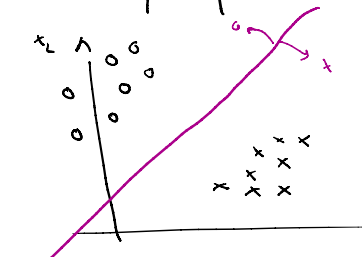


III Supervised Learning

A) The perception



separation → you want to find this separation automatically

I want to find \leftarrow
 get such that

$$f((x_1, x_2) \in x) \rightarrow x$$

$$f((x_1, x_2) \in o) \rightarrow o$$

in D dimensions, how do you define a line?

$$\rightarrow w_0 + \sum_{i=1}^D w_i x_i = 0 \quad \text{eq of a line}$$

\circledast
coordinates

parameters of the line: $\{w_0, \vec{w}\}$ w_0 : the bias
 (weights)

Let's define the model:
$$y = \text{sgn}\left(w_0 + \sum_{i=1}^D w_i x_i\right)$$

\downarrow \downarrow
 output input

We have to discretize the labels of the dataset

x: +1

o: -1

I want to find $\{w_0, \vec{w}\}$ s.t. $\forall i$ in my dataset

$$\tilde{y}_i = f(w_0 + \sum w_i x_i)$$

↳ label of data i

Perception algo: for each i , change the

weight \therefore if i is not correctly classified

$$\vec{w} \leftarrow \vec{w} + \tilde{y}_i \vec{x}_i$$

$$w_0 \leftarrow w_0 + \tilde{y}_i$$

you can show that this algo reaches a solution

\nexists : the dataset is linearly ~~separable~~ separable

B) Multiclass Perception

Notation: $\vec{x} = \begin{pmatrix} 1 \\ \vec{x} \end{pmatrix} \in \mathbb{R}^{D+1}$, $\vec{w} = \begin{pmatrix} w_0 \\ \vec{w} \end{pmatrix}$

pb with binary perception:

- sign jet cannot be derive
- no information on the probab of a data to belong to one cluster.

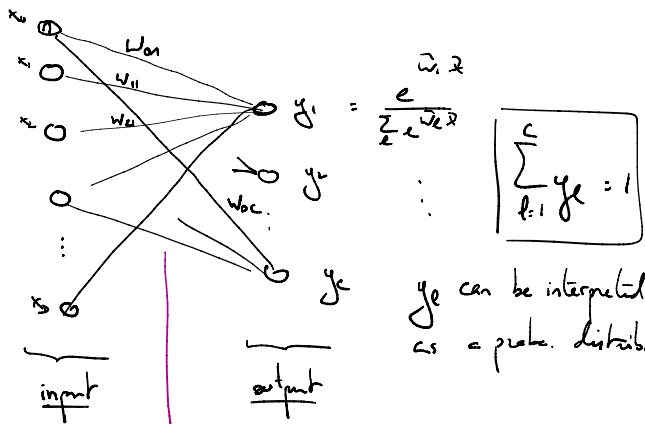
We will have C classes.

1) Let's define the model: Instead of 1 output I will have C of them

output n^k :
$$y_k = \frac{\exp(\sum_i w_{ki} x_i)}{\sum_{k=1}^C \exp(\sum_i w_{ki} x_i)} = \frac{\exp(\vec{w}_k \cdot \vec{x})}{\sum_{k=1}^C \exp(\vec{w}_k \cdot \vec{x})}$$

for each output, I have an associated set of weights \vec{w}_k
 $\vec{w}_k \in \mathbb{R}^{D+1}$

2) graphical representation



The edges represent the weights

How classes are separated?

$$\frac{y_h}{y_l} = \frac{e^{\vec{w}_h \cdot \vec{x}}}{e^{\vec{w}_l \cdot \vec{x}}} = \exp\left(\underbrace{(\vec{w}_h - \vec{w}_l) \cdot \vec{x}}\right)$$

it defines an hyperplan of dim $D-1$ between class $h \neq l$.

