Loss function
$$\ell : \mathbb{R}^p \times \mathbb{R}^p \to \overline{\mathbb{R}}_+$$

$$\ell(a,b) = \frac{1}{2} \| a - b \|^2 \qquad \forall \; a, b \in \mathbb{R}^p$$

In ML $\quad f : X \to Y \qquad \ell(\underbrace{f(x)}_{\substack{\uparrow \\ \text{prediction} \\ \text{on } x}}, \overset{\curvearrowleft \text{ lable}}{y})$

Def
A functional $\quad R : X \to \mathbb{R} \quad$ is a map from a space $X$ (usually $\infty$-dimensional) to real numbers

$$\pi \quad \longrightarrow \quad \underset{\Omega}{\int} F(z) \; \underbrace{\frac{d\pi(z)}{dz}}_{\searrow \; g(z)\,\underline{dz}} \qquad z \sim (x,y)$$

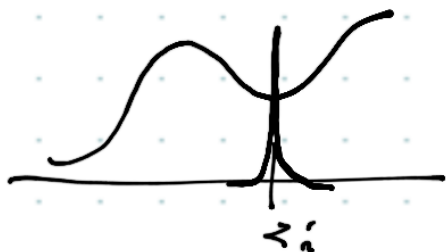If $\pi = \mathcal{L}$ (Lebesgue measure) $\quad d\pi(x,y) = dx\,dy$

Going from the Risk functional to the empirical risk functional

$$R(f) = \underset{\Omega}{\int} \ell(f(x), y) \; d\pi(x,y)$$

$$\pi \to \pi_n = \frac{1}{n} \sum_{i=1}^{n} \delta_{z_i} \qquad \overset{\delta(x-x_i)\,\delta(y-y_i)}{\underset{dx\,dy}{\uparrow}}$$

$$\frac{1}{n} \underset{\Omega}{\int} \ell(f(x),y) \sum_{i=1}^{n} d\delta_{z_i}(x,y) = \frac{1}{n} \sum_{i=1}^{n} \underset{\Omega}{\int} \ell(f(x),y)\, d\delta_{z_i}(x,y)$$

$$= \frac{1}{n} \sum_{i=1}^{n} \ell(f(x_i), y_i) = R_n(f)$$

Empirical risk
functional

$$f_n^* \in \boxed{\underset{f \in \mathcal{F}}{\operatorname{argmin}} \; R_n(f)}$$

↓
This is a set

$$f_n^* = \underset{f \in \mathcal{F}}{\operatorname{argmin}} \; R_n(f)$$

↙ function      ↘ set of functions

The set $X^Y = \{ f : X \to Y \}$

## Uniform Convergence

$$R(f_n^*) - \inf_{f \in \mathcal{F}} R(f) = \left( R(f_n^*) - R_n(f_n^*) \right) \quad ①$$

$$+ \left( R_n(f_n^*) - R_n(f^*) \right) \quad ②$$

$$+ \left( R_n(f^*) - R(f^*) \right) \quad ③$$

$$f^* \in \underset{f \in \mathcal{F}}{\operatorname{argmin}} \; R(f)$$

$$② = R_n(f_n^*) - R_n(f^*)$$

Finite-dimensional example    $R_n(f) = f^2 =$

$$f_n^* = 0 \qquad \mathcal{R}(f) = (f-1)^2$$



$$f^* = 1$$

$$\mathcal{R}_n(f_n^*) = 0$$

$$\mathcal{R}(f^*) = 0$$

$$\mathcal{R}_n(f^*) = 1$$

$$\mathcal{R}_n(f_n^*) - \mathcal{R}_n(f^*) = 0 - 1 = -1 \leq 0$$

We conclude that ② $\leq 0$

$$\mathcal{R}(f_n) - \inf_{f \in \mathcal{F}} \mathcal{R}(f) \leq \left( \mathcal{R}(f_n^*) - \mathcal{R}_n(f_n^*) \right)$$

$$+ \left( \mathcal{R}_n(f^*) - \mathcal{R}(f^*) \right)$$

The ③ term can be controlled using LLN

$$\eta \qquad (X_i) \qquad \frac{1}{n} \sum_{i=1}^{n} X_i \xrightarrow{P} \mu$$

are drawn from $\pi$

$$\mathcal{R}_n(f^*) = \frac{1}{n} \sum_{i=1}^{n} \ell(f^*(x_i), y_i)$$

$$\mathcal{R}(f^*) = \int_{\Omega} \ell(f^*(x), y) \, d\pi(x, y)$$

$$\left| R_n(f^*) - R(f^*) \right| \xrightarrow{n \to +\infty} 0 \quad \text{by LLN}$$

