

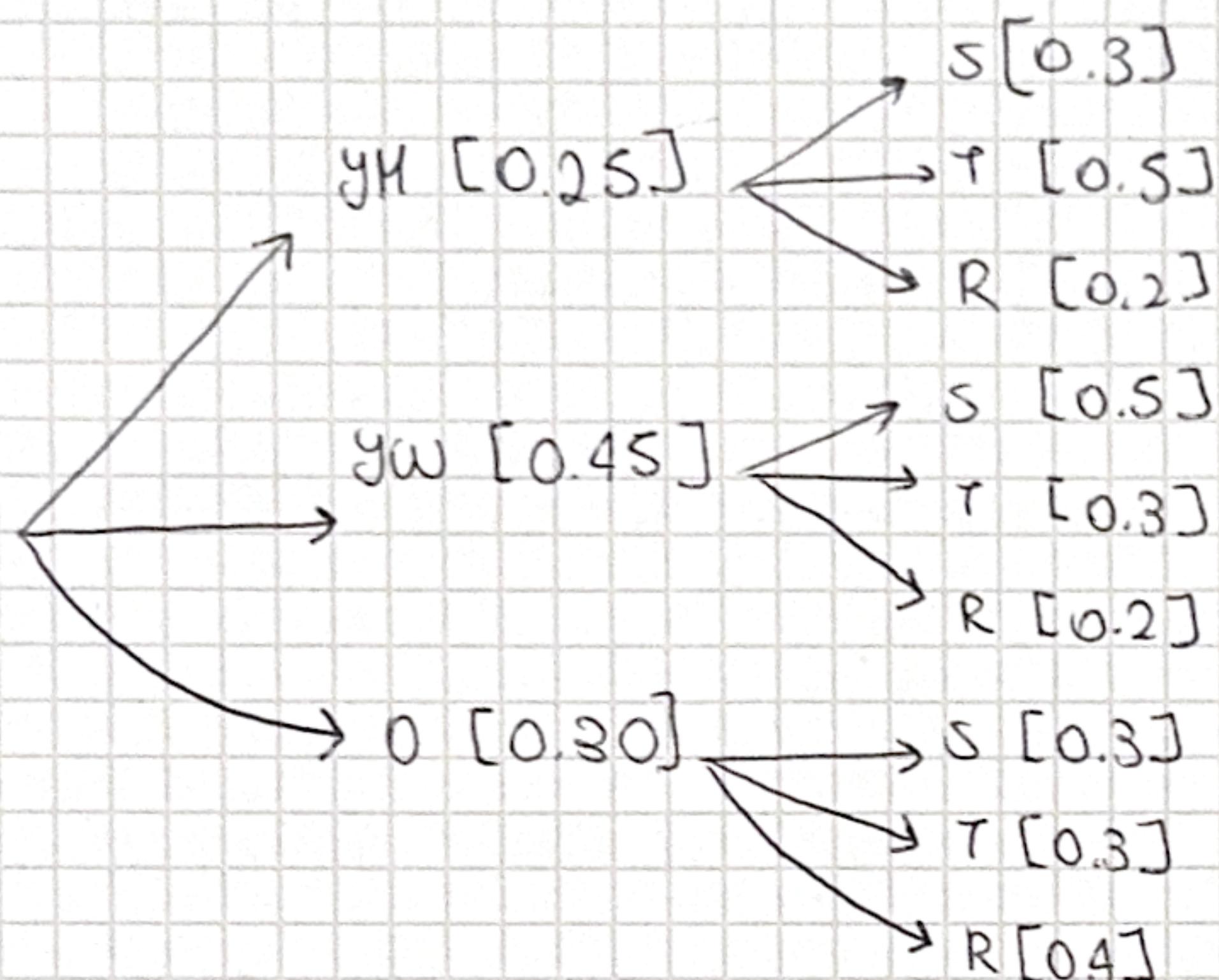
## EX. 1

$$P(Y_M) = 0.25 \quad P(Y_W) = 0.45 \quad P(O) = 0.30$$

$$P(Y_M|S) = 0.3 \quad P(Y_M|\tau) = 0.5 \quad P(Y_M|R) = 0.2$$

$$P(Y_W|S) = 0.5 \quad P(Y_W|\tau) = 0.3 \quad P(Y_W|R) = 0.2$$

$$P(O|S) = 0.3 \quad P(O|\tau) = 0.3 \quad P(O|R) = 0.4$$



$$V^* = \operatorname{argmax}_{c \in C} \{ p(\tau|c) p(c) \} = \underbrace{\{ (0.25 \cdot 0.5), (0.3 \cdot 0.45), (0.3 \cdot 0.3) \}}_{0.125} \leftarrow \operatorname{argmax}_c$$

