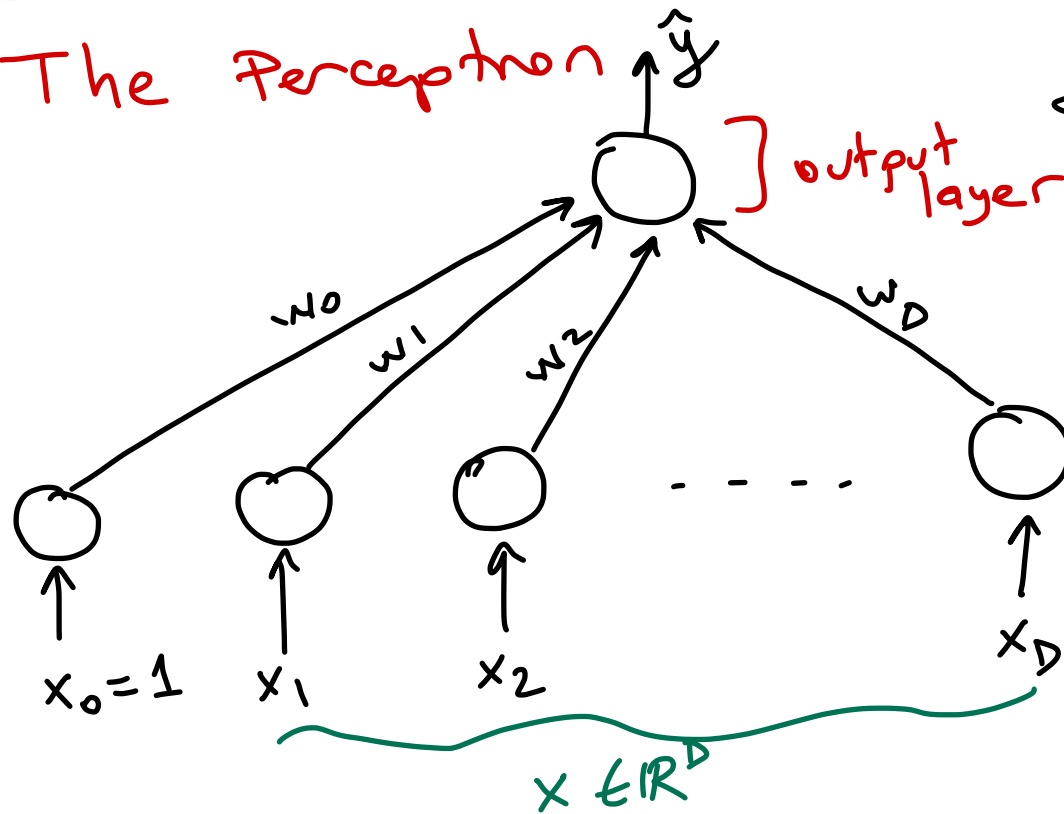


Multilayer Perceptrons

The Perceptron



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

$$= w_0 \cdot \overset{1}{x_0} + \sum_{d=1}^D w_d \cdot x_d$$

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

$$\Downarrow$$

$$W^T \cdot x$$

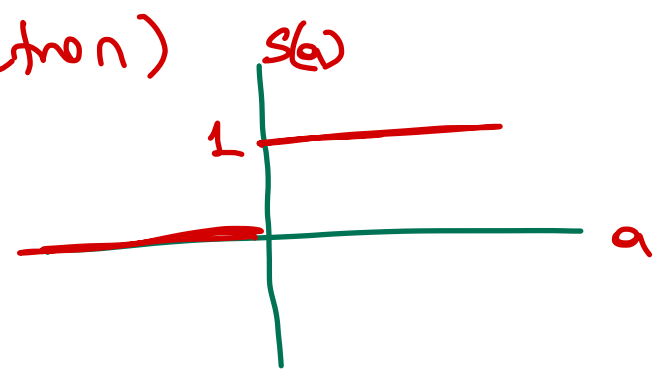
$$\hat{y} = [w_1 \ w_2 \ \dots \ w_D] \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_D \end{bmatrix} + w_0$$

