

## CONSISTENT

An hypothesis  $h$  is consistent with a set of training data  $D$  of target concept  $c$  if and only if  $h(x) = c(x)$  for each training samples  $(x, c(x))$  of  $D$

$$\text{consistent}(h, D) \equiv \forall x \in D \ h(x) = c(x)$$

## OVERFITTING

Hypothesis  $h$  overfits if exist another hypothesis  $h'$  given a data sample  $S$  if

- $\text{error}_S(h) < \text{error}_S(h')$  and  $\text{error}_D(h) > \text{error}_D(h')$

## CLASSIFICATION METRICS

$$\text{PRECISION} = \frac{\text{TP}}{\text{TP} + \text{FP}} \quad \text{RECALL} = \frac{\text{TP}}{\text{TP} + \text{FN}} \quad \text{F1-score} = \frac{2(\text{PRECISION} \cdot \text{RECALL})}{\text{PRECISION} + \text{RECALL}}$$

$$\text{ACCURACY} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{FP} + \text{FN} + \text{TN}}$$

## REGRESSION METRICS

$$\text{MAE} = \frac{1}{N} \sum_{i=1}^N |\hat{f}(x_i) - t_i| \quad \text{MSE} = \frac{1}{N} \sum_{i=1}^N (\hat{f}(x_i) - t_i)^2 \quad \text{RMSE} = \sqrt{\text{MSE}}$$

$$\text{ENTROPY} \rightarrow \text{entropy}(S) = - \sum_{i=1}^n p_i \log_2(p_i)$$

$$\text{GAIN} \rightarrow \text{entropy}(S) - \sum_{v \in \text{values}(A)} \left( \frac{|S_v|}{|S|} \right) \text{entropy}(S_v)$$

## REDUCE OVERFITTING (decision tree)

→ 1. REDUCED ERROR PRUNING: using validation dataset if a node doesn't reduce validation error prune

2. RULE POST PRUNE (C4.5): if accuracy is greater or equal when we have prune we have overfitting

GAIN NATION  $\rightarrow$  GAIN/SPLIT INFORMATION

SPLIT INFORMATION  $\rightarrow - \sum_{i=1}^N |S_i| / |S| \log_2 (|S_i| / |S|)$

$$\text{Likelihood} \\ \downarrow \\ P(h|D) = \frac{\underset{n}{\prod} P(D|h) P(h)}{P(D)}$$

MAXIMUM A POSTERIORI  $h_{MAP} \rightarrow h_{MAP} = \operatorname{argmax}(P(h|D)) =$

$$= \operatorname{argmax}((P(D|h) P(h)) / P(D)) = \operatorname{argmax}(P(D|h) P(h))$$

MAXIMUM LIKELIHOOD  $h_{ML} \rightarrow$  if we assume  $P(h_i) = P(h_j) \rightarrow$   
 $\rightarrow h_{ML} = \operatorname{argmax}(P(D|h))$

EXAMPLE Usefull

$$h_1(x_1) = + \quad h_2(x_1) = - \quad h_3(x_1) = - \quad P(D|h_1) = 0.4 \quad P(D|h_2) = 0.3 \quad P(D|h_3) = 0.3$$

$$h_{ML} = \operatorname{argmax} \{0.4, 0.3, 0.3\} = 0.4 \rightarrow \text{class } +$$

$$P(+|x, h_1) = 1 \quad P(+|x, h_2) = 0 \quad P(+|x, h_3) = 0$$

$$P(-|x, h_1) = 0 \quad P(-|x, h_2) = 1 \quad P(-|x, h_3) = 1$$

$$V_{OC} = \operatorname{argmax} \{ (0.4 \cdot 1 + 0 + 0), (0 + 0.3 \cdot 1 + 0.3 \cdot 1) \} = \text{class } -$$

BAYESIAN OPTIMAL CLASSIFIER  $\rightarrow V_{OB} = \operatorname{argmax} \sum_{v_j \in V} \sum_{h_i \in H} P(v_j|x, h_i) P(h_i|D)$

## NAIVE BAYSIAN CLASSIFIER

①  $x$  conditionally independent of  $y$  given  $z$

$$p(x, y | z) = p(x|y, z)p(y|z) = p(x|z)p(y|z)$$

assume a target function  $f: X \rightarrow V$  where each instance of  $x$  can be described  $\langle a_1, \dots, a_n \rangle$

$$\text{argmax}_{v_j \in V} p(v_j | x, D) = \text{argmax}_{v_j \in V} p(v_j | a_1 \dots a_n, D)$$

② NAIVE BAYES ASSUMPTION

$$p(a_1, \dots, a_n | v_j, D) = \prod_{i=1}^N p(a_i | v_i, D)$$

$$\text{VNB} = \text{argmax}_{v_j \in V} p(v_j | D) \prod_{i=1}^N p(a_i | v_i, D)$$

