

Comp Sci, January 24, 2023

NN architecture (model)

- # hidden layers
- # hidden neurons in a layer
- activation function (Sigmoid, tanh, ReLU, ELU, ...)

cost function & optimization

- Type of cost/loss
- Regularization, l_1 , l_2
- GD, stochastic, batch, learning rate, epochs
 - optimization schemes

Other specifics: pretraining, initialization, dropout, ...

Universal approximation theorem:

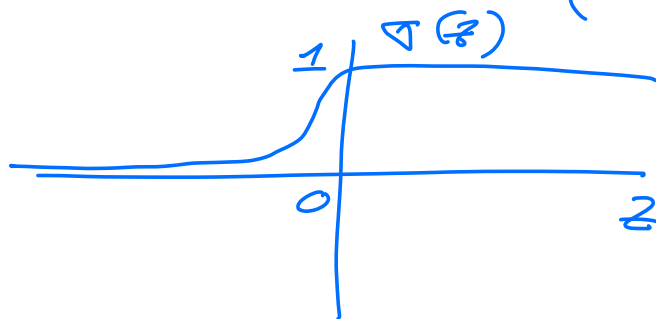
$y = F(x) \quad x \in [0, 1]^d$
 deterministic function assumed to be continuous

$$F \in \mathbb{C} [0, 1]^d$$

(Cybenko, 1989)

Let σ be any continuous sigmoidal function

$$\sigma(z) \rightarrow \begin{cases} 1 & \text{as } z \rightarrow \infty \\ 0 & \text{as } z \rightarrow -\infty \end{cases}$$



given a function $F \in \mathbb{C} [0, 1]^d$
 and $\epsilon > 0$, there is a one-
 layer NN $f(x; \Theta)$ of the
 form $W \in \mathbb{R}^{m \times n}$ and
 $b \in \mathbb{R}^m$ for which

