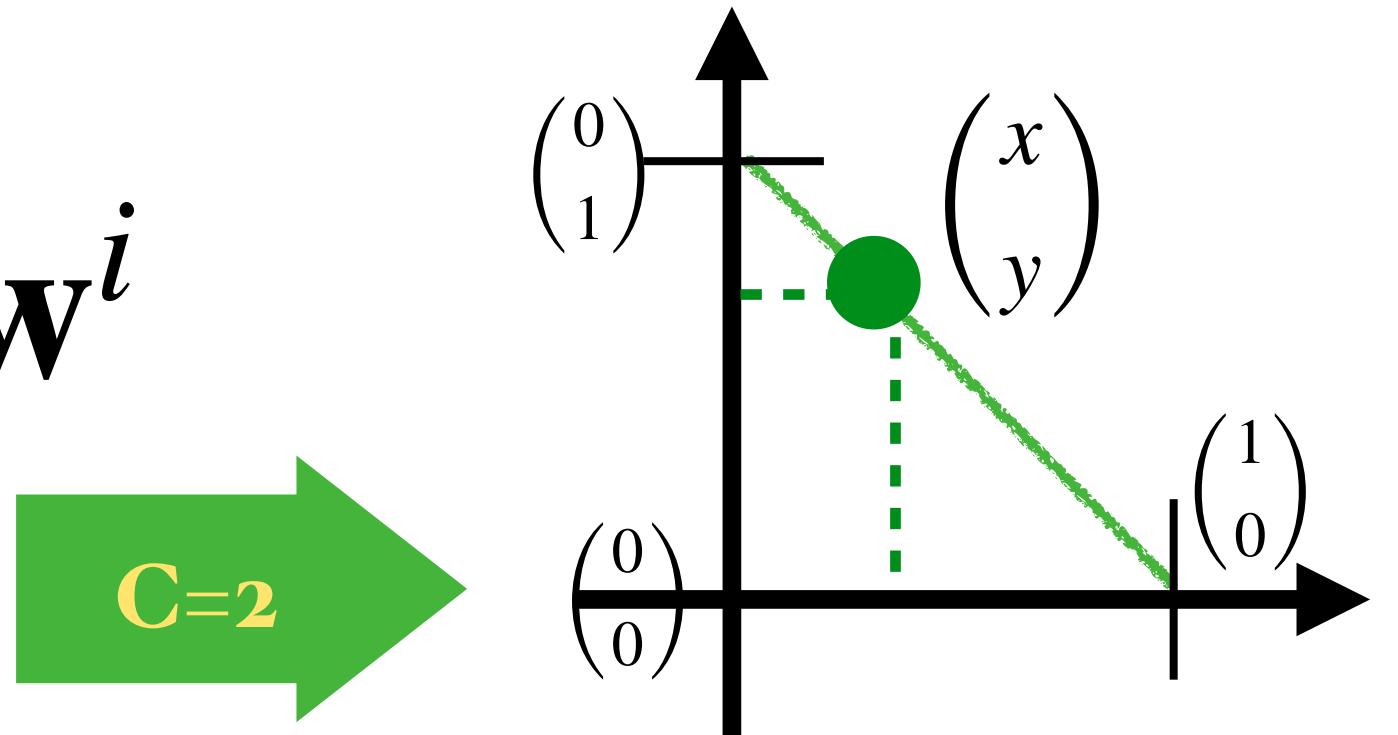
CAT  $\mapsto \mathbf{e}_4 =$

$$\begin{pmatrix} 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{pmatrix}$$

SHIP  $\mapsto \mathbf{e}_9 =$

$$\begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 1 \\ 0 \end{pmatrix}$$