$$= \operatorname{argmax} \{ 0.125, 0.09, 0.135 \} = Y_W$$

3. IN MAXIMUM POSTERIOR  $\forall h_{MAP} = \operatorname{argmax}_{h \in H} (p(O, h) p(h))$

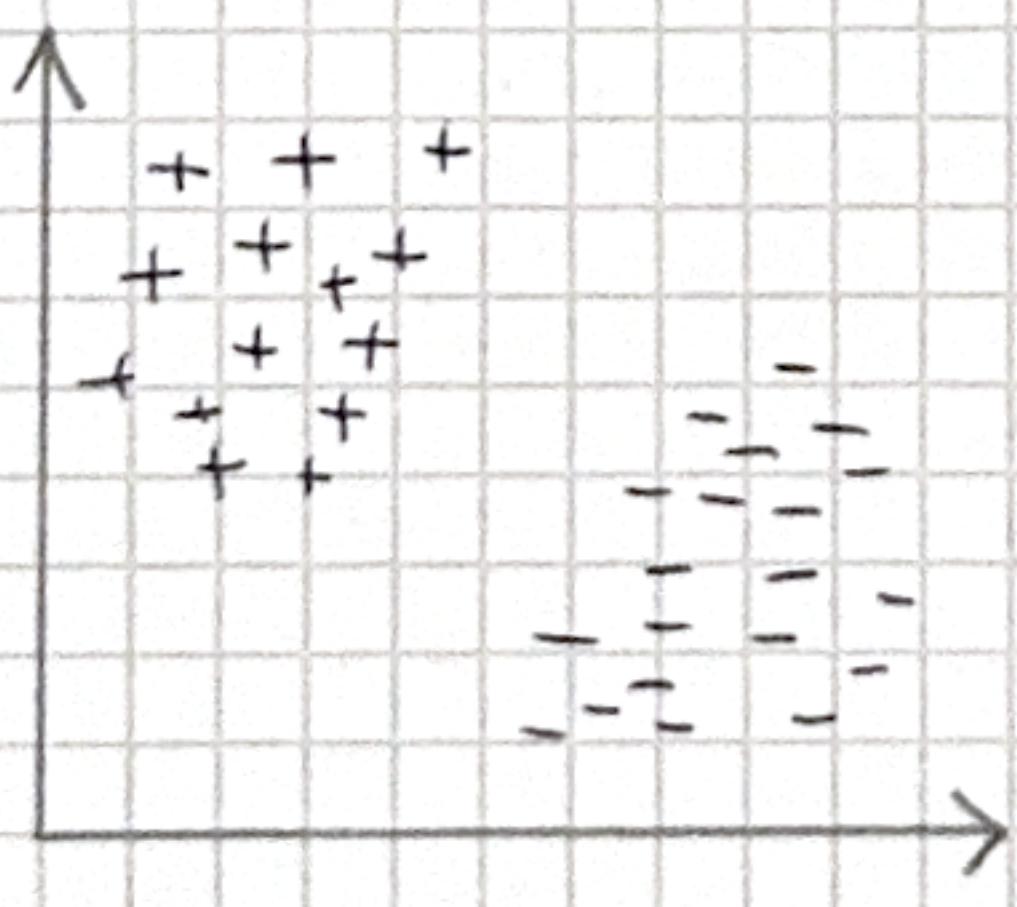
## EX. 2

1.

A dataset is LINEARLY SEPARABLE if exist an hyperplane (surface) that splits our instance into two regions much that different

