

# FORMULARIO ML:

## SLIDES 1.

SUPERVISED LEARNING  $D = \{(x_i, y_i)\}_{i=1}^M$

UNSUPERVISED LEARNING  $D = \{x_i\}_{i=1}^M$

REINFORCEMENT LEARNING  $\pi: S \rightarrow A \quad D = \{(\langle s_0, a_1, \pi_1, s_1, \dots, a_n, \pi_n, r_n \rangle)\}_{i=1}^n$

CONSISTENT  $(h, D) \equiv (\forall x \in D) h(x) = c(x)$

## SLIDES 2.

ERROR\_D(h)  $\equiv \Pr_{x \in D} [\hat{f}(x) \neq h(x)]$

ERROR\_S(h)  $\equiv \frac{1}{n} \sum_{x \in S} \delta(\hat{f}(x) \neq h(x))$

BIAIS  $\equiv E_S[\text{error}_S(h)] - \text{error}_D(h)$

## K FOLD CROSS VALIDATION:

for  $i = 1 \dots k$  do

$T_i \leftarrow D - S_i$

$h_i \leftarrow L(T_i)$

ACCURACY =  $1 - \frac{\text{error}_D}{k}$

$\delta_i \leftarrow \text{error}_{S_i}(h_i)$

RETURN  $\text{error}_{L,D} = \frac{1}{k} \sum_{i=1}^k \delta_i$

OVERRFIT:  $\text{error}_S(h) < \text{error}_S(h')$

$\text{error}_D(h) > \text{error}_D(h')$

RECALL =  $\frac{TP}{TP + FN}$

Precision =  $\frac{TP}{TP + FP}$

F1-SCORE =  $\frac{2 \cdot (\text{Precision} \cdot \text{Recall})}{\text{Precision} + \text{Recall}}$

} CLASSIFICATION

MAE =  $\frac{1}{M} \sum_{i=1}^M |\hat{f}(x_i) - t_i|$

MSE =  $\frac{1}{M} \sum_{i=1}^M (\hat{f}(x_i) - t_i)^2$

RMSE =  $\sqrt{MSE} = \sqrt{\frac{1}{M} \sum_{i=1}^M (\hat{f}(x_i) - t_i)^2}$

### SLIDES 3.

conclude binary classification (+, -)

$$\text{ENTROPY}(S) = -P_{(+)} \log_2 P_{(+)} - P_{(-)} \log_2 P_{(-)}$$

$$\text{ENTROPY}(S) = \sum_{i=1}^C -P_i \log_2 P_i$$

$$\text{GAIN}(S, A) = \text{ENTROPY}(S) - \sum_{v \in \text{Values}(A)} \frac{|S_v|}{|S|} \text{ENTROPY}(S_v)$$

$$\text{GAIN RATIO}(S, A) = \frac{\text{GAIN}(S, A)}{\text{SPLIT INFORMATION}(S, A)}$$

$$\text{SPLIT INFORMATION}(S, A) = -\sum_{i=1}^C \frac{|S_i|}{|S|} \log_2 \frac{|S_i|}{|S|}$$

~~$$\text{BAYES'S RULE } P(\text{CAUSE} | \text{EFFECT}) = P(\text{CAUSE})$$~~

$$\text{BAYES'S RULE} = P(\text{CAUSE} | \text{EFFECT}) = \frac{P(\text{EFFECT} | \text{CAUSE}) P(\text{CAUSE})}{P(\text{EFFECT})}$$

### SLIDES 4.

$$\text{THEOREM OF TOTAL PROBABILITY: } P(B) = \sum_{i=1}^N P(B|A_i) P(A_i)$$

LEARNING AS PROBABILISTIC ESTIMATION

$$P(h | D) = \frac{P(D|h) P(h)}{P(D)}$$

$P(h)$ : prior probability of hypothesis  $h$  (prior)

$P(D)$ : prior probability of training data  $D$  (normalization factor)

$P(h|D)$ : prior probability of  $h$  given  $D$  (posterior)

$P(D|h)$ : probability of  $D$  given  $h$  (likelihood)

$$\text{MAP hypothesis } h_{\text{MAP}} = \underset{h \in H}{\operatorname{argmax}} P(h|D) = \underset{h \in H}{\operatorname{argmax}} \left( \frac{P(D|h) P(h)}{P(D)} \right)$$

MAXIMUM APOSTERIORI HYPOTHESIS

$$\text{ML hypothesis } h_{\text{ML}} = \underset{h \in H}{\operatorname{argmax}} P(D|h)$$

BAYES OPTIMAL CLASSIFIER:

$$V_{OB} = \underset{v_j \in V}{\operatorname{argmax}} \sum_{h \in H} P(v_j | x, h) P(h|D)$$

+ target function  $f: X \rightarrow V$   $V = \{v_1, \dots, v_k\}$  disjoint  $D$

BERNOULLI DISTRIBUTION

$$P(X=k; \theta) = \theta^k (1-\theta)^{1-k}$$

## MULTIVARIATE BERNoulli DISTRIBUTION

$$P(X_1 = k_1, \dots, X_M = k_M; \theta_1, \dots, \theta_M) = \prod_{i=1}^M P(X_i = k_i; \theta_i) = \prod_{i=1}^M \theta_i^{k_i} (1 - \theta_i)^{1-k_i}$$

## BINOMIAL DISTRIBUTION

$$P(X = k; n, \theta) = \binom{n}{k} \theta^k (1 - \theta)^{n-k}$$

## MULTINOMIAL DISTRIBUTION

$$P(X_1 = k_1, \dots, X_d = k_d; n, \theta_1, \dots, \theta_d) = \frac{n!}{k_1! \dots k_d!} \theta_1^{k_1} \dots \theta_d^{k_d}$$

## NAIVE BAYES CLASSIFIER

~~x is conditionally independent to approximate the posterior to Y given Z~~

$$P(X, Y | Z) = P(X|Y, Z)P(Y|Z) = P(X|Z)P(Y|Z)$$

$$VNB = \underset{v_j \in V}{\operatorname{argmax}} P(v_j | D) \prod_i P(a_i | v_j, D) \text{ with assumption}$$

$$P(a_1, \dots, a_n | v_j, D) = \prod_i P(a_i | v_j, D)$$

## NAIVE BAYES ALGORITHM

target function  $f: X \rightarrow V \quad X = A_1 \times \dots \times A_m \quad V = \{v_1, \dots, v_n\}$

new instance  $x = (a_1, a_2, \dots, a_n)$

for each target value  $v_j \in V$

$$\hat{P}(v_j | D) \leftarrow \text{estimate } P(v_j | D)$$

for each attribute  $A_k$

for each attribute  $a_i \in A_k$

$$\hat{P}(a_i | v_j, D) \leftarrow \text{estimate } P(a_i | v_j, D)$$

$$VNB = \underset{v_j \in V}{\operatorname{argmax}} \hat{P}(v_j | D) \prod_{a_i \in x} \hat{P}(a_i | v_j, D)$$

## NAIVE BAYES ESTIMATION

$$\hat{P}(a_i | v_j, D) = \frac{\text{# of } a_i \text{ in } v_j}{\text{# of words in } v_j}$$

## WORD REPRESENTATION

1. BOOLEAN FEATURE VECTOR  $\{d_i = 1 \text{ if } w_i \text{ appears in doc} \quad 0 \text{ otherwise}\}$

2. ORDINAL FEATURE VECTOR  $d_i = k \text{ if } w_i \text{ appears } k \text{ times in doc}$   
(Multinomial Distribution)  $\rightarrow$  (Bag of Words)

$$3. \text{tf-idf} : d_i = \text{tf}(w_i, \text{doc}) \cdot \text{idf}(w_i, D)$$

term frequency inverse document frequency

MULTIVARIATE BERNOULLI DISTRIBUTION (NAIVE BAYES)

boolean feature vector  $d = [d_1, \dots, d_n]$  of generic doc  $\in \text{Docs}$

$$P(d | c_j, D) = \prod_{i=1}^m P(w_i | c_j, D)^{d_i} (1 - P(w_i | c_j, D))^{1-d_i}$$

MAXIMUM LIKELIHOOD SOLUTION

$$\hat{P}(w_i | c_j, D) = \frac{t_{ij} + 1}{t_j + 2}$$

$t_{c_j}$  = number of documents in  $D$  of class  $c_j$  containing word  $w_i$

$t_j$  = number of documents in  $D$  of class  $c_j$

$1, 2 = \text{LAPLACIAN TERM}$

MULTINOMIAL NAIVE BAYES DISTRIBUTION

$$P(d | c_j, D) = \frac{m!}{d_1! \dots d_m!} \prod_{i=1}^m P(w_i | c_j, D)^{d_i}$$

MAXIMUM LIKELIHOOD SOLUTION

$$\hat{P}(w_i | c_j, D) = \frac{\sum_{\text{doc} \in D} t_{f,i,j} + \alpha}{\sum_{\text{doc} \in D} t_{f,i,j} + \alpha |V|}$$

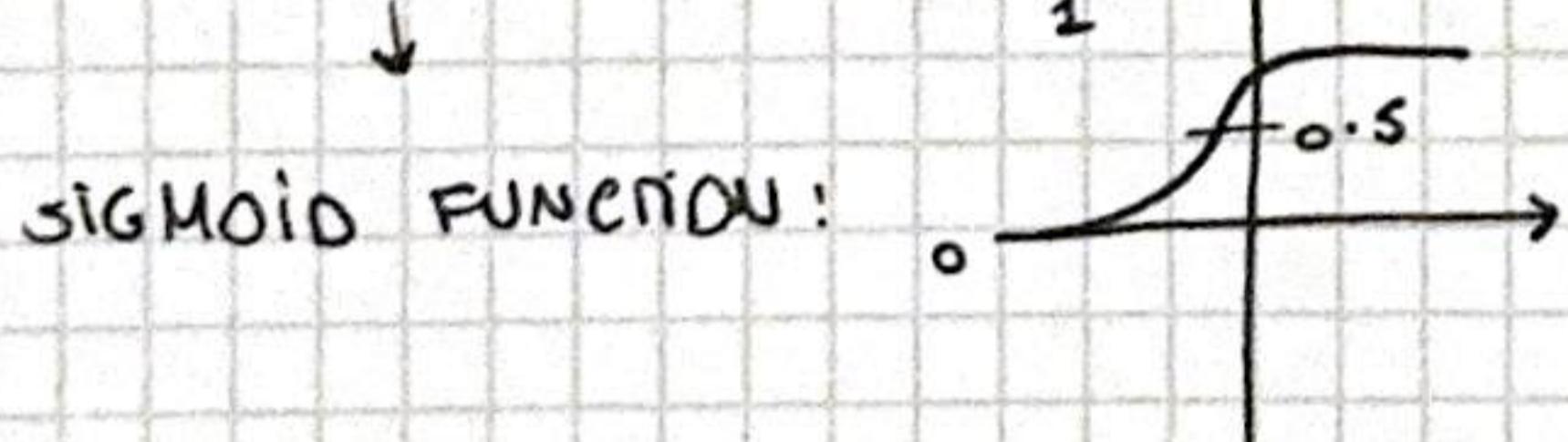
SIDE 6:

PROBABILISTIC GENERATIVE MODEL

$$P(c_1 | x) = \frac{P(x | c_1) P(c_1)}{P(x)} = \frac{P(x | c_1) P(c_1)}{P(x | c_1) P(c_1) + P(x | c_2) P(c_2)} =$$

$$= \frac{1}{1 + \exp(-\alpha)} = \delta(\alpha)$$

$$\alpha = \ln \left( \frac{P(x | c_1) P(c_1)}{P(x | c_2) P(c_2)} \right)$$



MAXIMUM LIKELIHOOD SOLUTION K CLASSES

$$P(c_k | x) = \frac{P(x | c_k) P(c_k)}{\sum_j P(x | c_j) P(c_j)} = \frac{\exp(\alpha_k)}{\sum_j \exp(\alpha_j)}$$

SOFT MAX

## NEWTON-RAPSON:

gradient of error with respect to  $\tilde{\omega}$

$$\nabla E(\tilde{\omega}) = \sum_{m=1}^N (y_m - t_m) \tilde{x}_m$$

gradient descent step

$$\tilde{\omega} \leftarrow \tilde{\omega} - H(\tilde{\omega})^{-1} \nabla E(\tilde{\omega})$$

$H(\tilde{\omega}) = \nabla \nabla E(\tilde{\omega})$  is the Hessian matrix of  $E(\tilde{\omega})$

## ITERATIVE METHOD

1. initialize  $\tilde{\omega}$

2. repeat until termination

$$\tilde{\omega} \leftarrow \tilde{\omega} - (\tilde{x}^\top R(\tilde{\omega}) \tilde{x})^{-1} \tilde{x}^\top (y(\tilde{\omega}) - t)$$

with:

$$\tilde{x} = \begin{pmatrix} \tilde{x}_1^\top \\ \vdots \\ \tilde{x}_N^\top \end{pmatrix} \quad t = \begin{pmatrix} t_1 \\ \vdots \\ t_N \end{pmatrix} \quad y(\tilde{\omega}) = (y_1 \dots y_N)^\top$$

$R(\tilde{\omega})$ : diagonal matrix with  $R_{NN} = y_m(1-y_m)$

## MULTICLASS LOGISTIC REGRESSION

$$P(\tau | \tilde{\omega}_1, \dots, \tilde{\omega}_K) = \prod_{m=1}^M \prod_{k=1}^K P(C_k | \tilde{x}_m)^{t_{mk}} = \prod_{m=1}^M \prod_{k=1}^K y_{mk}^{t_{mk}}$$

cross entropy error function

$$E(\tilde{\omega}_1, \dots, \tilde{\omega}_K) = -\ln(P(\tau | \tilde{\omega}_1, \dots, \tilde{\omega}_K)) = -\sum_{m=1}^M \sum_{k=1}^K t_{mk} \ln y_{mk}$$

## SLIDES 7:

### • LINEAR DISCRIMINANT FUNCTION

$$y: X \rightarrow \{C_1, \dots, C_K\}$$

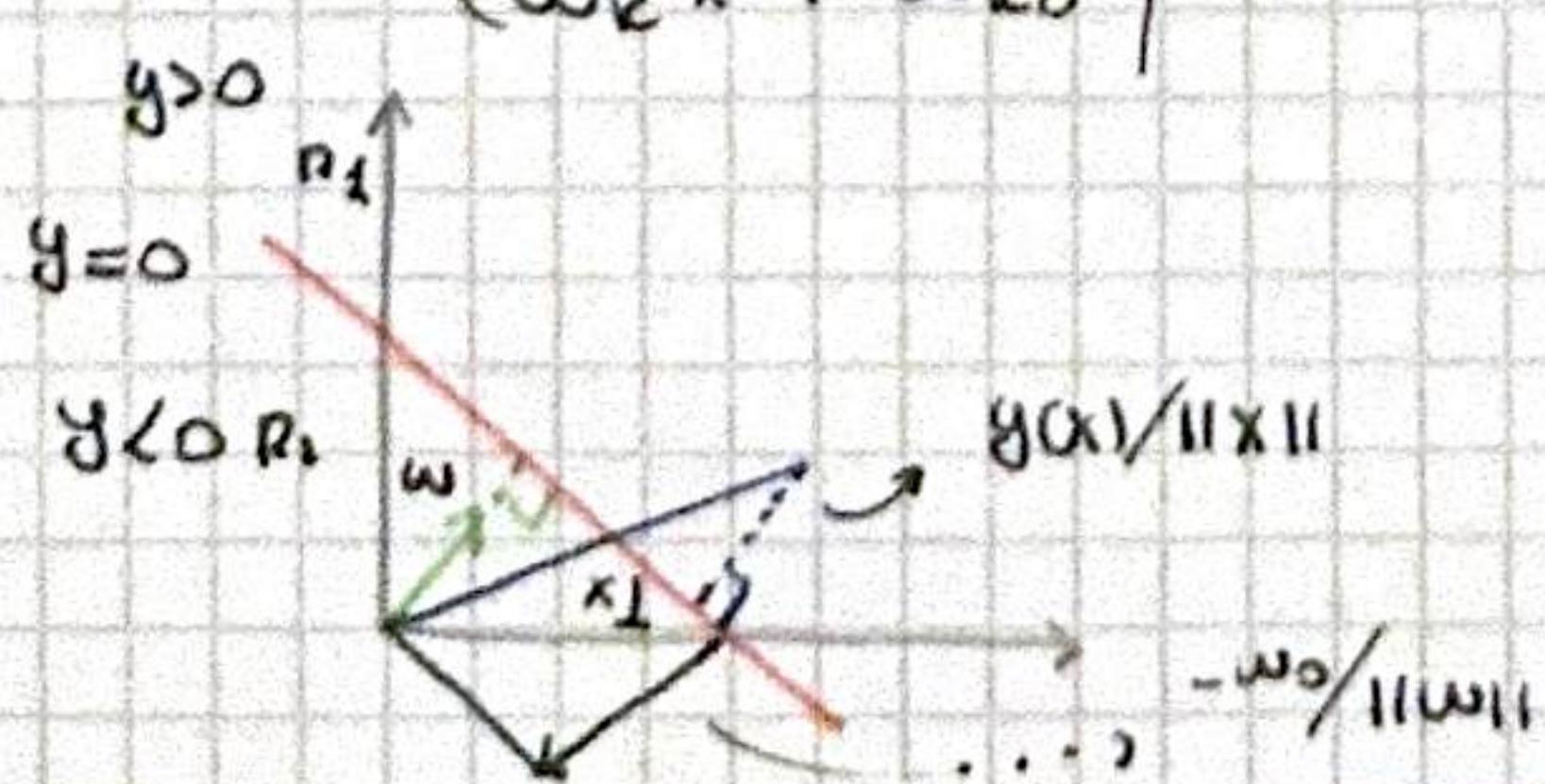
$$k \text{ classes} \quad y_1(x) = \tilde{w}_1^\top x + w_{10}$$

$$\vdots$$

$$y_K(x) = \tilde{w}_K^\top x + w_{K0}$$

### COMPACT NOTATION

$$y(x) = \begin{pmatrix} y_1(x) \\ \vdots \\ y_K(x) \end{pmatrix} = \begin{pmatrix} \tilde{w}_1^\top x + w_{10} \\ \vdots \\ \tilde{w}_K^\top x + w_{K0} \end{pmatrix} = \begin{pmatrix} \tilde{w}_1^\top \\ \vdots \\ \tilde{w}_K^\top \end{pmatrix} \tilde{x} = \tilde{w}^\top \tilde{x}$$



## GAUSSIAN NAIVE MODEL

$$P(C_k) = \pi_k \quad P(x|C_k) = N(x, \mu_k, \Sigma)$$

TRAINING  $\rightarrow$  estimating from D

$$\hat{\pi}_k = \frac{N_k}{N} \quad \hat{\mu}_k = \frac{1}{N_k} \sum_{n=1}^N t_{nk} x_n \quad \hat{\Sigma} = \sum_{k=1}^K \frac{N_k}{N} S_k$$

$$S_k = \frac{1}{N_k} \sum_{n=1}^N t_{nk} (x_n - \mu_k)(x_n - \mu_k)^T$$

PREDICTION

$$\underset{C \in C}{\operatorname{argmax}} \hat{P}(C_k | x) = \underset{C \in C}{\operatorname{argmax}} \frac{\exp(\alpha_k)}{\sum_j \exp(\alpha_j)} \quad \alpha_k = \ln(\hat{P}(x|C_k)P(C_k))$$

