§ Neural Network.

(*) Multiple layers



$\sum_{i=1}^{m}$   sign(·)

*(green)* difference
switch to other $\sigma(\cdot)$

- "Neurons" generalized a linear classifier.

$\begin{cases} \text{linear classifier}: \hat{d} = \text{sign}(\vec{x}\,w) \\ \text{★ Neuron}: \hat{d} = \underline{\sigma(\vec{x}\,w)} \quad \text{non-linear } \sigma(\cdot) \text{ function.} \end{cases}$
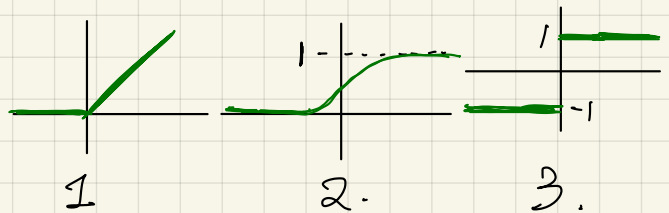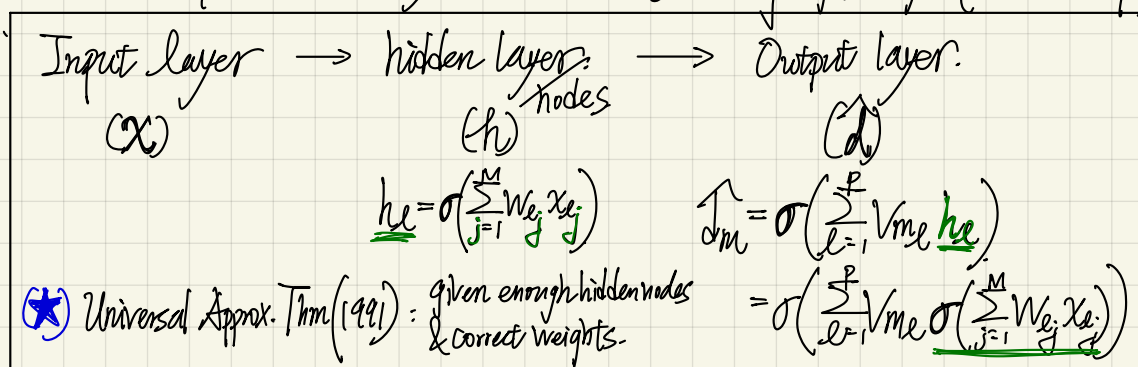
$\Rightarrow$ popular choices of $\sigma(\cdot)$:

1. ReLU: $\sigma(z) = \max\{0, z\}$
2. Logistic: $\sigma(z) = \frac{1}{1+e^{-z}} = (1+e^{-z})^{-1}$
3. Sign: $\sigma(z) = \text{sign}(z)$



1.        2.        3.

- A "neural network": 3 layer framework. $\begin{cases} \text{Need to decide \# of I/H/O layers.} \\ \text{Choosing right weights (nonconvex opt.)} \end{cases}$

Input layer $\longrightarrow$ hidden layer: $\longrightarrow$ Output layer.
(X)                    (h) nodes              (d)

$h_\ell = \sigma\left(\sum_{j=1}^{M} W_{\ell j}\, x_j\right)$     $\hat{d}_m = \sigma\left(\sum_{\ell=1}^{P} V_{m\ell}\, \underline{h_\ell}\right)$

$= \sigma\left(\sum_{\ell=1}^{P} V_{m\ell}\, \sigma\left(\sum_{j=1}^{M} W_{\ell j}\, x_j\right)\right)$

(★) Universal Approx. Thm (1991): given enough hidden nodes & correct weights.

* Multiple outputs to solve multiple problems
* "Deep learning" refers to "hidden layers" (how many we have?)
* Input is the "shallowest" layer; Output is the "deepest" layer
* Use SGD & Backpropagation

- Backpropagation: work back from deep (output) to shallow (input) layer.