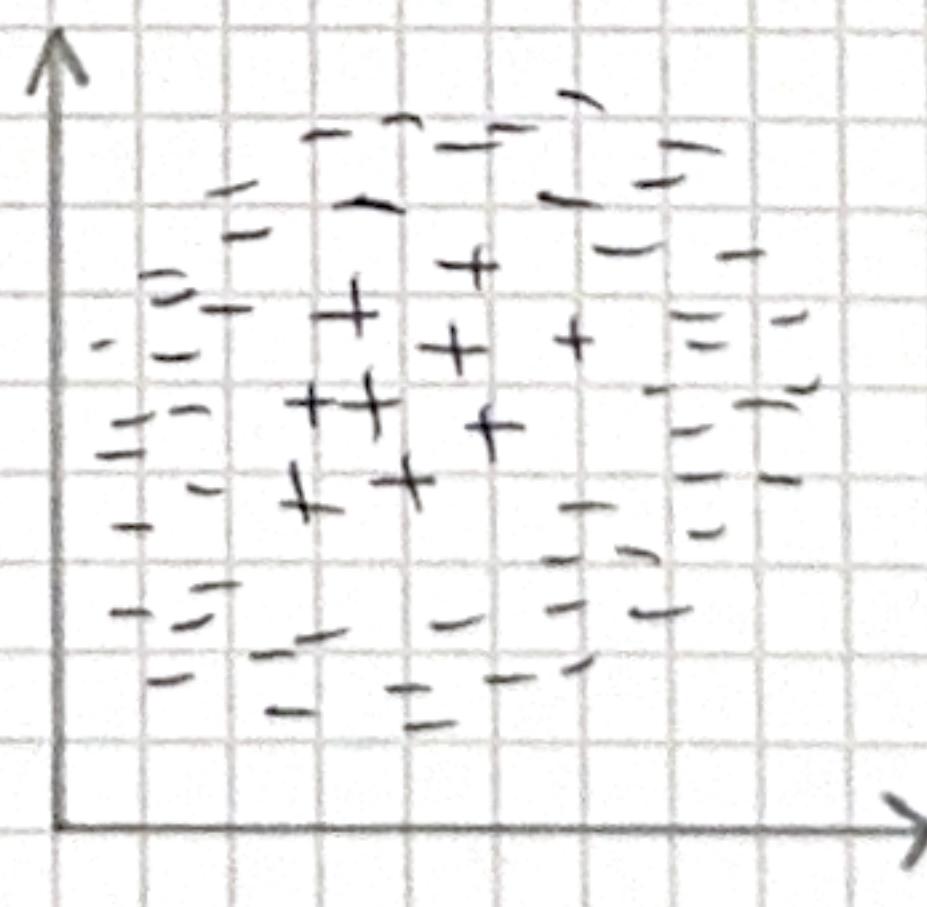
elamified and expanated

2.



SEPARABLE

3.



NOT SEPARABLE

4.

the first dataset is linearly separable and SVM finds at maximum margin providing a better accuracy  
 given a dataset D and an hyperplane h the margin is computed

$$\min_{m=1, \dots, N} \frac{1}{\|w\|} \cdot |y(x_m)| = \dots = \frac{1}{\|w\|} \min_{m=1, \dots, N} [t_m (w^T x_m + w_0)]$$

so given a dataset D and an hyperplane  $h^* = w^T x + w_0$  maximum margin is computed as:

$$w^*, w_0^* = \underset{w, w_0}{\text{argmax}} \frac{1}{\|w\|} \min_{m=1, \dots, N} [t_m (w^T x_m + w_0)]$$

a more stable version is obtain by averaging all support vector

$$w_0^* = \frac{1}{|SV|} \sum_{x_k \in SV} \left( t_k - \sum_{j \in S} \alpha_j^* t_j x_k^T x_j \right)$$

where  $\alpha_j^*$  one lagrange multiplier used to solve maximization problem write before

