Coding session will be exciting. We'll see how to get our machines to recognise _handwritten digits_!

$$\mathcal{X} = \text{array of "pixels"} \approx [0,1]^{m \times n} \approx \left\{ \frac{k}{N} : 0 \leq k \leq N \right\}^{n \times m}$$

$$\mathcal{Y} = \{0, \ldots, 9\}$$

$D = $ classified images of handwritten digits

_Recall_: We have data points $(x_i, y_i)_{i=1}^n = D$. We <u>assume</u> $y_i = f^*(x_i)$ (+ noise) for some $f: \mathcal{X} \to \mathcal{Y}$.

Our goal is to determine $f^*$.

① Consider a pool of candidate functions $\mathcal{F}$

Ex.: $\mathcal{X} = \mathcal{Y} = \mathbb{R}$, $\left\{ x \mapsto \beta_1 x + \beta_2 x^2 + \beta_{100} x^{100} : \beta_1, \beta_2, \beta_{100} \in \mathbb{R} \right\}$

② Use $D$ and a loss function $\ell(\cdot, \cdot)$ to define the empirical risk $\widehat{R}_D : \mathcal{F} \to [0, \infty]$

$$\widehat{R}_D(f) = \frac{1}{n} \sum_{i=1}^n \ell(y_i, f(x_i)). \qquad \text{Ex: } \ell(y, y') = \frac{1}{2} \| y - y' \|^2 \quad (\text{quadratic loss})$$

③ Use optimisation methods to find (an approximate) optimiser $\widehat{f}$ of $\widehat{R}_D$

Ex.: use GD, SGD, SGD with minibatching, ...

( Should be clear to everyone what ①, ②, ③ are in the case of linear regression )

A neural network (NN) is just <u>a specific choice of $\mathcal{F}$.</u>

- NN: originally introduced to model biological neurons (→ "perceptron") and their interaction, hence the name

- Allow for efficient implementation of (a version of) GD (via "backpropagation")

- Perform well in practice

We discuss the most basic case, fully-connected feedforward NN.

Let $\sigma : \mathbb{R} \to \mathbb{R}$ be a continuous function. Ex: $\sigma(x) = x^+$ (ReLU), $\sigma(x) = \dfrac{e^x}{1+e^x}$, ...

Notation: $x \in \mathbb{R}^n$, $\sigma(x) := (\sigma(x_1), \ldots, \sigma(x_n)) \in \mathbb{R}^n$.

$k \in \mathbb{N}$, $n_0 \quad n_k \in \mathbb{N}$, $W_i \in \mathbb{R}^{n_i \times n_{i-1}}$, $b_i \in \mathbb{R}^{n_i}$, $1 \leq i \leq k$,

$$f(x) = W_k \left( \ldots \sigma \left( W_2 \sigma \left( \underbrace{W_1 x + b_1}_{\in \mathbb{R}^{n_1}} \right) + b_2 \right) \ldots \right) + b_k, \quad x \in \mathbb{R}^{n_0}.$$

$W_1 \quad W_k \quad$ 'weights'

$b_1 \quad b_k \quad$ 'biases'

$\sigma \quad\quad$ 'activation fun.'

$\max\limits_{0 \leq i \leq k} n_i = $ 'width', $\quad k = $ 'depth', $\quad k - 2 = \#$ hidden layers

$$f : \mathbb{R}^{n_0} \xrightarrow{(W_1, b_1)} \mathbb{R}^{n_1} \xrightarrow{\sigma} \mathbb{R}^{n_1} \xrightarrow{(W_2, b_2)} \mathbb{R}^{n_2} \xrightarrow{\sigma} \mathbb{R}^{n_2} \ldots \xrightarrow{\sigma} \mathbb{R}^{n_{k-1}} \xrightarrow{(W_k, b_k)} \mathbb{R}^{n_k}$$
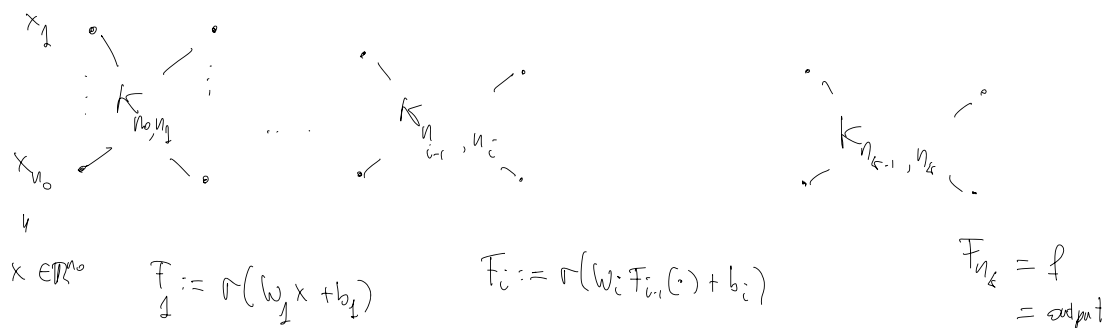
Sometimes, $b_k = 0$ is imposed.

