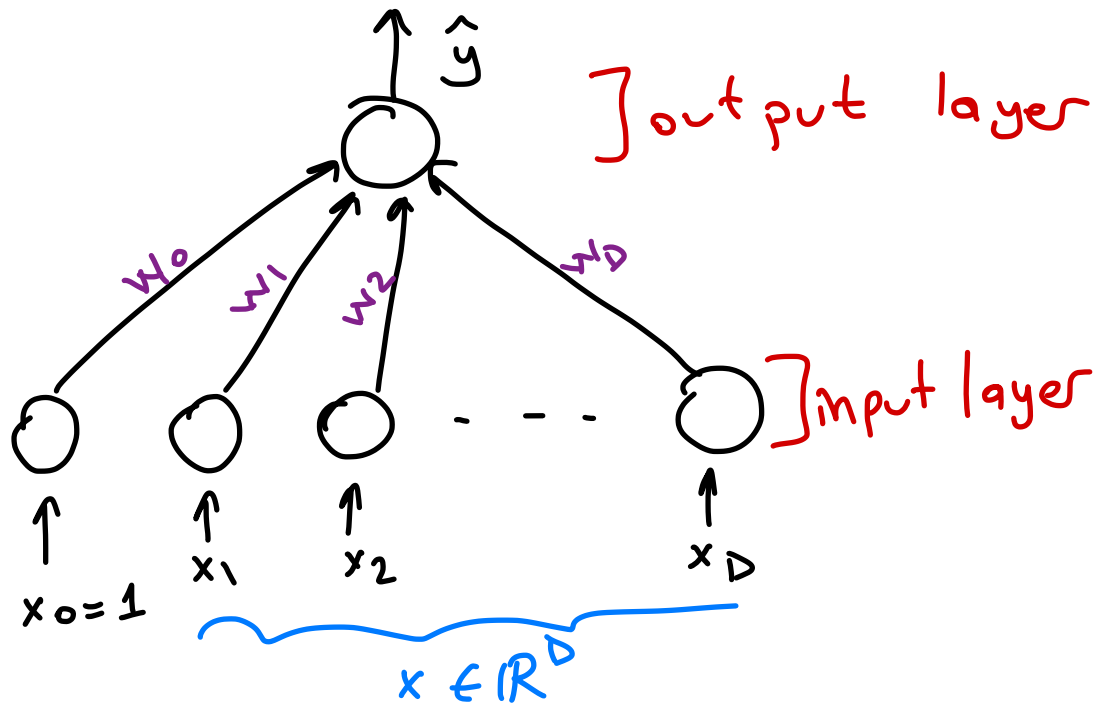


The Perceptron



$$\hat{y} = w_0 x_0 + w_1 x_1 + w_2 x_2 + \dots + w_D x_D$$

$$= w_0 x_0 + \sum_{d=1}^D w_d x_d \quad w^T \cdot x$$

$$= w^T \cdot x + w_0$$

total weighted signal received

$$w^T \cdot x + w_0 > 0 \Rightarrow \text{feel pain}$$

$$w^T \cdot x + w_0 < 0 \Rightarrow \text{does not feel}$$

$$w^T \cdot x > \boxed{-w_0} \Leftarrow \text{threshold}$$

$$\hat{y} = [w_1 \quad w_2 \quad \dots \quad w_D] \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_D \end{bmatrix} + w_0$$

$$= [w_0 \quad w_1 \quad w_2 \quad \dots \quad w_D] \begin{bmatrix} x_0=1 \\ x_1 \\ x_2 \\ \vdots \\ x_D \end{bmatrix} = \tilde{w}^T \tilde{x} \rightarrow (D+1) \times 1$$