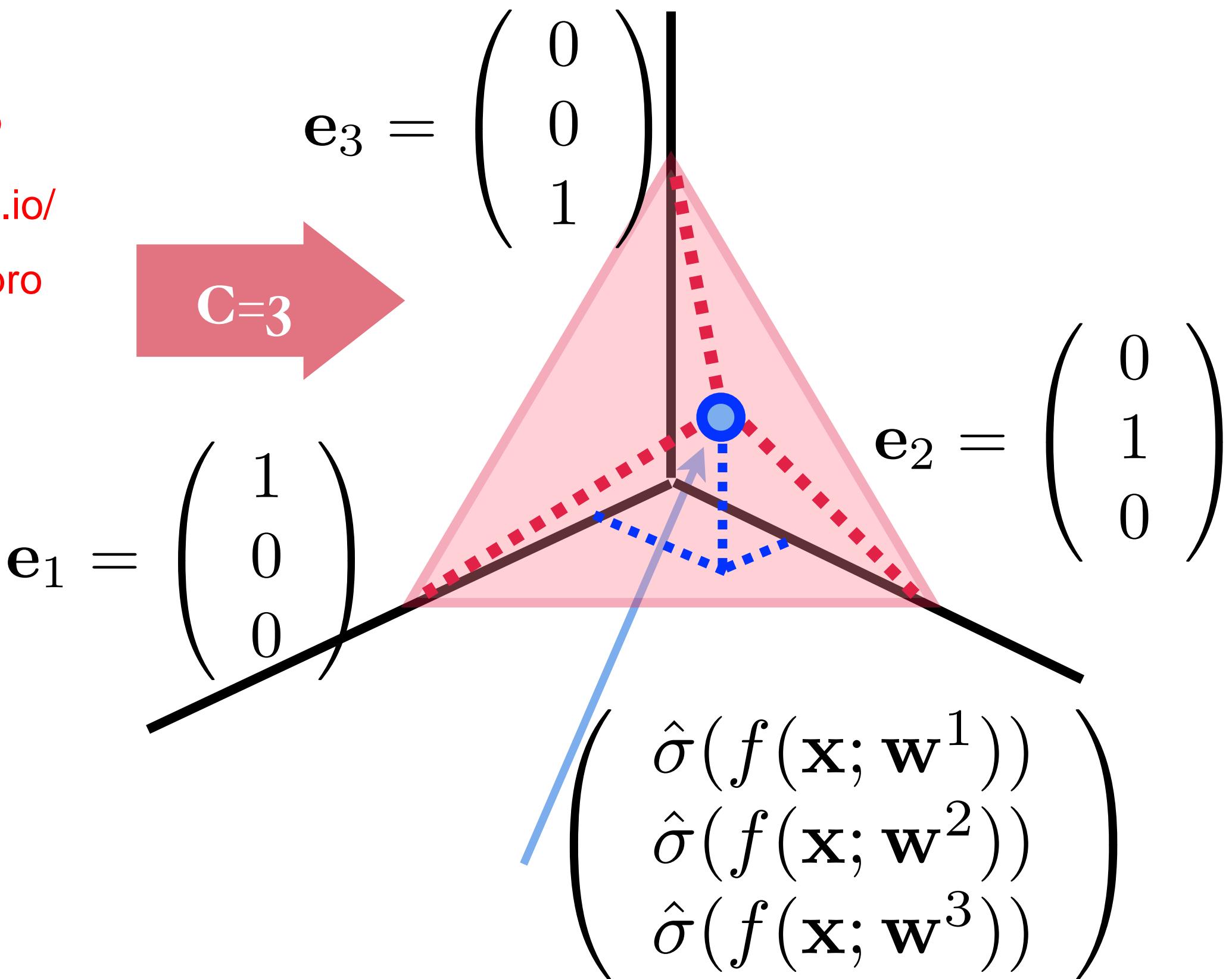
# C-classes C different weight vectors $\mathbf{w}^i$



$$x \geq 0, y \geq 0, x + y = 1$$

- For each input vector (say a representation of an image or sound file), produce an output on the  $(C-1)$ -dimensional surface embedded in  $C$ -dimensional Euclidean space
- Cost for input  $\mathbf{x}^{(n)}$  = measure of mismatch between  $n$  C-dimensional prediction  $\hat{\sigma}(f(\mathbf{x}^{(n)}; \mathbf{w}^i))$  and true label  $\mathbf{e}_i$
- Hat  $\hat{\cdot}$  on  $\hat{\sigma}$  indicates normalisation: entries add up to one:  $\hat{\sigma}(f(\mathbf{x}; \mathbf{w}^1)) + \hat{\sigma}(f(\mathbf{x}; \mathbf{w}^2)) + \hat{\sigma}(f(\mathbf{x}; \mathbf{w}^3)) = 1$

