# Comp Sci program, JAN 31, 2023

## FFNN

input            hidden 1



$$\sigma^1(z^1) \qquad \sigma^2(z^2)$$

$$w^1 \qquad w^1 \; b^1 \qquad w^2 \; b^2$$

fixed

$$z^1 = w^1 x + b^1$$

output

$$\sigma^{L-1}(z^{L-1}) \qquad \sigma^L(z^L)$$

$$w^L \; b^L$$

$$y = a^L(\Theta, x) = \sigma^L(z^L)$$

$$\Theta = \{ w^1, b^1, w^2, b^2 \; \cdots \; w^L, b^L \}$$

$$\sigma^L(z^L) = \sigma^L(\sigma^{L-1}(\sigma^{L-2}(\cdots \sigma^1(z^1)\cdots)$$

$$C(\Theta) = \| t - a^L(\Theta, x) \|_2^2$$

( MSE in regression )

$$\boxed{\frac{\partial C}{\partial \Theta^L} = 0}$$

$$\frac{\partial C}{\partial \Theta^{L-1}} = \frac{\partial C}{\partial a^L} \frac{\partial a^L}{\partial \Theta^{L-1}}$$

$$\frac{\partial C}{\partial \Theta^{L-2}} = \frac{\partial C}{\partial a^L} \frac{\partial a^L}{\partial a^{L-1}} \frac{\partial a^{L-1}}{\partial \Theta^{L-2}}$$

$$\frac{\partial C}{\partial \Theta^\ell} = \frac{\partial C}{\partial a^L} \frac{\partial a^L}{\partial a^{L-1}} \frac{\partial a^{L-1}}{\partial a^{L-2}} \cdots \frac{\partial a^{\ell+1}}{\partial \Theta^\ell}$$

$$\Theta^\ell_{(k+1)} \leftarrow \Theta^\ell_{(k)} - \mu_{(k)} \frac{\partial C}{\partial \Theta^\ell}\bigg|_{\Theta^\ell = \Theta^\ell_{(k)}}$$

## Automatic differantion

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

$x \cdot x = 1$ FLOP

$\exp(x^2) = \exp(x \cdot x) = 2$ FLOP

$x^2 \, \textcircled{+} \, \exp(x^2) = 1$ FLOP

$$\text{SQRT}(\cdot\cdot) = 1 \text{ FLOP}$$

---

5 FLOP

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

10 FLOPS

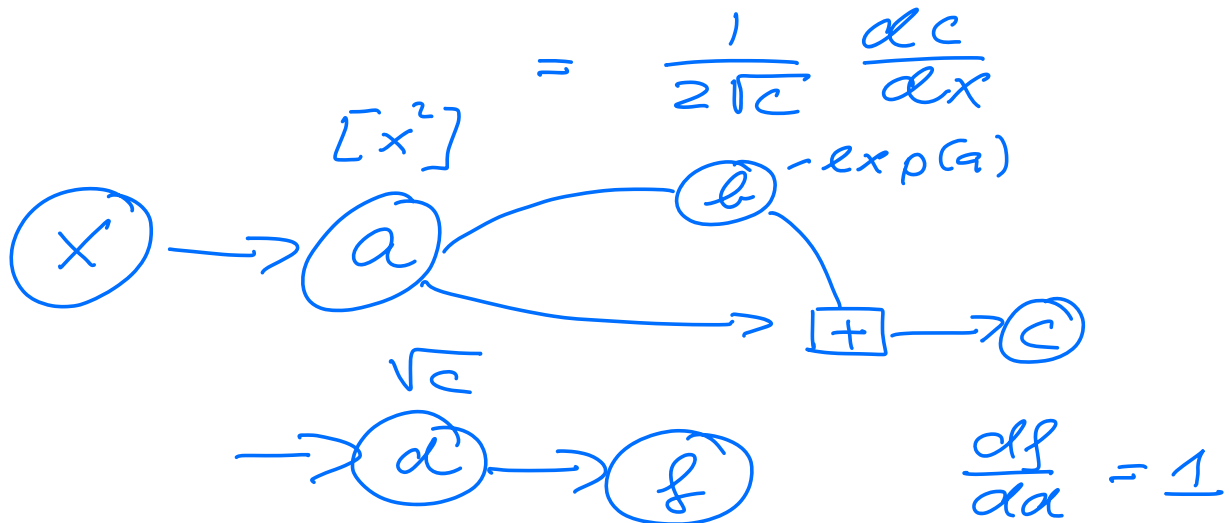automatic diff:

$$a = x^2$$

$$b = \exp(a)$$

$$c = a + b$$

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

$$\frac{da}{dx} = 2x \qquad \frac{db}{dx} = \frac{db}{da}\frac{da}{dx}$$

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

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

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

$$\frac{d\,d}{dc} = \frac{1}{2\sqrt{c}}$$

$$\frac{d\,d}{dx} = \frac{df}{dx} = \frac{d\,d}{dc}\frac{dc}{dx}$$

$$= \frac{1}{2\sqrt{c}}\frac{dc}{dx}$$

$[x^2]$

$-\exp(a)$

$\sqrt{c}$

$$\frac{df}{da} = 1$$

compute $\dfrac{df}{dx}$ in backward/ reverse mode

$$\frac{df}{dc} = \frac{df}{dd}\frac{dd}{dc} = \frac{1}{2\sqrt{c}}$$

$\underset{\searrow}{1}$

$$\frac{df}{db} = \frac{df}{dc}\frac{dc}{db} = \frac{1}{2\sqrt{c}}$$

$$c = a+b \qquad \frac{dc}{db} = 1$$

$$\frac{df}{da} = \frac{df}{db}\frac{db}{da} + \frac{df}{dc}\frac{dc}{da}$$

$$= \frac{1}{2\sqrt{c}}\left[1 + \exp(a)\right]$$

$$\frac{df}{dx} = \frac{df}{da}\frac{da}{dx}$$

$$= \frac{x\left(1 + \exp(x^2)\right)}{\sqrt{x^2 + \exp(x^2)}}$$

$$\frac{df}{dx} = \frac{x\left(1 + b\right)}{\sqrt{a + b}}$$

$b$ is calculated in $f$

numerator : 2 Flops
denominator : 2 Flops
+ Division : 5 Flops

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

$$\frac{df}{dx} = \frac{x(1+b)}{f(x)}$$

Formalization

assume we have $x_1, x_2 \ldots x_d$
input variables to $f$,

$x_{d+1}, x_{d+2} \ldots x_D$ intermediate
variables $x_D$ = output
                      variable

in previous example

$x_1 = x$       $d = \underline{1}$

$x_2 = a$       $x_3 = b$       $x_4 = c$

$x_D = d = f$

For $i = d+1, \ldots D$

$$x_i = g_i(x_{pa(x_i)})$$

$g_i$ are elementary functions
and $x_{pa(x_i)}$ are the parent
nodes of the variable $x_i$

$$g_2 = (x \cdot x)^2 = a$$

$$g_3 = \exp(\sqrt{\ }) = b$$

$$g_4 = c = a + b$$

$$g_5 = \sqrt{c} = d = f$$

By def $\quad \dfrac{df}{dx_D} = 1$

Reverse mode a Backprop

$$\frac{df}{dx_i'} = \sum_{\substack{x_j \\ x_i = Pa(x_j)}} \frac{df}{dx_j} \frac{dg_j}{dx_i'}$$

$Pa(x_j) = $ set of parent

nodes of $x_i'$

$$\frac{df}{dd} = 1$$

$$\frac{df}{dc} = \underbrace{\frac{df}{dd}}_{=1} \frac{dd}{dc} = \frac{1}{2\sqrt{c}}$$

$$\frac{df}{da} = \frac{df}{db}\frac{db}{da} + \frac{df}{dc}\frac{dc}{da}$$

$$\frac{df}{dx} = \frac{df}{da}\frac{da}{dx} = \frac{x(1+b)}{a}$$

Simple neural network
example :

---



$$a = f(x, w_1) \qquad y = g(a, w_2)$$

$$= g(f(x, w_1), w_2)$$

$$C = \frac{1}{2}(t-y)^2$$

$$w_1^{(k+1)} \leftarrow w_1^{(k)} - \delta \frac{\partial C}{\partial w_1}\Big|_{w_1 = w_1^{(k)}}$$

$$w_2^{(k+1)} \leftarrow w_2^{(k)} - \delta \frac{\partial C}{\partial w_2}\Big|_{w_2 = w_2^{(k)}}$$

linear activation function

$$f(x_1, w_1) = w_1 x = a$$

$$g(a, w_2) = w_2 \cdot a = w_2 f(x_1, w_1)$$

$$\frac{\partial C}{\partial w_1} = -(t-y) \frac{\partial y}{\partial a} \frac{\partial a}{\partial w_1}$$

$$\frac{\partial C}{\partial w_2} = -(t-y) \frac{\partial y}{\partial w_2}$$

$$\frac{\partial C}{\partial w_1} = -(t-y) x w_2$$

$$\frac{\partial C}{\partial w_2} = -(t-y) w_1 x$$

$$\frac{\partial C}{\partial w_1} = -(t-y) \frac{\partial y}{\partial a_2} \frac{\partial a}{\partial a_{L-1}} \cdots$$

$$\frac{\partial a_L}{\partial a_1} \frac{\partial a_1}{\partial w_1}$$

# Ingredients for an NN-code

- architecture (model)
  - \# layers
  - \# nodes
  - \# activation functions and their derivatives
- cost function
- regularization & optimization

  - regularization parameter $\lambda$ with $l_1$ or $l_2$

  - gradient descent methods (GD)
    - GD with momentum
    - SGD with & without momentum
    - learning rates
      - Adagrad
      - RMSprop
      - ADAM