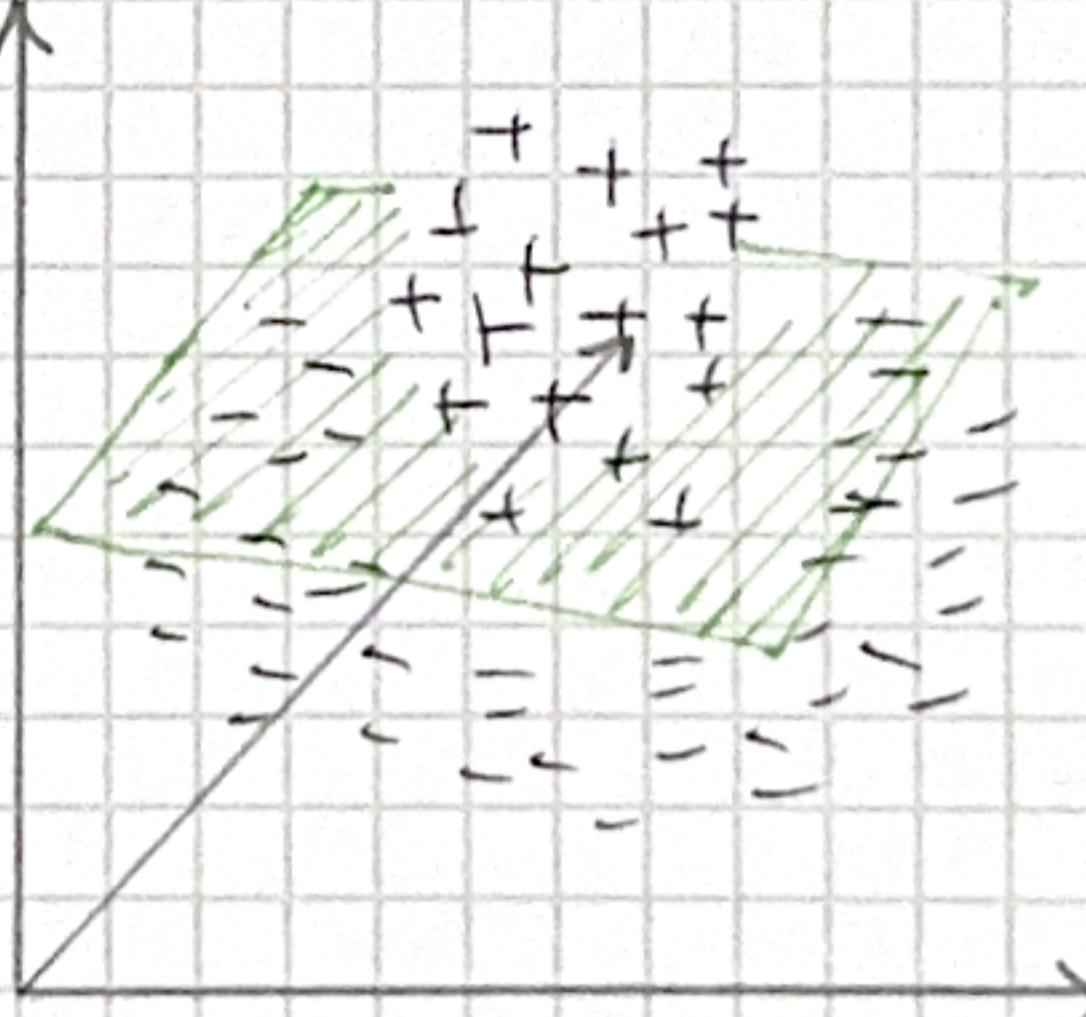
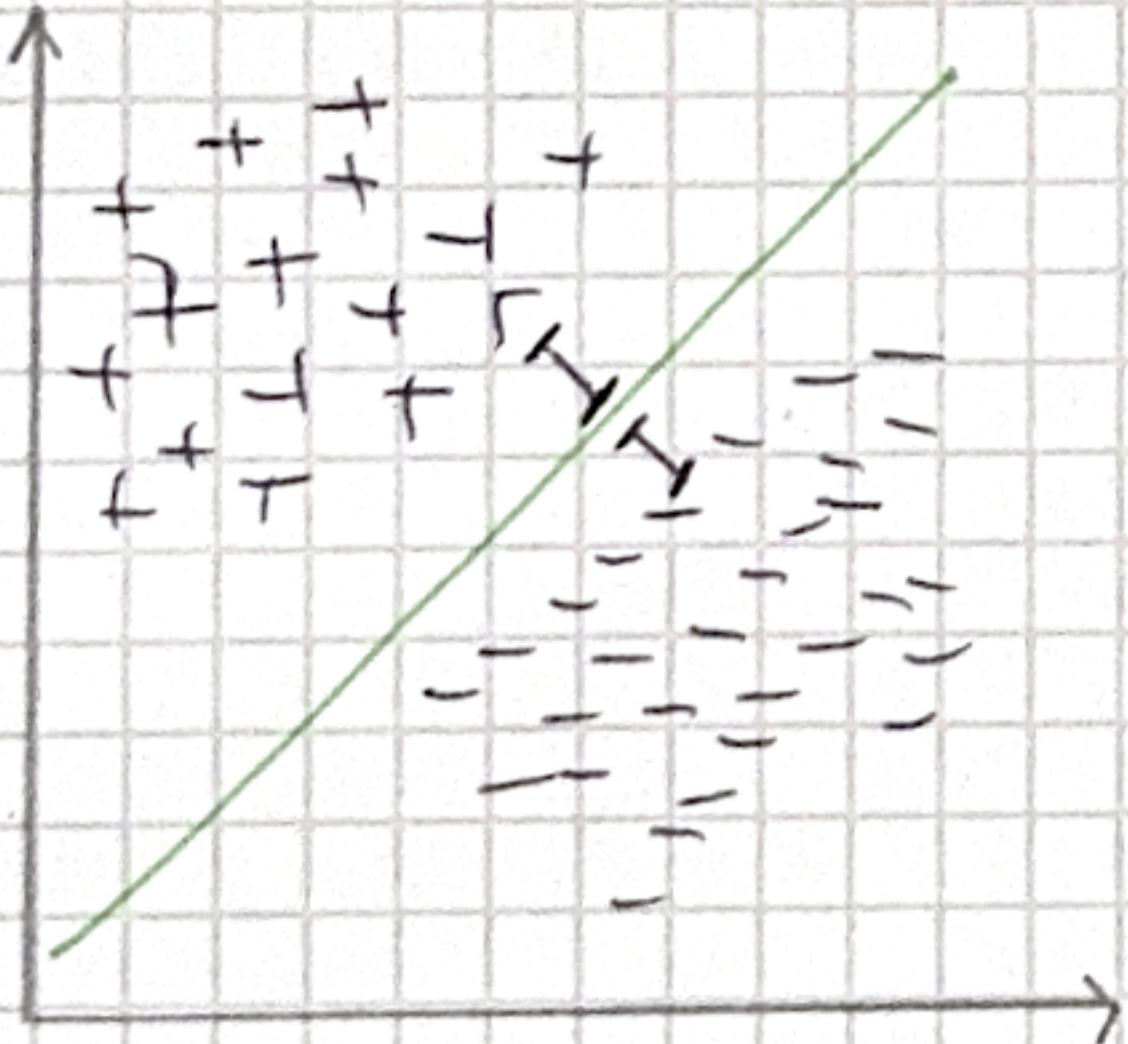
$$\tilde{L}(a) = \sum_{n=1}^N a_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^M a_n a_m t_n t_m x_n^T x_m$$