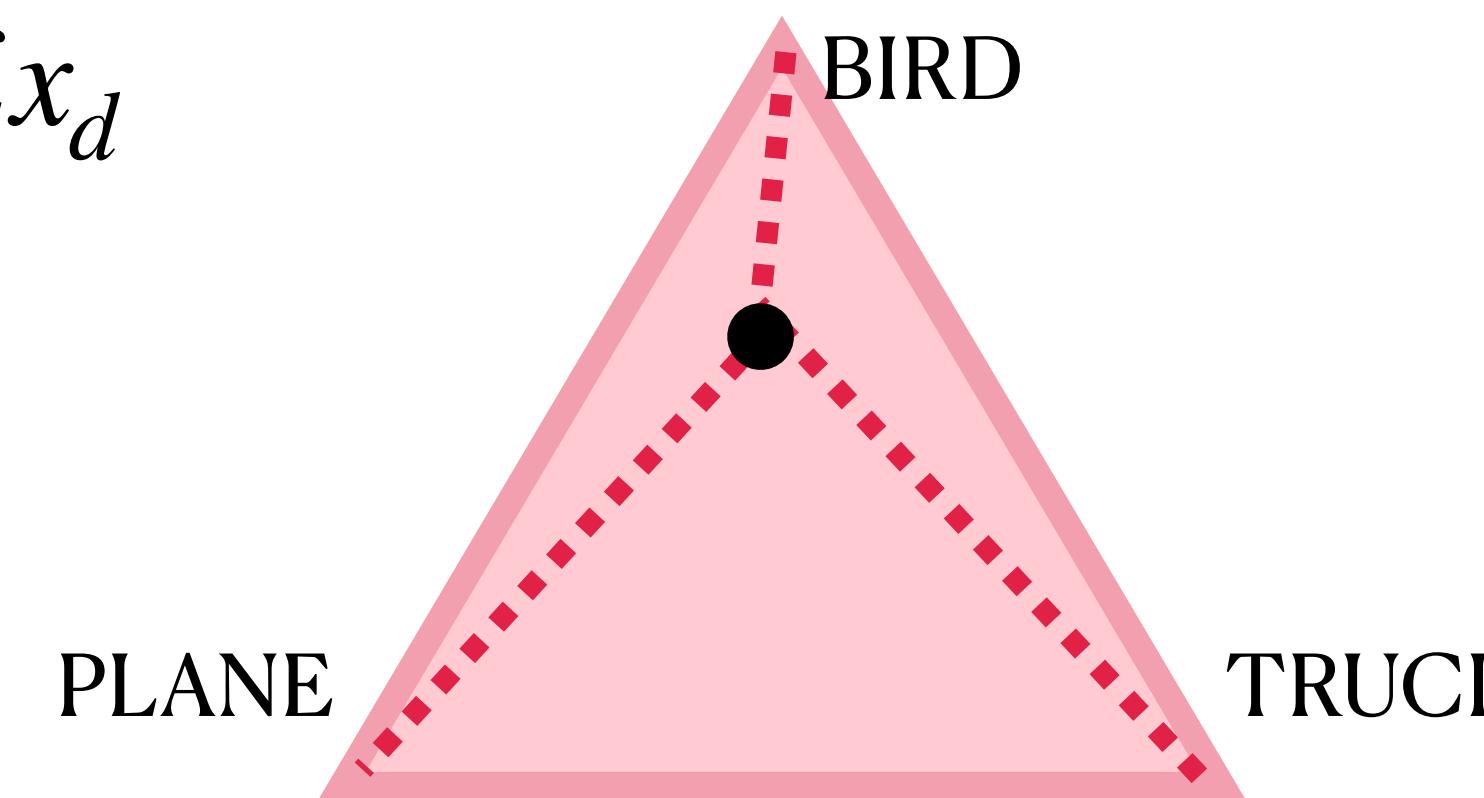
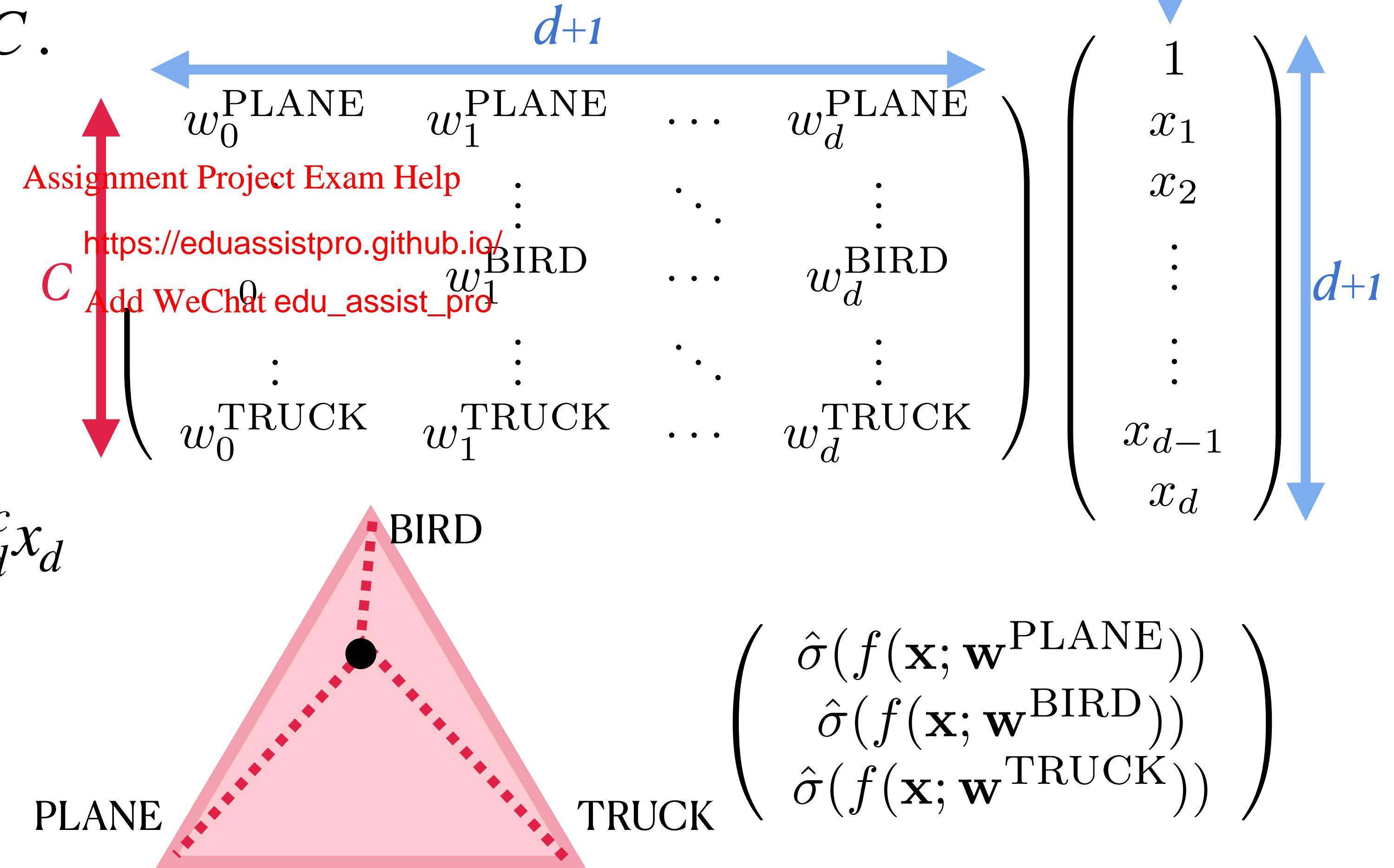
Assignment Project Exam Help  
<https://eduassistpro.github.io/>  
 Add WeChat edu\_assist\_pro