## POSTERIOR DISTRIBUTIONS WITH PARAMETERS MODELS

K CLASSES

$$P(C_i|x) = \frac{\exp(\alpha_i)}{\sum_j \exp(\alpha_j)} \quad \alpha_i = \tilde{w}^T x + w_0$$

## PROBABILISTIC DISCRIMINATIVE MODELS

$$\text{estimate directly } P(C_k|\tilde{x}, D) = \frac{\exp(\alpha_k)}{\sum_j \exp(\alpha_j)} \quad \alpha_k = \tilde{w}^T \tilde{x}$$

$$\text{MAXIMUM LIKELIHOOD } \tilde{w}^* = \underset{\tilde{w}}{\operatorname{argmax}} \ln(P(t|\tilde{w}, x))$$

LOGISTIC REGRESSION 2 classes dataset  $D = \{(x_m, t_m)\}_{m=1}^N$  with  $t \in \{0, 1\}$

likelihood:

$$p(t|\tilde{w}) = \prod_{m=1}^M y_m^{t_m} (1-y_m)^{1-t_m}$$

$$\text{with } y_m = p(C_1|x_m) = g(\tilde{w}^T x_m)$$

CROSS ENTROPY:

$$E(\tilde{w}) = -\ln(p(t|\tilde{w}))$$

$$\text{optimization} \rightarrow \underset{\tilde{w}}{\operatorname{argmin}} E(\tilde{w})$$

## MULTICLASS

$$y(x) = \begin{pmatrix} y_1(x) \\ \vdots \\ y_k(x) \end{pmatrix} = \begin{pmatrix} \tilde{w}_1^T \tilde{x} \\ \vdots \\ \tilde{w}_k^T \tilde{x} \end{pmatrix} = \tilde{w}^T \tilde{x}$$

classify  $x$  as  $C_k$  with  $k = \operatorname{argmax}_{j=1, \dots, k} \{y_j(x)\}$

## LEAST SQUARES

minimize sum of squares error function:

$$E(\tilde{w}) = \frac{1}{2} \operatorname{tr} \{ (\tilde{x} \tilde{w} - \tau)^T (\tilde{x} \tilde{w} - \tau) \}$$

closed-form solution:

$$\tilde{w} = \underbrace{(\tilde{x}^T \tilde{x})^{-1}}_{\tilde{X}} \tilde{x}^T \tau$$

learned model

$$y(x) = \tilde{w}^T \tilde{x} = \tau^T (\tilde{x}^T)^T \tilde{x}$$

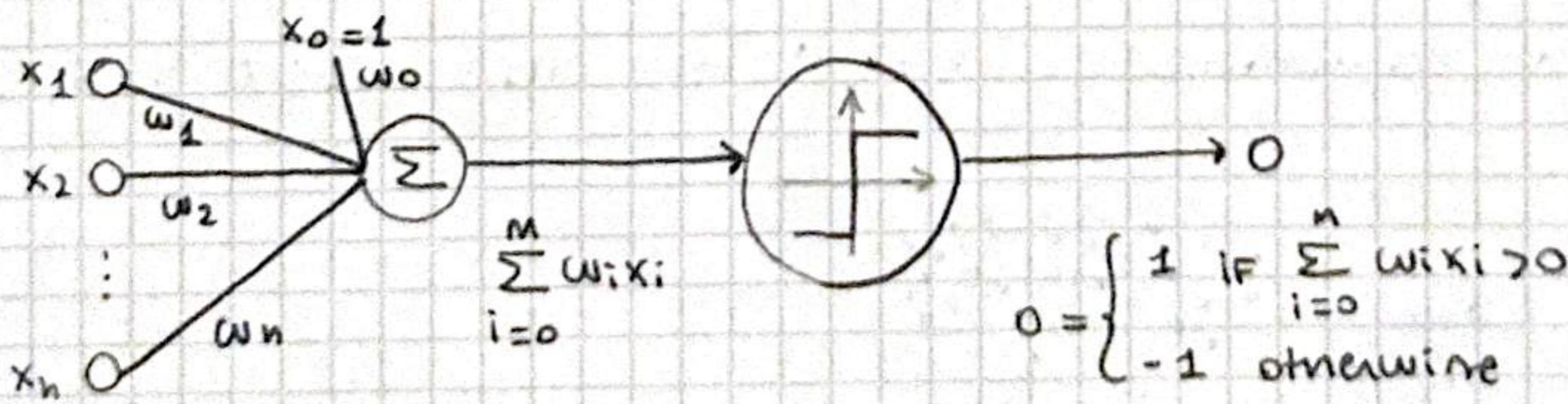
use  $\tilde{w}$  to compute  $y(x) = \tilde{w}^T \tilde{x} = \begin{pmatrix} y_1(x) \\ \vdots \\ y_k(x) \end{pmatrix}$

predict  $k = \operatorname{argmax}_{j \in \{1, \dots, k\}} \{y_j(x)\}$



## PERCEPTRON

$$f: \mathbb{R}^d \rightarrow \{-1, +1\}$$



$$o(x_1, \dots, x_d) = \begin{cases} 1 & \text{if } w_0 + w_1 x_1 + \dots + w_d x_d > 0 \\ -1 & \text{otherwise} \end{cases}$$

$$o(x) = \begin{cases} 1 & \text{if } w^T x > 0 \\ -1 & \text{otherwise} \end{cases} = \text{sign}(w^T x)$$

minimize squared error function:

$$E(w) = \frac{1}{2} \sum_{n=1}^N (t_n - o_m)^2 = \frac{1}{2} \sum_{n=1}^N (t_n - w^T x_n)^2$$

$$\frac{\partial E(w)}{\partial w} = \sum_{n=1}^N (t_n - w^T x_n) (-x_{i,m})$$

update of weights

$$w_i \leftarrow w_i + \Delta w_i$$

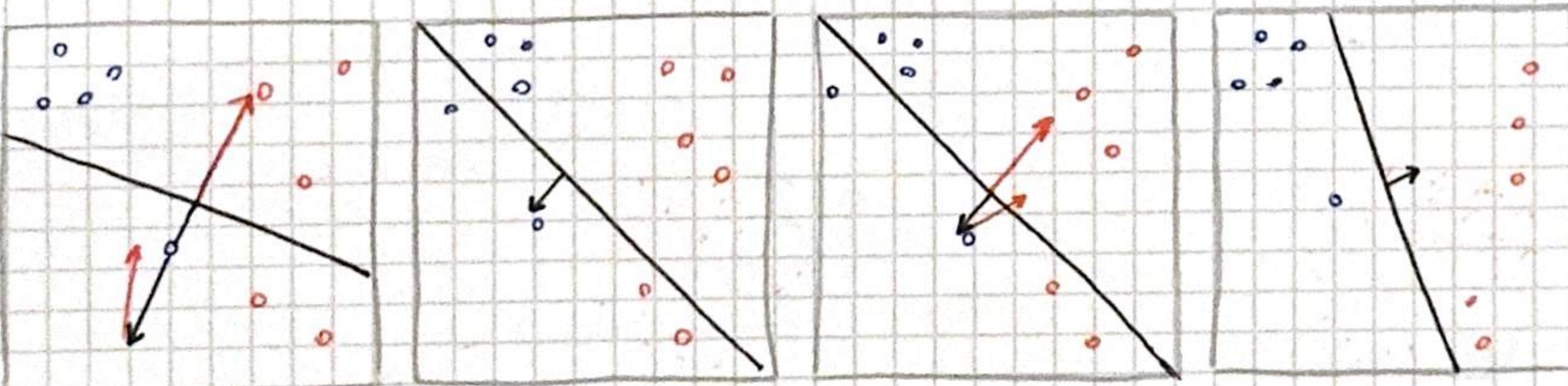
$$\Delta w_i = -\eta \frac{\partial E}{\partial w_i} = \eta \sum_{m=1}^N (t_n - w^T x_n) x_{i,m}$$

BATCH-MODE  $\rightarrow$  all dataset  $\Delta w_i = \eta \sum_{(x,t) \in D} (t - o(x)) x_i$

MINI BATCH-MODE  $\rightarrow$  choose small dataset  $\Delta w_i = \eta \sum_{(x,t) \in S} (t - o(x)) x_i$

INCREMENTAL-MODE  $\rightarrow$  choose one sample  $(x,t) \in D$

$$\Delta w_i = \eta (t - o(x)) x_i$$

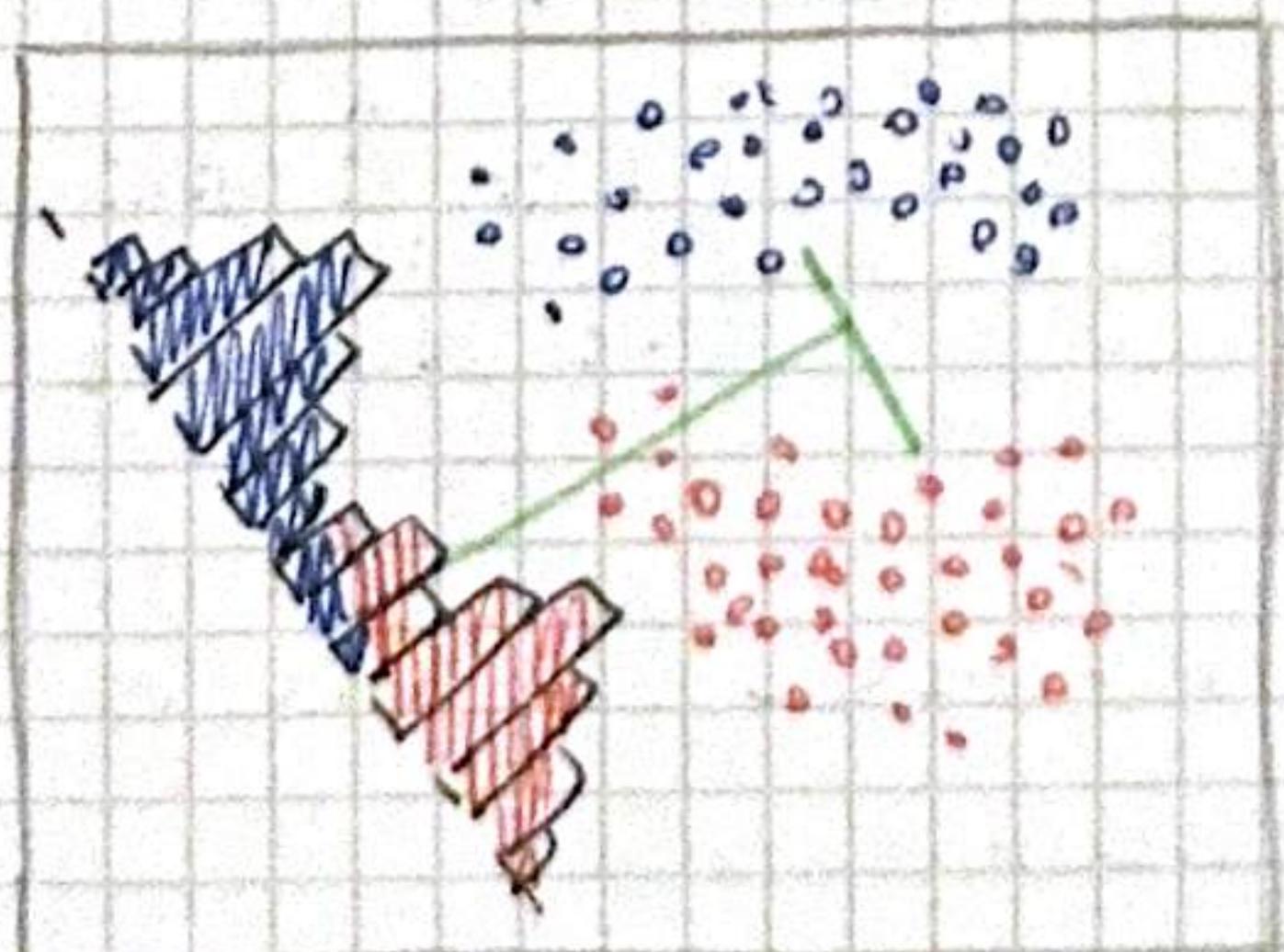


### FISHER LINEAR DISCRIMINANT

dataset  $N_1$  points in  $C_1$  and  $N_2$  points in  $C_2$

$$m_1 = \frac{1}{N_1} \sum_{n \in C_1} x_n \quad m_2 = \frac{1}{N_2} \sum_{n \in C_2} x_n$$

choose  $\vec{w}$  that maximizes  $J(w) = w^T (m_2 - m_1)$  subject to  $\|w\|=1$



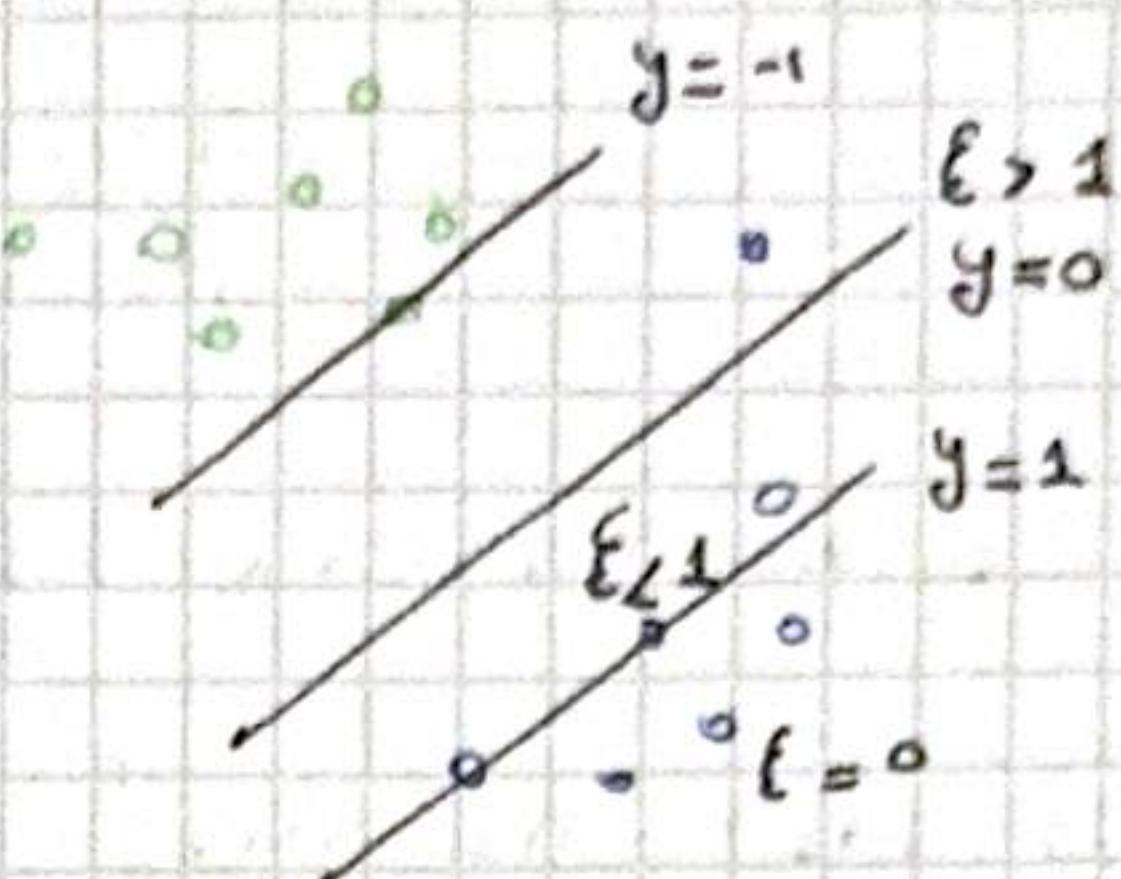
## SVM SOFT MARGIN CONSTRAINT

$\xi_n = 0$  ON OR INSIDE decision boundaries

$0 < \xi_n < 1$  INSIDE margin correct node

$\xi_n > 1$  wrong side margin

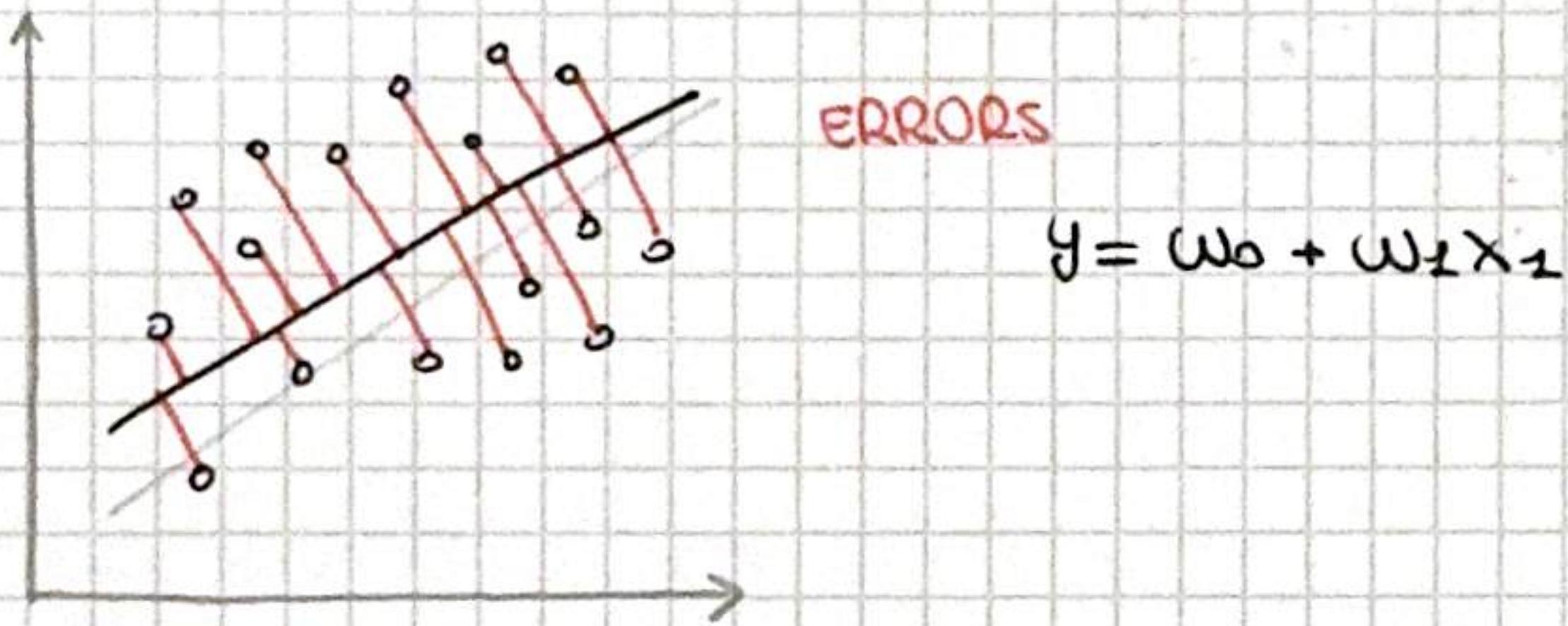
$$\omega^*, \omega_0^* = \operatorname{arg\min} \frac{1}{2} \|\omega\|^2 + C \sum_{n=1}^N \xi_n$$