and:

$$SV \in \{x_k \in D \mid t_k y(x_k) = 1\}$$

for the second dataset because is not linearly separable we need to use some basis function like polynomial no result

$$\text{will be } g(x) = w^T \phi(x) + w_0$$



## EX.3

1. in PCA we want maximize the data variance after the projection to some dimension  $u_1$  the projected points are  $x_n^T u_1$  variance of projected points

$$\frac{1}{N} \sum_{n=1}^N [u_1^T x_n - u_1^T \bar{x}]^2 = u_1^T S u_1$$

where  $S$  ( $d \times d$ ) is covariance matrix define:

$$S = \frac{1}{N} \sum_{n=1}^N (x_n - \bar{x})(x_n - \bar{x})^T = \frac{1}{2} x^T x$$

we want maximize  $\max_{u_1} u_1^T S u_1$  with a Lagrange multiplier

$$\max_{u_1} u_1^T S u_1 \Rightarrow \lambda (1 - u_1^T u_1)$$

solution  $\lambda_1 = u_1^T S u_1$

Variance is maximum when compressed us correspond to a largest eigenvector called principal component

2. yes because the intrinsic dimensionality of data are 3 so we can write images  $\mathbb{R}^{W \times H \times C} \rightarrow \mathbb{R}^3$  in a pairs represent in 3D dimension

## EX.4

1.  $W_{in} = 32 \quad h_{in} = 32 \quad d_{in} = 3$

~~$W_{out1} = \frac{32 - 5 + 2 \cdot 2}{1} + 1 = 34$~~

~~$H_{out1} = \frac{32 - 5 + 2 \cdot 2 + 1}{1} = 32$~~

\*parameters

\* parameters 1 =  $5 \times 5 \times 3 \times 16 + 16 = 1216$

\* parameters 2 =  $3 \times 3 \times 3 \times 32 + 32 = 896$

\* parameters 3 =  $5 \times 3 \times 64 + 64 = 2844$

\* fcs =  $200 \cdot 10 = 2000$

2. I will use cross-entropy loss

$$J(\theta) = E_{x, \text{cnn}} [-\ln(p(x, \theta))]$$

assuming additive noise

$$J(\theta) = \underset{x, t \text{ i.i.d.}}{\mathbb{E}} [\frac{1}{2} \|t - f(x, \theta)\|^2]$$