$$= [w_0 \ w_1 \ w_2 \ \dots \ w_D] \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \\ x_D \end{bmatrix} = \underset{\substack{\uparrow \\ 1 \times (D+1)}}{W^T} \cdot \underset{\substack{\rightarrow \\ (D+1) \times 1}}{x}$$

threshold function (activation function) $s(a)$

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

↑
received message



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

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

↳ sigmoid activation

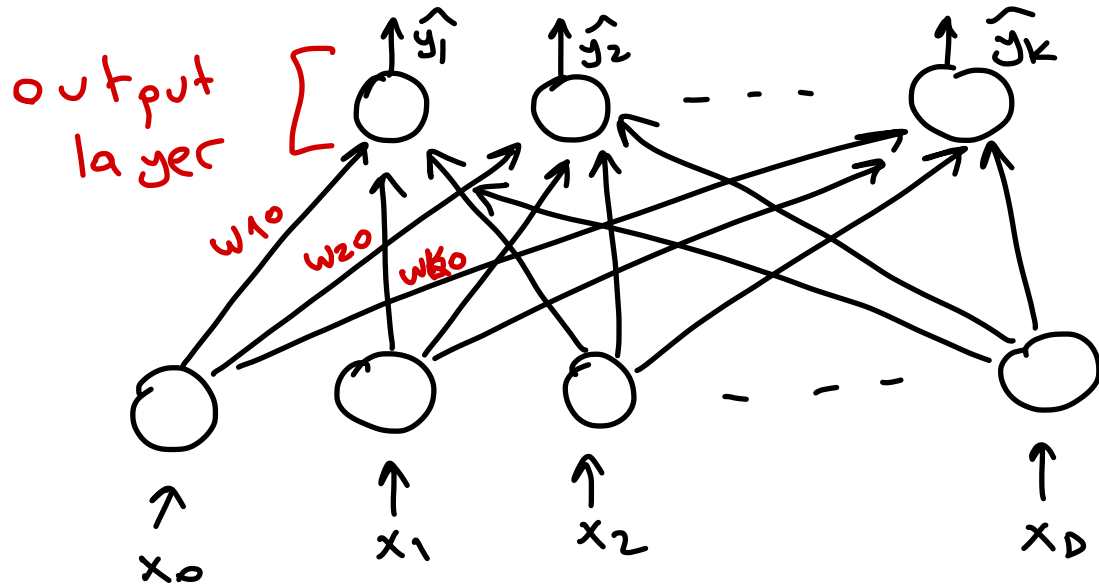
$$s(w^T \cdot x) = w^T \cdot x$$

↳ linear activation

} binary classification

} regression

$$s(a) = a$$



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

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

$$W_1 = \begin{bmatrix} w_{10} \\ w_{11} \\ w_{12} \\ \vdots \\ w_{1D} \end{bmatrix} \quad W_2 = \begin{bmatrix} w_{20} \\ w_{21} \\ w_{22} \\ \vdots \\ w_{2D} \end{bmatrix} \quad \dots \quad W_K = \begin{bmatrix} w_{K0} \\ w_{K1} \\ w_{K2} \\ \vdots \\ w_{KD} \end{bmatrix}$$