## SLIDES 8

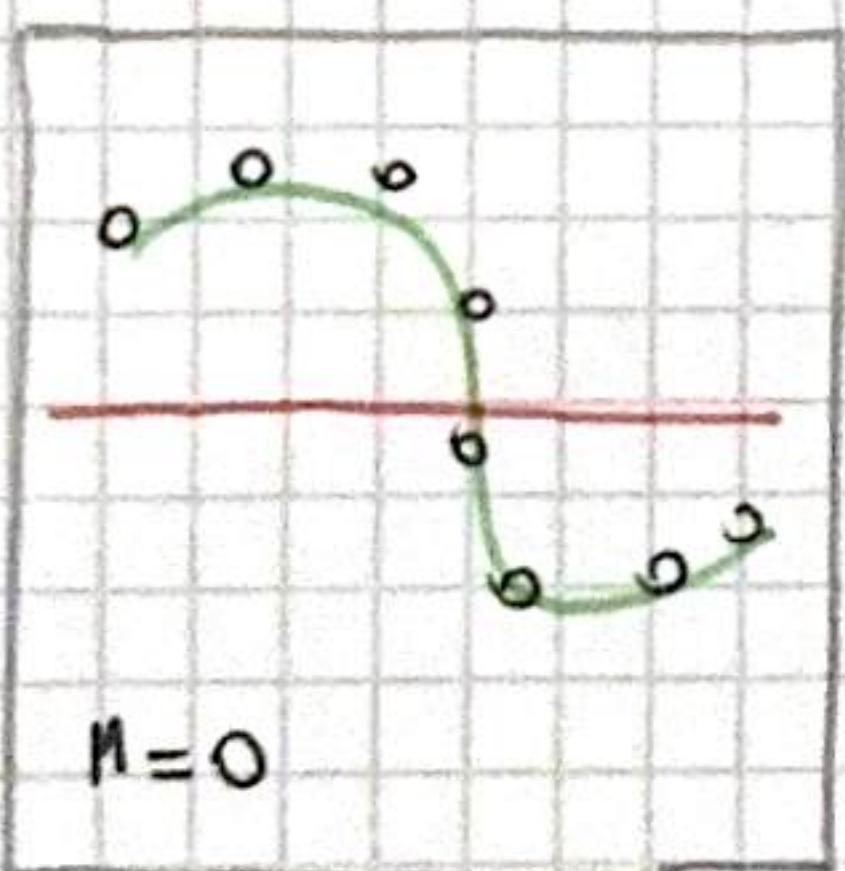
### LINEAR MODEL FOR REGRESSION:

$$y(x; \omega) = \omega_0 + \omega_1 x_1 + \dots + \omega_d x_d = \omega^T x \quad \text{with} \quad x = \begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_d \end{bmatrix} \quad \omega = \begin{bmatrix} \omega_0 \\ \omega_1 \\ \vdots \\ \omega_d \end{bmatrix}$$

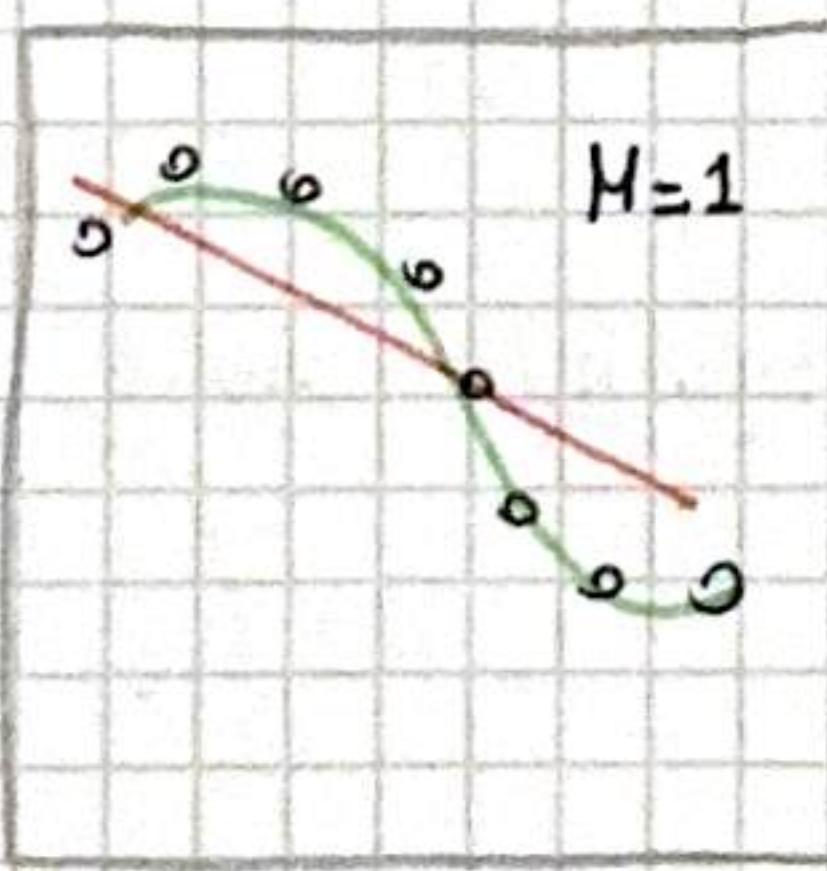


USING NON-LINEAR FUNCTIONS OF INPUT VARIABLES:  $y(x, \omega) = \sum_{j=0}^M \omega_j \phi_j(x) \Leftrightarrow \omega^T \phi(x)$

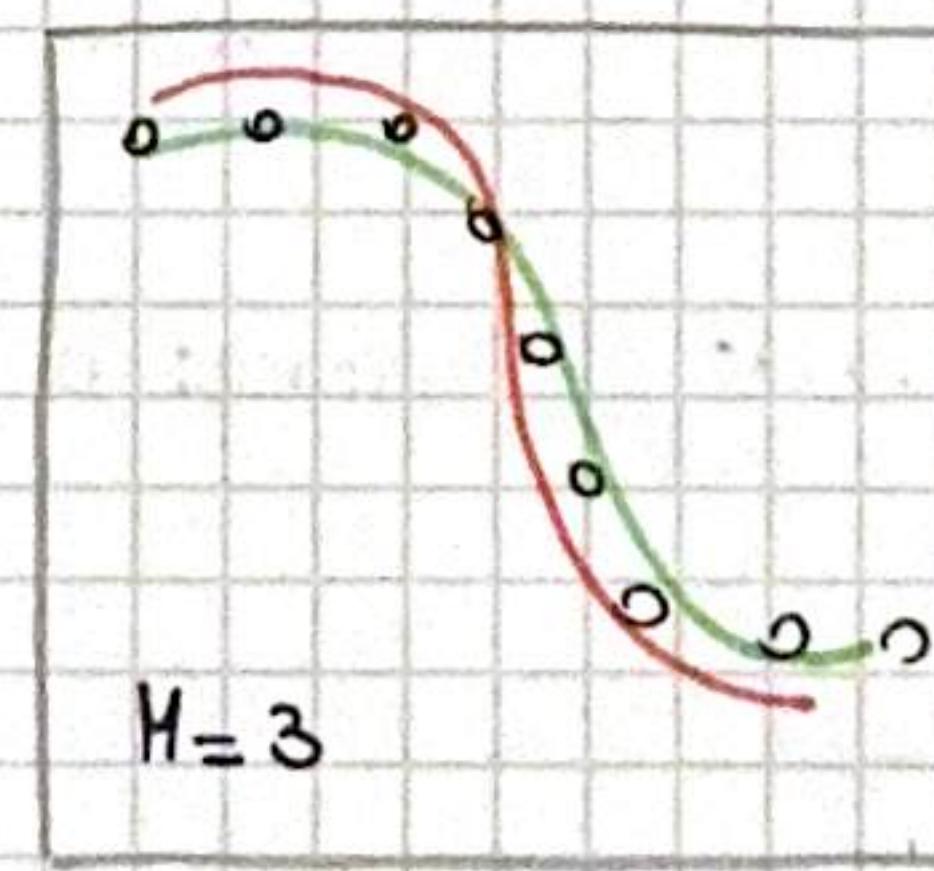
### POLYNOMIAL CURVE FITTING:



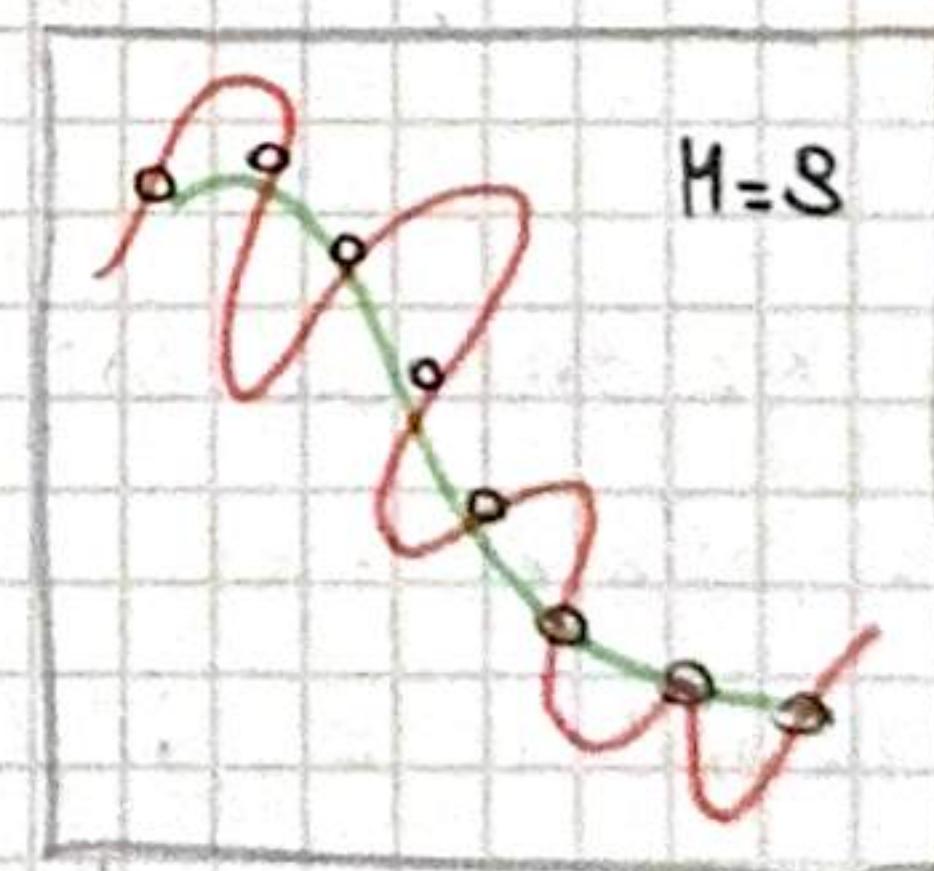
UNDERFITTING



H=1



H=3



H=8

OVERRFITTING

### MAXIMUM LIKELIHOOD (ZERO-MEAN GAUSSIAN NOISE ASSUMPTION)

$$\operatorname{arg\max}_{\omega} P(\{t_1, \dots, t_N\} | x_1, \dots, x_N, \omega, \beta)$$

CORRESPOND TO LEAST SQUARE ERROR MINIMIZATION

$$\operatorname{arg\min}_{\omega} E_D(\omega) = \operatorname{arg\min}_{\omega} \frac{1}{2} \sum_{m=1}^N [t_m - \omega^T \phi(x_m)]^2$$

$$E_D(\omega) = \frac{1}{2} (t - \phi \omega)^T (t - \phi \omega)$$

$$\omega_{ML} = \underbrace{(\phi^T \phi)^{-1}}_{\text{PSEUDO INVERSE}} \phi^T t$$

$$\phi = \begin{bmatrix} \phi(x_1) & \cdots & \phi_N(x_1) \\ \vdots & \ddots & \vdots \\ \phi(x_N) & \cdots & \phi_N(x_N) \end{bmatrix} \quad t = \begin{bmatrix} t_1 \\ \vdots \\ t_N \end{bmatrix}$$

## FISHER CRITERION

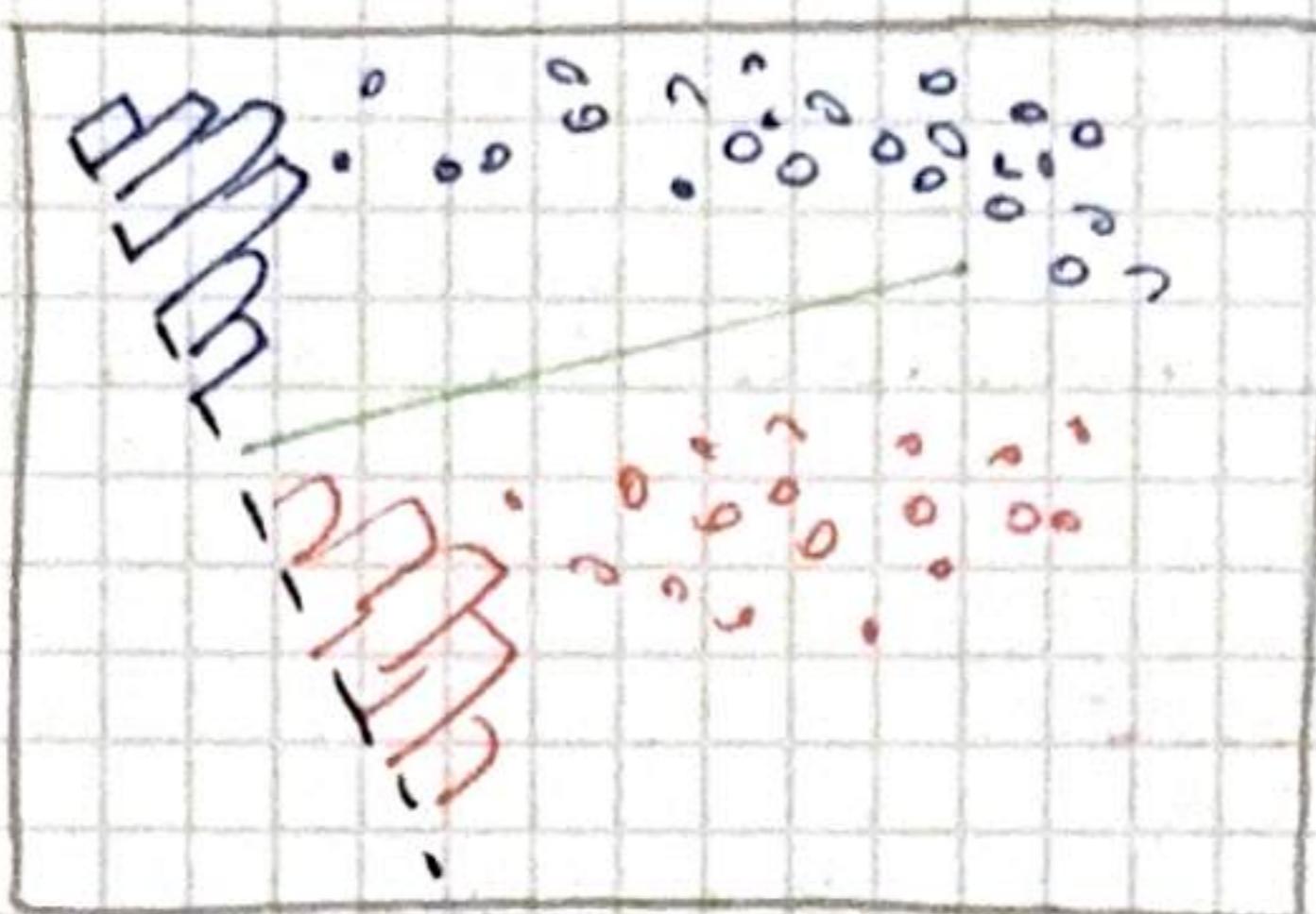
$$J(W) = \frac{W^T S_B W}{W^T S_W W}$$

$S_W$

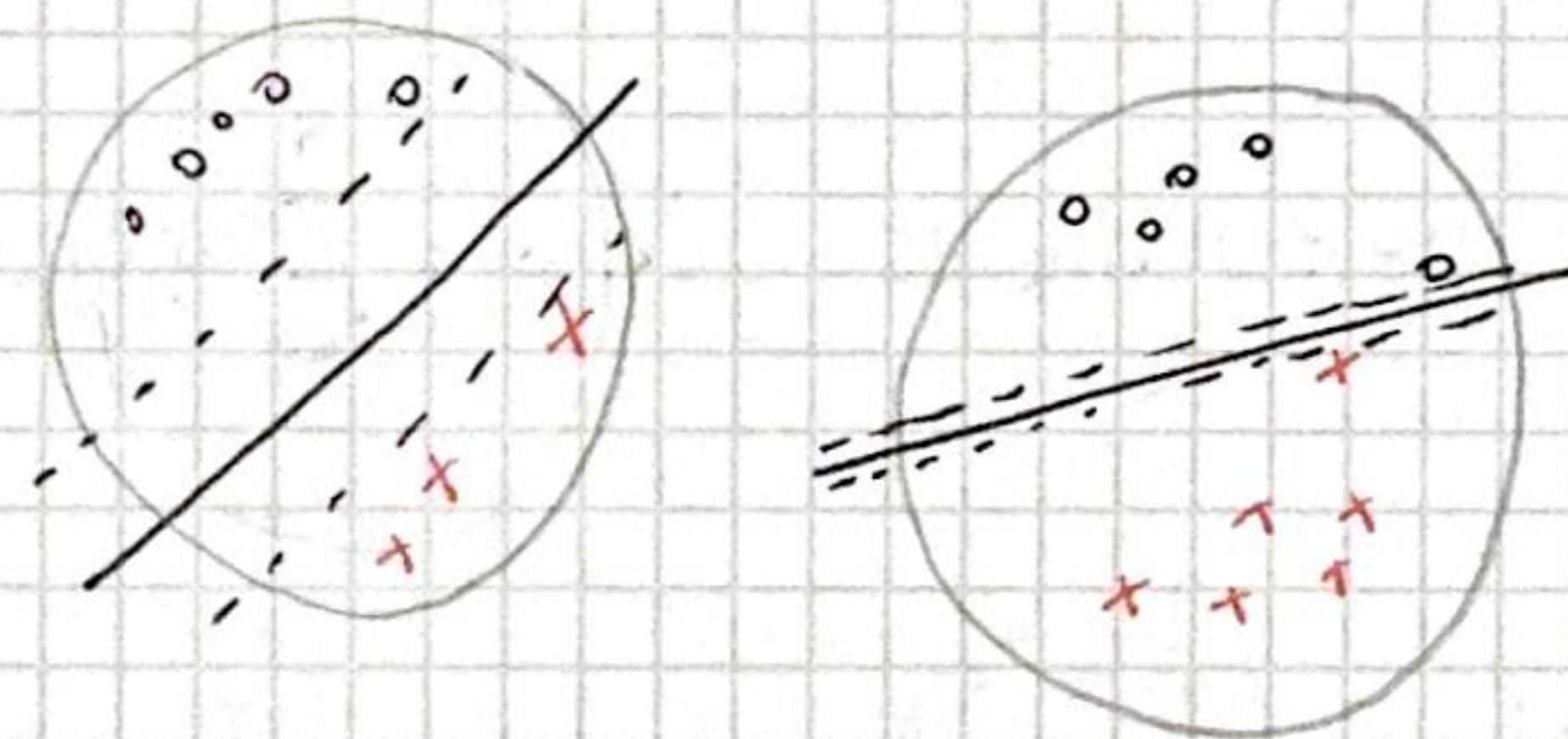
$$S_W = \sum_{m \in C_1} (x_n - m_1)(x_m - m_1)^T + \sum_{m \in C_2} (x_n - m_2)(x_m - m_2)^T$$

$$S_B = (m_2 - m_1)(m_2 - m_1)^T$$

maximite  $J(W)$  by solving  $\frac{d}{dW}(J(W)) = 0$



## SUPPORT VECTOR MACHINE



let  $x_k$  the closest point of dataset D to the hyperplane

$$h = W^T x + w_0 = 0$$

$$\text{MARGIN is } \frac{|y(x_k)|}{\|W\|}$$

given dataset D and hyperplane h the margin is computed as

$$\min_{m=1, \dots, N} \frac{|y(x_m)|}{\|W\|} = \dots = \frac{1}{\|W\|} \min_{m=1, \dots, N} [t_n (W^T x_n + w_0)]$$

given dataset D and hyperplane  $h^*: W^T w x + w_0$  with maximum margin as

$$W^*, w_0^* = \underset{w, w_0}{\text{anymax}} \frac{1}{\|W\|} \min_{m=1, \dots, N} [t_n (W^T x_n + w_0)]$$

quadratic problem solved with LAGRANGIAN METHOD

$$W^* = \sum_{m=1}^N \alpha_m^* t_m x_m$$

as in result of LAGRANGIAN METHOD optimization problem

$$\tilde{L}(\alpha) = \sum_{n=1}^N \alpha_n - \frac{1}{2} \sum_{m=1}^N \sum_{M=1}^N \alpha_m \alpha_M t_m t_M x_n^T x_m$$