## MULTINOMIAL NAIVE BAYES DISTRIBUTION

for ordinal feature vector  $d = \langle d_1, \dots, d_n \rangle$  of generic doc  $\in D$

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

maximum likelihood solution

$$\hat{p}(w_i | c_j, D) = \frac{\sum_{\text{doc} \in D} t_{fij} + \alpha}{\sum_{\text{doc} \in D} t_{fj} + \alpha |V|}$$

## GENERAL MODELS

$$P(c_1 | x) = \frac{P(x | c_1) P(c_1)}{P(x)} = \frac{P(x | c_1) P(c_1)}{P(x | c_1) P(c_1) + P(x | c_2) P(c_2)} = \frac{1}{1 - \exp(\alpha)}$$

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

consider now parametric model:

$$a = \ln \left( \frac{P(x|C_1)P(C_1)}{P(x|C_2)P(C_2)} \right) = \ln \left( \frac{N(x; \mu_1; \Sigma)}{N(x; \mu_2; \Sigma)} \right) = \dots = \omega^T x + \omega_0$$

$$\omega_0 = \Sigma^{-1}(\mu_1 - \mu_2) \text{ and } \omega = -\frac{1}{2}\mu_1^T \Sigma^{-1} \mu_2 + \frac{1}{2}\mu_2^T \Sigma^{-1} \mu_2 + \ln \left( \frac{P(C_2)}{P(C_1)} \right)$$

## LOGISTIC REGRESSION

given a dataset  $D = \{(x_n, t_n)\}_{n=1}^N$

LIKELIHOOD FUNCTION:  $P(t| \tilde{\omega}) = \prod_{n=1}^N y_n^{t_n} (1-y_n)^{1-t_n}$

CROSS ENTROPY ERROR FUNCTION (negative log LIKELIHOOD):

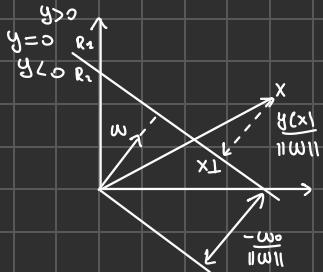
$$E(\tilde{\omega}) = -\log(P(t|\tilde{\omega})) = -\sum_{n=1}^N t_n \ln(y_n) + (1-t_n) \ln(1-y_n)$$

$$\tilde{\omega}^* = \operatorname{argmin} (E(\tilde{\omega}))$$

## LINEARLY SEPARABLE DATASET

a dataset is linearly separable if exist a surface that splits data in two regions such that different classified instances are separated

## LINEAR DISCRIMINANT FUNCTIONS



## LEAST SQUARES

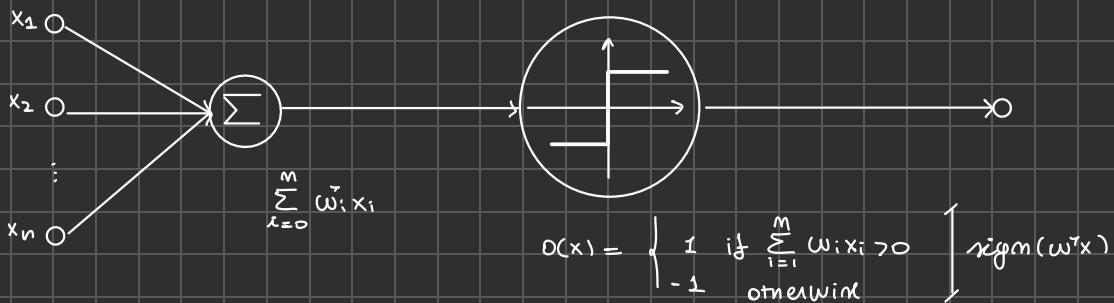
given a dataset  $D = \{(x_n, t_n)\}_{n=1}^N$  find linear discriminant  $y(x) = \tilde{\omega}^T x$

mimimite the sum of error functions:

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

with  $\tilde{\omega} = \underbrace{(\tilde{x}^T \tilde{x})^{-1}}_{\text{PSEUDO INVERSE}} \tilde{x}^T \tau$

### PERCEPTION



in error function we use SQUARED ERROR

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

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

### PERCEPTIONONE TRAIN RULE

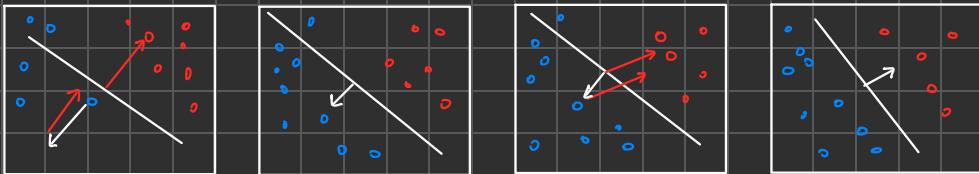
$$\omega_i \leftarrow \omega_i + \Delta \omega_i \quad \text{where} \quad \Delta \omega_i = \eta \frac{\partial E}{\partial \omega} = \eta \sum_{n=1}^N (t_n - \tilde{\omega}^T x_n) (x_{i,n})$$

