

Lecture 6: Bayesian learning

$P(h/D)$: probability of h given D (D is known)

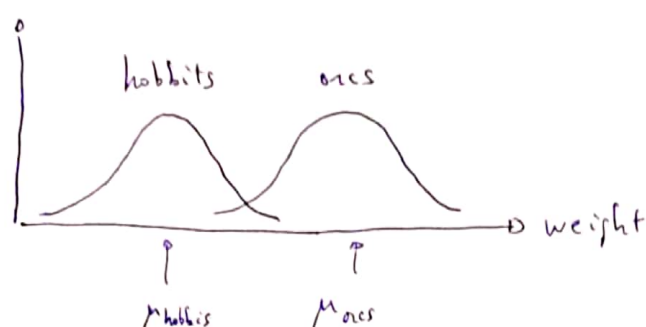
Example:

D : weight = 121 kg

h : orc or hobbit

Calculate $P(\text{hobbit} / 121 \text{ kg})$ vs. $P(\text{orc} / 121 \text{ kg})$ and select the highest

$P(D/h)$ read probability of 121 kg given orc/hobbit



With μ and σ you can use in Matlab the function `normpdf(weight, μ , σ)` to get the probability value

$P(h)$ prior probability \rightarrow knowledge without observation
how many orcs on average?

$P(D)$ probability of observation \rightarrow probability of someone that gets 121 kg

$\sum P(D/h_i) P(h_i)$ makes sure that $P(h_i/D)$ sum to 1

Bayesian probability = Posterior probability

h_{MAP} : maximum A posteriori hypothesis \rightarrow no other hypothesis is better

$$= \arg \max_{h \in H} P(h/D) = \arg \max_{h \in H} \frac{P(D/h) P(h)}{P(D)} = \arg \max_{h \in H} P(D/h) P(h)$$

\hookrightarrow denominator is a constant

Bayesian

Maximum likelihood

} doctors

\hookrightarrow standard in medicine and gambling

$$h_{MAP} = \arg \max_{h \in H} P(h_i/D)$$

$$h_{ML} = \arg \max_{h \in H} P(D/h_i)$$

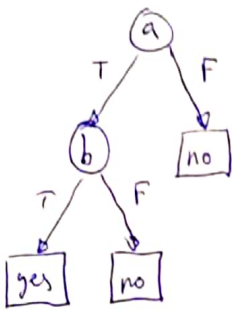
\rightarrow changes of tumor given headache

\rightarrow chances of headache given tumor

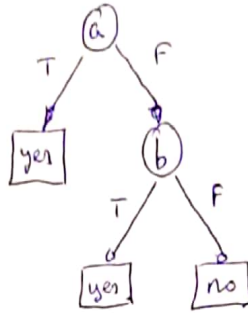
\hookrightarrow standard in engineering
 \hookrightarrow types of decision

Homework 3

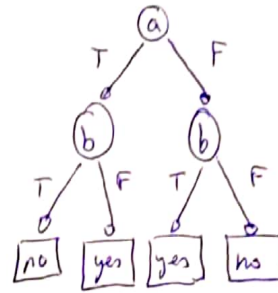
a) $A \wedge B$



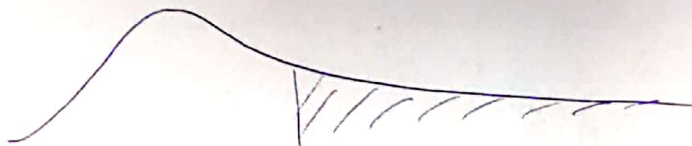
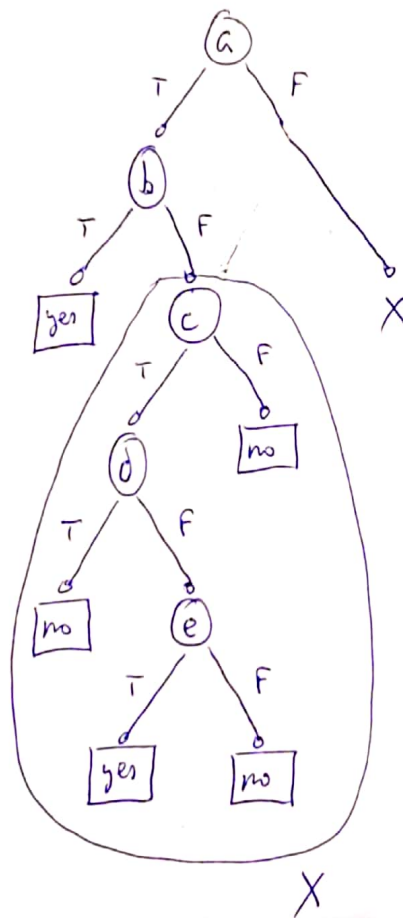
b) $A \vee B$



c) $A \oplus B$

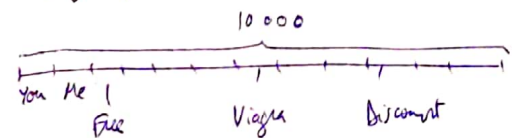


d) $(A \wedge B) \vee (C \wedge \neg D \wedge E)$



Spam classification

Representation of a text document
"bag of words"