SEQUENTIAL LEARNING:  $\hat{\omega} \leftarrow \hat{\omega} - \eta E_n$

REGULARIZATION TECHNIQUE  $\underset{\omega}{\operatorname{argmin}} E_0(\omega) + \lambda E_{\omega}(\omega)$  w/ with  $\lambda$  regularization factor

SLIDES 8

• LINEAR MODELS GRAM-MATRIX

consider a linear model  $y(x, \omega) = \omega^T x$  with dataset  $D = \{(x_n, t_n)\}_{n=1}^N$

$$\textcircled{1} \quad E(\omega) = \frac{1}{2} \sum_{m=1}^N (\omega^T x_m - t_m)^2 + \lambda \omega^T \omega$$

$$\textcircled{2} \quad E(\omega) = (t - x\omega)^T (t - x\omega) + \lambda \|\omega\|^2$$

optimal solution  $\nabla E(\omega) = 0 \quad \textcircled{1}$

$$\omega^* = -\frac{1}{\lambda} \sum_{m=1}^N (\omega^T x_m - t_m) x_m \quad \alpha_m = -\frac{1}{\lambda} (\omega^T x_m - t_m)$$

$$\omega^* = \sum_{n=1}^N \alpha_n x_n$$

optimal solution  $\nabla E(\omega) = 0 \quad \textcircled{2}$

$$\omega^* = (x^T x + \lambda I_N)^{-1} x^T t = x^T (x x^T + \lambda I_N)^{-1} t$$

$$\alpha = (x x^T + \lambda I_N)^{-1} t$$

$$\text{then } y(\omega, x) = \sum_{m=1}^N \alpha_m x_m^T x \quad \alpha = (K + \lambda I_N)^{-1} t$$

$$\bullet \text{ GRAM-MATRIX} \rightarrow K = \begin{bmatrix} x_1^T x_1 & \dots & x_1^T x_N \\ \vdots & \ddots & \vdots \\ x_N^T x_1 & \dots & x_N^T x_N \end{bmatrix}$$

• KERNEL FUNCTION

similarity function  $K: X \times X \rightarrow \mathbb{R}$

come prima ma con  $K$

### KERNEL TRICK

replace inner product  $x^T x'$  with  $K(x, x')$

### KERNEL FUNCTION

a real valued function  $K(x, x') \in \mathbb{R}$  for  $x, x' \in X$  where  $X$  is some abstract space, usually:

• SYMMETRIC

• NON NEGATIVE

## NORMALIZATION:

- min MAX  $\rightarrow \bar{x} = \frac{x - \text{min}}{\text{MAX} - \text{min}}$
- standardization  $\rightarrow \bar{x} = \frac{x - \mu}{\sigma}$

## KERNEL FAMILIES:

- LINEAR  $k(x, x') = x^T x'$
- POLYNOMIAL  $k(x, x') = (\beta x^T x' + \gamma)^d \quad d \in \{2, 3, \dots\}$
- RADIAL BASIS FUNCTION (RBF)  $k(x, x') = e^{-\beta \|x - x'\|^2}$
- SIGMOID  $k(x, x') = \tanh(\beta x^T x' + \gamma)$

SVM kernelized  $\rightarrow$  classification

$$y(x; \alpha) = \text{sign}\left(w_0 + \sum_{n=1}^N \alpha_n x_n^T x\right) \rightarrow \text{sign}\left(w_0 + \sum_{n=1}^N \alpha_n k(x_n, x)\right)$$

## REGRESSION

$$E(w) = \sum_{n=1}^N E(y_n, t_n) + \lambda \|w\|^2$$

$$y(x; w^*) = \sum_{n=1}^N \alpha_n k(x_n, x) \quad \alpha = (k + \lambda I_N)^{-1} t$$

## SLIDES 10

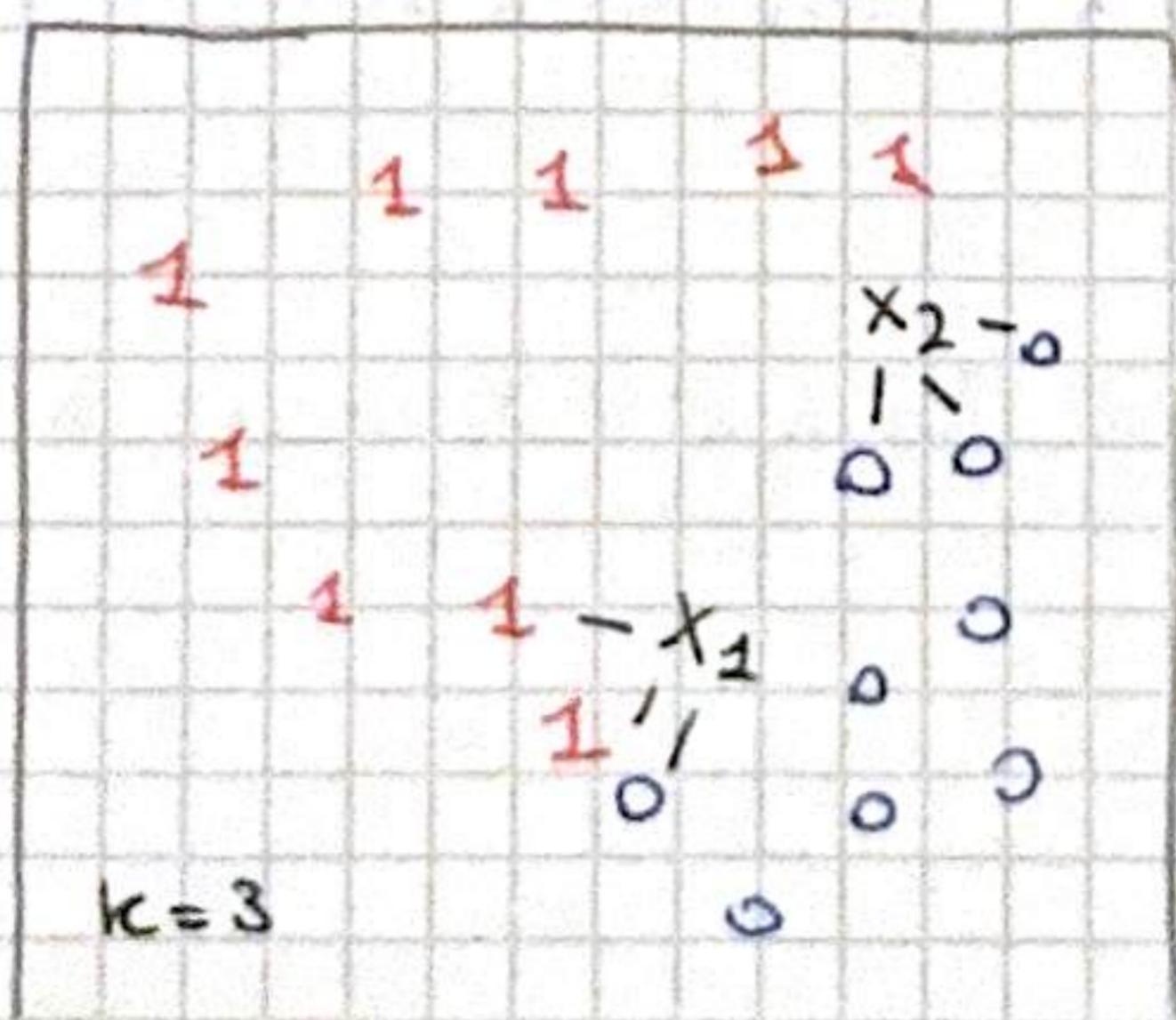
### K-NEAREST NEIGHBORS

- find k nearest neighbors of new instance x
- Assign to x the most common label among the majority of neighbors

LIKELIHOOD OF CLASS C FOR NEW INSTANCES X

$$P(C|x, D, k) = \frac{1}{k} \sum_{x_n \in N_k(x_m, D)} \mathbb{I}(t_n = c)$$

$N_k(x_m, D)$  the k nearest points to  $x_m$  and  $\mathbb{I}(e) = \begin{cases} 1 & \text{if } e = \text{true} \\ 0 & \text{otherwise} \end{cases}$



## SLIDES 11

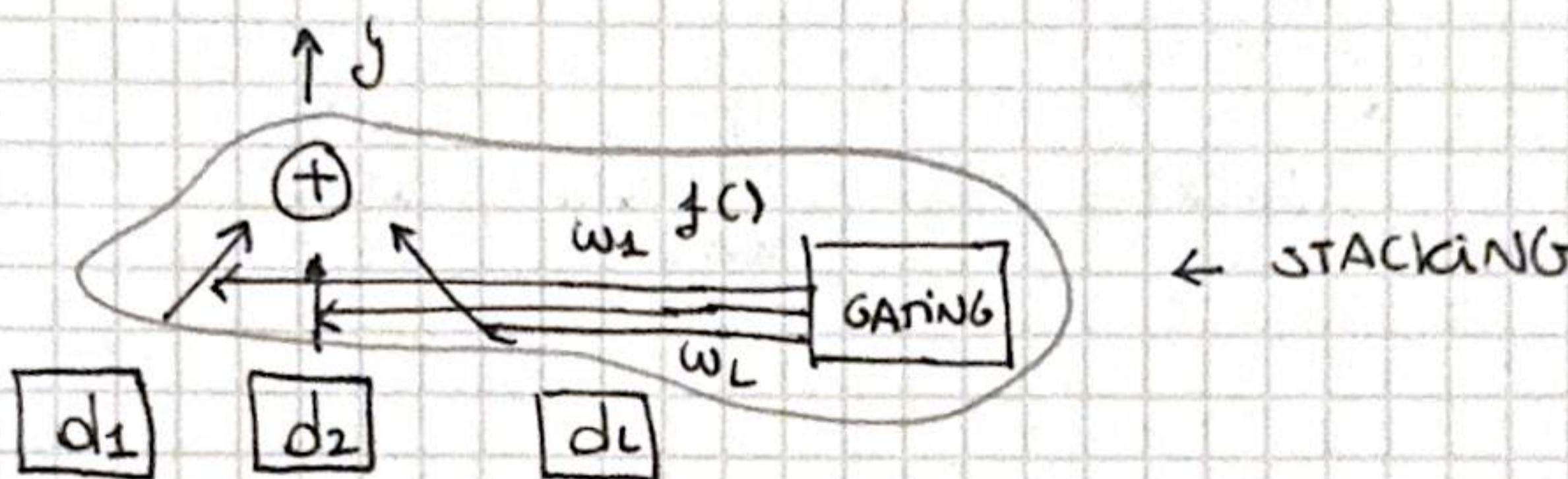
### VOTING:

1. dataset D to train models  $y_m(x)$

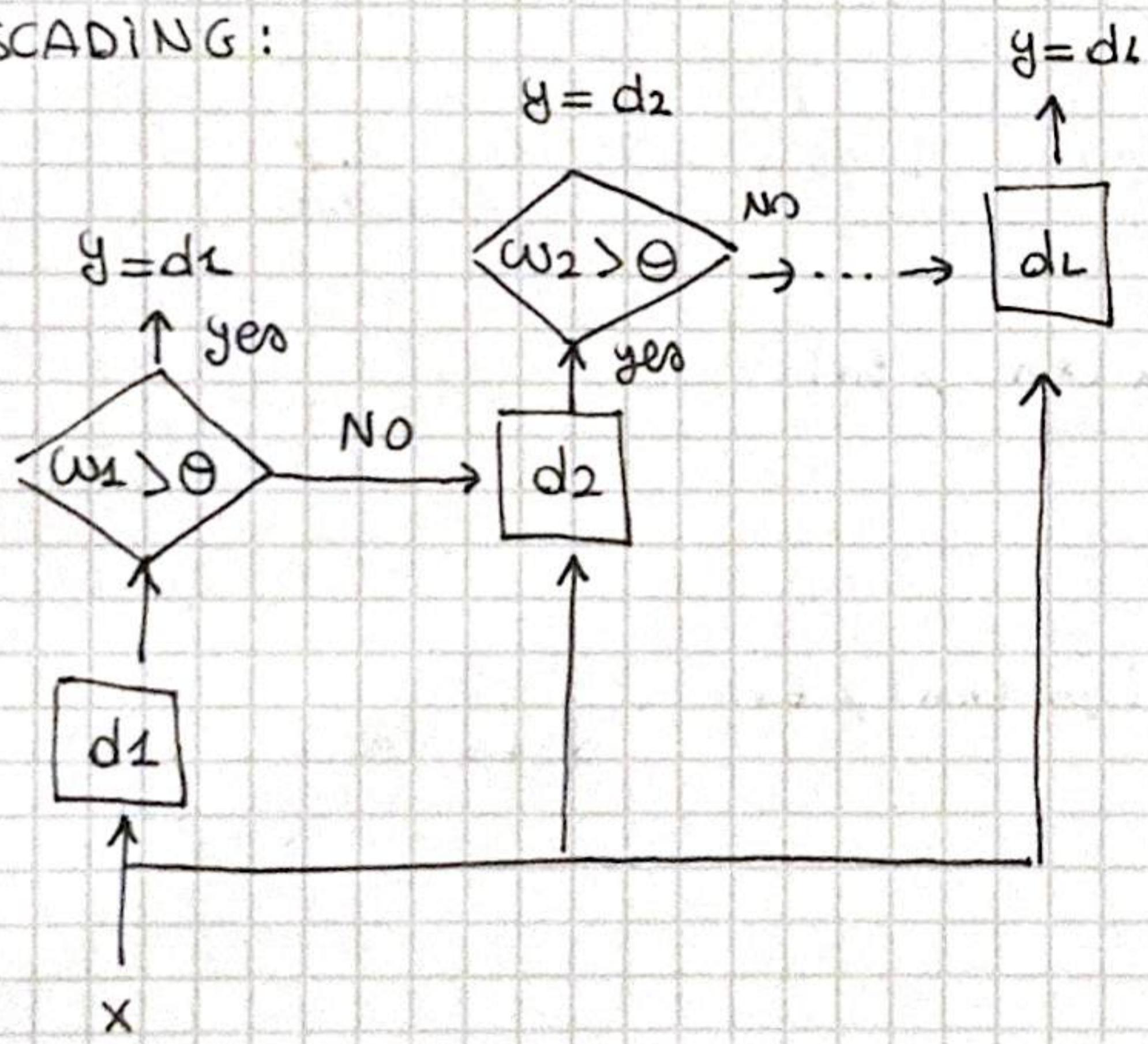
2. make predictions with:

$$y_{\text{VOTING}}(x) = \sum_{m=1}^M w_m y_m(x) \quad \text{REGRESSION}$$

$$y_{\text{VOTING}}(x) = \operatorname{argmax}_c \sum_{m=1}^M w_m I(y_m(x) = c) \quad \text{weighted majority (classification)}$$



### CASCADING:



### BAGGING

1. generate M bootstraps datasets  $D_1, \dots, D_M$

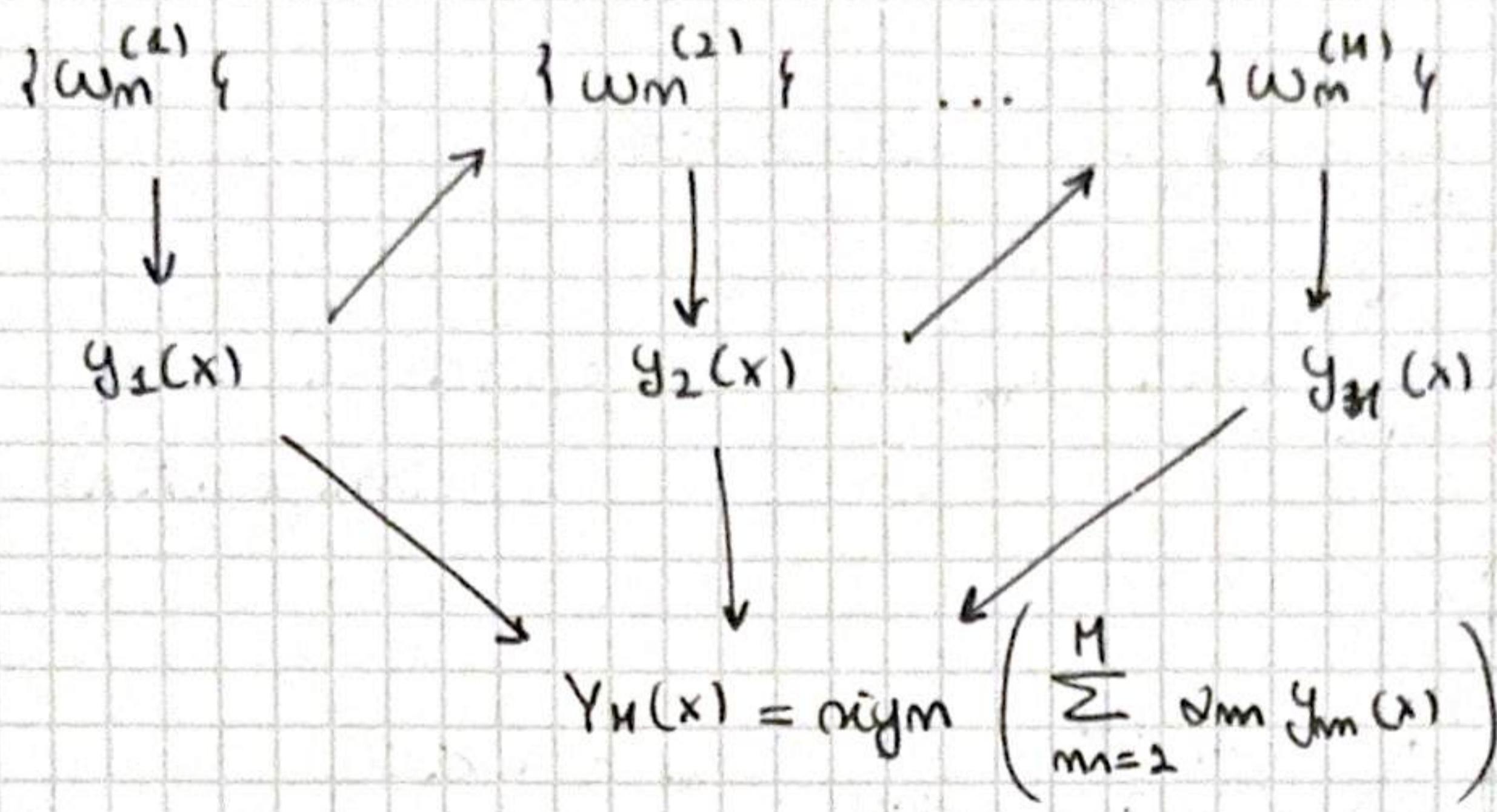
2. use them to train  $y_m(x)$  model

3. make prediction with voting scheme

$$y_{\text{bagging}}(x) = \frac{1}{M} \sum_{m=1}^M y_m(x)$$

## BAGGING:

bootstrapping train the first model and use those weights for train the next model