BATCH MODE  $\rightarrow \Delta \omega_i = \eta \sum_{(x,t) \in D} (t - O(x)) x_i$  high computational cost

MINIBATCH MODE  $\rightarrow \Delta \omega_i = \eta \sum_{(x,t) \in S} (t - O(x)) x_i$  best one

INCREMENTAL MODE  $\rightarrow \Delta \omega_i = \eta (t - O(x)) x_i$  high variance

## PERCEPTION TRAIN RULE



## FIGUER DISCRIMINANT RULE

consider two classes can determine  $y = \omega^T x$  and classify  $x \in C_1$  if  $y \geq -w_0$   
 $x \in C_2$  otherwise

compute normalised the projection on a line determine by  $\omega$

we want adjoint  $\omega$  in order to find direction that maximise separation

$$m_1 = \frac{1}{N_1} \sum_{n \in C_1} x_n$$

$$m_2 = \frac{1}{N_2} \sum_{n \in C_2} x_n$$

choose  $\omega$  that maximise  $J(\omega) = \omega^T (m_2 - m_1)$  subject to  $\|\omega\| = 1$

FISHER CRITERION:

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

with  $S_B = (m_2 - m_1)^T (m_2 - m_1)$   $S_W = \sum_{n \in C_2} (x_n - m_1)(x_n - m_1)^T + \sum_{n \in C_1} (x_n - m_2)(x_n - m_2)^T$

we need to choose  $\omega$  that maximise  $J(\omega)$  so

$$\frac{d J(\omega)}{d \omega} = 0 \Rightarrow \omega^* \propto S_W^{-1} (m_2 - m_1)$$

MULTIPLE CLASSES  $J(\omega) = \text{Tr} \{ (\omega S_w \omega^*)^{-1} (\omega S_B \bar{w}) \}$

## SUPPORT VECTOR MACHINE SVM

SVM for classification aims at maximum margin providing for better accuracy

MARGIN (normalised distance from  $x_k$  to  $\bar{w}$ ) is defined  $|y(x_k)| / \|\omega\|$

$$\text{argmax} [ \min (|y(x_k)| / \|\omega\|) ] = \text{argmax} [ \frac{1}{\|\omega\|} \min |y(x_k)| ] =$$

$$= \text{argmax} \left[ \frac{1}{\|w\|_1} \min (\text{tr}(w^T x + w_0)) \right] \text{ assuming } |y(x_k)| = \text{tr}(y(x_k))$$

so we want find

$$w^*, w_0^* = \text{argmax} \left[ \frac{1}{\|w\|_1} \min (\text{tr}(w^T x + w_0)) \right]$$

for solve this optimization problems we can use LAGRAGIAN MULTIPLIERS  
and for better estimation averaging over all support vector SV

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

final we have:

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

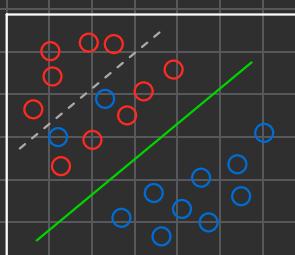
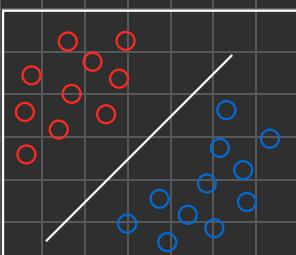
$$y(x, w) = \min \left( \sum_{x_k \in SV} \alpha_k^* t_k x_k^T x + w_0^* \right)$$

### SOFT MARGIN CONSTRAINTS

- $\xi_n = 0$  if point is ON or INSIDE the correct margin boundary
- $0 < \xi_n \leq 1$  if point is INSIDE the margin but comet side
- $\xi_n > 1$  if point on wrong side of boundary

### SVM WITH SOFT MARGIN CONSTRAINTS

$$w^*, w_0^* = \text{argmax} \frac{1}{\|w\|^2} + C \sum_{n=1}^N \xi_n$$



## LINEAR REGRESSION

$$g(x, \omega) = \sum_{j=0}^M \omega_j \phi_j(x) \quad \text{is linear in } M$$

target value is effected by additive noise  $\epsilon \rightarrow t = g(x, \omega) + \epsilon$

assuming gaussian noise  $P(\epsilon | \beta) = N(\epsilon | 0, \beta^{-1})$  so we have

$$P(t | x, \omega, \beta) = N(t | g(x, \omega), \beta^{-1})$$

we need maximum likelihood function:

$$P(t_1, \dots, t_N | x_1, \dots, x_N, \beta) = -\frac{\beta}{2} \sum_{n=1}^N (t_n - \omega^\top \phi(x_n))^2 - \frac{N}{2} \ln(2\pi\beta^{-1})$$

conventional least square minimization

$$\text{argmin}_{\omega} \frac{1}{2} \sum_{n=1}^N (t_n - \omega^\top \phi(x_n))^2$$