\downarrow
 $1 \times (D+1)$

threshold function (activation function)

$$s(a) = \begin{cases} 1 & \text{if } a > 0 \\ 0 & \text{otherwise} \end{cases}$$



$$s(w^T \cdot x) = \begin{cases} 1 & \text{if } w^T \cdot x > 0 \\ 0 & \text{otherwise} \end{cases}$$

$$s(w^T \cdot x) = \frac{1}{1 + \exp[-w^T \cdot x]}$$

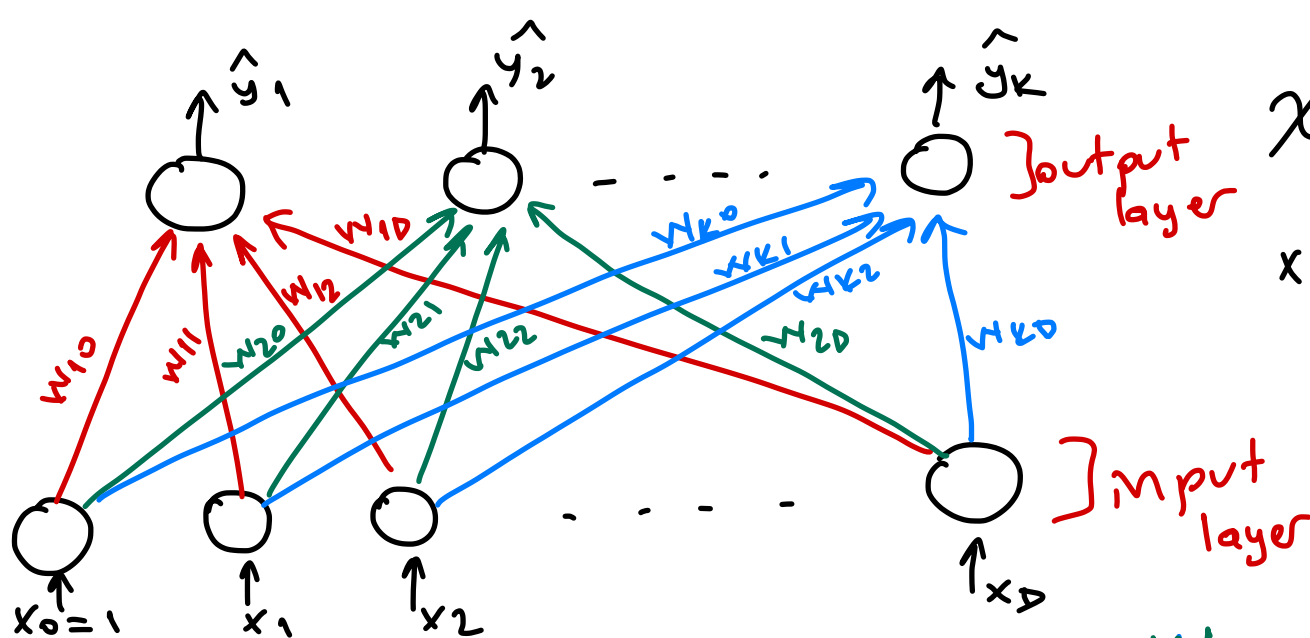
↳ sigmoid activation

} binary classification

$$s(w^T \cdot x) = w^T \cdot x$$

↳ linear activation

} regression



$$\mathcal{X} = \{ (x_i, y_i) \}_{i=1}^N$$

$$x_i \in \mathbb{R}^D \quad y_i \in \{1, 2, \dots, k\}$$

$$\Rightarrow W_1 = \begin{bmatrix} w_{10} \\ w_{11} \\ w_{12} \\ \vdots \\ w_{1D} \end{bmatrix} \quad W_2 = \begin{bmatrix} w_{20} \\ w_{21} \\ w_{22} \\ \vdots \\ w_{2D} \end{bmatrix} \quad \dots \quad W_k = \begin{bmatrix} w_{k0} \\ w_{k1} \\ w_{k2} \\ \vdots \\ w_{kD} \end{bmatrix}$$

$$\hat{y}_c = \sum_{d=1}^D w_{cd} \cdot x_d + w_{c0} = W_c^T \cdot X$$

$$\begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_k \end{bmatrix} \begin{matrix} 1st \\ 2nd \\ kth \end{matrix} = \begin{bmatrix} w_{10} & w_{11} & w_{12} & \dots & w_{1D} \\ w_{20} & w_{21} & w_{22} & \dots & w_{2D} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ w_{k0} & w_{k1} & w_{k2} & \dots & w_{kD} \end{bmatrix} \begin{matrix} K \times 1 \\ K \times (D+1) \end{matrix}$$

$$\begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_D \end{bmatrix} \begin{matrix} (D+1) \times 1 \end{matrix}$$

$$\hat{y} = W \cdot X$$

$$\begin{matrix} K \times 1 & K \times (D+1) & (D+1) \times 1 \\ \downarrow & \downarrow & \downarrow \end{matrix}$$