## ADA BOOST

1. train a weak learner  $y_m(x)$

$$\delta_m = \sum_{m=1}^N w_m^{(m)} I(y_m(x_m) \neq t_n) \quad I(c) = \begin{cases} 1 & \text{if } c = \text{true} \\ 0 & \text{otherwise} \end{cases}$$

$$2. \text{ evaluate } E_m = \frac{\sum_{i=1}^N w_i^{(m)} I(y_m(x_m) \neq t_n)}{\sum_{i=1}^N w_i^{(m)}}$$

$$3. \text{ update } w_m^{(m+1)} = w_m^{(m)} e^{(\alpha_m I(y_m(x_m) \neq t_n))} \quad \alpha_m = \ln \left( \frac{1 - E_m}{E_m} \right)$$

$$4. \quad Y_M(x) = \sum_{m=1}^M \alpha_m y_m(x)$$

## SLIDES 12

Why:

linear models cannot interaction between input variables

Kernell method → require suitable methods

FNN → complex combination of many parametric functions

## REGRESSION

linear units:  $y = W^T h + b$

Gaussian distribution noise model:  $P(t|x) = N(t|\mu, \sigma^2)$

LOSS FUNCTION → MEAN SQUARED ERROR

OUTPUT UNITS:  $y = \delta(w^T h + b)$

Loss function:  $J(\theta) = \text{E}_{\text{BINARY}} [-\ln(p(t|x))]$

} binary classification

OUTPUTS UNITS:  $y_i = \text{softmax}(\alpha^{(i)}) = \frac{e^{\alpha_i}}{\sum_j e^{\alpha_j}}$

Loss function: CATEGORICAL CROSS ENTROPY

$$\hookrightarrow J(\theta) = \text{E}_{\text{CROSS}} [-\ln(\text{softmax}(\alpha^{(i)}))]$$

CHAIN RULE:

$$g: \mathbb{R}^m \rightarrow \mathbb{R}^m \text{ and } f: \mathbb{R}^m \rightarrow \mathbb{R} \quad \frac{\partial z}{\partial x_i} = \sum_j \frac{\partial z}{\partial y_j} \frac{\partial y_j}{\partial x_i}$$

FORWARD STEP

R: network step  $e$

R:  $w^{(i)}$ ,  $i \in \{1, \dots, e\}$  weighted matrices

R:  $b^{(i)}$ ,  $i \in \{1, \dots, e\}$  bias parameter

R:  $x$  input value

R:  $t$  target value

$$h^{(0)} = x$$

for  $k = 1, \dots, e$  DO

$$\alpha^{(k)} = b^{(k)} + w^{(k)} h^{(k-1)}$$

$$h^{(k)} = f(\alpha^{(k)})$$

end for

$$y = h^{(e)}$$

$$J = L(t, y)$$

BACKWARD STEP

$$g \leftarrow \nabla_y J = \nabla_y L(t, y)$$

for  $k = e, e-1, \dots, 1$  DO

propagate gradients to the pre- nonlinearity activations

$$g \leftarrow \nabla_{\alpha^{(k)}} J = g \odot f'(\alpha^{(k)}) \quad (\odot \text{ element wise product})$$

$$\nabla_{b^{(k)}} J = g$$

$$\nabla_{w^{(k)}} J = g (h^{(k-1)})^T$$

propagate gradients to the next lower-level hidden layer

$$g \leftarrow \nabla_{h^{(k-1)}} J = (W^{(k)})^T g$$

end for

## STOCHASTIC GRADIENT DESCENT

$\eta$ : learning rate  $\eta \geq 0$

$\theta$ : initial value of  $\theta^{(1)}$

$k \leftarrow 1$

while stopping criterion not met do

SAMPLE ~~DATA~~ SUBSET (mini batch)  $\{x_1^{(i)}, \dots, x_m^{(i)}\}$  of  $m$  examples from  $D$

compute gradient estimate:  $\hat{g} = \frac{1}{m} \nabla_{\theta} \sum_i L(f(x_i^{(i)}; \theta^{(k)}), t_i^{(i)})$

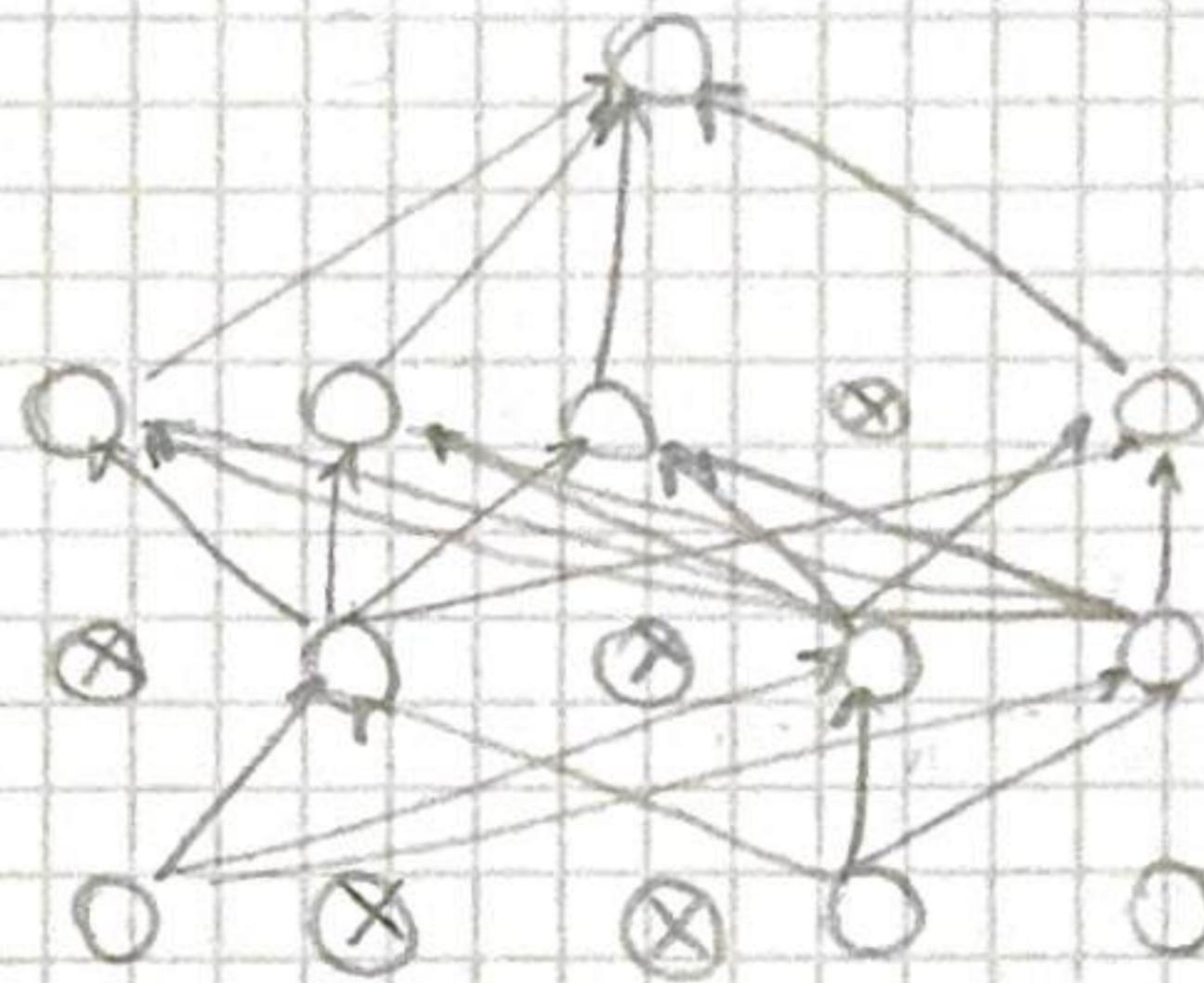
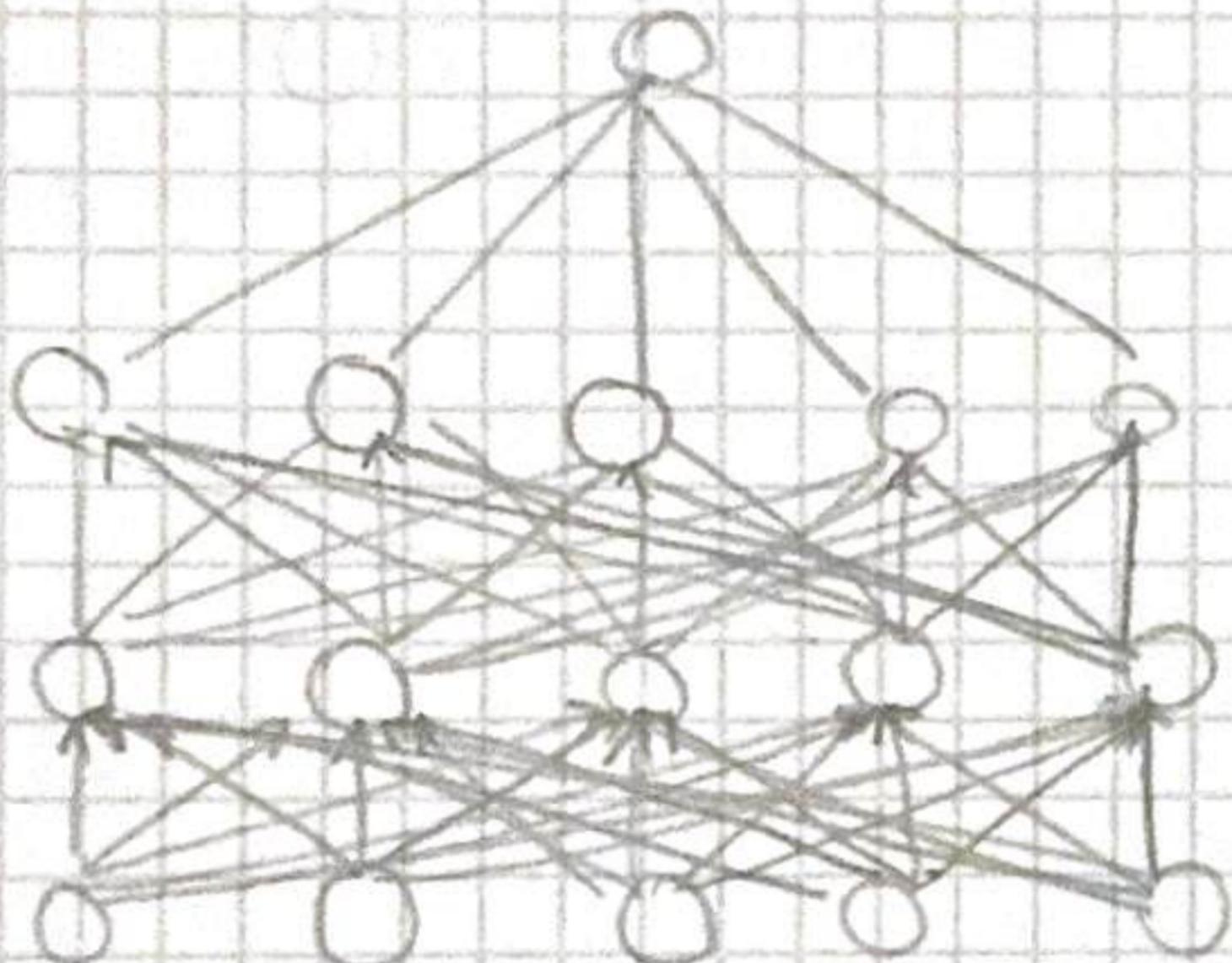
apply update  $\theta^{(k+1)} \leftarrow \theta^{(k)} - \eta \hat{g}$

$k \leftarrow k + 1$

backprop

end while

DROP-OUT



## SLIDES 13

CONV

CONVOLUTION:  $(x * w)(t) = \int_{a=-\infty}^{\infty} x(a)w(t-a) d(a)$  continuous

$(x * w)(t) = \sum_{a=-\infty}^{\infty} x(a)w(t-a)$  discrete

2D convolution  $(I * k)(i, j) = \sum_{m \in S_1} \sum_{n \in S_2} I(m, n) k(i-m, j-n)$

cross correlation  $(I * k)(i, j) = \sum_{m \in S_1} \sum_{n \in S_2} I(m+i, n+j) k(m, n)$

GENERAL

input:  $w_{in} \times h_{in} \times d_{in}$

kernel:  $d_{out}$  of size  $w_k \times h_k \times d_k$   $d_k = d_{in}$

output:  $w_{out} \times h_{out} \times d_{out}$

## NUMBERS OF PARAMETERS

consider input of size  $W_{in} \times h_{in} \times d_{in}$ , filter kernel of size  $W_{fc} \times h_{fc} \times d_{in}$

dimensions of output feature map:

$$W_{out} = \frac{W_{in} - W_{fc} + 2P}{S} + 1$$

$$h_{out} = \frac{h_{in} - h_{fc} + 2P}{S} + 1$$

number of trainable parameters:

$$|\theta| = \underbrace{W_{fc} \cdot h_{fc} \cdot d_{in} \cdot d_{out}}_{\text{kernels weights}} + \underbrace{d_{out}}_{\text{bias}}$$

## SLIDES 14

INPUT DATA available  $D = d \times n_f$

$$\text{GAUSSIAN MIXTURE MODEL (GMM)}: P(x) = \sum_{k=1}^K \pi_k N(x, \mu_k, \Sigma)$$

where  $\pi_k$ : prior probability

$\Sigma$ : covariance matrix

$\mu_k$ : Mean

## K-MEANS

STEP. 1  $\rightarrow$  decision value  $k = \text{number of cluster}$

STEP. 2  $\rightarrow$  put initial partition that classifies the data into  $k$  clusters

1. take training samples as single-element cluster

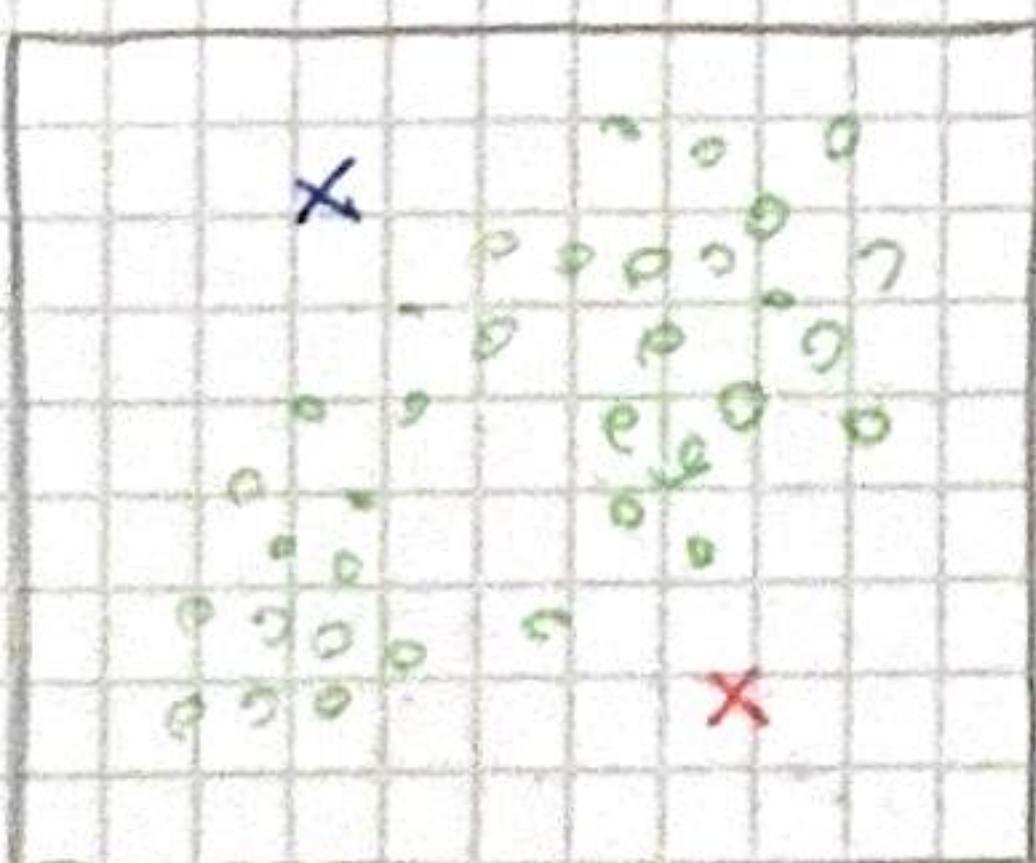
2. assign  $N-k$  to the cluster with nearest centroid

after each assignment recompute centroid

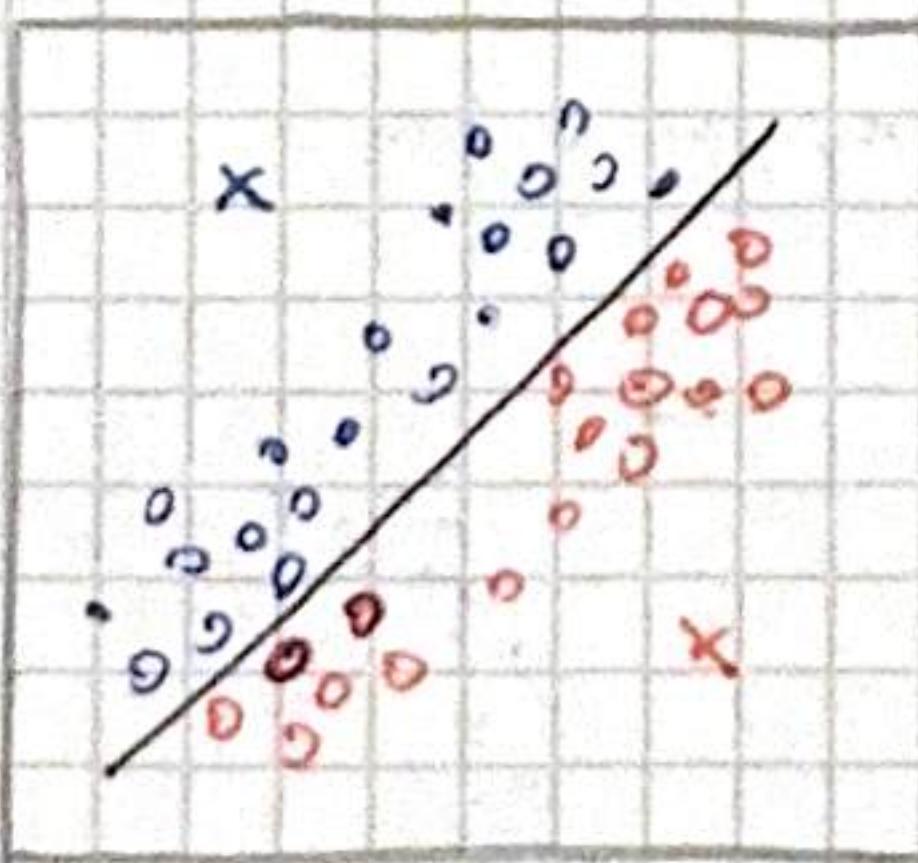
STEP 3  $\rightarrow$  take each sample in sequence and compute its distance from centroid of each cluster