REGULARIZATION TECHNIQUE:  $\text{argmin}_{\omega} E_\lambda(\omega) + \lambda E_\omega(\omega)$

## KERNEL MODELS

consider a linear model

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

$$E(\omega) = \frac{1}{2} (t - \omega^\top x)(t - \omega^\top x)^\top + \lambda \|\omega\|^2$$

where  $x = [x_1, \dots, x_N]^\top$  is design matrix  $t = [t_1, \dots, t_N]$  output vector

the optimal solution is given by solving  $\nabla E(\omega) = 0$

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

We can write  $\alpha = (X^T X + \gamma I_N)^{-1}$  and we will have  $y(x, w^*) = \sum_{n=1}^N \alpha_n x_n^T x$

$$K = x^T x \text{ in gram matrix} = \begin{bmatrix} x_1^T x_1 & \dots & x_1^T x_N \\ \vdots & \ddots & \vdots \\ x_N^T x_1 & \dots & x_N^T x_N \end{bmatrix}$$

KERNEL TRICK if in algorithm appear an INNER PRODUCT  $x^T x$  we can substitute them with kernel function  $K(x, x')$

$$\text{gram matrix will become} = \begin{bmatrix} K(x_1, x_1) & \dots & K(x_1, x_N) \\ \vdots & \ddots & \vdots \\ K(x_N, x_1) & \dots & K(x_N, x_N) \end{bmatrix}$$

Kernel families: 1. linear  $K(x, x') = x^T x'$

$$2. \text{ polynomial } K(x, x') = (Bx^T x' + \gamma)^d$$

$$3. \text{ RBF } K(x, x') = \exp(-\beta \|x - x'\|^2)$$

$$4. \text{ sigmoid } K(x, x') = \tanh(\beta x^T x' + \gamma)$$

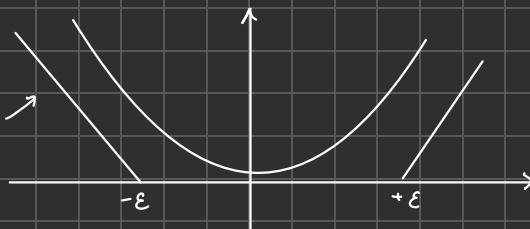
## KERNELIZED VERSION OF SVM

$$y(x, \alpha) = \operatorname{sign} \left( \sum_{n=1}^N \alpha_n y_n K(x_n, x) + w_0 \right)$$

## KERNELIZED VERSION SVM REGRESSION

$$E_D(w) = \frac{1}{2} \|w\|^2 + C \sum_{n=1}^N E_\epsilon(y_n - t_n)$$

doesn't penalize  
too much outliers



is not differentiable  $\rightarrow$  difficult to solve

## KNN

is an untrained based model. considering a classification problem  $f: X \rightarrow C$   
 where  $C = \{+, -\}$  and a dataset  $D = \{(x_n, t_n)\}_{n=1}^N$

classification with KNN:

- ① find  $k$  nearest neighbor of new instance  $x$
- ② assign to new instance  $x$  the most common label among the majority of neighbors

likelihood function

$$P(C|x, D, k) = \frac{1}{k} \sum_{x_n \in N_k(x, D)} \mathbb{I}(t_n = C) \quad \text{with } \mathbb{I}(e) = \begin{cases} 1 & e = \text{true} \\ 0 & \text{otherwise} \end{cases}$$

## MULTIPLE LEARNERS

VOTING:

given a dataset  $D$

1. use  $D$  for train a set of models  $y_m(x)$
2. make predictions

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

REGRESSION  
WEIGHTED MAJORITY

$$y_{\text{VOTING}}(x) = \arg \max \sum_{m=1}^M w_m \mathbb{I}(y_m(x) = C)$$

OF CLASS

BAGGING  $\rightarrow$  parallel methods

1. generate  $M$  bootstrap dataset from  $D$
2. train each model with  $M$  bootstrap dataset
3. make prediction using voting method

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

## ADABOOST

Given a dataset  $D = \{(x_n, t_n)\}_{n=1}^N$

1. initialize  $w_m = 1/N$ ,  $m = 1, \dots, M$

2. for  $m = 1, \dots, M$

train a weak learner by minimizing the weighted error function

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

$$\cdot \text{evaluate } \epsilon_m = \frac{\sum_{n=1}^N w_n \mathbb{I}(y_m(x_n) \neq t_n)}{\sum_{n=1}^M w_n} \quad \text{and} \quad \alpha_m = \ln \left( \frac{1 - \epsilon_m}{\epsilon_m} \right)$$

$$\cdot \text{update coefficients } w_n^{(m+1)} \leftarrow w_n^{(m)} \exp [\alpha_m \mathbb{I}(y_m(x_n) \neq t_n)]$$