# Multiclass classification

**Weight vectors for each class**

- $\mathbf{w}^c = (w_0^c, w_1^c, \dots, w_d^c), c = 1, \dots, C.$
- $C$  - number of classes, 10 for CIFAR-10
- $d$  - dimensionality of data,  $\mathbf{x} = (x_1, x_2, \dots, x_d)$
- $f(\mathbf{x}; \mathbf{w}^c) = w_0^c \cdot 1 + w_1^c x_1 + \dots + w_d^c x_d$   
: for each input data point, compute output for all classes



# Set up gradient descent of loss for classification

## Re-phrasing what has been done

- For each class each data point  $\mathbf{x}^{(n)}$  is assigned a score  $s_c^{(n)} = f(\mathbf{x}^{(n)}; \mathbf{w}^c)$ ,  $c = 1, \dots, C$

- Choose the largest of the  $C$  scores as the predicted class for  $\mathbf{x}^{(n)}$

Assignment Project Exam Help

$$\bullet \quad c^* = \arg \max_{c \in \{1, \dots, C\}} s_c^{(n)}$$

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

- Replace max by softmax:  $\max(s_1, s_2, s_3) \longrightarrow \text{softmax}(s_1, s_2, s_3) = \ln(e^{s_1} + e^{s_2} + e^{s_3})$

- Exponential function: monotonic in argument ( $x \nearrow \implies e^x \nearrow$ )

- Normalise exponential scores:  $s_c^{(n)} \mapsto \frac{e^{s_c^{(n)}}}{e^{s_1^{(n)}} + e^{s_2^{(n)}} + \dots + e^{s_C^{(n)}}} =: \hat{\sigma}(s_c^{(n)}) = [\hat{y}^{(n)}]_c$