The difficult term is ①

$$\left( R(f_n^*) - R_n(f_n^*) \right) \le \sup_{f \in \mathcal{F}} \left| R(f) - R_n(f) \right|$$

If the loss is "well behaved" we can prove

$$\sup_{f \in \mathcal{F}} \left| R(f) - R_n(f) \right| \to 0$$

Uniform law of large numbers



error

test

training

Early stopping

Regularization

Gradient flow is a solution to the ODE

$$(*) \quad \begin{cases} w'(t) = -\nabla \phi(w(t)) \leftarrow & \phi \text{ any } \quad C^{1,1}(\mathbb{R}^N ; \mathbb{R}) \\ w(0) = w^0 \end{cases}$$

$$t \mapsto w(t) \in \mathbb{R}^N$$

We are interested in the case $\phi = \mathbb{R}_n$

Example $N=1$
$$\phi(x) = \frac{x^2}{2} \qquad \nabla \phi(x) = x$$
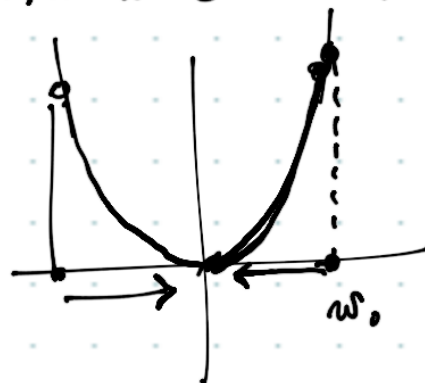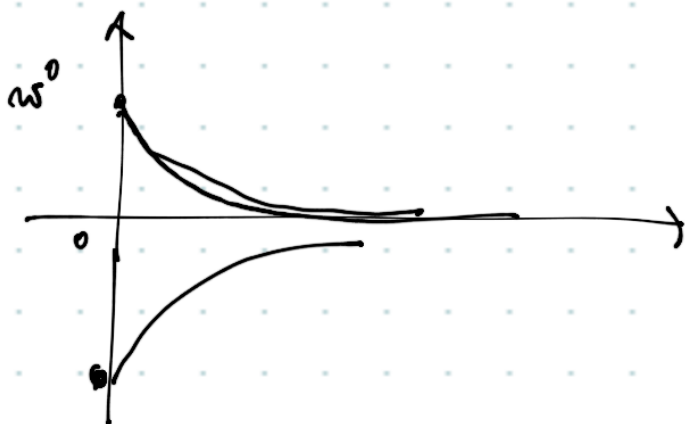$$\nabla \phi(w(t)) = w(t)$$

Problem $(*)$ becomes

$$\begin{cases} w'(t) = -w(t) \\ w(0) = w^0 \end{cases}$$

$$w(t) = w^0 e^{-t}$$
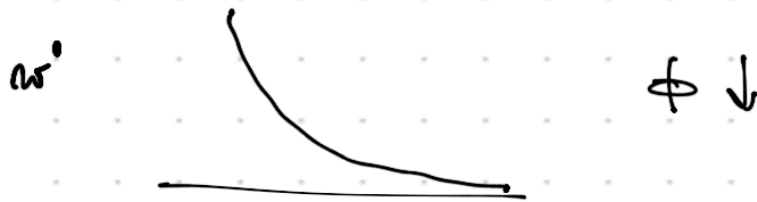$$w'(t) = -w^0 e^{-t} = -w(t)$$
$$w(0) = w^0 e^{-0} = w^0$$



"$\phi$ decreases along the gradient flow"   $\phi : \mathbb{R}^N \to \mathbb{R}$

$$\frac{d}{dt} \phi(w(t)) = \underset{\uparrow}{\nabla \phi(w(t))} \cdot w'(t) = -\underset{\uparrow}{w'(t) \cdot w'(t)}$$

Solves $(*)$     dot prod in $\mathbb{R}^N$     Because $w' = -\nabla \phi$

$$= -\| w' \|^2 \quad , \quad \frac{d}{dt} \phi(w(t)) = -\| w' \|^2 \leq 0$$

$w'$



$\phi \downarrow$

**Exercise**   Find the gradient flow when $N = 1$

$$\phi(x) = (x^2 - 1)^2 \quad i.e. \quad \text{solve}$$

$$\nabla \phi(x) = 2(x^2 - 1) \, 2x = 4x(x^2 - 1)$$

$$\begin{cases} w'(t) = -4 \, w(t) \, (|w(t)|^2 - 1) \\ w(0) = w^0 \end{cases}$$
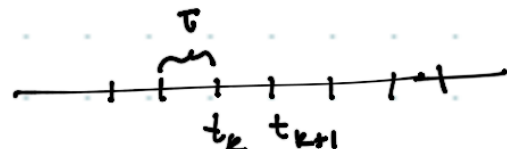