3. make prediction

$$y_M = \operatorname{sign} \left( \sum_{m=1}^M \alpha_m y_m(x) \right)$$

## ANN

• REGRESSION  $\rightarrow$  doesn't saturate

LINEAR UNITS  $\rightarrow$  identity activation function  $y = w^T h + b$

use a Gaussian distribution model  $p(t|x) = N(t|y, \sigma^2)$

LOSS-FUNCTION  $\rightarrow$  MSE mean squared error  $\frac{1}{N} \sum_{n=1}^N (f(x_n) - t_n)^2$

• BINARY CLASSIFICATION  $\rightarrow$  SATURATE only correct answer

OUTPUTS UNITS  $\rightarrow$  sigmoid activation function  $y = \sigma(w^T h + b)$

LOSS-FUNCTIONS  $\rightarrow$  Binary cross entropy  $J(\theta) = -\frac{1}{N} \sum_{n=1}^N [t_n \ln(y_n) + (1-t_n) \ln(1-y_n)]$

• MULTICLASS CLASSIFICATION  $\rightarrow$  SATURATE minimal errors

OUTPUTS UNITS  $\rightarrow$  softmax activation function  $y_i = \text{softmax}(z_i) = \frac{\exp(z_i)}{\sum_j \exp(z_j)}$

LOSS-FUNCTION  $\rightarrow$  categorical cross entropy  $J(\theta) = -\frac{1}{N} \sum_{n=1}^N \sum_{i=1}^C t_{ni} \ln(y_{ni})$

## BACKPROPAGATION

• is not a learning algorithm • is only compute gradient

### FORWARD PASS

repune network depth  $l$

repune  $W^{(k)}$ ,  $i \in \{1, \dots, l\}$  weight matrices

repune  $b^{(k)}$ ,  $i \in \{1, \dots, l\}$  bias parameters

repune  $x$  input values

repune  $t$  target values

$h^{(0)} = x$

for  $k = 1, \dots, l$  do

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

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

end for

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

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

### BACKWARD PASS

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

for  $k = l, l-1, \dots, 1$  do

$$g \leftarrow \nabla_{\alpha^{(k)}} J = g \odot f'(\alpha^{(k)})$$

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

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

propogate gradient to the next lower layer

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

end for

## STOCHASTIC GRADIENT DESCENT

repune learning rate  $\eta > 0$

repune initial values of  $\theta^{(1)}$

$$k \leftarrow 1$$

while stop criterion not met do

sample a subset (minibatch) from  $D$

compute the gradient estimate  $g = \frac{1}{m} \nabla_{\theta} L(f(x^{(i)}, \theta^{(k)}), t^{(i)})$

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

$$k \leftarrow k+1$$

end while

## CNN

DIMENSIONS OF OUTPUT feature map:

$$w_{out} = \frac{w_{in} - w_k + 2p}{s} + 1 \quad h_{out} = \frac{h_{in} - h_k + 2p}{s} + 1$$

NUMBER OF TRAINABLE PARAMS:

$$|\theta| = w_k \cdot h_k \cdot \text{din} \cdot d_{out} + \underbrace{d_{out}}_{\text{BIAS}}$$

## GAUSSIAN MIXTURE MODELS

mixed probability distribution  $p$  formed by  $k$  different Gaussian distributions

$$p(x) = \sum_{k=1}^K \pi_k N(x, \mu_k, \Sigma_k)$$

$\pi_k$  = prior probability    $\mu_k$  = mean    $\Sigma_k$  = covariance matrix

