

3. take each sample in sequence and compute its centroid from centroid.

If sample is not in the current centroid switch and update the centroid of two cluster

4. repeat all the steps until convergence

II.

termination condition

1. for each switch the sum of conditions distance as each training samples decrease to centroid decrease
2. the are only many finite partitions of the training examples into K cluster

ML EXAMS 12/02/2018

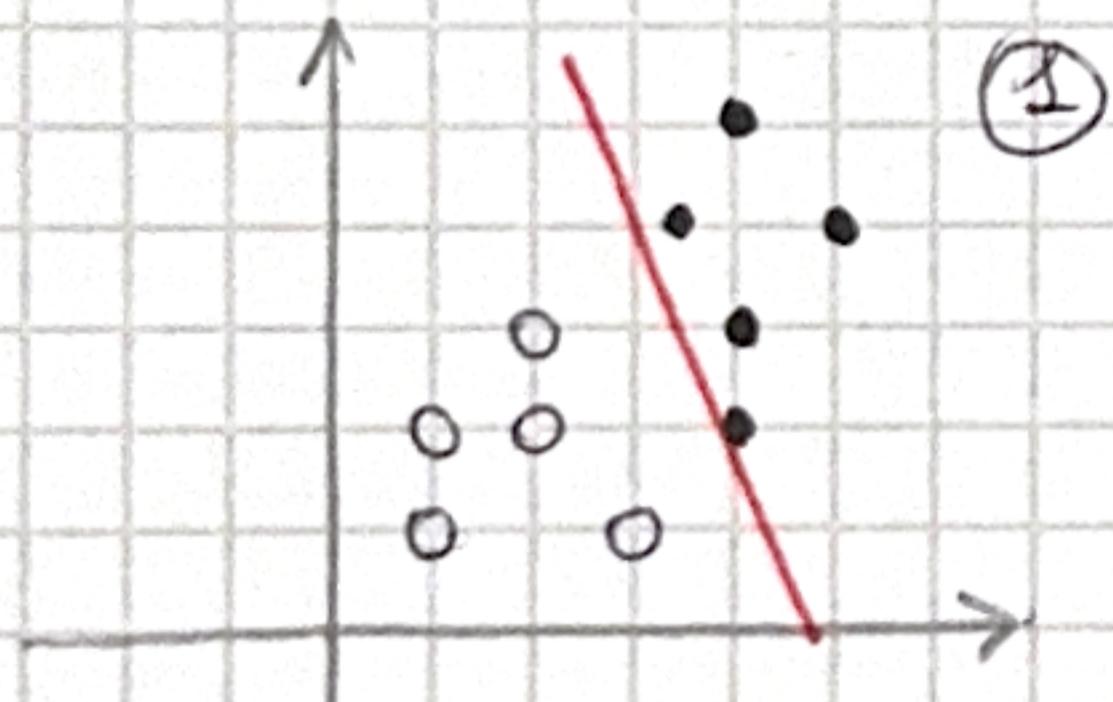
EXERCISE 1

supervised learning are machine problem where the model that will solve it, are trained with a dataset composed by x and y