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

$$\begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_K \end{bmatrix}_{K \times 1} = \begin{bmatrix} w_{10} & w_{11} & w_{12} & \dots & w_{1D} \\ w_{20} & w_{21} & w_{22} & \dots & w_{2D} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ w_{K0} & w_{K1} & w_{K2} & \dots & w_{KD} \end{bmatrix}_{K \times (D+1)} \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_D \end{bmatrix}_{(D+1) \times 1}$$

$$\Rightarrow \hat{\mathbf{y}} = \mathbf{W} \cdot \mathbf{x}$$

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

$$\hat{y}_c = \frac{\exp[w_c^T \cdot x]}{\sum_{k=1}^K \exp[w_k^T \cdot x]}$$

} softmax activation

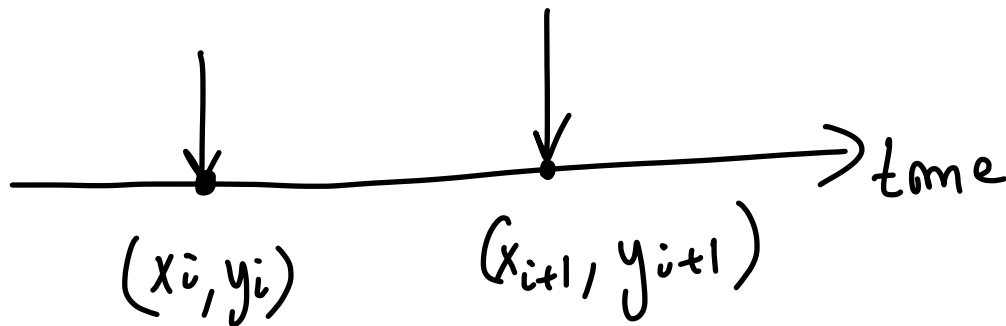
a new data point x^* \Rightarrow choose $\hat{y}^* = \arg \max_c \hat{y}_c$

LEARNING

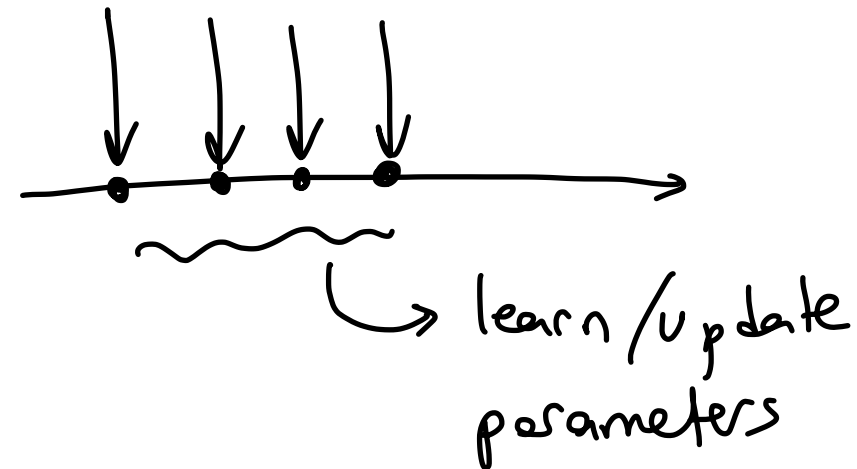
Online Learning

vs

Batch Learning



- samples are coming one by one



Regression 1

$$\text{Error}_i(w|x_i, y_i) = \frac{1}{2} (y_i - \hat{y}_i)^2 \quad \leftarrow \text{squared error}$$

$$= \frac{1}{2} [y_i - s(w^T \cdot x_i)]^2$$

$$= \frac{1}{2} [y_i - \underline{w^T \cdot x_i}]^2$$

$$\frac{\partial \text{Error}_i}{\partial w} = \frac{1}{2} \cdot 2 \cdot [y_i - w^T \cdot x_i] \cdot \frac{\partial [y_i - w^T \cdot x_i]}{\partial w}$$

$$= \underbrace{[y_i - w^T \cdot x_i]}_{1 \times 1} \underbrace{(-x_i)}_{(D+1) \times 1} = - \underbrace{(y_i - \hat{y}_i)}_{y_i - \hat{y}_i} \cdot x_i$$

$$\frac{1}{2} \sum_{i=1}^N (y_i - \hat{y}_i)^2 = \frac{1}{2} (y_1 - \hat{y}_1)^2 + \frac{1}{2} (y_2 - \hat{y}_2)^2 + \dots + \frac{1}{2} (y_N - \hat{y}_N)^2$$

$$\frac{\partial w^T \cdot x_i}{\partial w} = \frac{\partial x_i^T \cdot w}{\partial w} = x_i$$