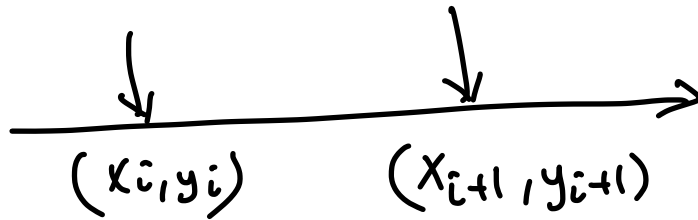
$$\hat{y}_c = \frac{\exp(w_c^T \cdot x)}{\sum_{d=1}^K \exp(w_d^T \cdot x)}$$

} softmax activation

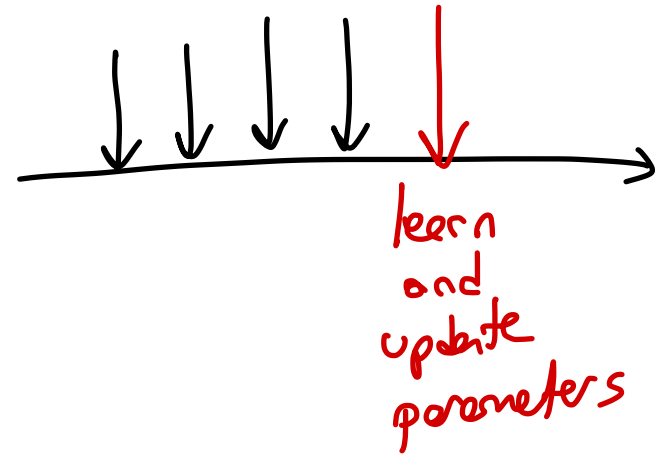
a new data point x^* \Rightarrow choose $y^* = \arg \max_c \hat{y}_c$

LEARNING

Online Learning versus Batch Learning



\rightarrow samples are coming one by one



Regression
 $\{(x_i, y_i)\}_{i=1}^N$

$y_i \in \mathbb{R}$

$$\begin{aligned} \text{Error}_i(w | x_i, y_i) &= \frac{1}{2} (y_i - \hat{y}_i)^2 \\ &= \frac{1}{2} (y_i - s(w^T \cdot x_i))^2 \\ &= \frac{1}{2} (y_i - w^T \cdot x_i)^2 \end{aligned}$$

$$\begin{aligned} \frac{\partial \text{Error}_i}{\partial w} &= \frac{1}{2} \cancel{2} (y_i - \cancel{w^T \cdot x_i}) \cdot (-x_i) \\ &= -(y_i - \underbrace{w^T \cdot x_i}_{\hat{y}_i}) x_i = -(y_i - \hat{y}_i) \cdot x_i \end{aligned}$$

$$\Delta w = -\eta \cdot \frac{\partial \text{Error}_i}{\partial w} = \eta \cdot \underline{(y_i - \hat{y}_i)} \cdot x_i$$

Binary Classification

$$\{x_i, y_i\}_{i=1}^N \Rightarrow y_i \in \{0, 1\} \quad \text{Error}_i(w | x_i, y_i) = -[y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i)]$$

$$\text{Hint: } \frac{\partial \log(\hat{y}_i)}{\partial w}$$

$$\frac{\partial \log(f(w))}{\partial w} = \frac{1}{f(w)} \frac{\partial f(w)}{\partial w}$$

$$\hat{y}_i = s(w^T x_i) = \frac{1}{1 + \exp[-w^T x_i]}$$

$$= - \left[y_i \log \left[\frac{1}{1 + \exp(-w^T x_i)} \right] + (1 - y_i) \log \left[1 - \frac{1}{1 + \exp(-w^T x_i)} \right] \right]$$

$$\frac{\partial \text{Error}_i(w | x_i, y_i)}{\partial w} = -(y_i - \hat{y}_i) x_i$$

$$\Delta w = -\eta \frac{\partial \text{Error}_i}{\partial w} = \eta (y_i - \hat{y}_i) \cdot x_i$$