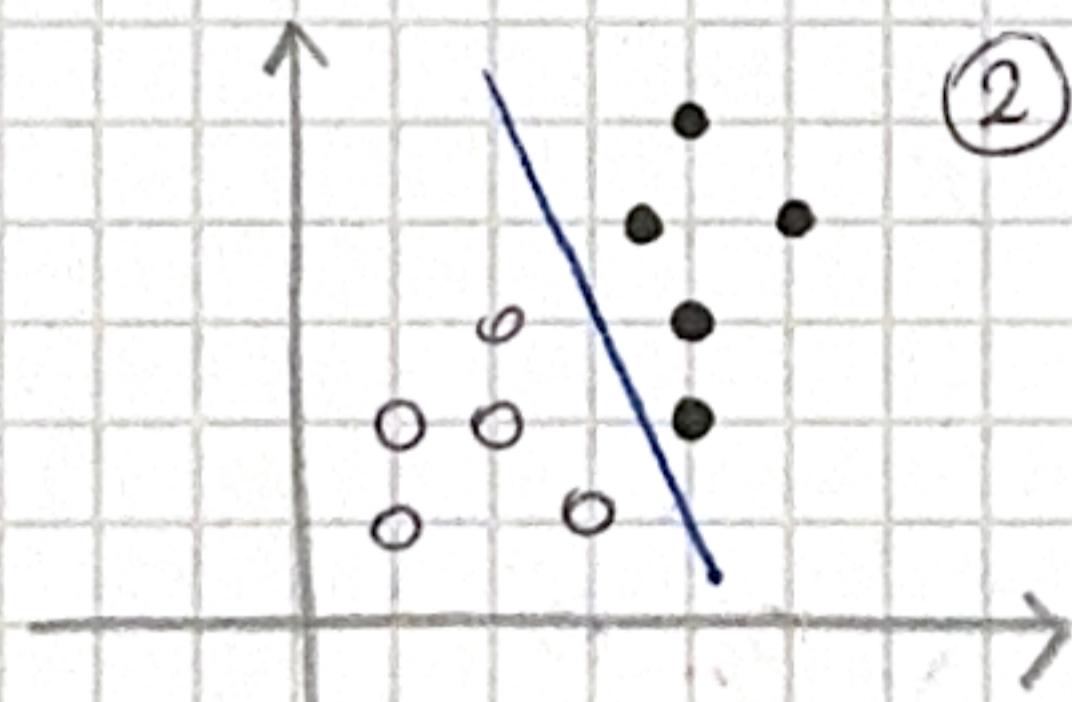
$$D = \{(x_i, y_i)\}_{i=1}^N$$

different in the case of unsupervised learning where the model used for solve them task are trained with dataset $D = \{x_i\}_{i=1}^N$ the label are NOT used in this case

EXERCISE 2



PERCEPTRON



SVM

the difference between perceptrone and SVM is perceptron is sequential algorithm based that depend on learning rate η while SVM try to maximize the margin so in this case due the fact data is linearly separable I will use SVM

EXERCISE 3

1. if $C = c_1$ and $B = b_1$ then NO ①
- if $C = c_1$ and $B = b_2$ then YES ②
- if $C = c_2$ and $A = a_1$ then YES ③
- if $C = c_2$ and $A = a_2$ and $B = b_1$ YES ④
- if $C = c_2$ and $A = a_2$ and $B = b_2$ NO ⑤
- if $C = c_2$ and $A = a_3$ then NO ⑥
- if $C = c_3$ then NO ⑦

2.

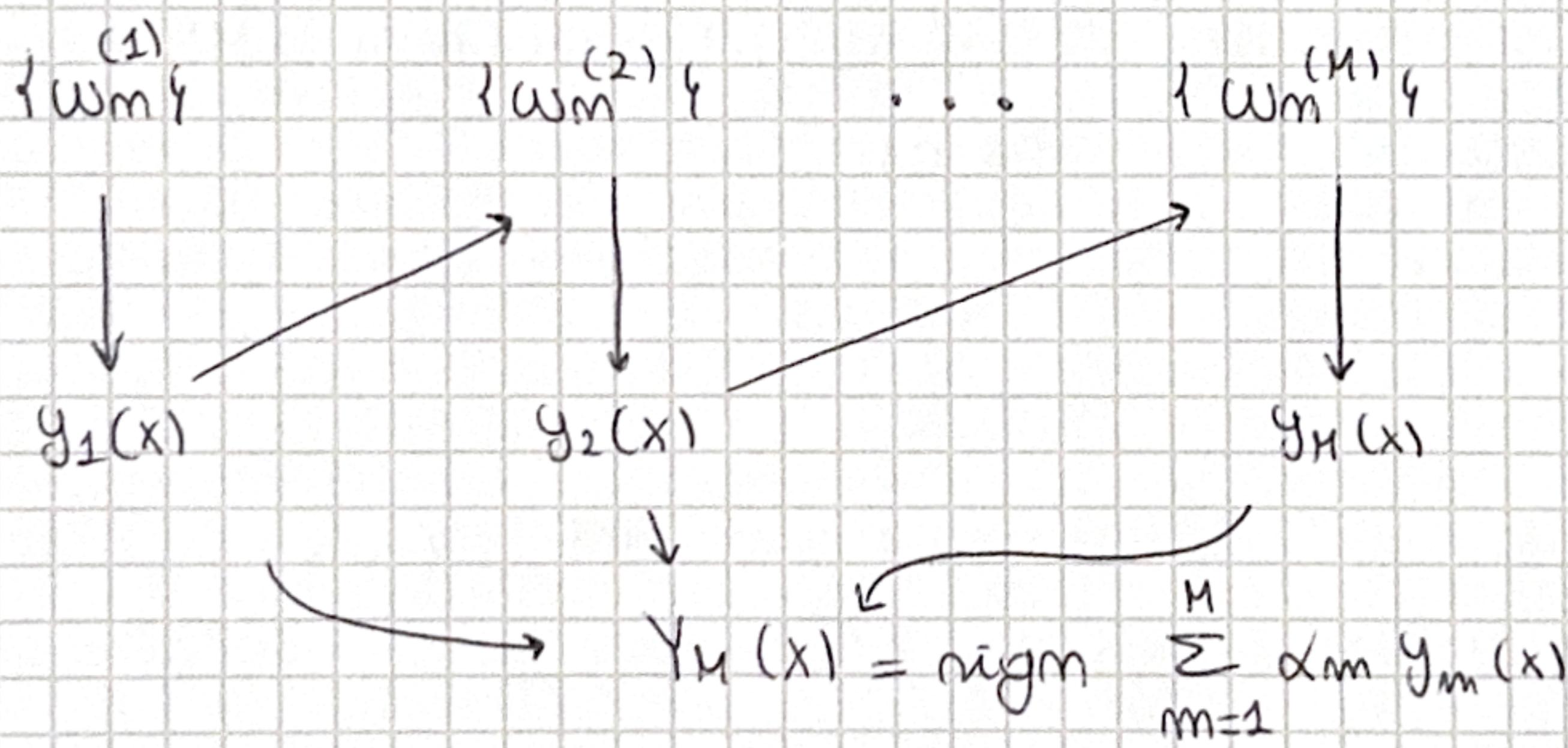
this is consistent with S_1 because of ①

