

* **Activation (Transformation)** for Deep Learning

$$S_j^{(L)} = \sum_{i=0}^{d^{(L-1)}} w_{ij}^{(L)} x_i^{(L-1)}$$

$$x_j^{(L)} = \phi_j^{(L)}(S_j^{(L)})$$

* previously

$\phi_j^{(L)}$ for Last layer : depend on desired output

$\phi_j^{(L)}(s) = s \leftarrow$ regression

e.g.

$\phi_j^{(L)}(s) = \frac{1}{1+e^{-s}} \leftarrow$ logistic regression (for bin. classification)

$\phi_j^{(L)}$ for hidden layer : "soft perceptron"

e.g. $\phi_j^{(L)}(s) = \tanh(s)$

* scalar activation \Rightarrow

Vector activation (joint)

$\phi(s)$ $\phi(\underline{s})$

used when needing to consider a vector "jointly" of scores

e.g. converting scores of multi-classes to prob. estimates

that (Σ to 1) Δ non-negative

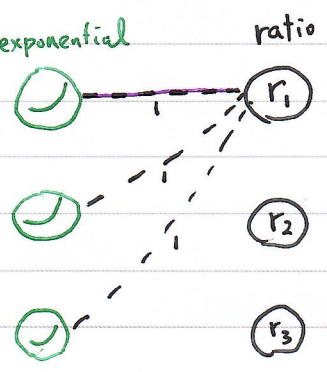
* softmax activation

$$\phi(\underline{s}) = \left[\frac{e^{s_1}}{\sum e^{s_j}}, \frac{e^{s_2}}{\sum e^{s_j}}, \dots, \frac{e^{s_d}}{\sum e^{s_j}} \right]$$

extension of $\frac{1}{1+e^{-s}}$ for multiclass classification

(often for $\phi^{(L)}$)

equivalently exponential



$$r_i(\underline{x}) = \frac{x_i}{\sum x_j}$$

$$r_{\bar{i}}(\underline{x}) = \frac{x_{\bar{i}}}{\sum x_j}$$

backprop (chain rule) still applicable after modifying

* $\phi_j^{(l)}$: the bottleneck of being deep

$$\begin{aligned} \nabla_{i,j}^{(l)} &= \delta_j^{(l)} x_i^{(l-1)} \\ \frac{\partial e}{\partial s_j^{(l)}} &= \sum_k \delta_k^{(l+1)} (w_{jk}^{(l+1)}) \phi'(s_j^{(l)}) \\ &= \sum_k \sum_m \delta_m^{(l+2)} (w_{km}^{(l+2)} w_{jk}^{(l+1)}) \phi'(s_k^{(l+1)}) \phi'(s_j^{(l)}) \end{aligned}$$

early layers : lots of $\begin{pmatrix} \phi' \\ \underline{w} \end{pmatrix}$ in $\nabla^{(l)}$ on top of $\underline{\delta}^{(l)}$

* traditional NNet

$\phi(s) = \tanh(s)$

- $x_j^{(l)} \in (-1, 1)$
- $\phi'(s) \in (0, 1)$

if $w_{ij}^{(l)}$ large $\phi'(s) \approx 0 \rightarrow$ saturation

if $w_{ij}^{(l)}$ small $\rightarrow \delta_i^{(l-1)}$ small $\rightarrow \nabla_{ki}^{(l-1)}$ small) get smaller in earlier layers

\downarrow vanishing gradient

* modern deep learning

$\phi(s) = \max(s, 0)$

rectified linear unit (ReLU)

- $x_j^{(l)} \in [0, \infty)$
- $\phi'(s) = \begin{cases} 1 : s \geq 0 \\ 0 : s < 0 \end{cases}$

(most) if $w_{ij}^{(l)}$ too negative $\phi'(s) = 0, \phi(s) = 0$

\rightarrow dead neuron

- less issues about diminishing gradient (0)
- more efficient arithmetic operations (0)
- benefits of sparsity (0)
- some points not differentiable (Δ nowadays)

* ReLU variants

$\phi(s) = \max(0, s) + \underbrace{0.1 \min(0, s)}_{\text{wake-up}}$ leaky ReLU

$\phi(s) = \max(0, s) + \underbrace{"W" \min(0, s)}_{\text{learned wake-up}}$ parameterized ReLU

"usually" ReLU & leaky ReLU suffice

* Initialization for Deep Learning

- activation \Rightarrow bad $W_{ij}^{(l)}$ properties

w.r.t. $x_i^{(l)}$

e.g. tanh: $|W_{ij}^{(l)}|$ too big \rightarrow saturation

$W_{ij}^{(l)} = 0, W_{ij}^{(l)} = \text{constant}$ \rightarrow saddle, symmetry, etc.

ReLU: $W_{ij}^{(l)}$ too negative for $x_i^{(l)} \geq 0$

* want: randomly small