□ convexity property: let's have two pts $A \neq B$
in class h according to my model.

$$\Rightarrow \bar{w}_h \cdot \bar{x}_A > \bar{w}_j \cdot \bar{x}_A, \forall j \neq h$$

$$\bar{w}_h \cdot \bar{x}_B > \bar{w}_j \cdot \bar{x}_B, \text{ ---}$$

hence $\forall t$ defined as $\bar{x} = \lambda \bar{x}_A + (1-\lambda) \bar{x}_B, \lambda \in [0,1]$

$$\text{we have: } \bar{w}_h \cdot \bar{x} = \lambda \bar{w}_h \cdot \bar{x}_A + (1-\lambda) \bar{w}_h \cdot \bar{x}_B$$

$$> \lambda \bar{w}_j \cdot \bar{x}_A + (1-\lambda) \bar{w}_j \cdot \bar{x}_B, \forall j \neq h$$

ii) Learning of the model

we have to define a $\left\{ \begin{array}{l} \text{loss} \\ \text{objective} \end{array} \right.$ function. and $\left\{ \begin{array}{l} \text{minimize} \\ \text{maximize} \end{array} \right.$ it

Since we have an interpretation in terms of probab. we can define a likelihood.

we said: $\text{probab}(\bar{x} \text{ is in class } h) = y_h(\bar{x})$

\rightarrow we want to maximize the probab that \bar{x} is in the correct class

Dataset: $(\bar{x}_d, \tilde{y}_d), d=1, \dots, D$

\hookrightarrow label of state d

we will transform it into a one-hot vector.

instead of having $\tilde{y}_d = z \rightarrow \tilde{y}_d = \begin{pmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix} = \begin{cases} 1 & \text{on coord } z \\ 0 & \text{elsewhere} \end{cases}$

For data (\bar{x}, \tilde{y}) I define the likelihood: $p_c(\bar{x} \in \text{class } \tilde{y})$

$$= \prod_{h=1}^c (y_h(\bar{x}))^{\tilde{y}_h}$$

\rightarrow log-likelihood: $\mathcal{L}(\bar{x}) = \sum_h \tilde{y}_h \log(y_h(\bar{x}))$ (cross-entropy)

we will make a gradient ascent of \mathcal{L}

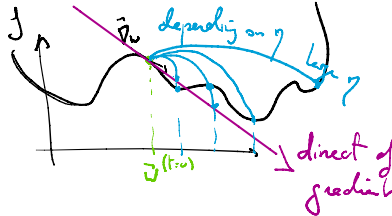
we will make a gradient descent of L

Compute: $\frac{\partial L}{\partial w_i} = \kappa_i (\tilde{y}_e - y_e(\tilde{x}_i))$

strongly suggest to derive it

iteration rule: $w_i \leftarrow w_i - \eta \cdot \kappa_i (\tilde{y}_e - y_e(\tilde{x}_i))$

In practice: η ? the learning rate



• Mini-batch: The dataset is partitioned into small groups of mb data

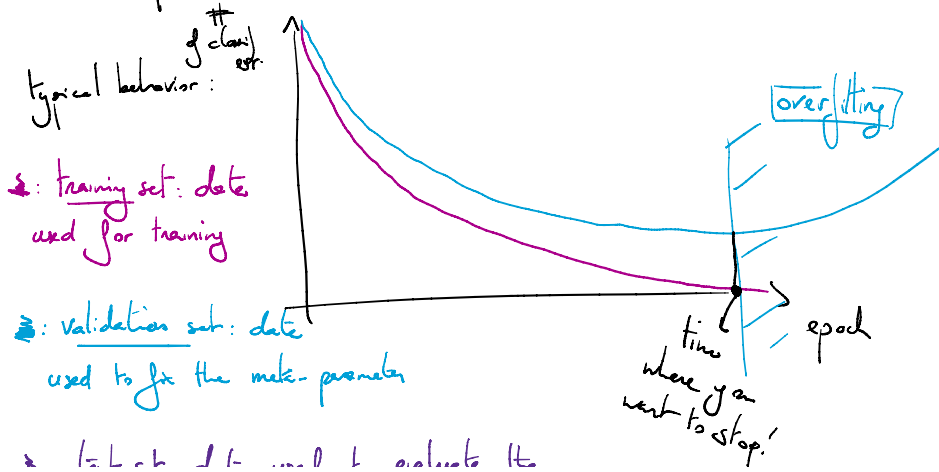
- $\frac{D}{mb}$ groups / mini-batches
- each group/mini-batch contains mb data points
- a data-pt belong to only one group/mini-batch

typically $mb \in [8, 200]$

the update becomes: $w_i \leftarrow w_i - \eta \frac{1}{mb} \sum_{d \in \text{batch}(b)} \kappa_i (\tilde{y}_e - y_e(\tilde{x}_i))$