IF NOT cluster  $\rightarrow$  switch with that cluster and update centroid

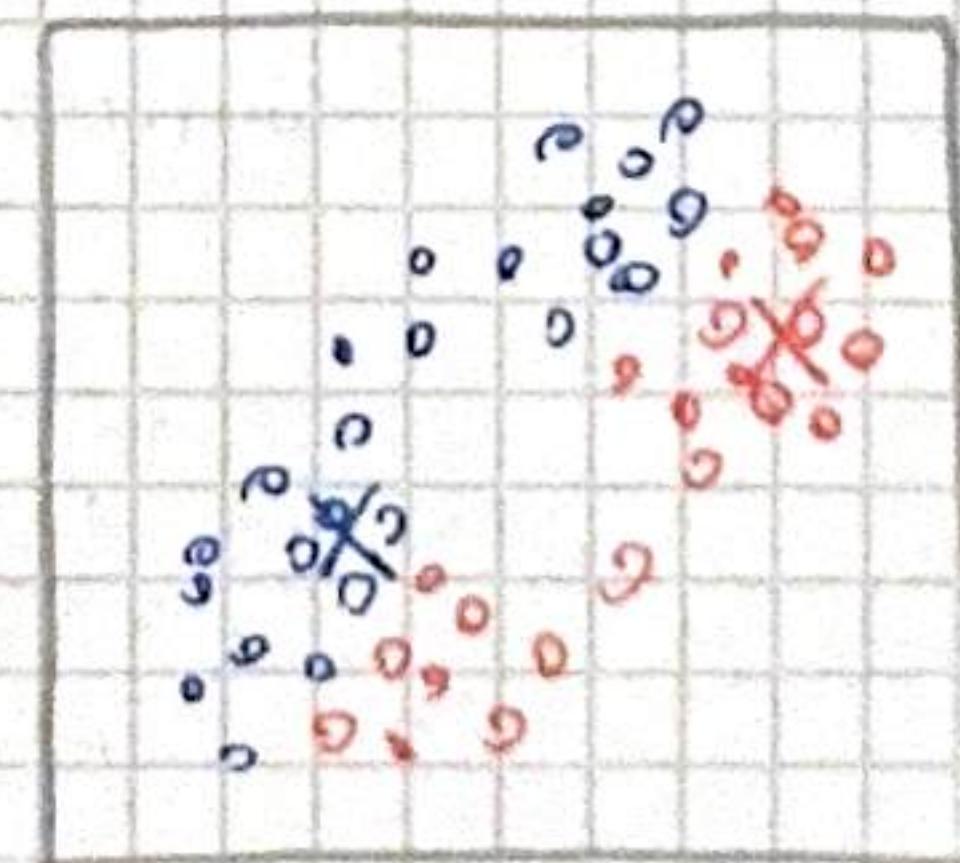
STEP 4 repeat until convergence



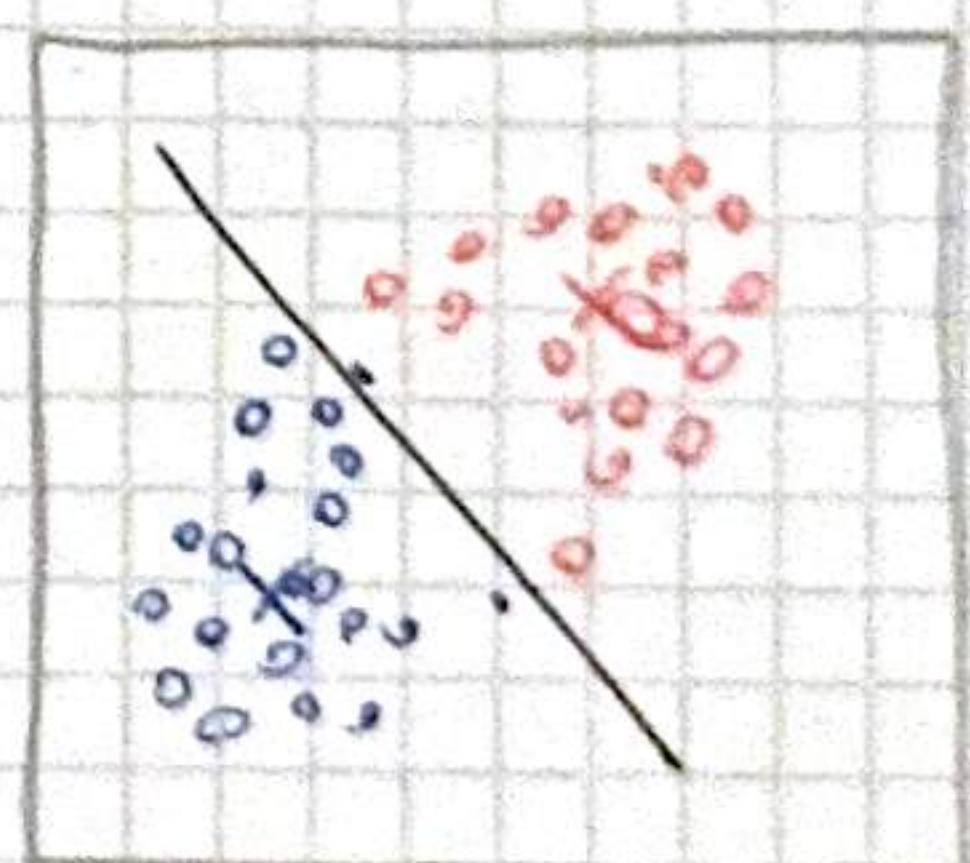
INITIALIZATION



CREATE CLUSTER GIVEN THE MEAN



COMPUTE THE MEAN GIVEN THE CLUSTER



REPEAT UNTIL CONVERGENCE

## LATENT VARIABLE

introduce  $z_k \in \{0, 1\}$  with  $z = (z_1, \dots, z_K)^T$

$$P(z_k=1) = \pi_k \text{ thus } P(z) = \prod_{k=1}^K \pi_k^{z_k}$$

$$P(x|z) = \prod_{k=1}^K N(x; \mu_k, \Sigma_k)^{\pi_k}$$

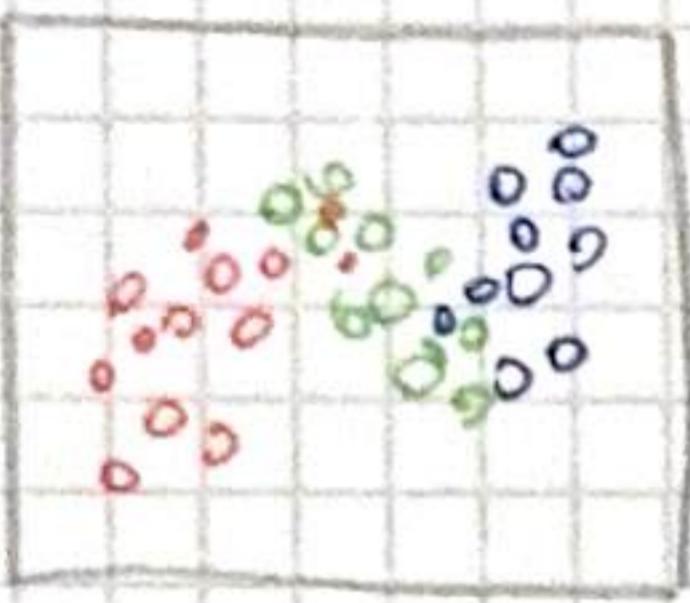
$$P(x) = \sum_z P(z) P(x|z) = \sum_{k=1}^K \pi_k N(x; \mu_k, \Sigma_k)$$

$$\gamma(z_k) = \frac{\pi_k N(x; \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j N(x; \mu_j, \Sigma_j)}$$

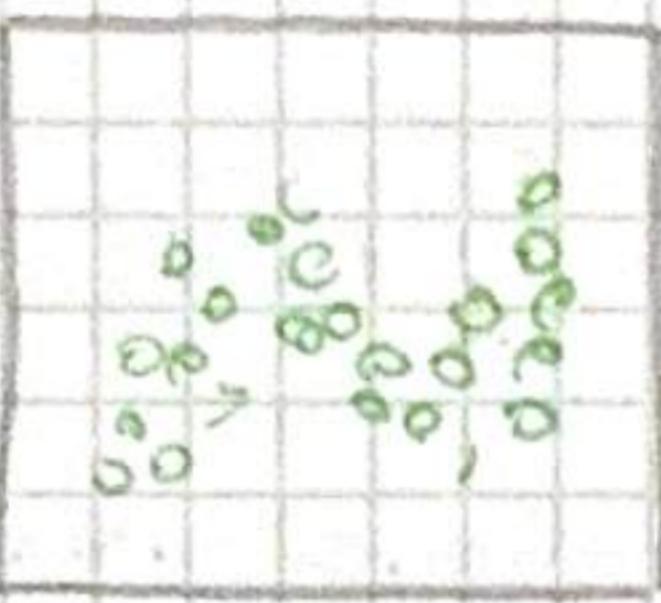
$\pi_k$ : prior probability

$\gamma(z_k)$ : posterior probability

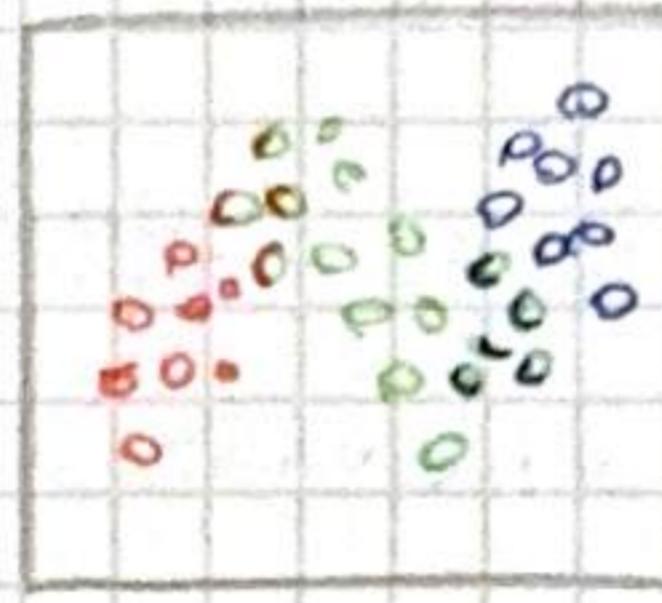
after observation of  $x$



$P(x, z)$  with 3 latent variable



$P(x)$  marginalized distribution



$\gamma(z_n, k)$  posterior distribution

## EXPECTATION MAXIMIZATION (EM)

MAXIMUM LIKELIHOOD:  $\underset{\pi, \mu, \Sigma}{\operatorname{argmax}} \ln(P(x|\pi, \mu, \Sigma))$

$$\mu_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) x_n$$

$$\pi_k = \frac{N_k}{N} \text{ with } N_k = \sum_{n=1}^N \gamma(z_{nk})$$

$$\Sigma_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) (x_n - \mu_k) (x_n - \mu_k)^T$$

### E STEP

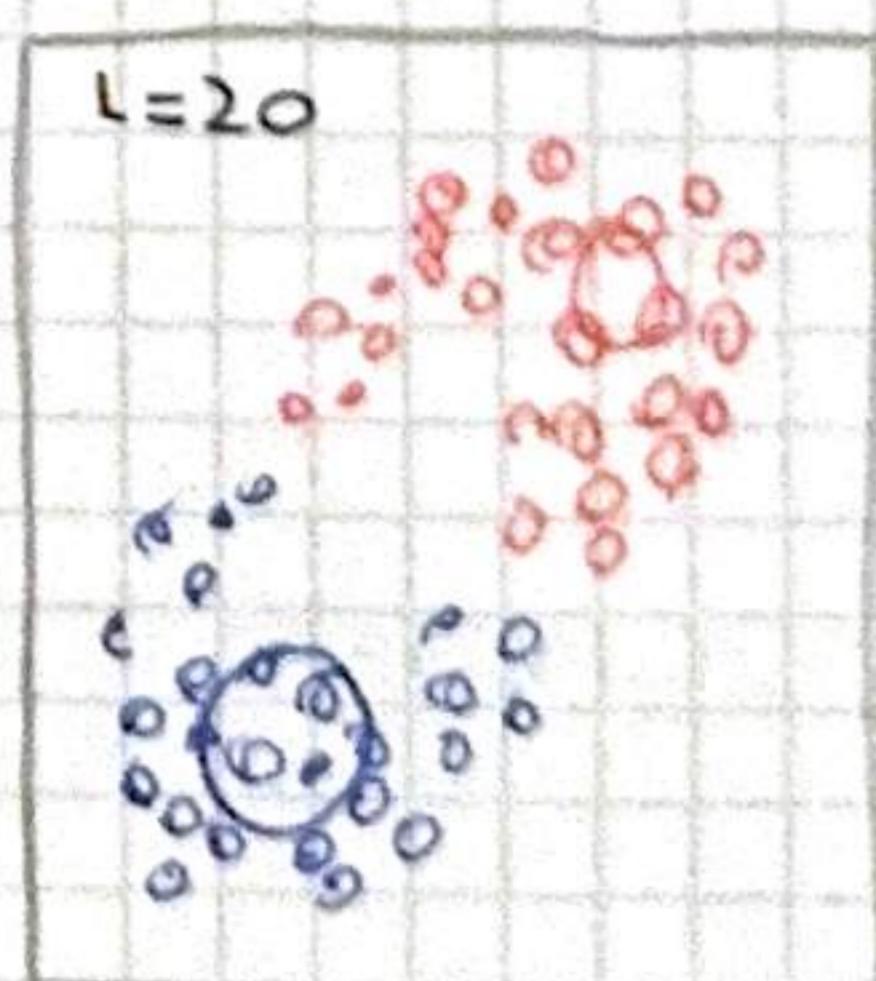
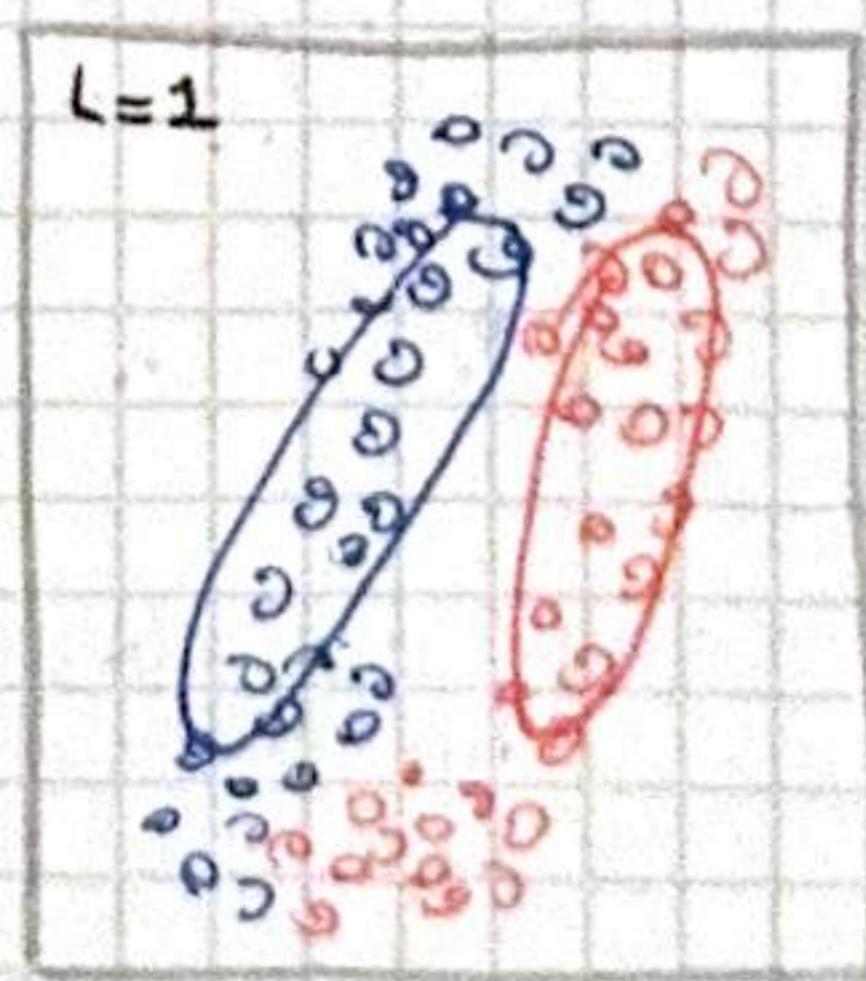
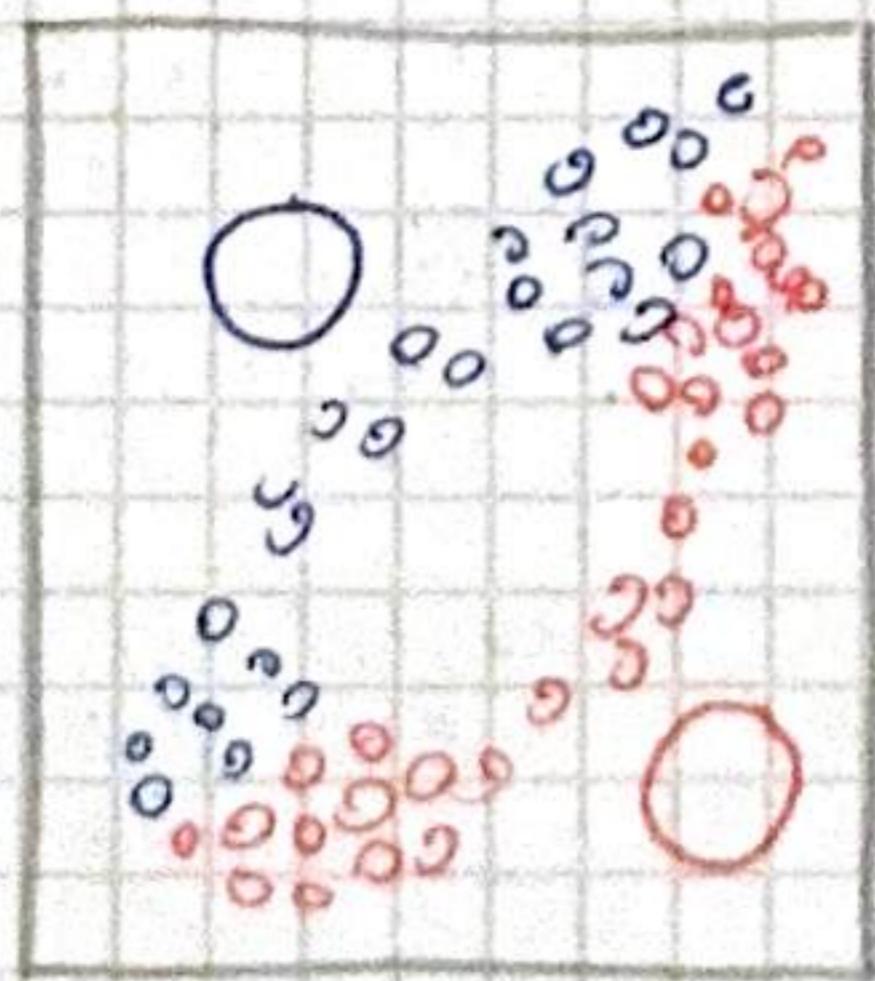
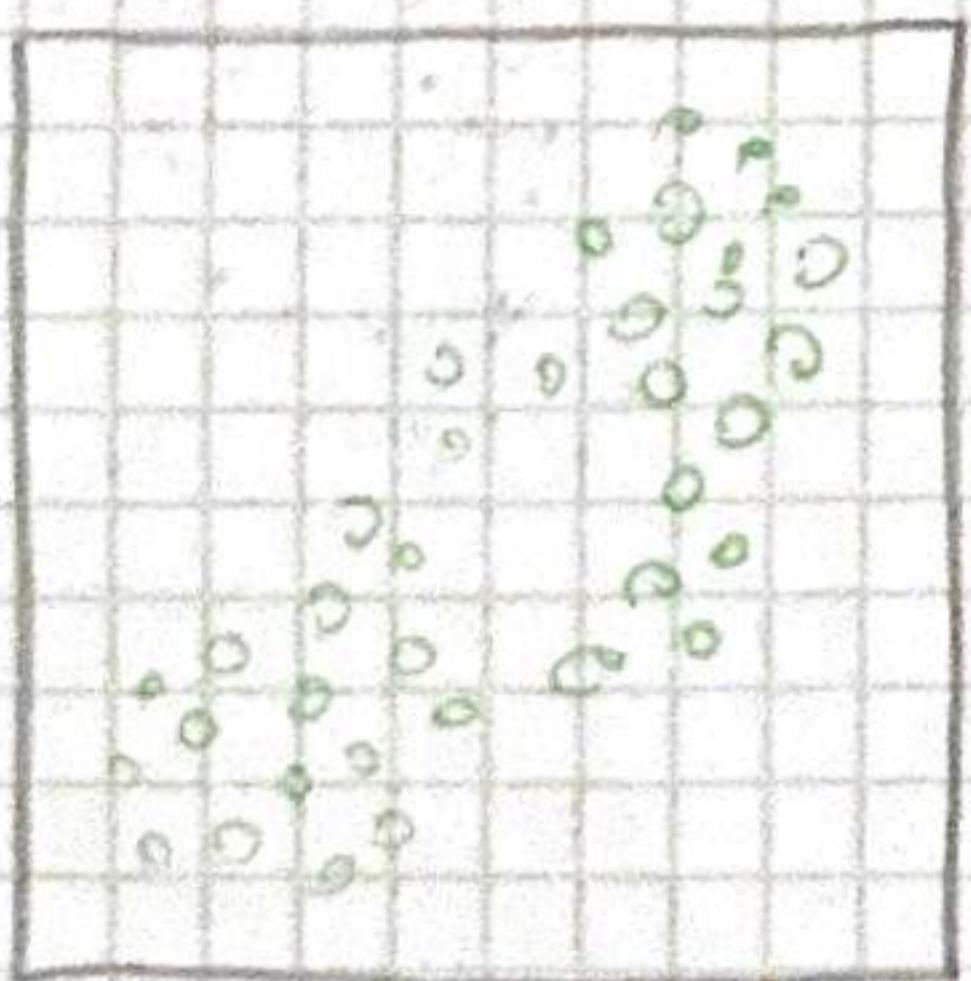
given  $\pi_k, \mu_k, \Sigma_k$