$\Rightarrow$ for $N$ training samples (input) & $Q$ outputs: $\min\limits_{W_{ej}, V_{ej}} \sum\limits_{i=1}^{N} \sum\limits_{q=1}^{Q} \frac{1}{2} \left( \hat{d}_{i,q} - d_{i,q} \right)^2$

backward.

$1°$. Initialisation: guess $W_{ej}$, $V_{ej}$

$2°$ SGD process: choose inputs $i$ randomly.

$3°$ Calculate $h_{i,p}$, $\hat{d}_{i,q}$ for data $i$.

$4°$ GD updates: first $V_{ej}$, then $W_{ej}$. (deep back to shallow)

# 5. Backpropagation for N.N.
(BPPG)

- **BPPG uses SGD to train weights** $(W \& V)$

$$\hat{d_q} = \sigma\left(\sum_{i=1}^{P} V_{qn} h_n\right) \text{ where } h_n = \sigma\left(\sum_{\ell=1}^{M} W_{n\ell} x_\ell\right)$$

We take $\sigma(z) = \left(1 + e^{-z}\right)^{-1} = \frac{1}{1+e^{-z}}$ $\quad$ ≪ Activation function ≫

$N$ training samples (inputs) $(x_1^i, x_2^i \cdots x_M^i ; d^i)$ $i = 1, \cdots, M$

$1°$ Initialization: Set $W_{n\ell}$, $V_{qn}$

$2°$ For $t = 1, 2, 3, \cdots$ Select $i_t \in \{1, \cdots M\}$ randomly.

$\quad\quad$ Compute $h_j^{i_t}$, $\hat{d_k}^{i_t}$.

$\Rightarrow$ Calculate: $V_{k\ell}^{(t+1)} = V_{k\ell}^{(t)} - \alpha_t \left(\dfrac{d f^{i_t}}{V_{k\ell}}\right)$

$\quad\quad\quad\quad W_{mj}^{(t+1)} = W_{mj}^{(t)} - \alpha_t \left(\dfrac{d f^{i_t}}{d W_{mj}}\right)$

where $\alpha_t$ is step size.

$f^{i_t}$: loss function.

- **Gradients from Chain Rule.**

Consider "sq error loss": $f(V, W) = \sum_{i=1}^{N} \underbrace{\left(\frac{1}{2} \sum_{q=1}^{Q} (\hat{d_q^i} - d_q^i)^2\right)}_{f^i(V,W)}$

$*$ $\quad \dfrac{\partial f^i}{\partial V_{k,\ell}} = \left(\dfrac{\partial f^i}{\partial \hat{d_k^i}}\right) \dfrac{\partial \hat{d_k^i}}{\partial V_{k,\ell}}$

$\sigma(z) = \left(1 + e^{-z}\right)^{-1}$

$\Rightarrow \sigma'(z) = \sigma(z)(1 - \sigma(z))$

① $\quad \dfrac{\partial f^i}{\partial V_{k,\ell}} = \left(\underbrace{(\hat{d_k^i} - d_k^i)}_{\text{relevant term}} \cdot \underbrace{\sigma'\left(\sum_{n=1}^{P} h_n^i V_{k,n}\right)}_{\substack{\text{chain rule.} \\ d = \sigma(\cdot)}}\right) \cdot \underbrace{h_\ell^i}_{(\partial \hat{d_k^i}/V_{k,\ell})}$

$(= 0 \text{ when correct label})$

$\begin{cases} \delta_k^i : \text{error related to output.} \\ h_\ell^i : \text{layer input.} \end{cases}$

$= \underbrace{(\hat{d_k^i} - d_k^i)\, \hat{d_k^i}(1 - \hat{d_k^i})}_{= \delta_k^i} h_\ell^i \equiv \delta_k^i h_\ell^i$ where

$\langle \text{update} \rangle$

$\Rightarrow V_{k\ell}^{(t+1)} = V_{k\ell}^{(t)} - \alpha_t\, \delta_k^i h_\ell^i$