• An epoch: it is when all the data of the dataset have been used one time (only one) to update the parameters

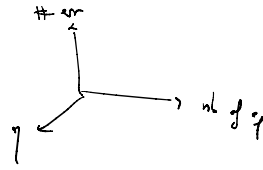
• Meta-parameters & evaluation



used to fix the meta-parameters

• test set: data used to evaluate the performance of your network

more you want to stop!



- Other meta-parameters:
- nb of epoch
 - learning rate η
 - L_2 regularization λ
 - L_1 _____ λ_1

L_2 -reg: you add $-\frac{\lambda}{2} \sum w_{ii}^2$ to the log-likelihood
 \sim gaussian prior on the weights

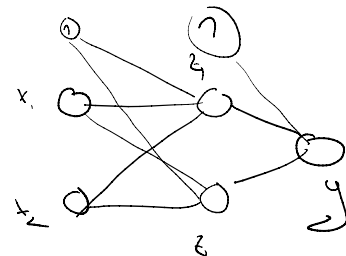
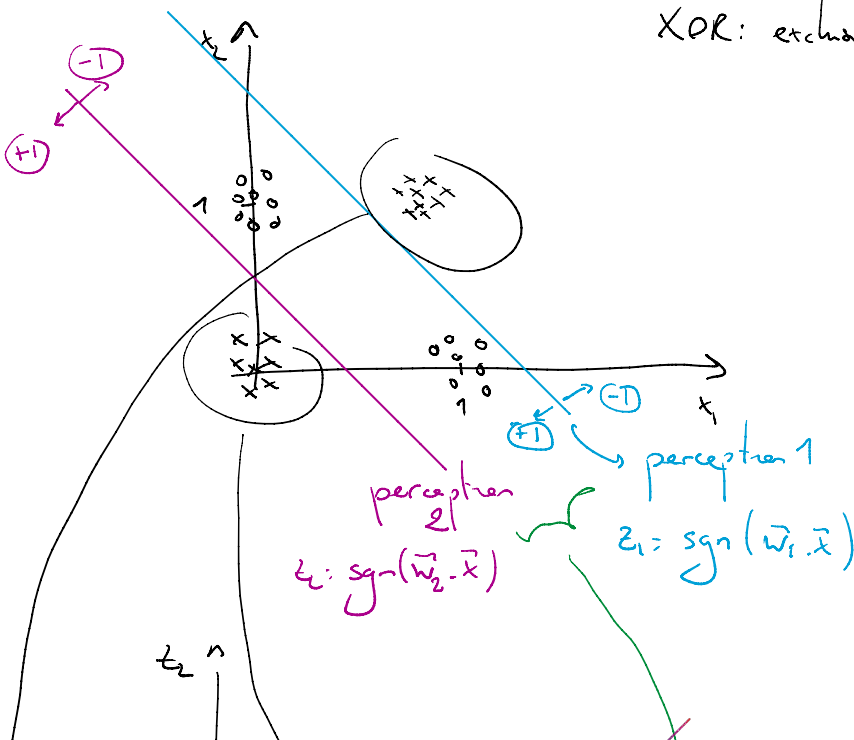
L_1 -reg: $-\lambda \sum |w_{ii}|$ to the log-l.
 \sim Laplace prior on the weights