EX.5

1.

$$\text{the model can be defined as } y(x, \omega) = \sum_{n=1}^N \omega_n \phi_n(x)$$

we know that target value are  $t = y(x, \omega) + \epsilon$  with noise  $\epsilon$

$$\text{if noise is gaussian } P(\epsilon | \beta) = P(\epsilon | 0, \beta^{-1})$$

$$P(\epsilon | x_1, \dots, x_N, \beta) = N(t | y(x, \omega), \beta^{-1})$$

now if we assume iid hypothesis  $P(t_1, \dots, t_N | x_1, \dots, x_N, \omega, \beta) =$

$$= \prod_{n=1}^N N(t_n | \omega^\top \phi(x_n), \beta^{-1}) = -\beta \underbrace{\frac{1}{2} \sum_{n=1}^N (t_n - \omega^\top \phi(x_n))^2}_{E_D(\omega)} - \frac{N}{2} \ln(2\pi\beta^{-1})$$

2.

we consider error function

$$E_D(\omega)$$

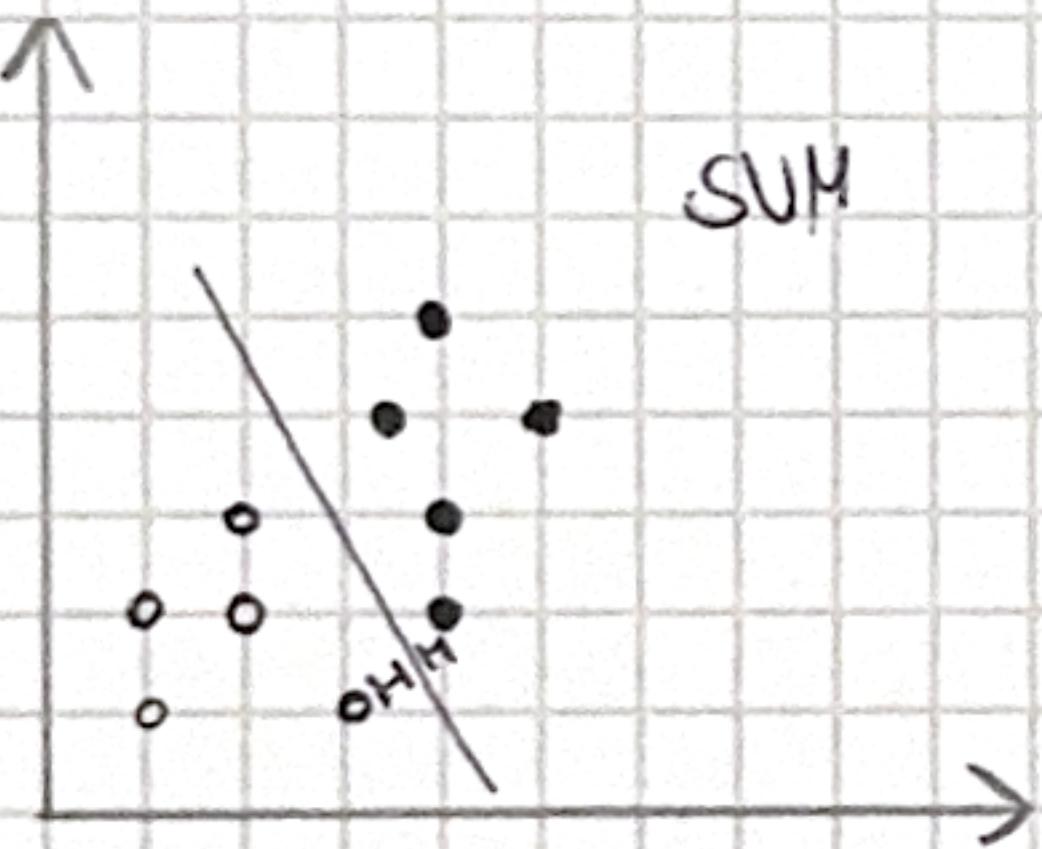
$$E_D(\omega) = \frac{1}{2} \sum_{n=1}^N (t_n - \omega^\top x_n)^2$$