- Treat component  $c$  of  $[\hat{y}^{(n)}]_c = \hat{\sigma}(s_c^{(n)})$  as probability that  $\mathbf{x}^{(n)}$  belongs to class  $c$ :  $P(c | \mathbf{x}^{(n)})$

# Multi-class loss function: cross entropy

**Measures information about label distribution from input data and choice of weights**

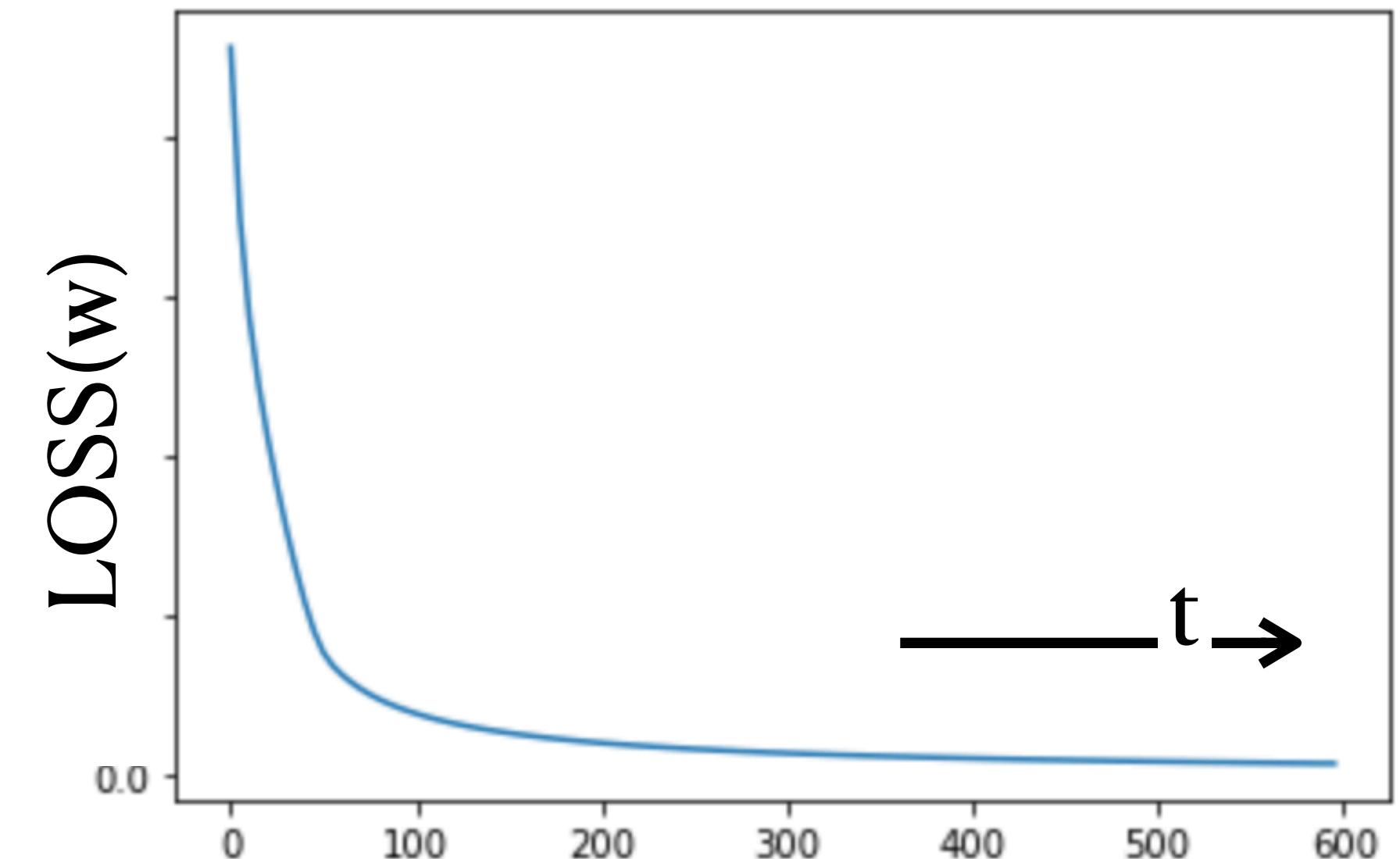
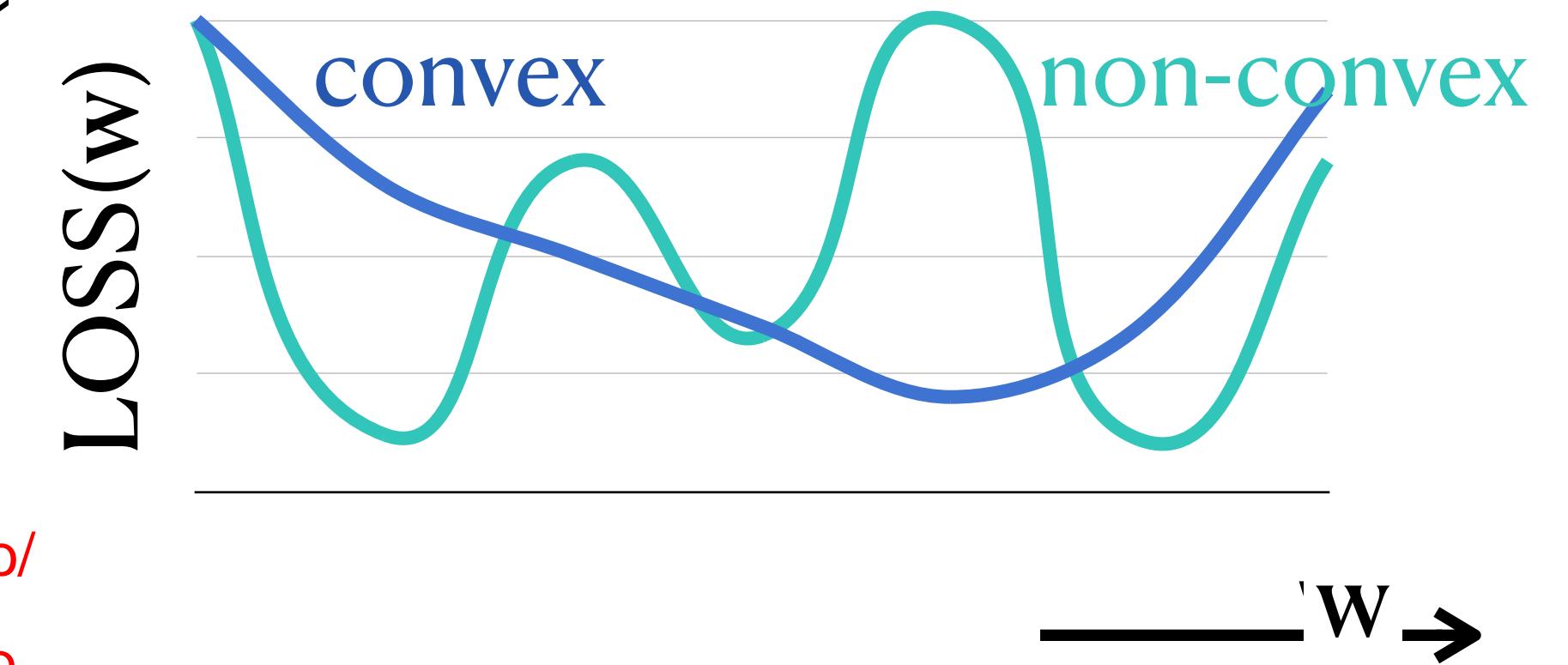
- For each data point  $\mathbf{x}^{(n)}$  sum over costs  $-\sum_{c=1}^C y_c^{(n)} \ln \hat{y}_c^{(n)}$  for all classes
- Sum costs over all data points  $L(\mathbf{W}) := L(\{\mathbf{w}^{(n)}\}_{n=1}^N) = -\sum_{n=1}^N \sum_{c=1}^C y_c^{(n)} \ln \hat{y}_c^{(n)}$ , called cross-entropy.  
Assignment Project Exam Help  
<https://eduassistpro.github.io/>  
Add WeChat edu\_assist\_pro
- Eg: target  $y^{(n)} = \begin{pmatrix} 0 \\ 1 \\ 0 \\ 0 \end{pmatrix}$  prediction  $\hat{y}^{(n)} = \begin{pmatrix} \hat{y}_1^{(n)} \\ \hat{y}_2^{(n)} \\ \hat{y}_3^{(n)} \\ \hat{y}_4^{(n)} \end{pmatrix}$ :  $-(0 \cdot \ln \hat{y}_1^{(n)} + 1 \cdot \ln \hat{y}_2^{(n)} + 0 \cdot \ln \hat{y}_3^{(n)} + 0 \cdot \ln \hat{y}_4^{(n)}) = -\ln \hat{y}_2^{(n)}$
- $L(\mathbf{W}) = -\ln (\hat{y}_{c_1}^{(1)} \cdot \hat{y}_{c_2}^{(2)} \cdots \hat{y}_{c_N}^{(N)}) = -\sum_{n=1}^N \ln \hat{y}_{c_n}^{(n)}$ : reduce negative of log(predicted probabilities)