- " " " " " " " " of ④
- " " " " " " " " // ~~②~~
- " " " " " " " " // ②

EXERCISE 4

Boosting is an ensemble method where instead of training one single complex classifier we train different classifiers and we combine their results.

Base classifiers are trained in sequence using a weighted data set where weights are based on performance of previous classifiers



2.

$$J_m = \sum_{n=1}^N w_n^{(m)} I(y_m(x_n) \neq t_n)$$

$$\alpha_m = \frac{1}{2} \ln \left[\frac{1 - e_m}{e_m} \right]$$

~~error function:~~

$$\delta_m = \sum_{n=1}^N w_n^{(m)} I(y_m(x_n) \neq t_n) \quad \begin{cases} 1 & \text{if } e = \text{true} \\ 0 & \text{otherwise} \end{cases}$$

$$E = \sum_{m=1}^M \exp [-t_m f_M(x_m)]$$

$$\text{where } f_M(x) = \frac{1}{2} \sum_{m=1}^M \alpha_m y_m(x) \quad t_n \in [-1, +1]$$

instead of minimize E regressively globally we can minimize it regressively:

$$\begin{aligned} E &= \sum_{m=1}^N e^{-t_m f_{M-1}(x_m)} = \\ &= \sum_{m=1}^N w_m^{(M)} e^{-\frac{1}{2} t_m \alpha_m y_m(x_m)} \text{ where } w_m^{(M)} = e^{-t_m f_{M-1}(x_m)} \\ \text{Note } y_m(x) &= \alpha_m \left(\sum_{m=1}^M \alpha_m y_m(x) \right) \end{aligned}$$

EXERCISE 5

1. this is regression task we can use a feed forward network $f(x, \theta) = f^{(3)}(f^{(2)}(f^{(1)}(x, \theta^{(1)}); \theta^{(2)}); \theta^{(3)})$

2. for hidden units we can use ReLU function for the output we can use identity activation function $y = w^T h + b$

3. the error function that we can use is Mean Square Error (MSE)

$$E(\theta) = \frac{1}{N} \sum_{n=1}^N (t_n - f(x, \theta))^2$$

EXERCISE 6

top

1. the GRAM MATRIX is $\mathbf{V} = \mathbf{X}\mathbf{X}^T$. if we have a model $y(x, \alpha) = \sum_{m=2}^N \alpha_m x_n^T x_m$ with a kernel function $k(x, x') = x^T x'$ the gram will be

$$K = \begin{bmatrix} k(x_1, x_1) & k(x_1, x_M) \\ \vdots & \vdots \\ k(x_N, x_1) & \cdots & k(x_N, x_M) \end{bmatrix}$$

2. the error function becomes $E(w) = \sum_{m=1}^N E(y_m, t_m) + \lambda \|w\|^2$

applying kernel trick $y(x, w^*) = \sum_{n=1}^N \alpha_n k(x_n, x)$ with $\alpha = (K + \lambda I_N)^{-1}$