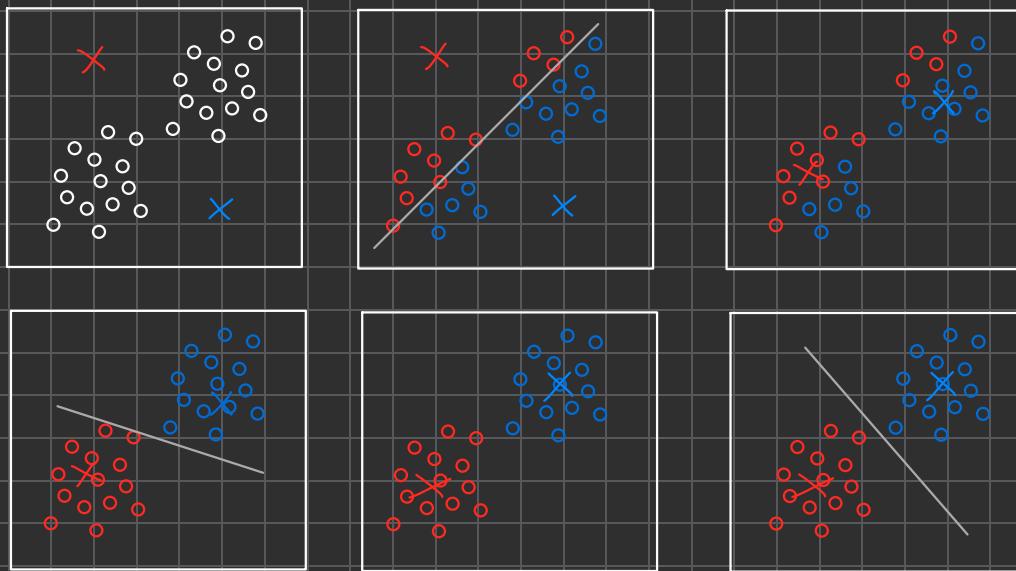
## K-MEANS

is an unsupervised learning algorithm

1. begin with a selection value  $k$  = number of clusters
2. put any partition that classifies the data into  $k$ -clusters. this can be done randomly or systematically
  - o consider  $k$  training samples as centre element of clusters
  - o assign to remain  $N-k$  training samples to the cluster with nearest centroid and after each assign recompute the centroid of clusters
3. take each sample in sequence and compute the distance from centroid of each cluster if sample is not in cluster with nearest centroid switch and update the centroid
4. repeat step 3 until convergence

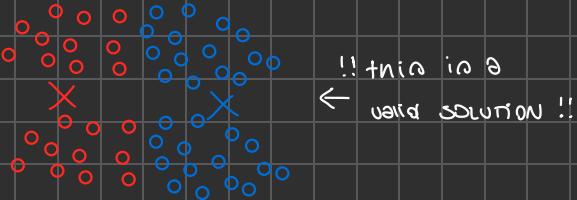
## CONVERGENCE

1. for each cluster in step 2 the sum of distances for each training sample to centroid decreased
2. there are only finitely many partitions of training samples into  $k$  clusters



### PROBLEMS with K-MEAN

- NOT robust to outliers
- SENSITIVE TO INITIAL CONDITION



### EXPECTATION MAXIMIZATION

$$\gamma(z_{nk}) = \frac{\pi_k N(x, \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j N(x, \mu, \Sigma)} \rightarrow \text{POSTERIOR}$$

we have two steps : 1. E-STEP  $\rightarrow$  given  $\pi_k, \mu_k$  and  $\Sigma_k$  compute  $\gamma(z_{nk})$   
 2. M-STEP  $\rightarrow$  given  $\gamma(z_{nk})$  compute  $\pi_k, \mu_k, \Sigma_k$

LATENT VARIABLE

### PCA Principal Component Analysis

is widely technique used for various task:

1. dimensionality reduction
2. data compression
3. feature extraction
4. data visualization

## VARIANCE MAXIMIZATION

goal maximize data variance after projection to norme direction  $u_1$   
 projected points  $\vec{x}_n^T u_1$

$$\text{mean value of data points } \bar{x} = \frac{1}{N} \sum_{n=1}^N x_n$$

$$\text{variance of projected points } \frac{1}{N} \sum_{n=1}^N [\vec{u}_1^T x_n - \vec{u}_1^T \bar{x}]^2 = \vec{u}_1^T S \vec{u}_1$$

$$\text{covariance matrix of dataset } S = \frac{1}{N} \sum_{n=1}^N (x_n - \bar{x})(x_n - \bar{x})^T = \frac{1}{N} \vec{x} \vec{x}^T$$

we want maximize the projected variance:  $\max_{\vec{u}_1} \vec{u}_1^T S \vec{u}_1$

is equivalent to unconstrained maximization with lagrange multiplier:

$$\max_{\vec{u}_1} \vec{u}_1^T S \vec{u}_1 + \lambda_1 (1 - \vec{u}_1^T \vec{u}_1)$$

solution derive w.r.t.  $\vec{u}_1 \Rightarrow \vec{u}_1^T S \vec{u}_1 = \lambda_1 \rightarrow$  variance is maximal when  $\vec{u}_1$  is proportional to the largest eigenvector

CALLED FIRST PRINCIPAL COMPONENT

## ERROR MINIMIZATION

GOAL: approximate  $x_n$  using a lower dimensional representation:

$$\hat{x}_n = \sum_{i=1}^m z_{ni} u_i + \sum_{i=m+1}^d b_i u_i$$

↑   ↑  
 HIGHEST EIGENVALUE    LOWEST EIGENVALUE

we will lose norme information for minimization we compute the error we commit doing this approximation