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

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

Binary Classification

$$\text{Error}_i(w | x_i, y_i) = - \underbrace{\left[y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i) \right]}_{\text{log-likelihood}}$$

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

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

$\nearrow f(w)$

Hint: $\frac{\partial \log(\hat{y}_i)}{\partial w} \Rightarrow \frac{\partial \log[f(w)]}{\partial w} = \frac{1}{f(w)} \cdot \frac{\partial f(w)}{\partial w}$

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

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

Multiclass Classification

$$\text{Error}_i(\{w_c\}_{c=1}^K | x_i, y_i) = - \underbrace{\sum_{c=1}^K y_{ic} \log(\hat{y}_{ic})}_{\text{log-likelihood}} \quad \rightarrow f(w)$$

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

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

$$\frac{\partial \text{Error}_i(\{w_k\}_{k=1}^K | x_i, y_i)}{\partial w_c} = -(y_{ic} - \hat{y}_{ic}) \cdot x_i$$

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

$$\text{Update} = (\text{Learning Factor}) \times (\text{True Output} - \text{Predicted output}) \times (\text{Input})$$

$$[\eta] \times [y_i - \hat{y}_i] \times [x_i] \Rightarrow \text{Regression}$$

$$[\eta] \times [y_i - \hat{y}_i] \times [x_i] \Rightarrow \text{Binary Classification}$$

$$[\eta] \times [y_{ic} - \hat{y}_{ic}] \times [x_i] \Rightarrow \text{multiclass Classification}$$

$$f(x) = 2x$$

\hookrightarrow linear

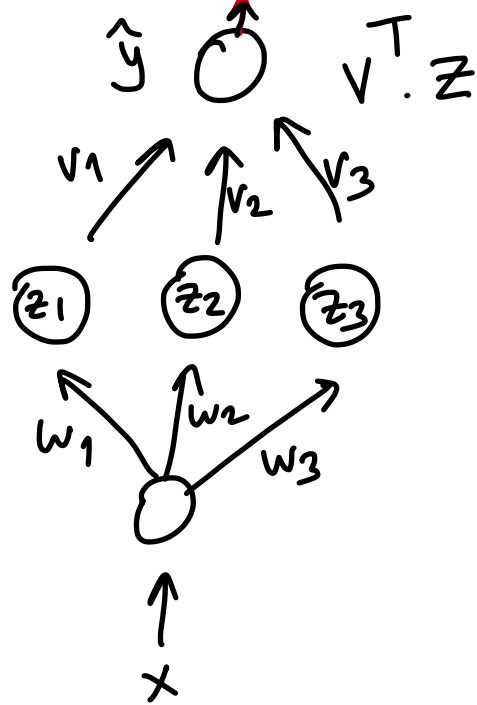
$$g(x) = 3x$$

\hookrightarrow linear

$$f \circ g(x) = 6x$$

\hookrightarrow linear

at least $f(x)$
or $g(x)$
should be
nonlinear



$$\underbrace{\text{Scalar } V_1 \cdot \text{vector } w_1^T}_{z_1} x + \underbrace{\text{Scalar } V_2 \cdot \text{vector } w_2^T}_{z_2} x + \underbrace{\text{Scalar } V_3 \cdot \text{vector } w_3^T}_{z_3} x$$

$$\tilde{w}_1^T = V_1 \cdot w_1^T$$

$$\tilde{w}_2^T = V_2 \cdot w_2^T$$

$$\tilde{w}_3^T = V_3 \cdot w_3^T$$

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