$$|f(x; \epsilon) - F(x)| < \epsilon$$

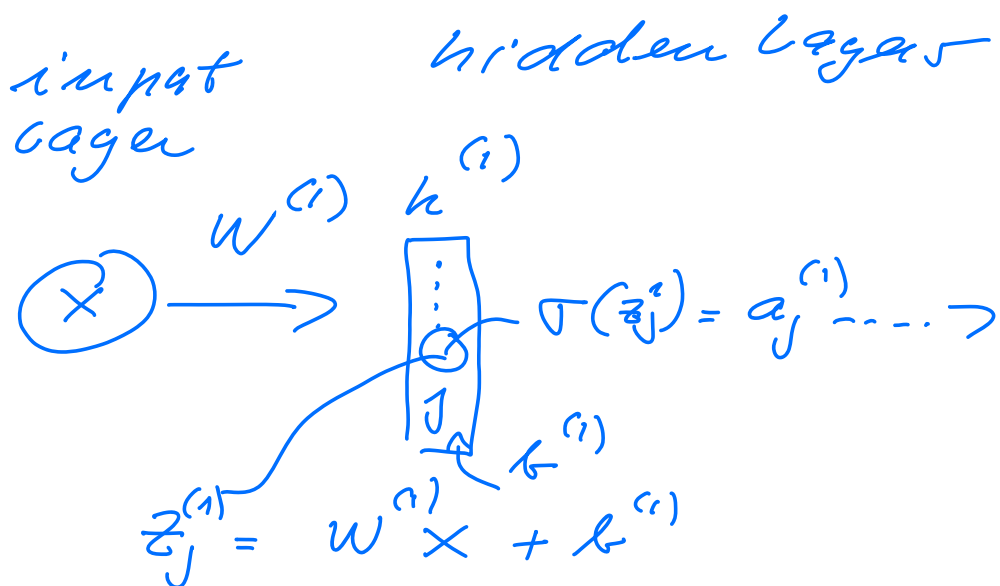
$$\epsilon = \{w, b\}$$

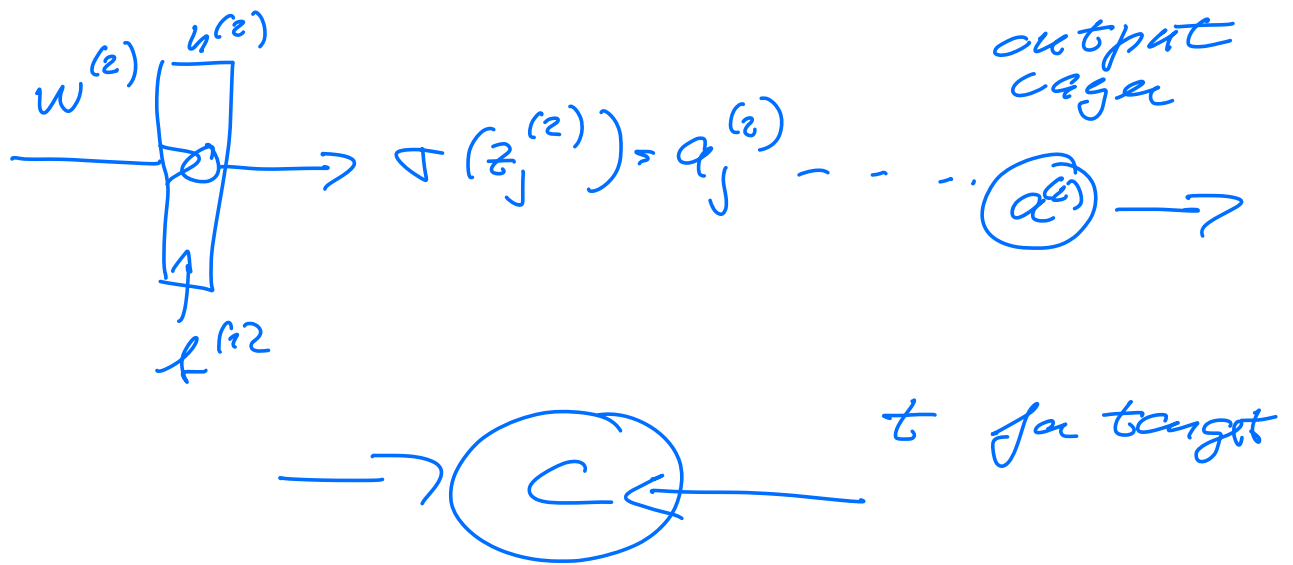
for all $x \in [a, b]^d$

Hornik (1991) extended the theorem to apply to any non-constant, bounded activation function.

Basics of ANN

— Feed Forward stage





- Back propagation stage

Final output layer

$$a^{(L)}(x; \Theta)$$

$$\Theta = \left\{ w^{(1)}, b^{(1)}, w^{(2)}, b^{(2)}, \dots, w^{(L)}, b^{(L)} \right\}$$

$$\text{MSE}$$

$$C(\Theta) = \frac{1}{n} \| (t - a^{(L)}(x; \Theta)) \|_2^2$$

$$a^{(L)}(x; \Theta) = \sigma^{(L)}(\sigma^{(L-1)}(\sigma^{(L-2)} \dots \sigma^{(1)}(x w^{(1)} + b^{(1)})) \dots)$$

Back propagation algo

- problems with vanishing gradient (or exploding)

consider a simple NN in which x, w, b are all scalars and reals.

$$L = 2$$

$$f(x; \theta) = \sigma_2(w_2 \sigma_1(w_1 x + b_1) + b_2)$$

$$\partial_{w_1} f(x; \theta) = \sigma_2' (w_2 \sigma_1(w_1 x + b_1) + b_2) \times w_2 \sigma_1'(w_1 x + b_1) x$$

with L -layers

$$\left[\prod_{l=2}^L w_l \right] \times \left[\prod_{l=1}^L \sigma_l'(z_l) \right] x$$

$$z_l = A_l [\nabla_{l-1} (A_{l-1} (\dots \nabla_1 (A_1(x)) \dots))]$$

if ∇_l is a sigmoid function (or tanh)

then $\nabla_l'(z_l)$ will be small when $|z_l| \gg 0$

Derivation of Back prop algo

— analytical expressions

for $\frac{\partial C}{\partial w^{(l)}} \wedge \frac{\partial C}{\partial b^{(l)}}$