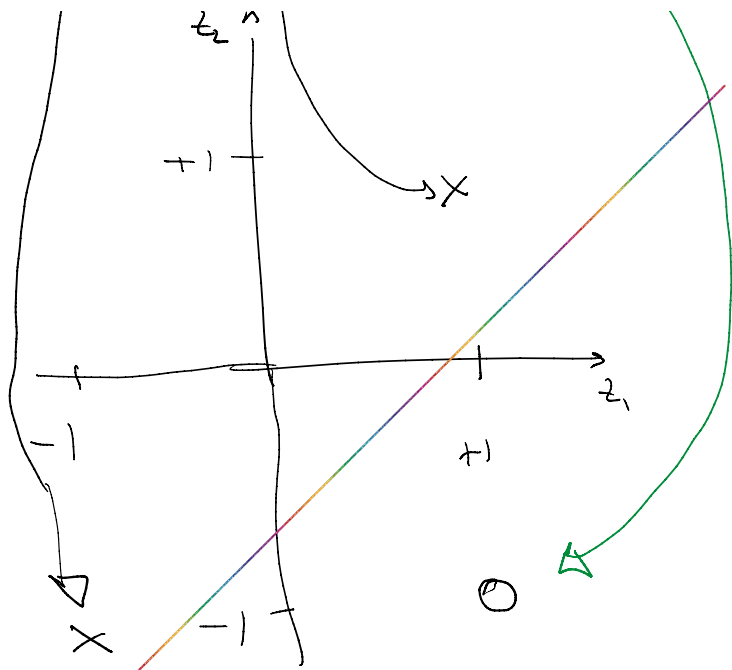
C] Multilayer Perceptron

1) Why multilayer is good?

XOR: exclusive OR

a	b	a XOR b
0	0	0
0	1	1
1	0	1
1	1	0

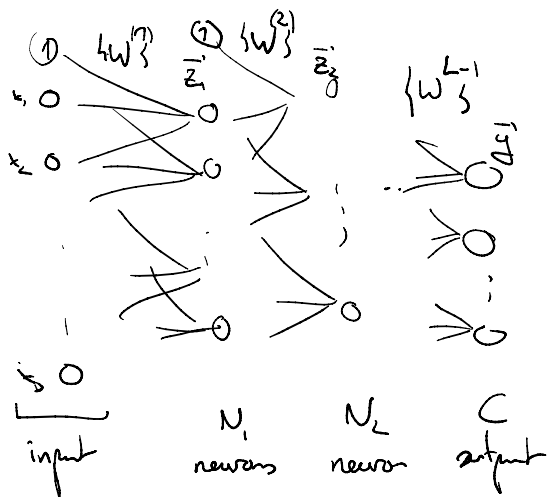




perceptron $y = \text{Sgn} \left(z_1 w'_1 + z_2 w'_2 + w'_0 \right)$

Multi-layer perceptron is a gd model for complex dataset.

ii) Definition:



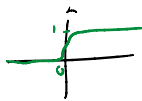
$$z_i^{(1)} = f_1 \left(\sum_{i=1}^{N_1} \vec{w}_i^{(1)} \cdot \vec{x} \right), \quad f_i: \text{an activation function}$$

$$z_i^{(2)} = f_2 \left(\sum_{i=1}^{N_2} \vec{w}_i^{(2)} \cdot \vec{z}^{(1)} \right), \quad i=1, \dots, N_2$$


$$\vdots$$

$$\dots \exp \left(\frac{-z^{(L-1)} + z^{(L-2)}}{w_h \cdot z} \right)$$

$$y_k = \frac{\exp(\vec{w}_k^{(L-1)} \cdot \vec{z}^{(L-2)})}{\sum_l \exp(\vec{w}_l^{(L-1)} \cdot \vec{z}^{(L-2)})}$$

Typical activation fct: - sigmoid: $\sigma(x) = \frac{1}{1+e^{-x}}$ 

- ReLU: rectified linear unit

$$\text{relu}(x) = \max(0, x)$$


possible output: for class: softmax $\frac{e^{\vec{w}_k \cdot \vec{z}}}{\sum_l e^{\vec{w}_l \cdot \vec{z}}}$

• regression: linear act (identity)

• regression in $[0,1]$: sigmoid

iii) Learning: maximizing the likelihood

$p(\vec{x} \in \text{class } y)$

$$\mathcal{L}(w) = \sum_k \tilde{y}_k \log(y_k(\vec{x}))$$

$$\frac{\partial \mathcal{L}}{\partial w^{(L)}}; \frac{\partial \mathcal{L}}{\partial w^{(L-1)}}; \dots; \frac{\partial \mathcal{L}}{\partial w^{(L-1)}}$$

$\vec{z}_i^{(L-2)} (\tilde{y}_i - y_i(\vec{x}))$

chain rule