$(x^2-1)^2$

these solution reaches the minima of $\phi$ <u>asymptotically</u>

## Gradient Descent

Gradient descent can be interpreted as an Euler method to solve the gradient flow (*)

$$w'(t) = f(w(t), t) \qquad$$



$\tau$

$t_k \quad t_{k+1}$

$$w(t_k) = w^k$$

$$\frac{w^{k+1} - w^k}{\tau} = \begin{cases} f(w^k, t_k) & \text{explicit Euler method} \\ f(w^{k+1}, t_{k+1}) & \text{implicit Euler method} \end{cases}$$

If you apply this to (*)

$$\frac{w^{k+1} - w^k}{\tau} = -\nabla\phi(w^k) \qquad \text{explicit}$$

$$\frac{w^{k+1} - w^k}{\tau} = -\nabla\phi(w^{k+1}) \qquad \text{implicit}$$

Which are usually written as

Learning rate

$$w^{k+1} = w^k - \tau\,\nabla\phi(w^k) \longleftarrow$$

$$w^{k+1} = w^k - \tau\,\nabla\phi(w^{k+1}) \longleftarrow$$

A "better" way to define GD

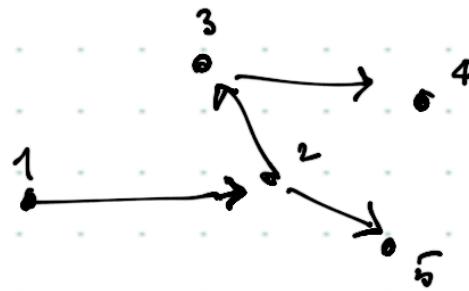(1) $w^{k+1} \in \underset{s \in \mathbb{R}^N}{\text{argmin}} \quad \phi(s) + \frac{1}{2\tau}\|s - w^k\|^2 \longleftarrow$

(2) $w^{k+1} \in \underset{s \in \mathbb{R}^N}{\text{argmin}} \quad \phi(w^k) + \nabla\phi(w^k)\cdot(s - w^k) + \frac{1}{2\tau}\|s - w^k\|^2$

<u>Exercise</u> Prove (1) is equivalent to the implicit GD

(2) " " , " explicit GD

Back prop. on a dag

$$G = (\underset{\uparrow}{V}, A)$$

set of vertices



1 → 2

is an arc of the dag.

$A = \{ 1 \to 2,$
$\qquad 2 \to 3,$
$\qquad 3 \to 4,$
$\qquad \vdots$
$\qquad \}$

    These are not Dags

given $i \in V$         $pa(i) = \{ j \in V : j \to i \in A \}$



$j \in pa(i)$

$$ch(i) = \{ j \in V : i \to j \in A \}.$$

$j \in ch(i)$



What is a FNN? In general is a dag + 

$i \mapsto x^i$

$j \to i \to w_{ij}$



$$x^3 = \sigma(w_{31} \, x_1)$$

With the following computational rule

$$z^i = \begin{cases} \sigma\left( \underbrace{\sum_{j \in pa(i)} w_{ij} z^j}_{a^i} \right) & \text{if } i \in H \cup O \\ y^i & \text{if } i \in I \end{cases}$$
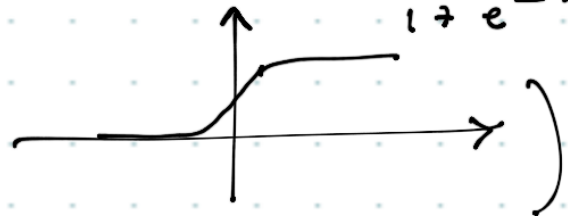
$$I = \{ i \in V : pa(i) = \emptyset \}$$

$$O = \{ i \in V : ch(i) = \emptyset \}$$

$$H = V \setminus \{ O \cup I \}.$$

$\sigma \to$ activation function $\left( \text{for instance } \sigma(z) = \dfrac{1}{1 + e^{-z}} \right.$



$z$ is the input to
the network

With this computational rules a network describes a
function

$$z \to f(z, w) \quad \left( \text{which are the values} \atop \text{of } z^i \text{ on } O \right)$$

What is Backprop?
It is a clever (optimal) way to compute

$$\nabla R_n(w)$$

which is the key ingredient for GD on $\mathbb{R}_n$

when $\quad R_n(w) = \frac{1}{n} \sum\limits_{i \geq 1}^{n} \ell\big(f(z_i, w), y_i\big)$

$\qquad\qquad\qquad\qquad\qquad\qquad\uparrow$

$\qquad\qquad\qquad\qquad\qquad$ this is a NN

$\dfrac{\partial}{\partial w_{ij}} \ell = \dfrac{\partial \ell}{\partial a^i} \boxed{\dfrac{\partial a^i}{\partial w_{ij}}} \qquad a^i = \sum\limits_{j \in pa(i)} w_{ij} z^j$

$\dfrac{\partial a^i}{\partial w_{ij}} = \dfrac{\partial}{\partial w_{ij}} \sum\limits_{k \in pa(i)} w_{ik} z^k$

$\qquad\qquad = \sum\limits_{k \in pa(i)} \delta_{kj} z^k = z^j$

$\qquad\qquad\qquad\qquad\qquad\uparrow$ Kronecker $\delta$

$\dfrac{\partial \ell}{\partial a^i} = \delta_i \qquad \forall i \in H \cup O$

$\dfrac{\partial \ell}{\partial a^i} = \sum\limits_{k \in ch(i)} \dfrac{\partial \ell}{\partial a^k} \dfrac{\partial a^k}{\partial a^i}$

$\rule{8cm}{0.4pt}$

$\qquad\qquad = \sum\limits_{k \in ch(i)} \delta_k \dfrac{\partial a^k}{\partial a^i}$

$a^k = \sum\limits_{j \in pa(k)} w_{kj} z^j = \sum\limits_{j \in pa(k)} w_{kj} \, \sigma(a^j)$

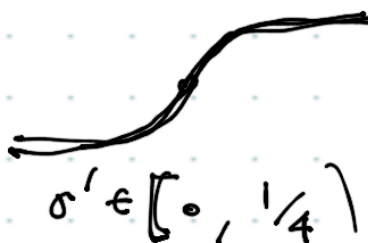$\dfrac{\partial a^k}{\partial a^i} = w_{ki} \, \sigma'(a^i)$

$\boxed{\delta_i := \sigma'(a_i) \cdot \sum\limits_{k \in ch(i)} \delta_k \, w_{ki}}$

$\dfrac{\gamma \ell}{\gamma w}$

$\sigma' \qquad \sigma'$

$\sigma' \in \left[0, \tfrac{1}{4}\right)$

$\left(\tfrac{1}{4}\right)^{10}$

ReLu
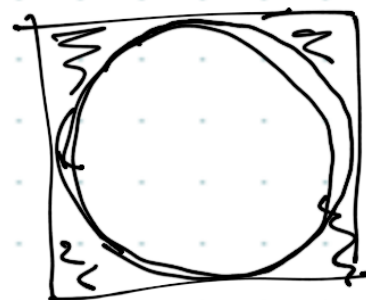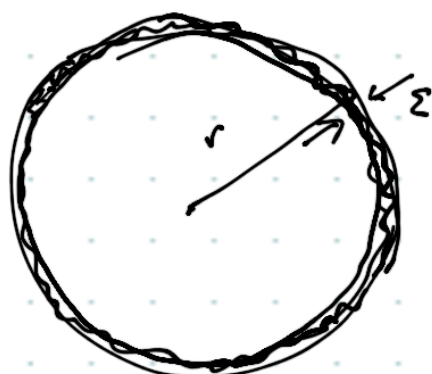
Vanishing gradient problem

---

Exercise



$\dfrac{V_S}{V_C} \xrightarrow{n \to +\infty} 0$

All the volume lives on the surface as $n \to +\infty$

$$\overbrace{\dfrac{\text{vol}(B_r) - \text{vol}(B_{r-\varepsilon})}{\text{vol}(B_r)}}$$

$$\text{vol}(B_r) = \underset{\substack{\uparrow \\ \text{volume of the unit Ball in } \mathbb{R}^n}}{\omega_n \, r^n}$$

$0 < \dfrac{\varepsilon}{r} < 1$

$$\dfrac{\omega_n r^n - \omega_n (r-\varepsilon)^n}{\omega_n r^n} = 1 - \overbrace{\left(1 - \dfrac{\varepsilon}{r}\right)^n} \to \begin{array}{l} 1 - 0 \\ = 1 \end{array}$$