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

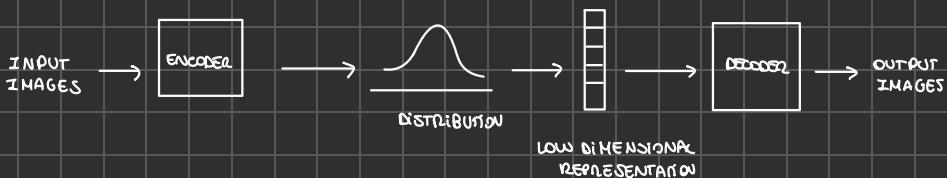
minimize the approximation error will become:  $\hat{S} = \sum_{i=m+1}^N \vec{u}_i^T S \vec{u}_i + \lambda_1 (1 - \vec{u}_1^T \vec{u}_1)$

Setting all weights of a  $W_i$  to zero we have  $S_{W_i} = \pi_{W_i}$

the approximation error is then given by  $J = \sum_{i=1}^N \pi_{W_i}$

VAE Variational Auto-Encoder → focusing on learning latent variable space

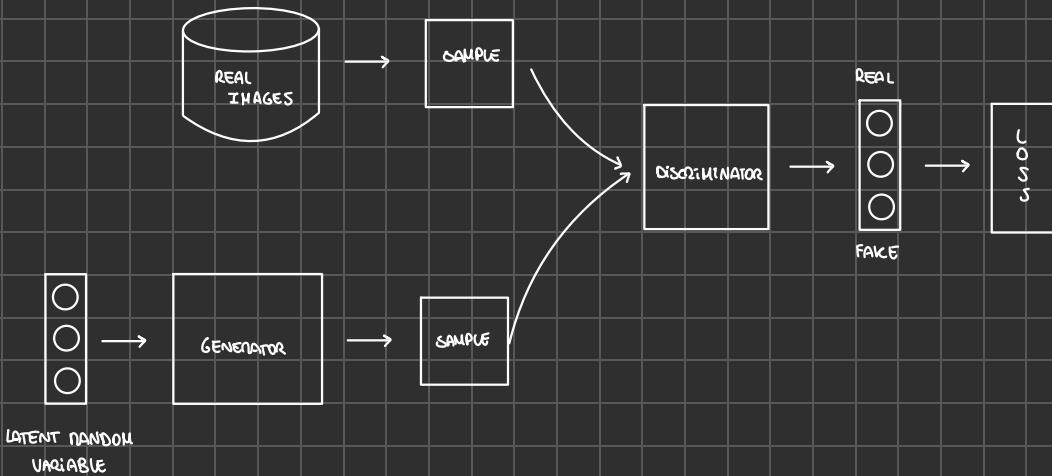
- modify data in specific dimensions
- identifies meaningful dimensions in latent space



now produce distribution → using parametric distribution

SAMPLING IS NOT DIFFERENTIABLE → re-parametrization

GAN Generative Adversarial Network



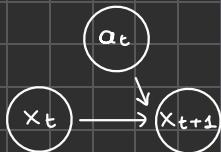
"the loss of generation and the loss of discrimination compete each other" the meaning is when one loss goes down the other goes up and vice versa the goal is when both stabilize

generators try to fool discriminators in believing that the images are real and discriminators try to discriminate which images are real and which no

LOSS OF GENERATOR AND DISCRIMINATOR DOESN'T DECREASE MONOTONIC

### HARKON PROPERTY

1. Once the current state is known the evolution of the dynamic system doesn't depend on the history of states, actions and observations
2. current state contains all information needed to predict the future
3. future one conditionally independent of past states and observations given the current state



### DETERMINISTIC TRANSITIONS

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

$X$ : finite set of states

$A$ : finite set of actions

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

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

### NON-DETERMINISTIC TRANSITIONS

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

$X$ : finite set of states

$A$ : finite set of actions

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

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

### STOCHASTIC TRANSITIONS

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

$X$ : finite set of states

$A$ : finite set of actions

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

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

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

if  $r(a_i)$  is deterministic and know

$$\pi^*(x_0) = \operatorname{argmax}_{a \in A} r(a)$$

if  $r(a_i)$  is deterministic and unknown

for each  $a_i \in A$

- execute  $a_i$  and collect reward  $r(a_i)$

$$\pi^*(x_0) = a_i \text{ with } i = \operatorname{argmax}_{a_i} r(a_i)$$

+ ... |A|

if  $r(a_i)$  non-deterministic and known

$$\pi^*(x_0) = \operatorname{argmax} (\mathbb{E}(r_i))$$

if  $r(a_i)$  non-deterministic and un-known

I. initialize data structure  $\Theta$

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

- choose an action  $a_{t,i} \in A$
- execute  $a_{t,i}$  and collect reward  $r_{t,i}$
- update data structure  $\Theta$

