$$\min_{w} \sum_{i=1}^{n} L\left(w^T x_i, y_i\right) + R(w)$$

$\underbrace{\phantom{\sum_{i=1}^{n} L\left(w^T x_i, y_i\right)}}$ ↑ LOSS

$R(w)$ ⟶ Regularizer

---

## NEURAL NETWORKS

$x \in \mathbb{R}^d$ $\qquad$ $\text{sign}\left(w^T x\right)$

$[x_1 \ x_2 \ \cdots \ x_d]$



$x_1 \rightarrow \bigcirc \ w_1$

$x_2 \rightarrow \bigcirc \ w_2$

$\vdots$

$x_d \rightarrow \bigcirc \ w_d$ $\quad \rightarrow w^T x$

---

PARAMETERS

$\left\{ w^{(1)}, \cdots, w^{(k)} \right\}$ $\quad w^{(i)} \in \mathbb{R}^d$

$w^{out} \in \mathbb{R}^k$

Input Layer $\qquad$ Hidden Layer $\qquad$ output Layer



$a\left(w^{(1)T} x\right) \cdot w_1^{out}$

$a\left(w^{(2)T} x\right)$

$a\left(w^{(k)T} x\right) \ w_k^{out}$

$w^{out^T} \left[ a\left(w^{(1)T} x\right) \cdots a\left(w^{(k)T} x\right) \right]$

# hidden layer nodes $\qquad$ activation function

$$\hat{y} = \sum_{i=1}^{k} w_i^{out} \ a\left(w^{(i)T} x\right)$$

---

Examples of activation functions / non-linearities



$\bullet \quad a(z) = \dfrac{1}{1 + e^{-z}}$ $\qquad$ [SIGMOID]



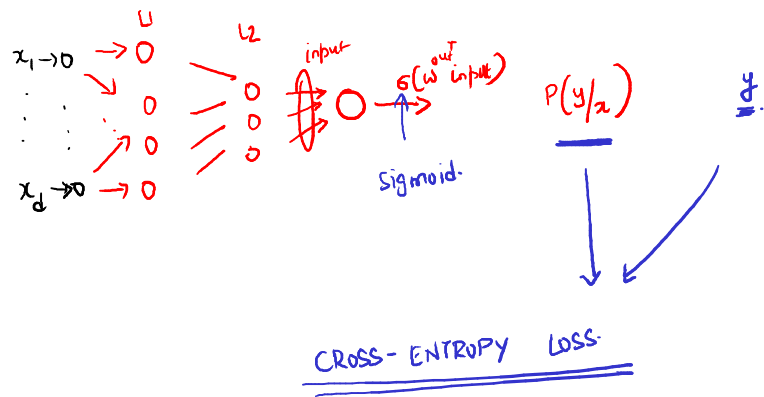$\bullet \quad a(z) = \max(0, z)$ $\qquad$ [Rectified Linear unit]

$$\{w^{(1)}, \ldots, w^{(L)}, w^{out}\}$$

Regression $\qquad \mathcal{L}\left(NN(x_i, \theta), \; y_i\right)$

$$= \sum_{i=1}^{n} \left(NN(x_i; \theta) - y_i\right)^2 \leftarrow$$

$$\equiv \underline{w}\, x_i$$

▸ Learn $\theta^*$ using __Gradient descent__



Feature mapping

Input

Layer 1    Layer 2   ····   Layer L

output Layer

$\in \mathbb{R}^d$         $\in \mathbb{R}^M$    $\in \mathbb{R}$

- Gradient computed taking advantage of chain rule $\rightarrow$ __BACK-PROPAGATION__

- Converges to local minima!



$x_1 \rightarrow 0$   L1   L2   input   $\sigma(w^{out^T} input)$   $P(y/x)$   $y$

$\vdots$

$x_d \rightarrow 0$     Sigmoid.

CROSS-ENTROPY LOSS

Conclusion

     NN

CNN, RNN,    • Learn's local minima of non-concave functions

LSTM, Attentions

Transformers,    • Typically works very well in practice.

             ↳ especially for unstructured data.