Multiclass Classification

$$\text{Error}_i(\{w_c\}_{c=1}^k | x_i, y_i) = - \sum_{c=1}^k y_{ic} \log(\hat{y}_{ic})$$

$$= - \sum_{c=1}^k y_{ic} \log \left[\frac{\exp[w_c^T \cdot x_i]}{\sum_{d=1}^k \exp[w_d^T \cdot x_i]} \right]$$

$$\hat{y}_{ic} = \frac{\exp[w_c^T \cdot x_i]}{\sum_{d=1}^k \exp[w_d^T \cdot x_i]}$$

$$\frac{\partial \text{Error}_i}{\partial w_c} = - (y_{ic} - \hat{y}_{ic}) \cdot x_i$$

$$\Delta w_c = - \eta \frac{\text{Error}_i}{\partial w_c} = \eta (y_{ic} - \hat{y}_{ic}) \cdot x_i$$

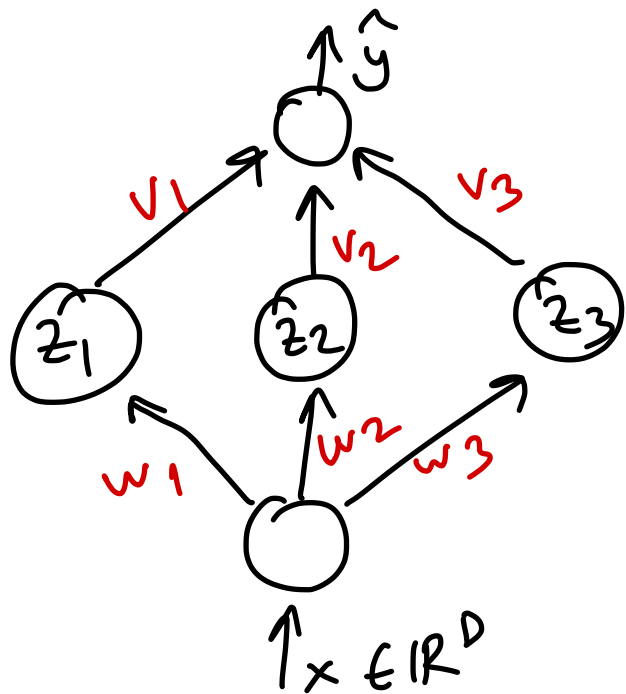
$$\text{Update} = (\text{Learning Factor}) \times \left(\begin{matrix} \text{True} \\ \text{output} \end{matrix} - \begin{matrix} \text{Predicted} \\ \text{output} \end{matrix} \right) \times (\text{Input})$$

$$f(x) = 2x \quad \hookrightarrow \text{linear func} \quad g(x) = 3x \quad \hookrightarrow \text{linear func.}$$

$$f \circ g(x) = 6x \quad \hookrightarrow \text{linear func.}$$

$$f(x) = 2x \quad \hookrightarrow \text{linear} \quad g(x) = \log(x) \quad \hookrightarrow \text{nonlinear}$$

$$f \circ g(x) = 2 \log(x) \quad \hookrightarrow \text{nonlinear}$$



$$z_1 = w_1^T \cdot x$$

$$z_2 = w_2^T \cdot x$$

$$z_3 = w_3^T \cdot x$$

$$\hat{y} = v_1 z_1 + v_2 z_2 + v_3 z_3$$

$$\hat{y} = \underbrace{v_1 \cdot w_1^T}_{\tilde{w}_1^T} x + \underbrace{v_2 w_2^T}_{\tilde{w}_2^T} x + \underbrace{v_3 w_3^T}_{\tilde{w}_3^T} x$$

$$\hat{y} = \tilde{w}_1^T \cdot x + \tilde{w}_2^T \cdot x + \tilde{w}_3^T \cdot x$$

$$= (\tilde{w}_1 + \tilde{w}_2 + \tilde{w}_3)^T \cdot x$$