Ex.: $x \in \mathbb{R}^{n_0}$, $x = $ handwritten digit, $f(x) = (f_0(x), \ldots, f_9(x))$,

$f_i(x) = \sim$ 'probability that $x$ represents the digit $i$' (after normalisation, eg with softmax)

'Pictorial' representation of $f$: $\qquad$ ($K_{s,t} = $ complete $(s,t)$-bipartite graph, $s,t \in \mathbb{N}$)



$x \in \mathbb{R}^{n_0}$ $\qquad$ $F_1 := \sigma(W_1 x + b_1)$ $\qquad$ $F_i := \sigma(W_i F_{i-1}(\cdot) + b_i)$ $\qquad$ $F_{n_k} = f$ = output

(There are other ways to connect neurons, see eg. convolutional NN.)

In practice, one fixes $k \in \mathbb{N}$, $n_0 \ldots n_k \in \mathbb{N}$ (and $\sigma$) ('hyperparameters') at the beginning, and then optimises over $W_i, b_i$ ('training').

$\mathcal{X} \subseteq \mathbb{R}^n, \ \mathcal{Y} \subseteq \mathbb{R}^m.$

We take $n_0 = n$ and $n_L = m$.

$$\mathcal{F}(n_1 \ \dots \ n_{L-1}; \sigma) := \text{functions as above}$$

$$\mathcal{F}(\sigma) := \bigcup_{L \in \mathbb{N}_{\geq 2}} \ \bigcup_{n_1 \dots n_{L-1} \in \mathbb{N}} \mathcal{F}(n_1 \ \dots \ n_{L-1}; \sigma).$$

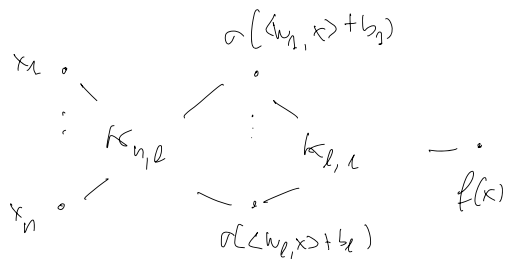Recall steps ①, ②, ③. Each one introduces an error :

① no approximation error, ② no estimation error, ③ no optimisation error.

We will not discuss ③ and only focus on $L = 2, m = n_2 = 1, b_2 = 0,$ and arbitrary $n_1 \in \mathbb{N}$, i.e.

functions of the form

$$f(x) = \sum_{j=1}^{\ell} \eta_j \ \underbrace{\sigma\left( \langle w_j, x \rangle + b_j \right)}_{\text{scalar product}}, \qquad x \in \mathbb{R}^n,$$

where $\ell \in \mathbb{N}$ and $\eta_j \in \mathbb{R}, \ (w_j, b_j) \in \mathbb{R}^n \times \mathbb{R}$ for $1 \leq j \leq \ell.$



Gather all such functions in the set $\mathcal{G}(n; \sigma), \ \mathcal{G}(n, \sigma) = $ single hidden layer of arbitrary width.

(1) Which functions can we approximate with elements of $\mathcal{G}(n,\sigma)$?

$\sigma \in C(\mathbb{R}) \Rightarrow \mathcal{G}(n,\sigma) \subseteq C(\mathbb{R}^n)$ ~~ Is $\mathcal{G}(n,\sigma)$ dense?

$C(\mathbb{R}^n)$ endowed with the (metrisable) top of unif. conv on cpt sets, ie.

$$f_n, f \in C(\mathbb{R}^n), \ f_n \to f \iff \forall K \subseteq \mathbb{R}^n \text{ cpt}, \quad \max_{x \in K} \left| f_n(x) - f(x) \right| \underset{n \to \infty}{\to} 0.$$

Thm (Pinkus, '99) Let $\sigma \in C(\mathbb{R})$. TFAE

- $\mathcal{G}(1,\sigma)$ is dense in $C(\mathbb{R})$
- $\mathcal{G}(n,\sigma)$ is dense in $C(\mathbb{R}^n)$ $\forall n$
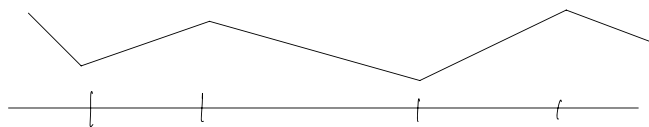- $\sigma$ is not a polynomial

In part., $\sigma = \text{ReLU}$ or $\sigma = $ sigmoid are okey.

Problem: No information on $\ell \in \mathbb{N}$ nor the size of $(w_j, b_j)_{j=1}^{\ell}$.

Could be useless in practice: We may be unable to construct a NN of the required size or to find the correct weights & biases.

It is possible to obtain more 'quantitative' approx. results, eg. with some smoothness assumption on the fun we want to approx. More technical. See eg. Weinon E.

<u>Ex.1</u>  $\sigma = ReLU$, $n=1$, $f : \mathbb{R} \to \mathbb{R}$ piecewise affine.



② Estimation error.

<u>Prop</u> (9.1) : Let $\sigma = ReLU$, $D, R > 0$ and

$$\widetilde{\mathcal{G}} = \left\{ f \in \mathcal{G}(n, \sigma) : \|m\|_1 \leq D, \|w_j\|_2^2 + \left(\frac{b_j}{R}\right)^2 = 1 \ \forall j \right\}.$$

If the loss fun. $\ell(\cdot, \cdot)$ is $G$-LIP in the second variable, then

$$\mathbb{E}\left[ \hat{R}_N(\hat{f}) - \inf_{f \in \widetilde{\mathcal{G}}} R(f) \right] \leq \frac{16 \, GDR}{\sqrt{N}} \qquad (1),$$

$N =$ sample size.

In particular, the RHS of (1) does <u>not</u> depend on the number of parameters, but only on their <u>size</u>.