

\* activation (transformation)

initialization

optimization

regularization

"pre-training"

L2, early stopping

\* activation

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

$$x_j^{(l)} = \phi_j^{(l)}(s_j^{(l)})$$

scalar

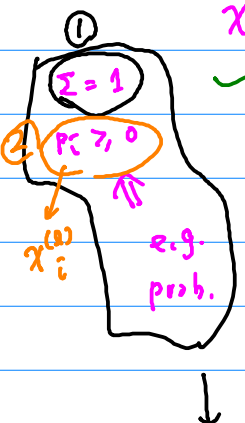
vector

$$\underline{x}^{(l)} = \underline{\phi}(\underline{s}^{(l)})$$

activation property

real values

(joint)



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

regression

(logistic)

bin. classification

$\phi_j^{(l)}$  hidden layer: "soft perc."

tanh

multi class

\* softmax activation

output often

exponential

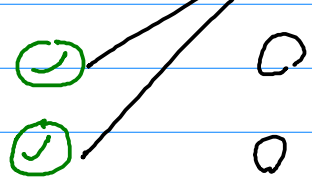
$$\underline{\phi}(\underline{s}) = \left( \frac{e^{s_1}}{\sum e^{s_i}}, \frac{e^{s_2}}{\sum e^{s_i}}, \frac{e^{s_3}}{\sum e^{s_i}}, \dots, \frac{e^{s_d}}{\sum e^{s_i}} \right)$$

ratio (normalized)

$\approx 0$

$\approx 1$

$$\frac{a}{a+b+c}$$



"general" backprop

\* bottleneck of being deep:

activation

$$s_j^{(l)} = \sum_{i=0}^{d^{(l-1)}} W_{ij}^{(l)} x_i^{(l-1)}$$

$$x_j^{(l)} = \phi_j^{(l)}(s_j^{(l)})$$

$$\frac{\partial e}{\partial s_j^{(l)}} = \sum_k \delta_k^{(l+1)} W_{jk}^{(l+1)} \phi'(s_j^{(l)})$$

$$= \sum_k \sum_m \delta_m^{(l+2)} (W_{km}^{(l+2)} W_{jk}^{(l+1)}) \phi'(s_k^{(l+1)}) \phi'(s_j^{(l)})$$

$$s_j^{(l)} = \text{many } W \cdot \text{many } \phi' \delta_j^{(l)}$$

\* traditional NNet

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

$$x_j^{(l)} \in (-1, 1)$$

$$\phi'(s) \in (0, 1)$$

①  $W_{ij}^{(l)}$  large  $\rightarrow s$  large

$$\phi(s) \approx 1$$

$s$

$-1$

$$\phi'(s) \approx 0$$

saturation

②  $W_{ij}^{(l)}$  small

$\rightarrow s_j^{(l)}$  small

$\phi'$  "shrink"

earlier layer small gradient  
vanishy gradient

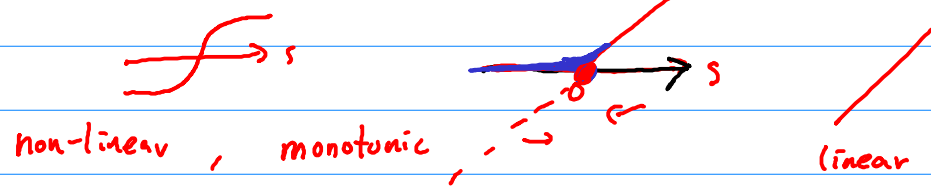
\* modern deep learning

rectified linear unit (ReLU)

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

$$x_j^{(l)} \in [0, \infty)$$

$$\phi'(s) = \begin{cases} 1 & s \geq 0 \\ 0 & s < 0 \end{cases}$$



$$\phi' \phi' \phi' \dots$$

dead

- o (o) - less issues on gradient vanishing
- x (o) - more efficient operations (no exp, log --)
- (o) - sparsity
- e (Δ) - not fully differentiable  $\phi'$