compute  $\gamma(z_{nk})$

### M STEP

given  $\gamma(z_{nk})$

compute  $\pi_k, \mu_k, \Sigma_k$



$L=1$

$L=20$

## SLIDES 15

MEAN OF PROJECTED POINTS:  $\bar{x}^T u_1$

VARIANCE OF PROJECTED POINTS:  $\frac{1}{N} \sum_{m=1}^N [\mu_1^T x_m - \bar{x}^T \bar{x}]^2 = \mu_1^T S \mu_1$

with  $S$  ( $d \times d$ ) covariance matrix of the dataset

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

### PCA VARIANCE MAXIMIZATION

$$\max_{u_1} u_1^T S u_1$$

equivalent to unconstrained maximization with LAGRANGE MULTIPLIER

$$\max_{u_1} u_1^T S u_1 + \lambda_1 (1 - u_1^T u_1)$$

solution:

$$S u_1 = \lambda_1 u_1 \quad u_1 \text{ must be an eigenvectors of } S$$

$$\text{left multiplying } u_1^T \text{ we have } u_1^T S u_1 = \lambda_1$$

variance is maximal when  $u_1$  is the eigenvector corresponding to the largest eigenvalue  $\lambda_1$  (called principal component)

### PCA ERROR MINIMIZATION

$$x_n = \sum_{i=1}^d \alpha_{ni} u_i \quad \text{obt} \quad \alpha_{ni} = x_n^T u_i \quad x_n = \sum_{i=1}^d (x_n^T u_i) u_i$$

goal: approximate  $x_n$  using a lower dimensional representation

$$\tilde{x} = \underbrace{\sum_{i=1}^m z_{ni} u_i}_{\text{HIGHEST EIGENVALUE}} + \underbrace{\sum_{i=m+1}^d b_i u_i}_{\text{LOWEST EIGENVALUE}}$$

evaluate Betrag des MSE

$$J = \frac{1}{N} \sum_{n=1}^N \|x_n - \tilde{x}_n\|^2$$

mimimizing  $z_{ni}$  and  $b_i$

$$z_{ni} = x_n^T u_i \quad b_i = \tilde{x}^T u_i$$

using this representation:

$$x_n - \tilde{x}_n = \sum_{i=m+1}^d [(x_n - \bar{x})^T u_i] u_i$$

the overall approximation:

$$J = \frac{1}{N} \sum_{m=1}^N \sum_{i=m+1}^d (x_m^\top u_i - \bar{x}^\top u_i)^2 = \sum_{i=m+1}^d \bar{u}_i^\top S u_i$$

$$\tilde{J} = \sum_{i=m+1}^d \bar{u}_i^\top S u_i + \lambda_i (1 - \bar{u}_i^\top u_i)$$

setting derivative of  $u_i = 0 \Rightarrow S u_i = \lambda_i u_i$

approximation error  $J = \frac{1}{N} \sum_{i=m+1}^d \lambda_i$

HIGH DIMENSIONAL DATA:

- centered data matrix  $X$
- compute the  $N-1$  eigenvalue  $\lambda_i$  and eigenvectors  $v_i$  of  $X^\top X$
- let  $u_i = \frac{1}{\sqrt{N\lambda_i}} X^\top v_i$

PROBABILISTIC PCA

$$P(x) = \int P(x|z) P(z) dz = N(x, \mu, C) \quad C = WW^\top + \sigma^2 I$$

posterior distribution

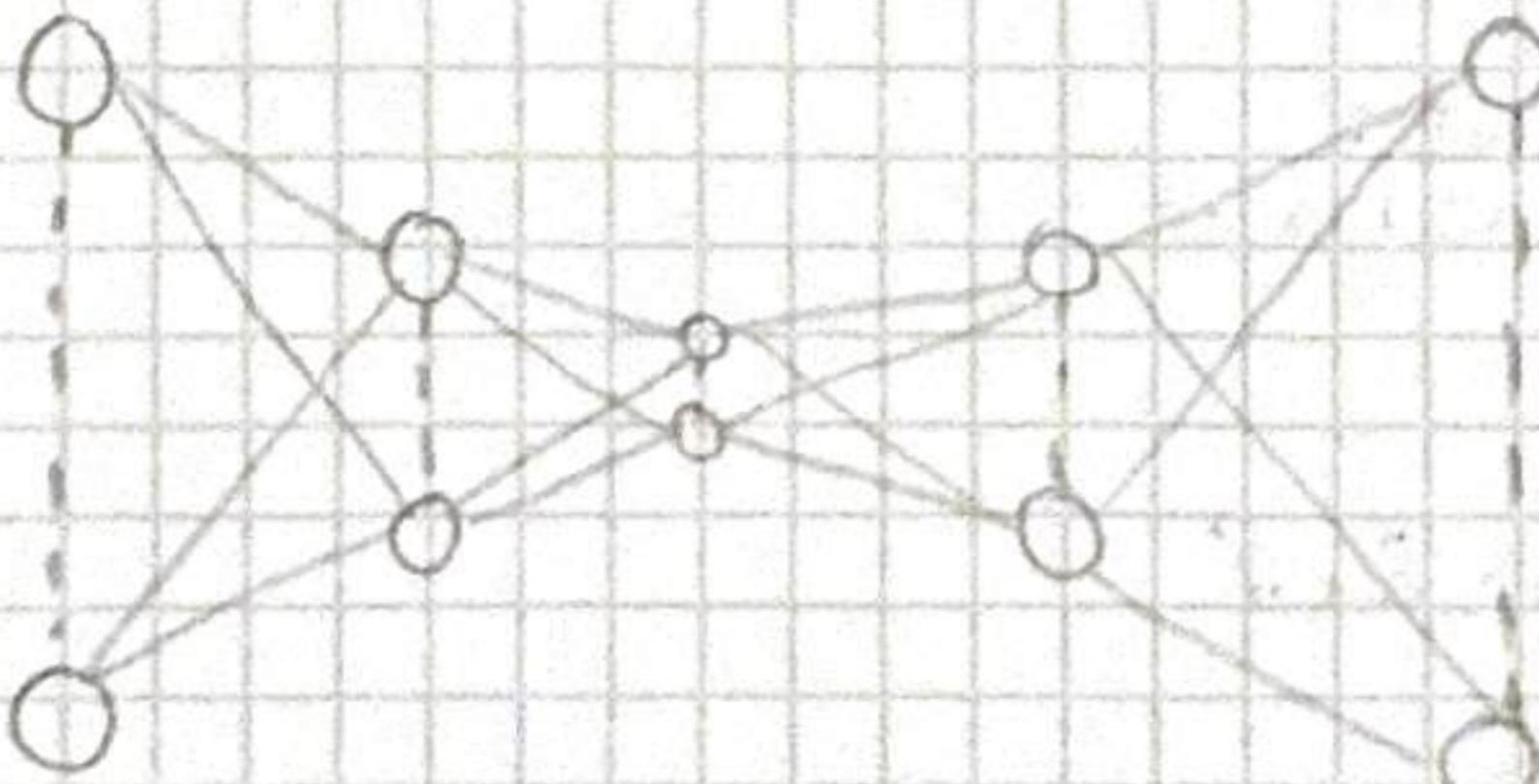
$$P(z|x) = N(z, M^{-1}W^\top(x - \mu), \sigma^2 M) \quad M = W^\top W + \sigma^2 I$$

maximum likelihood:

$$\underset{W, \mu, \sigma}{\operatorname{argmax}} \ln (P(x|W, \mu, \sigma^2)) = \sum_{n=1}^N \ln (P(x_n|W, \mu, \sigma^2))$$

AUTODECODERS

reduce size of hidden layers which learn to reconstruct their input by minimizing loss-function



ANOMALY DETECTION:

LEARN  $f: x \rightarrow \{m, a\}$  with dataset  $D = \{(x_n, m)\}$  (labels are irrelevant so we can consider an UNSUPERVISED DATASET)

GIVEN TEST  $x \notin D$  predict  $\hat{f}(x) \in \{m, a\}$  normal or abnormal

SOLUTION:

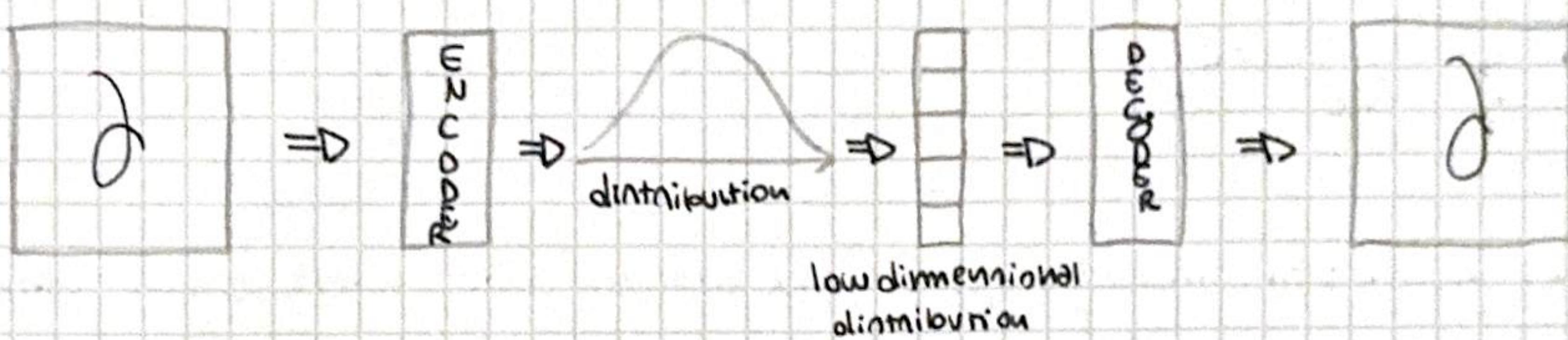
1. train AE with  $\{x_n\}$  compute final train loss
2. determine threshold  $f = \text{mean}(\text{loss}) + \text{std}(\text{loss})$
3. given  $x' \in D$  reconstruct  $x'$  with AE and compute  $\text{loss}'$
4. if  $\text{loss}' < f \rightarrow \text{normal}$  otherwise  $\text{abnormal}$

## GENERATIVE MODELS

- Variational AutoEncoders VAE  $\rightarrow$  focus on learning latent space structure
  - modify data in a specific direction
  - identify meaningful directions in latent space

ENCODER  $\rightarrow$  produces distribution instead of vector

DECODER  $\rightarrow$  operate on sample on this distribution



PROBLEM  $\rightarrow$  sampling operation not differentiable

$\hookrightarrow$  SOLUTION  $\rightarrow$  RE-PARAMETERIZATION

## ◦ GAN Generative Adversarial Network

- sample from the input data distribution  $X$

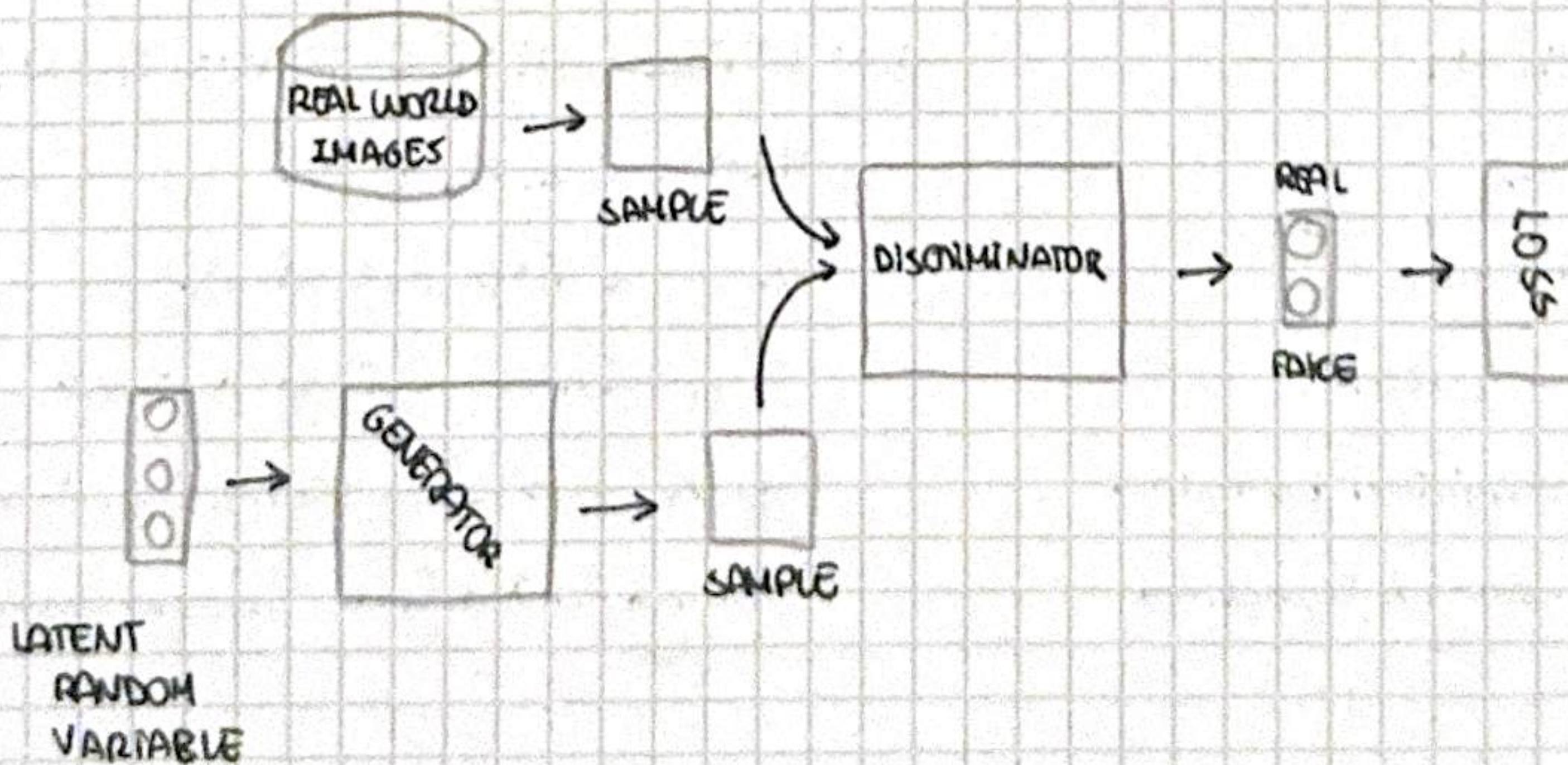
IDEA  $\rightarrow$  invert convolutional neural network via adversarial training

