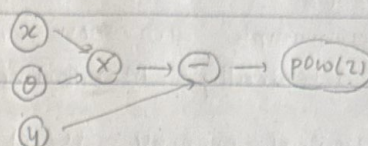
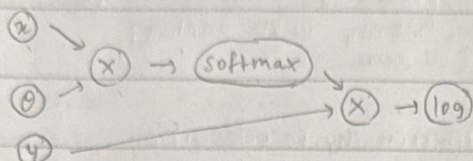


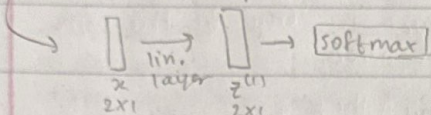
Computation Graph - analogous to a "MATH DAG". NN is one large comp. graph.

↳ los. Reg.

↳ Lin. Reg.

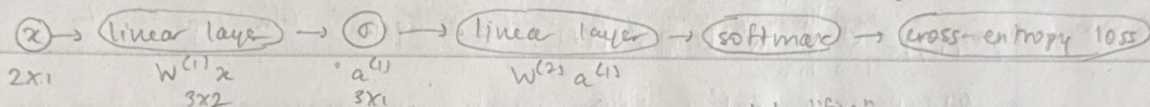


data (x, y); parameters (θ) - usually used once.



which layer

How to represent learned features - represent as logistic regression output ($\sigma(W^{(1)}x)$)



how much to "fire"
based on incoming signals

We can use non-linear functions (activation functions) in slw. layers (lin. als. lasagna)

Chain rule: $\frac{d}{dx} f(g(x)) = \sum_{j=1}^m \frac{dy_j}{dx} \frac{dz}{dy_j}$ (sum over all dimensions of y)

$$\frac{d}{dx} f(g(x)) = \frac{dy}{dx} \frac{dz}{dy}$$

$\begin{matrix} \xrightarrow{1 \times m} & & \xleftarrow{m \times 1} \end{matrix}$

$\left(\frac{dy}{dx}\right)_{ij} = \frac{dy_j}{dx_i}$

compute Jacobians right-to-left: saves compute time by factor of n (backprop).

Backpropagation algorithm: ("classic" version)

→ forward pass: calculate each $a^{(i)}$ and $z^{(i)}$

→ backward pass:

initialize $\delta = \frac{df}{dz^{(L)}}$

for each l w/ input x_l & params θ_l from end to start:

$$\frac{df}{d\theta_l} \leftarrow \frac{df}{dz^{(l)}} \delta$$

$$\delta \leftarrow \frac{df}{dx_l} \delta$$

NN arch details:

1. How many layers?
2. How big are the layers?
3. What type of activ. func.? Popular choice: ReLU (efficient).

$$z^{(l+1)} = W^{(l)} a^{(l)} + b^{(l)} \leftarrow \text{additional param in each lin. layer.}$$

$$\frac{dz}{dw} \delta = \sum_i \frac{dz_i}{dw} \delta_i$$

$$\frac{dz_i}{dw_{jk}} = \begin{cases} 0 & \text{if } j \neq i \\ \delta a^j & \text{otherwise} \end{cases}$$

$$\frac{dz}{da} \delta = W^T \delta$$