III. optimal policy  $\pi^*(x_0) = \dots$  depending with data structure  $\Theta$

assume  $r(a_i) = \mathcal{N}(\mu_i, \sigma_i)$

1. initialize  $\Theta(0)[i] \leftarrow 0$  and  $c[i] \leftarrow 0 \quad i = 1, \dots, |A|$

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

- choose an index  $i$  for action  $a_{t,i} = a_i \in A$
- execute  $a_{t,i}$  and collect reward  $r_{t,i}$
- increment  $c[i]$
- update  $\Theta(t)[i] \leftarrow \frac{1}{c[i]} (r_{t,i} + (c[i]-1)\Theta(t-1)[i])$

3. optimal policy:  $\pi^*(x_0) = a_i \text{ with } i = \operatorname{argmax}_{i=1, \dots, |A|} \Theta(t)[i]$

EXAMPLE K-ARM BANDIT  $\rightarrow r(a_i) = \mathcal{N}(\mu_i, \sigma_i)$

$$Q_n(a_i) \leftarrow Q_{n-1}(a_i) + \alpha [r_i - Q_{n-1}(a_i)] \quad \text{and} \quad \alpha = 1 / 1 + \sum_{a_j} (a_j)$$

## VALUE ITERATION

estimate value function and then compute  $\pi^*$

$$\pi^* = \operatorname{argmax}_{a \in A} (\mathbb{E}(r(x, a) + \gamma V^*(d(x, a)))$$

we don't know

## RECURSIVE FORMULA OF $Q(x, a)$

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

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

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

## Q-LEARNING ALGORITHM $\rightarrow$ OFF policy

for each  $x, a$  initialize table entry  $\hat{Q}_{(0)}(x, a) \leftarrow 0$

observe current state  $x$

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

- choose an action  $a$
- execute action  $a$
- observe new state  $x'$
- collect immediate reward  $r_t$
- update the table entry for  $\hat{Q}(x, a)$

$$\hat{Q}_{(t)}(x, a) \leftarrow r_t + \gamma \max_{a' \in A} \hat{Q}_{(t-1)}(x', a')$$

$x \leftarrow x'$

$$\text{optimal policy } \pi^*(x) = \arg \max_{a \in A} \hat{Q}_{(T)}(x, a)$$

## NON-DETERMINISTIC Q-LEARNING

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

{ } } }  
 WHAT I CAN COMPUTE WHEN      THE KNOWLEDGE  
 I EXECUTE AN ACTION      I HAVE BEFORE

## TEMPORAL DIFFERENCE LEARNING

Q-LEARNING reduce discrepancy between successive Q estimate

$$Q^\lambda(x_t, a_t) = r_t + \gamma [(1-\lambda) \max_{a' \in A} \hat{Q}(x_t, a') + \lambda Q^\lambda(x_{t+1}, a_{t+1})]$$

SARSA → ON POLICY → accordingly with current policy

is based on tuple  $\langle r, a, \pi, s', a' \rangle$

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

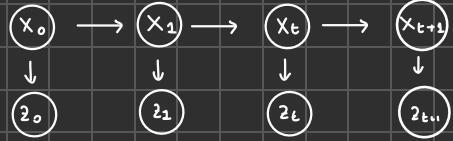
POLICY ITERATION → we don't estimate Q function we have an parametric representation of policy  $\pi_\theta(x)$  we estimate the policy gradient through experiment until terminal condition

### HIDDEN MARKOV CHAIN



future depends only on the current state

### UNKNOWN HIDDEN MARKOV MODELS



state  $x$  are discrete and non observable

$\langle x, z, \pi_0 \rangle$

transitions model  $P(x_t | x_{t-1})$

observation model  $P(z_t | x_t)$

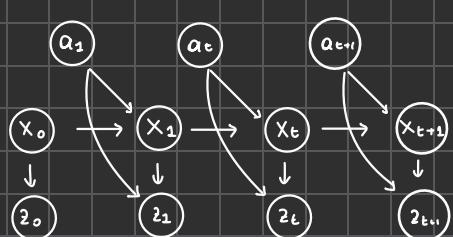
FILTERING → problem of estimating the state of a system at time  $t$  given all the past observations

$$P(x_t = k | z_{1:T}) = \alpha_t^k / \sum_j \alpha_t^j$$

SMOOTHING → given observation so far we want estimate what was the value of state at the past

$$P(x_t = k | z_{1:T}) = \alpha_t^k \beta_t^k / \sum_j \alpha_t^j \beta_t^j$$

POMDP = HDP + HMM



$\langle x, A, z, \delta, r, o \rangle$

$x$ : states

$o(x', a, z') = P(z' | x', a)$ : prob.

$A$ : actions

distribution over observations

$P(x_0)$ : prob. initial state

$\delta(x, a, x') = P(x' | x, a)$ : prob dist. over trans. function

$r(x, a)$ : reward function

we want final policy but we don't know rate  $\rightarrow$  beliefs( $x$ ) = prob. distribution over  
rates