— then we can update

$$w^{(l)} \leftarrow w^{(l)} - \eta^{(l)} \frac{\partial C}{\partial w^{(l)}}$$

$$b^{(l)} \leftarrow b^{(l)} - \eta^{(l)} \frac{\partial C}{\partial b^{(l)}}$$

input to node j in layer- l

$$z_j^l = \sum_{i=1}^{N_{l-1}} w_{ji}^l a_i^{l-1} + b_j^l$$

$$a_i^{l-1} = \sigma(z_i^{l-1})$$

$$\frac{\partial z_j^l}{\partial w_{ij}^l} = a_i^{l-1}$$

$$\frac{\partial z_j^l}{\partial a_i^{l-1}} = w_{ij}^l$$

$$\frac{\partial a_j^l}{\partial z_j^l} = ?$$

$$\text{assume } a_j^l = \sigma^l(z_j^l)$$

$$= \frac{1}{1 + e^{-z_j^l}}$$

$$\frac{\partial a_j^l}{\partial z_j^l} = \sigma(z_j^l) (1 - \sigma(z_j^l))$$

$$= a_j^L (1 - a_j^L)$$

Computational note:

these derivatives will change when we change activation function

Example of $C(\theta)$

$$C(\theta) = \frac{1}{2} \sum_i (a_i^L - t_i)^2$$

$$\frac{\partial C(\theta)}{\partial w_{jk}^L} \quad \wedge \quad \frac{\partial C}{\partial b^L}$$

$$\frac{\partial C}{\partial w_{jk}^L} = (a_j^L - t_i) \frac{\partial a_j^L}{\partial w_{jk}^L}$$

Chain rule:

$$\frac{\partial a_j^L}{\partial w_{jk}^L} = \frac{\partial a_j^L}{\partial z_j^L} \frac{\partial z_j^L}{\partial w_{jk}^L}$$

with sigmoid for $\sigma(z)$

$$\frac{\partial a_j^L}{\partial w_{jk}^L} = \underbrace{a_j^L (1 - a_j^L)}_{\sigma'(z)} a_k^{L-1}$$

$$\begin{aligned} \frac{\partial C}{\partial w_{jk}^L} &= (a_j^L - t_j) a_j^L (1 - a_j^L) a_k^{L-1} \\ &= (a_j^L - t_j) \sigma'(z_j^L) a_k^{L-1} \end{aligned}$$

$$\begin{aligned} \delta_j^L &= a_j^L (1 - a_j^L) (a_j^L - t_j) \\ &= \sigma'(z_j^L) \frac{\partial C}{\partial a_j^L} \end{aligned}$$

$$\frac{\partial C}{\partial w_{jk}^L} = \delta_j^L a_k^{L-1}$$

$$\delta_j^L = \frac{\partial C}{\partial z_j^L} = \frac{\partial C}{\partial a_j^L} \frac{\partial a_j^L}{\partial z_j^L}$$

$$\delta_j^L = \frac{\partial C}{\partial b_j^L} \frac{\partial b_j^L}{\partial z_j^L} = \frac{\partial C}{\partial b_j^L}$$

Collect for $l = L$

$$\frac{\partial C}{\partial w_{jk}^L} = \delta_j^L a_k^{L-1}$$

$$\delta_j^L = \nabla'(z_j^L) \frac{\partial C}{\partial a_j^L}$$

$$\delta_j^L = \frac{\partial C}{\partial b_j^L}$$

$L \rightarrow l \ (\neq L)$

$$\delta_j^l = \frac{\partial C}{\partial z_j^l} = \sum_k \frac{\partial C}{\partial z_k^{l+1}} \frac{\partial z_k^{l+1}}{\partial z_j^l}$$

$$= \sum_k \delta_k^{l+1} \frac{\partial z_k^{l+1}}{\partial z_j^l}$$

$$z_j^{l+1} = \sum_{i=1}^{N_l} w_{ij}^{l+1} \frac{a_i^l}{\Delta(z_j^l)} + b_j^{l+1}$$

$$\frac{\partial z_k^{l+1}}{\partial z_j^l} = w_{kj}^{l+1} \Delta'(z_j^l)$$

$$\delta_j^l = \sum_k \delta_k^{l+1} w_{kj}^{l+1} \Delta'(z_j^l)$$

Final algo

- (i) initialize $\Theta = \{W, b\}$
- (ii) Perform the 1st Feed forward pass
- (iii) a_i^L is used to compute $C(\theta)$
- (iv) Perform Back prop