# Gradient descent on cross-entropy finds optimal weights

For linear maps  $f$ , cross-entropy is convex

- Learning: Reduce  $L(\mathbf{W})$  by changing weights
- $\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - \eta \nabla_{\mathbf{w}} L(\mathbf{W})$
- All weights are contained in  $\mathbf{w}$
- Jupyter notebook

Assignment Project Exam Help  
<https://eduassistpro.github.io/>  
Add WeChat edu\_assist\_pro

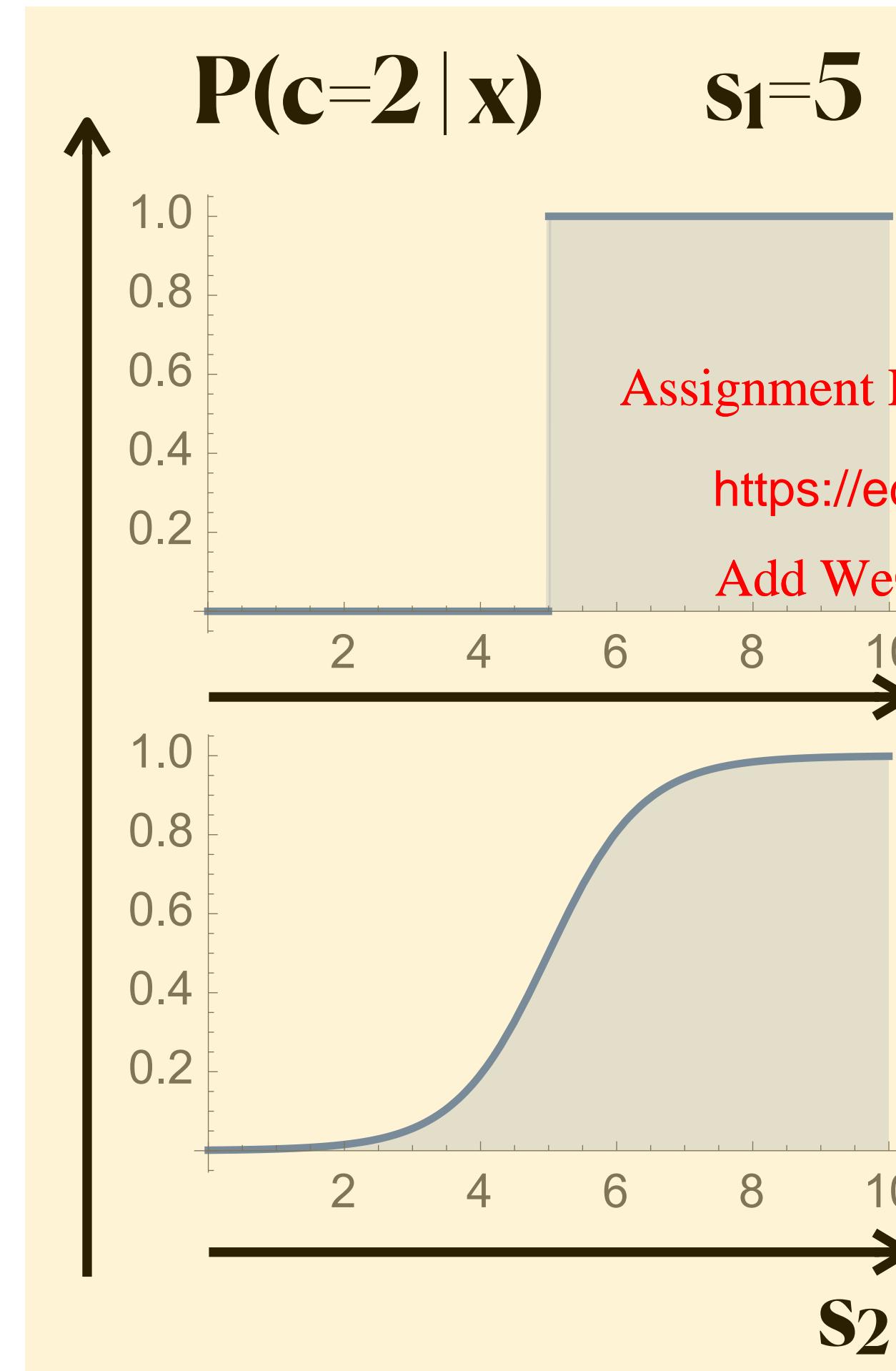


# Example: data x; 2-class problem

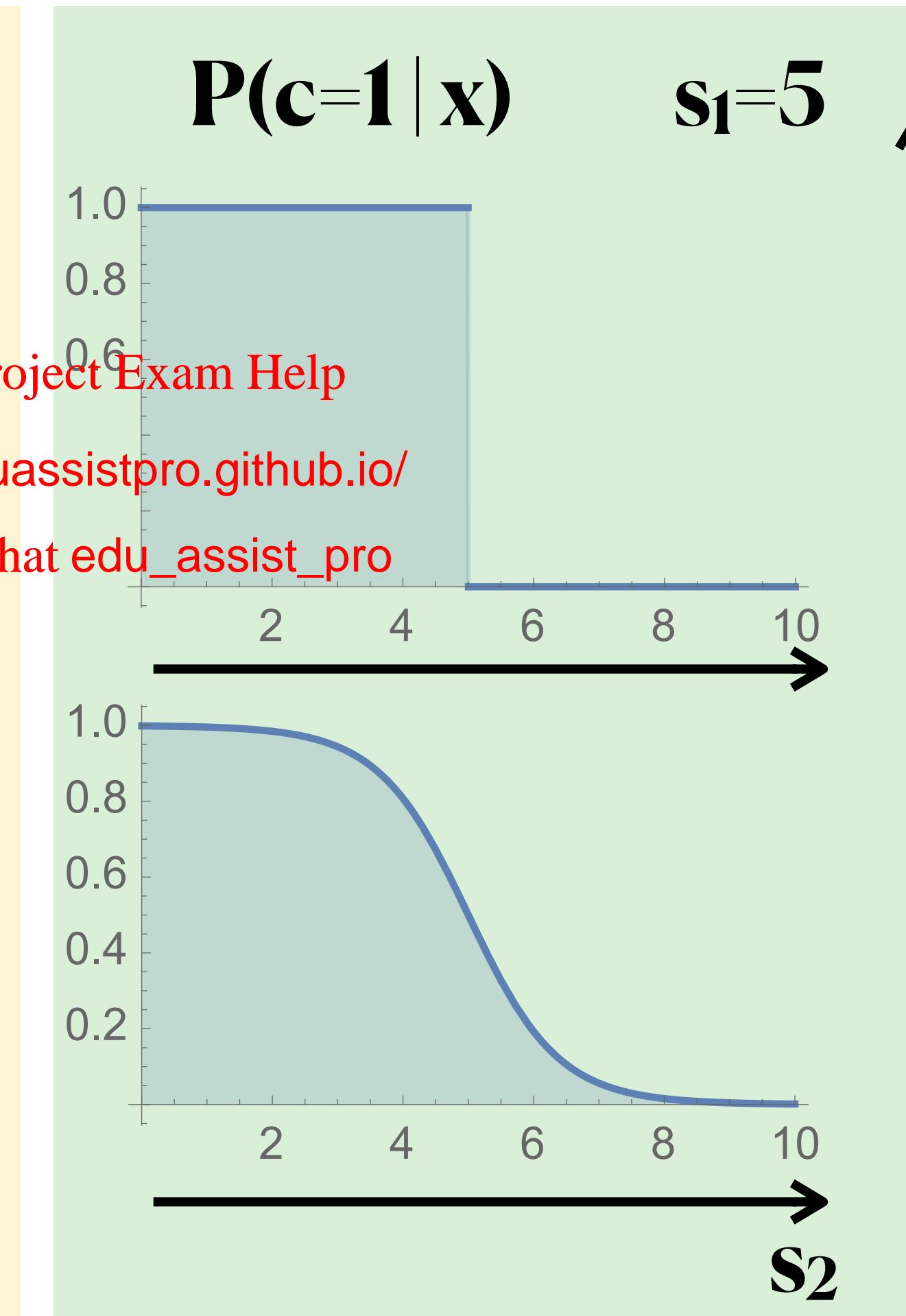
**Compare probability assignment for arg max with arg softmax**

$$\frac{f[0,(s - 5)]}{f[0,(s - 5)] + f[0, -(s - 5)]}$$

- $f = \text{Max}$



- $f = \text{Softmax}$



Assignment Project Exam Help  
<https://eduassistpro.github.io/>  
Add WeChat edu\_assist\_pro

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro