# Neural Networks & Gradient Descent

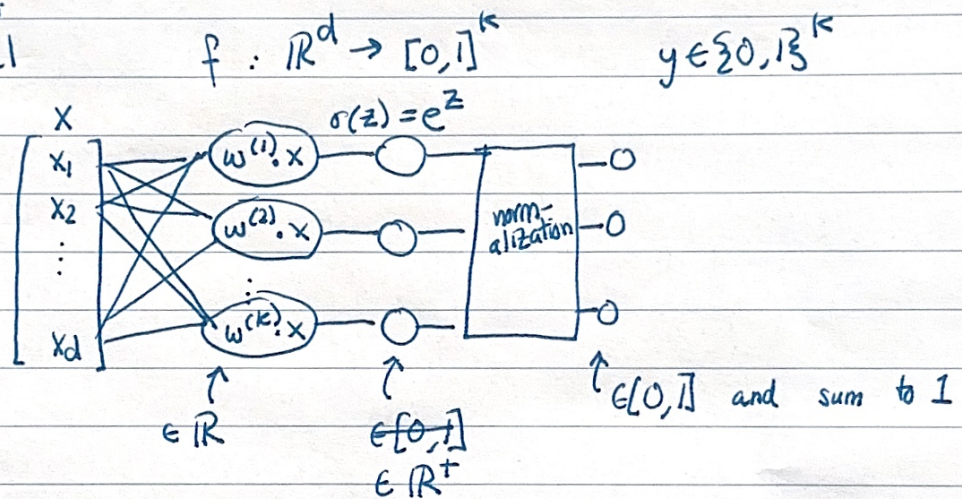**Plan:**
- Review
- Neural Networks
- Gradient Descent
- Back propagation

**Logistics:**
- zoom :)
- scribed notes (moved by 2)

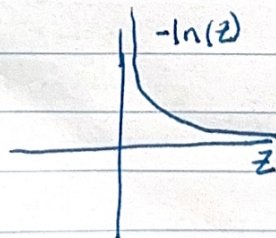## [Review]:

① Model $\qquad f : \mathbb{R}^d \to [0,1]^k \qquad\qquad y \in \{0,1\}^k$



$\sigma(z) = e^z$

norm-alization

$\in \mathbb{R}$

$\in [0,1]$
$\in \mathbb{R}^+$

$\in [0,1]$ and sum to 1

② Loss : Cross Entropy

$$H(y, f(x)) = -\sum_{i=1}^{k} y_i \log f(x)_i = -\log f(x)_{i*}$$

\+ efficient $\qquad\qquad$ + big gradient if $f(x)_{i*} \approx 0$



$-\ln(z)$

$z$

# Neural Networks:



$x$ — weight layer — activation — layer — activation

$$\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix}$$

nodes: $x \cdot w^{(1)}$, $x \cdot w^{(2)}$, $x \cdot w^{(k)}$

## Activation functions

| | |
|---|---|
| sigmoid | $\frac{1}{1+e^{-z}}$ |
| ReLU | $\max(0, z)$ |
| step | $\mathbb{1}[z > 0]$ |
| sign | $\mathbb{1}[z > 0] - \mathbb{1}[z \leq 0]$ |

$$\underset{k \times d}{\underbrace{\begin{bmatrix} \ \end{bmatrix}}} \ \underset{d \times 1}{\underbrace{\begin{bmatrix} \ \end{bmatrix}}} = \begin{bmatrix} w^{(1)} x \\ w^{(2)} x \\ \vdots \\ w^{(k)} x \end{bmatrix}$$

$W \quad x$

## Layers

↳ Dense / fully connected

↳ Convolutional : images or spatial data

↳ Recurrent : sequential data

↳ Residual : processing technique

↳ Attention : sequential data with short + long dependencies

## Gradient Descent:



$w^{(0)}$ ........ $w$

- Initalize $w^{(0)}$
- Iteratively improve loss

$$w^{(t+1)} \leftarrow w^{(t)} - \alpha \nabla_w \mathcal{L}(w)$$

- small $\alpha$ "guarantees" improvement but slowly
- big $\alpha$ moves quickly but can overshoot

$\boxed{\text{Activity:}}$

Gradient descent on linear regression

$$w^{(t+1)} \leftarrow \underline{\phantom{xxxxxxxxx}}$$

Time complexity of optimal solution:

Time complexity of one update:

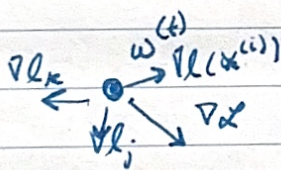$\boxed{\text{Stochastic Gradient Descent:}}$

$$\mathcal{L}(w) = \frac{1}{n} \sum_{i=1}^{n} \ell_i(x^{(i)}, y^{(i)}, w)$$

expensive! fit in memory?

$$S \subseteq [n]$$

$$w^{(t+1)} \leftarrow w^{(t)} - \alpha \frac{1}{|S|} \sum_{i \in S} \nabla_w \ell_i(x^{(i)}, y^{(i)}, w)$$

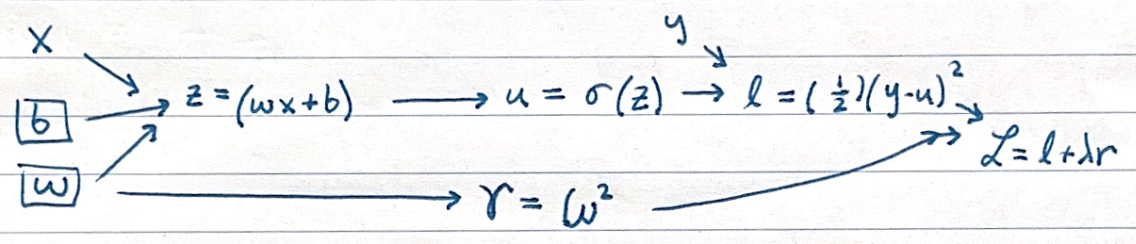Trade off: speed for "representativeness"

## Back propagation

$$\mathcal{L}(w, b) = \tfrac{1}{2}\left(y - \sigma(wx+b)\right)^2 + \lambda w^2$$

regularization on "size of $w$" $\propto$ "complexity"

$$\frac{\partial \mathcal{L}}{\partial w} = (y - \sigma(wx+b)) \cdot - \sigma'(wx+b) \cdot x + 2\lambda w$$

$$\frac{\partial \mathcal{L}}{\partial b} = (y - \sigma(wx+b)) \cdot - \sigma'(wx+b) \cdot 1$$
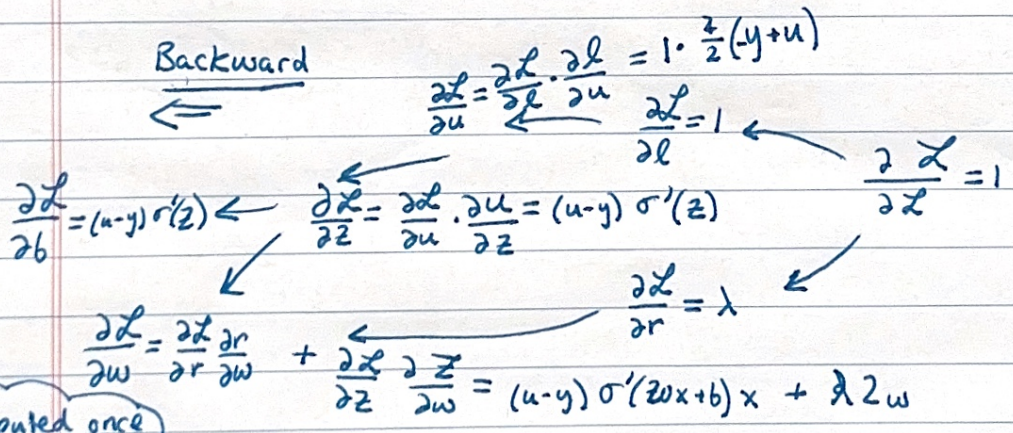
Redundancy!  Complicated!

$x$

$\boxed{b}$

$\boxed{w}$

$z = (wx+b) \longrightarrow u = \sigma(z) \to l = (\tfrac{1}{2})(y-u)^2 \qquad y$

$r = w^2$

$\mathcal{L} = l + \lambda r$

**Forward** $\Rightarrow$

for $i \in \{1, \ldots, N\}$:
compute $v_i$ as function of parents

**Backward** $\Leftarrow$

$$\frac{\partial \mathcal{L}}{\partial u} = \frac{\partial \mathcal{L}}{\partial l} \cdot \frac{\partial l}{\partial u} = 1 \cdot \tfrac{2}{2}(-y+u)$$

$$\frac{\partial \mathcal{L}}{\partial l} = 1$$

$$\frac{\partial \mathcal{L}}{\partial \mathcal{L}} = 1$$

$$\frac{\partial \mathcal{L}}{\partial z} = \frac{\partial \mathcal{L}}{\partial u} \cdot \frac{\partial u}{\partial z} = (u-y)\,\sigma'(z)$$

$$\frac{\partial \mathcal{L}}{\partial b} = (u-y)\,\sigma'(z)$$

$$\frac{\partial \mathcal{L}}{\partial r} = \lambda$$

$$\frac{\partial \mathcal{L}}{\partial w} = \frac{\partial \mathcal{L}}{\partial r}\frac{\partial r}{\partial w} + \frac{\partial \mathcal{L}}{\partial z}\frac{\partial z}{\partial w} = (u-y)\,\sigma'(wx+b)\,x + \lambda 2w$$

+ computed once
+ structured
+ modular

for $i \in \{N, \ldots, 1\}$:

compute $\dfrac{\partial \mathcal{L}}{\partial v_i} = \sum_{j \in \text{children}(v_i)} \dfrac{\partial \mathcal{L}}{\partial v_j} \dfrac{\partial v_j}{\partial v_i}$