for $l = L-1, L-2, \dots, 2$

$$\delta_j^l = \sum_k \delta_k^{l+1} w_{kj}^{l+1} \sigma'(z_j^l)$$

$$w_{jk}^l \leftarrow w_{jk}^l - \eta \delta_j^l a_k^{l-1}$$

$$b_j^l \leftarrow b_j^l - \eta \frac{\partial \mathcal{L}}{\partial b_j^l}$$
$$= b_j^l - \eta \delta_j^l$$

learning

end for

(V) repeat (ii) - (iv) till

$\mathcal{L}(\theta)$ stabilizes.

output $\theta = \{W, b\}$,
which gives us the optimal
description of $F(x)$

Automatic differentiation

JAX replaces autograd.

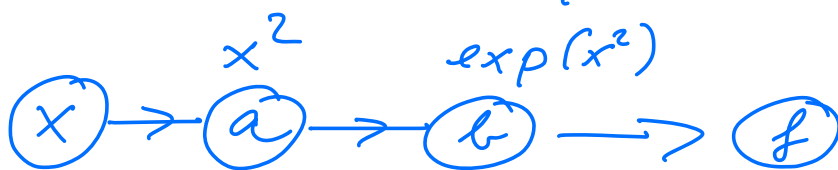
Example

$$f(x) = \exp(x^2)$$

$$f'(x) = 2x \exp(x^2)$$

Def $a = x^2$

$$b = \exp(a) = f(x)$$



$$\frac{df}{dx} = \left[\frac{df}{db} \frac{db}{da} \right] \frac{da}{dx}$$

$\stackrel{1}{\text{Reverse mode}}$

$$= \frac{df}{db} \left[\frac{db}{da} \frac{da}{dx} \right]$$

Forward mode

New example

$$f(x) = \sqrt{x^2 + \exp(x^2)}$$

$$a = x^2; \quad b = \exp(a)$$

$$c = a + b$$

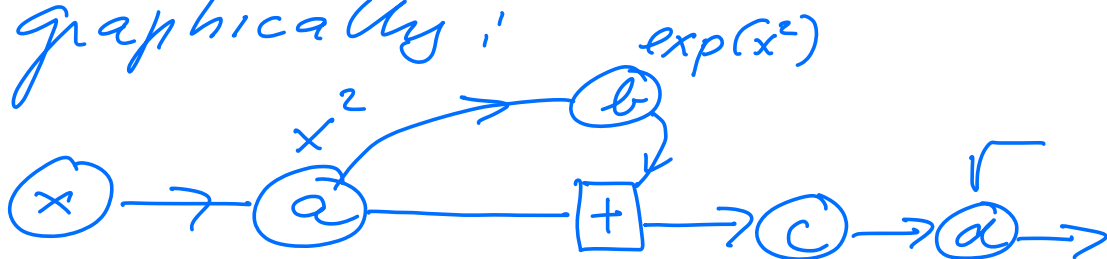
$$d = \sqrt{c} = f(x)$$

$$\frac{df}{dx} = \frac{x(1 + \exp(x^2))}{\sqrt{x^2 + \exp(x^2)}}$$

$$= \frac{x(1+b)}{\sqrt{a+b}} = \frac{x(1+b)}{\sqrt{c}}$$

$$= \frac{x(1+b)}{d}$$

graphically:



→ (f)

$$f(x) = d$$

$$\frac{\partial a}{\partial x} = 2x$$

$$\frac{\partial b}{\partial x} = \frac{\partial b}{\partial a} \frac{\partial a}{\partial x}$$

$$\frac{\partial c}{\partial x} = \left[\frac{\partial c}{\partial a} \frac{\partial a}{\partial x} + \frac{\partial c}{\partial b} \frac{\partial b}{\partial x} \right]$$

$$\frac{\partial d}{\partial c} = \frac{1}{2\sqrt{c}}$$

$$\frac{\partial d}{\partial x} = \frac{\partial d}{\partial c} \frac{\partial c}{\partial x} = \frac{\partial f}{\partial x}$$

Computation in reverse
mode : More next
week.