② for $W_{mj}$.

$$\frac{\partial f^i}{\partial W_{mj}} = \sum_{q=1}^{Q} \left(\frac{\partial f^i}{\partial \hat{d_q^i}}\right)\left(\frac{\partial \hat{d_q^i}}{\partial h_m^i}\right)\left(\frac{\partial h_m^i}{\partial W_{mj}}\right)$$

<span style="color:gray">loss/output    output/hidden    hidden/input.</span>

(i) loss/output: $\quad \frac{\partial f^i}{\partial \hat{d_q^i}} = \hat{d_q^i} - d_q^i$.

(ii) output/hidden: $\quad \frac{\partial \hat{d_q^i}}{\partial h_m^i} = \frac{\partial}{\partial h_m^i} \sigma\left(\sum_{r=1}^{P} h_r^i V_{qr}^i\right) = V_{qm}^i \cdot \sigma'\left(\sum_{r=1}^{P} h_r^i V_{qr}^i\right)$

$$= V_{qm}^i \cdot \left[\hat{d_q^i}(1-\hat{d_q^i})\right].$$

(iii) hidden/input: $\quad \frac{\partial h_m^i}{\partial W_{mj}} = \frac{\partial}{\partial W_{mj}} \sigma\left(\sum_{\ell=1}^{M} W_{m,\ell} X_\ell^i\right) = W_{mj} \cdot \sigma'\left(\sum_{\ell=1}^{M} W_{m\ell} X_\ell^i\right)$

$$= X_j^i \cdot \left[h_m^i(1-h_m^i)\right]$$

$$\Rightarrow \text{So}, \frac{\partial f^i}{\partial W_{mj}} = \sum_{q=1}^{Q} \left[(\hat{d_q^i}-d_q^i) V_{qm}^i \cdot \hat{d_q^i}(1-\hat{d_q^i}) \cdot h_m^i(1-h_m^i)\right] X_j^i$$

<span style="color:blue">$$= \gamma_m^i = \beta_q^i V_{qm}^i h_m^i(1-h_m^i)$$</span>

$$\equiv \gamma_m^i X_j^i \qquad \text{(another weighted sum of } X_j^i \text{)}$$

$$\Rightarrow \langle \text{Update} \rangle \quad W_{mj}^{(t+1)} = W_{mj}^{(t)} - \alpha_t \, \gamma_m^i X_j^i.$$

⊛ How to "Backpropagation"?

$1°$ Initialize $W_{mj}, V_{k\ell}$.

$2°$ Iteration; for $t = 0, 1, 2 \cdots$

     $\langle SGD \rangle$ Choose $i_t \in \{1 \cdots, N\}$. randomly.

     $\langle$ Forward thru NN $\rangle$ Compute $h_m^{i_t}$ by $(x_j^{i_t}, W_{mj}^t)$

                         $\hat{d}_q^{i_t}$ by $(h_m^{i_t}, V_{k\ell}^t)$

    $\langle$ Backward updates $\rangle$ $\quad \delta_k^{i_t} = (\hat{d}_k^{i_t} - d_k^{i_t}) \hat{d}_k^{i_t}(1 - \hat{d}_k^{i_t})$

             $\Rightarrow V_{k\ell}^{(t+1)} = V_{k\ell}^{(t)} - \alpha_t \, \delta_k^{i_t} \cdot h_\ell^{i_t}$

             $\Rightarrow \gamma_m^{i_t} = \sum_{q=1}^{Q} \delta_q^{i_t} V_{qm}^t \, h_m^{i_t}(1 - h_m^{i_t})$

             $\Rightarrow W_{mj}^{(t+1)} = W_{mj}^{(t)} - \alpha_t \, \gamma_m^{i_t} \cdot x_j^{i_t}$

⊛ Cost/Loss func. is non convex function of weights!

    * May converge to local minima.

    * Some empirical tricks: "Batch SGD", "Normalization"

    * Can add regularization.