### ENCODER

- process an image and produce a vector

### DECODER

- receives a code and produces an image

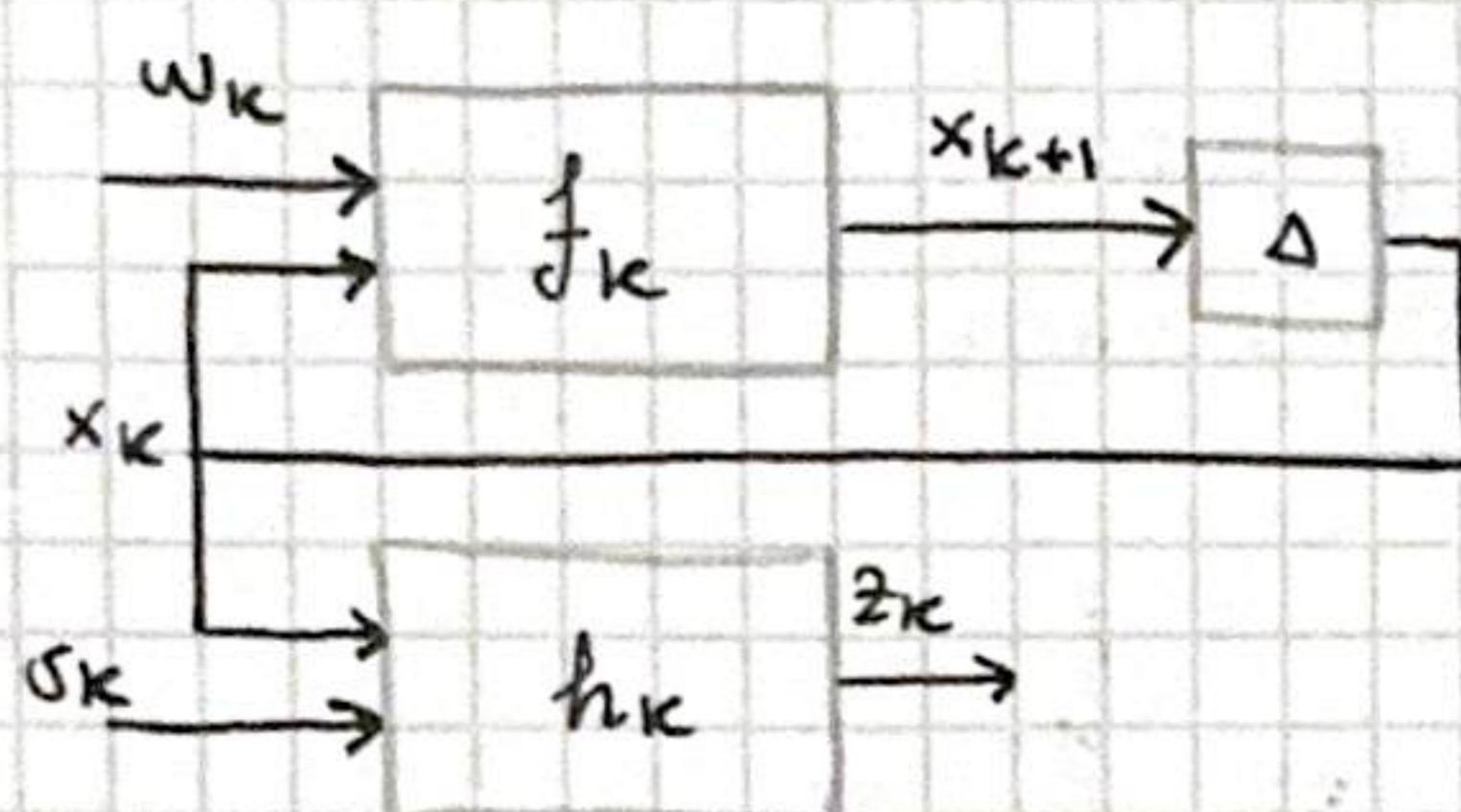


## TRAIN:

- train the discriminator with a batch of data  $\{(x_n, \text{Real}) \cup (x'_n, \text{Fake})\}$  where  $x_n$  comes from the dataset while  $x'_n$  is generated by GENERATOR with random value of the latent variable
- train the GENERATOR by unifying the entire model (generator + discriminator) with fixed layers (not trainable)  $\{z_k, \text{Real}\}$  where  $z_k$  are random values of LATENT VARIABLE

## SLIDES 16

### DYNAMIC SYSTEM



$x$ : state       $f$ : state transmission model  
 $z$ : observation       $h$ : observation model  
 $w, s$ : noise

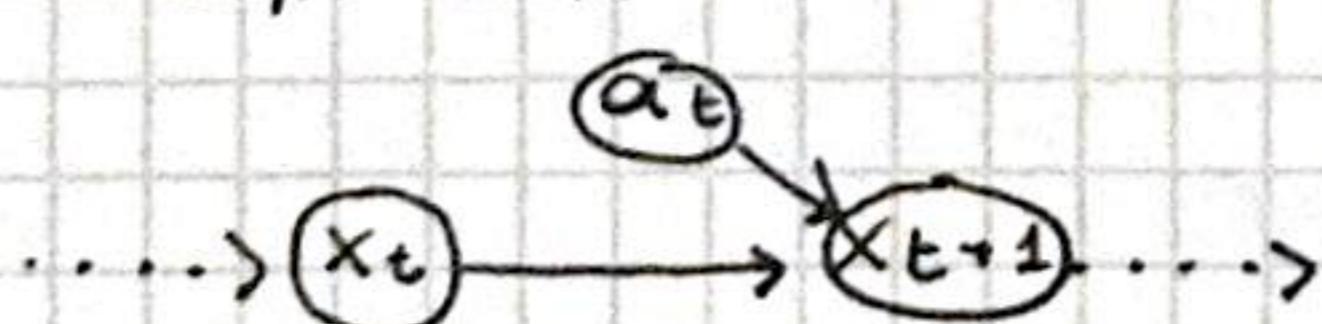
REASONING: given  $\theta$  the model ( $f, h$ ) and current state predict the future

LEARNING: given past experience  $(z_0, x)$  determine the model

STATE FULLY OBSERVABLE  $\rightarrow T: X \rightarrow A$  which action execute

### MARCOV PROPERTY:

- once the current state is known, the evolution of dynamic system doesn't depend on the history of states, actions and observations
- current state contains all information needed to predict the future
- future state are conditionally independent of past states and past observations given the current state
- the knowledge of the current state makes past, present and future observations statistically independent



### DETERMINISTIC TRANSITION:

$$\text{MDP} = \langle X, A, S, r \rangle$$

$X$  is finite set of states

$A$  is finite set of actions

$\delta: X \times A \rightarrow X$  is transition function

$r: X \times A \rightarrow \mathbb{R}$  is reward function

### NON DETERMINISTIC TRANSITION:

$$\text{MDP} = \langle X, A, \delta, r \rangle$$

$X$  is finite set of states

$A$  is finite set of actions

$\delta: X \times A \rightarrow 2^X$  is transition function

$r: X \times A \times X \rightarrow \mathbb{R}$  is reward function

## STOCHASTIC TRANSITIONS:

MDP  $\langle X, A, \delta, r \rangle$

$X$ : num of states finite

$A$ : finite set of actions

$P(x'|x, a)$  probability distributions over transitions

$r: X \times A \times X \rightarrow \mathbb{R}$  reward function

## OPTIMAL POLICY

$\pi: X \rightarrow A$  for each  $x \in X$   $\pi(x) \in A$  is the optimal action to be executed in such state

OPTIMALITY = MAXIMIZE the CUMULATIVE REWARD

EXPECTED VALUE OF THE CUMULATIVE REWARD

$\gamma$  = discount factor

$$V^\pi(x_1) = E[\bar{r}_1 + \gamma \bar{r}_2 + \gamma^2 \bar{r}_3 + \dots]$$

NON deterministic case

$$V^\pi(x_1) = \bar{r}_1 + \gamma \bar{r}_2 + \gamma^2 \bar{r}_3 + \dots \rightarrow \text{deterministic case}$$

OPTIMAL POLICY  $V^{\pi^*}(x_i) > V^\pi(x_i) \quad \forall i = \{1, \dots\}$

$\pi(a_i)$ : DETERMINISTIC KNOWN

OP:  $\underset{a_i \in A}{\text{argmax}} \pi(a_i) = \pi^*(x_0)$

$\pi(a_i)$ : DETERMINISTIC UNKNOWN

ALGORITHM:

1. for each  $a_i \in A$

- execute  $a_i$  and collect reward  $r_i$

$|A|$  iterations are needed

2. OP:  $\pi^*(x_0) = a_i$  with  $i = \underset{i=1 \dots |A|}{\text{argmax}} \pi(a_i)$

$\pi(a_i)$ : NON-DETERMINISTIC KNOW

OP:  $\pi^*(x_0) = \underset{a_i \in A}{\text{argmax}} E[\pi(a_i)]$

$\pi(a_i)$ : NON DETERMINISTIC UNKNOWN

ALGORITHM:

1. initialize structure  $\Theta$

3. optimal policy  $\pi^*(x_0) = \dots$  according

2. for each time  $t = 1, \dots, T$

with data structure

- choose an action  $a(t) \in A$

- execute  $a(t)$  and collect reward  $r(t)$

- update data structure  $\Theta$

$T \gg |A|$  iterations

## $\epsilon$ -GREEDY STRATEGY:

given  $0 \leq \epsilon \leq 1$

select an random action with probability  $\epsilon$

select best action with probability  $1-\epsilon$

$\epsilon$  decrease over time (from exploration THEN exploitation)

## SOFTMAX STRATEGIES:

actions with higher  $\hat{Q}$  values are assigned higher probabilities but every action is assigned a non-zero probabilities

$$P(a_i | k) = \frac{k^{\hat{Q}(x, a_i)}}{\sum_j k^{\hat{Q}(x, a_j)}}$$

$k > 0$  set how strongly

the selection favors actions with high  $\hat{Q}$

## EXAMPLE KARMEAD bandit

stochastic chance  $\pi(a_i) = N(\mu_i, \sigma_i)$

### $\epsilon$ -greedy strategies

training rule:

$$Q_n(a_i) \leftarrow Q_{n-1}(a_i) + \alpha [\bar{r} - Q_{n-1}(a_i)] \quad \alpha = \frac{1}{1 + V_{n-1}(a_i)}$$

$V_{n-1}(a_i)$  = number of execution

NOW consider also states  $X$  and state of evolution

$$\tilde{Q}^*(x, a) = r(x, a) + \gamma V^*(\delta(x, a))$$

$$Q(x, a) = r(x, a) + \gamma V^*(\delta(x, a))$$

$$\text{OP: } \pi^*(x) = \underset{a \in A}{\text{argmax}} Q(x, a)$$

## RECURSIVE FORMULA:

$$V^*(x) = \max_{a \in A} \{r(x, a) + \gamma V^*(\delta(x, a))\} = \max_{a \in A} Q(x, a)$$

$$Q(x, a) = r(x, a) + \gamma \max_{a' \in A} Q(\delta(x, a), a')$$

training rule

$$\hat{Q}(x, a) \leftarrow \bar{r} + \gamma \max_{a'} Q'(x'; a')$$

$\bar{r}$  is the immediate reward and  $x'$  is the state resulting from applying action  $a$  in  $x$

## NON DETERMINISTIC CASE

$$Q(x, a) = E[r(x, a) + \gamma V^*(\delta(x, a))]$$

$$Q(x, a) = E[r(x, a)] + \gamma \sum_{x'} P(x' | x, a) \max_{a'} Q(x', a')$$

training rule:

$$\hat{Q}(x, a) \leftarrow \hat{Q}_{n-1}(x, a) + \alpha [r + \gamma \max_{a'} (\hat{Q}_{n-1}(x', a') - \hat{Q}_{n-1}(x, a))]$$

$$\alpha = \alpha_{n-1}(x, a) = \frac{1}{1 + \text{visits}_{n-1}(x, a)}$$

TD TEMPORAL DIFFERENCE  $\rightarrow$  OFF POLICY

M-step time difference:

$$Q^{(n)}(x_t, a_t) = r_t + \gamma r_{t+1} + \dots + \gamma^{(n-1)} r_{t+n-1} + \gamma^n \max_a \hat{Q}(x_{t+n}, a)$$

using  $0 \leq \lambda \leq 1$

$$Q^\lambda(x_t, a_t) = (1-\lambda)[Q^{(1)}(x_t, a_t) + \lambda Q^{(2)}(x_t, a_t) + \lambda^2 Q^{(3)}(x_t, a_t) + \dots]$$

if  $\lambda=0$  Q-LEARNING

if  $\lambda > 0$  emphasis on discrepancy based on more distant look aheads

if  $\lambda=1$  only observed  $r_{t+i}$  are considered

SARSA  $\rightarrow$  ON POLICY

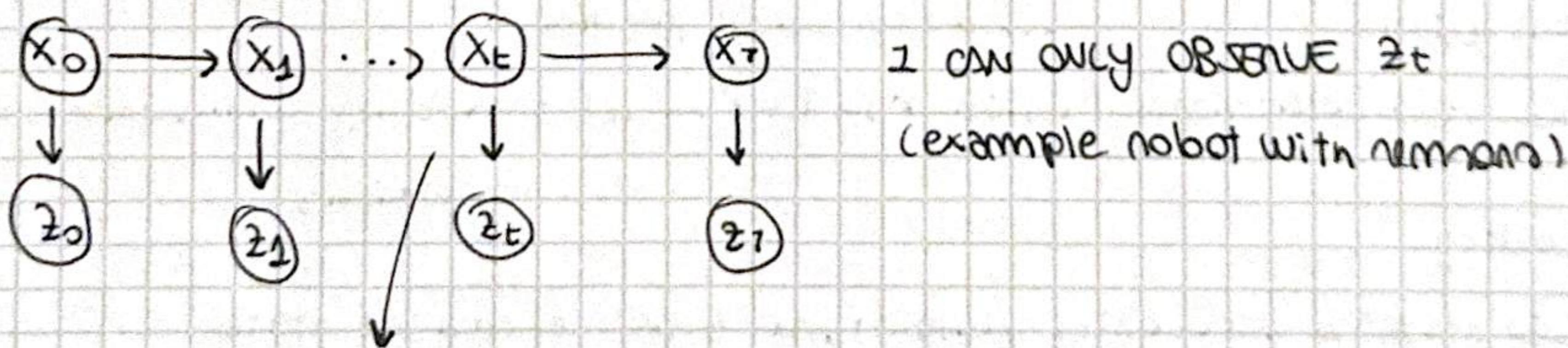
is based on tuple  $(s, a, r, s', a')$

$$\hat{Q}(x, a) = \hat{Q}_{n-1}(x, a) + \alpha [r + \gamma \hat{Q}_{n-1}(x', a') - \hat{Q}_{n-1}(x, a)]$$

$a'$  is chosen based on current policy

## SLIDES 17

STATE ~~OF~~ OF WORLD OF DYNAMIC SYSTEM IS NOT OBSERVABLE



$P(x_t = \underbrace{\dots}_{m} | x_{t-1} = \underbrace{\dots}_{m}) \rightarrow$  can be discrete way  $\rightarrow$  form a MATRIX

## HMM Hidden Markov Model

- FILTERING: problem of estimating the state of system at time  $t$  given all the past observations

CURRENT STATE

$$P(x_t = k | z_1:T) = \frac{\alpha_t^k}{\sum_j \alpha_t^j}$$

- SMOOTHING: given all the observations given so far we want estimate what was the value of any state in the past

$$P(x_t = k | z_1:T) = \frac{\alpha_t^k \beta_t^k}{\sum_j \alpha_t^j \beta_t^j}$$

## LEARNING in HMM

- states can be observed at training time

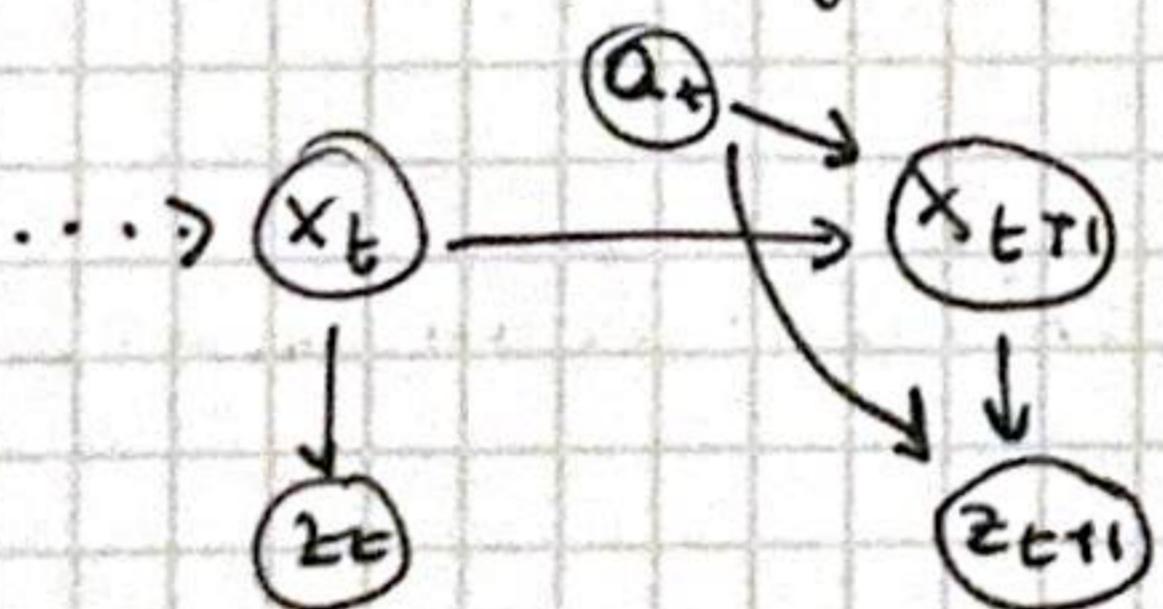
transitions and observations models can be estimated with statistical analysis

- states cannot be observed at training time

compute local maximum likelihood with EM

## POMDP Partially Observable Markov Decision Process

combining decision making of MDP with HMM



POMDP  $\langle X, A, S, R, Z, O \rangle$

$X$ : finite states

$\delta(x, a, x') = P(x' | x, a)$  is probability distribution over transitions

$A$ : finite actions

$r(x, a)$  is reward function

$Z$ : finite set of observations

$\delta(x', a, z') = P(z' | x', a)$  probability distribution over observations

we need policy

BELIEF STATES = probability distribution over the current state