we can use sequential algorithm to find  $\omega^* = \arg \min E_D(\omega)$

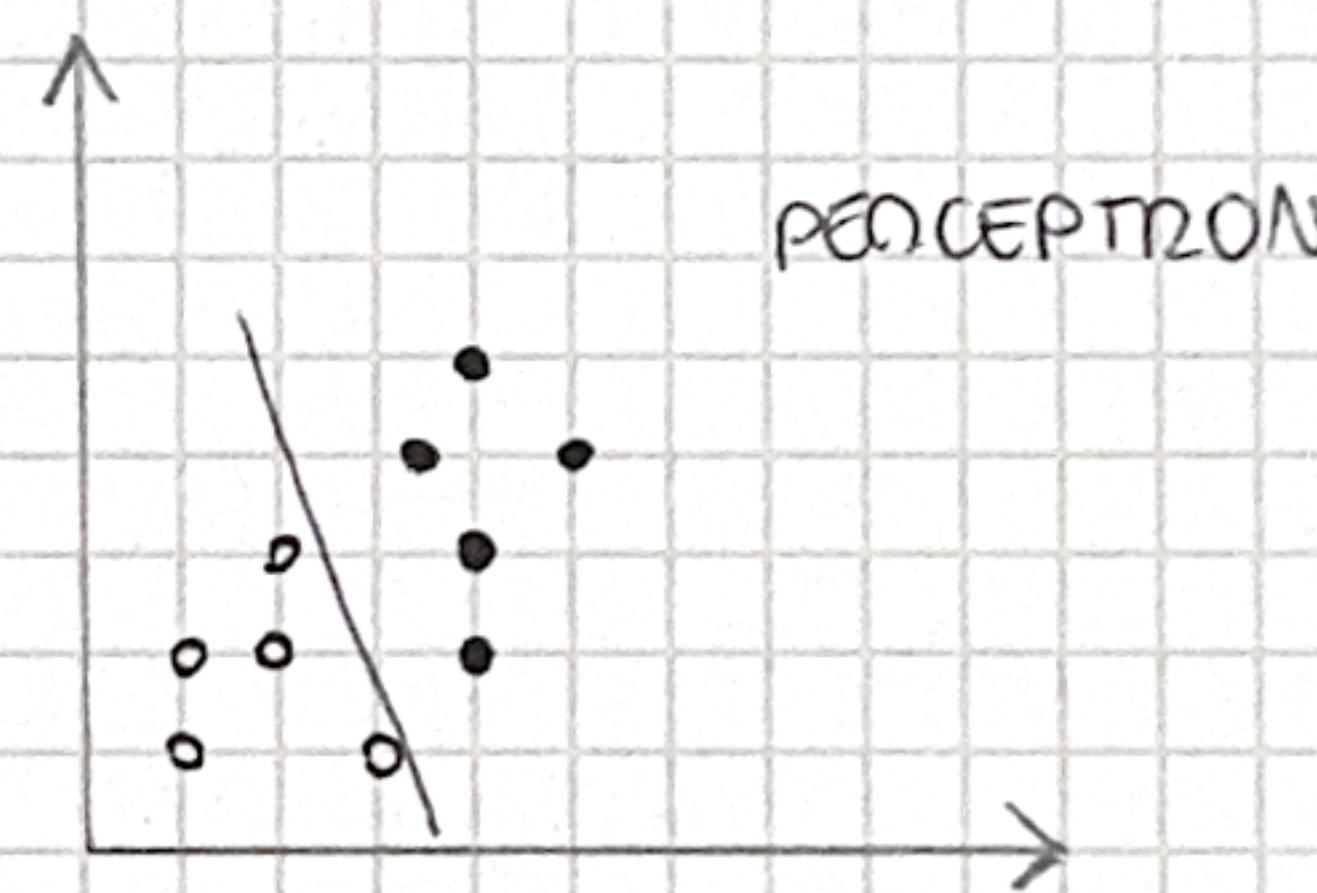
$$\hat{\omega} \leftarrow \hat{\omega} + \eta [t_n - \omega^\top \phi(x_n)] \phi(x_n)$$

EX. 6

1.



SVM



PERCEPTRON

2

the main difference is that SVM try to ~~maximize~~ aims to maximum with the better accuracy while perceptron no.

Perceptron loop iterations when finds a possible separation surface simply summing the vector related to the separation margin with the feature vector that misclassified

3. I prefer SVM because with linear separable dataset due the fact aims to maximum with better accuracy obtain better result while with perceptron we need to find the  $\eta$  connect us for have a convergence.