$P(\text{spam}) < 1, 0, 0, \dots, 0, 0, 1 >$

$P(\neg \text{spam}) < 1, 0, 0, \dots, 0, 0, 1 >$

If no examples of a specific combination both probabilities are zero

↳ Naïve Bayes approach → values closer to 1 or 0

The maximum likelihood approach gives for granted to have an enough quantity of samples

↳ in real world → sometimes it's not possible

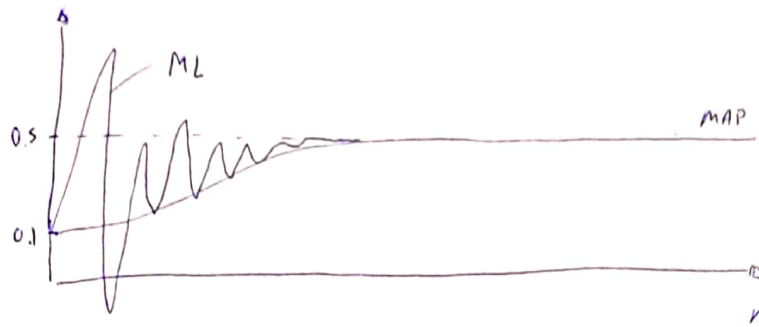
Tossing a coin → example

$P(\text{heads} | n \text{ coin tosses})$

$$\text{heads_prob}_{ML} = \frac{\# \text{ heads}}{n} \quad \rightarrow \text{with few samples, it's very inaccurate}$$

↳ use prior information

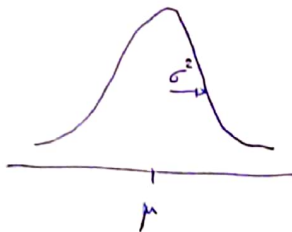
$$\text{heads_prob}_{MAP} = \frac{\# \text{ heads} + m \cdot p}{n + m}$$



Lecture 7 : Bayesian Line Fitting (regression)

Whatever is my observation, don't take the decision only on the observation, keep on with the prior

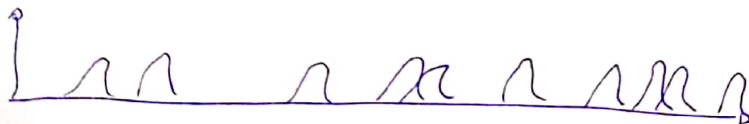
Gaussian :



$$\mu = \frac{1}{N} \sum_i x_i \quad \sigma^2 = \frac{1}{N} \sum_i (x_i - \mu)^2$$

$$f(x; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

CIFAR-10



Homework 4

A patient takes a lab test and the result comes back positive. The test returns a correct positive result in only 98% of the cases in which the disease is actually present, and a correct negative result in only 97% of the cases in which the disease is not present. Furthermore, .008 of the entire population have this cancer.

$$P(\text{cancer}) = 0.008$$

$$P(\neg \text{cancer}) = 0.992$$

$$P(+ | \text{cancer}) = 0.98$$

$$P(+ | \neg \text{cancer}) = 0.03$$

$$h_{\text{map}} \rightarrow \begin{cases} P(+ | \text{cancer}) \cdot P(\text{cancer}) = 0.98 \cdot 0.008 = 0.00784 \\ P(+ | \neg \text{cancer}) \cdot P(\neg \text{cancer}) = 0.03 \cdot 0.992 = 0.02976 \end{cases} \left\{ \begin{array}{l} h_{\text{map}} = \text{no cancer} \\ \downarrow \\ \text{no cancer} \end{array} \right.$$

Homework 5

Based on the table of positive and negative training examples for Play Tennis. Compute posterior probabilities $P(\text{Play Tennis} = \text{yes} | \text{data})$ vs

$P(\text{Play Tennis} = \text{no} | \text{data})$ where data is an observation vector

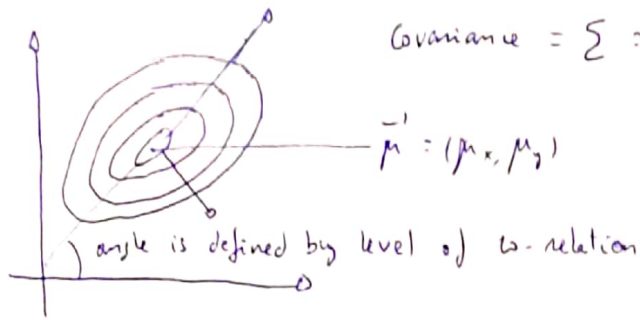
$$S = \langle \text{Outlook} = \text{Sunny}, \text{Temp} = \text{Cool}, \text{Humidity} = \text{High}, \text{Wind} = \text{Strong} \rangle$$

$$\begin{aligned} P(\text{yes} | S) &= P(\text{sunny} | \text{yes}) P(\text{cool} | \text{yes}) P(\text{high} | \text{yes}) P(\text{strong} | \text{yes}) P(\text{yes}) = \\ &= \frac{2}{9} \cdot \frac{3}{9} \cdot \frac{3}{9} \cdot \frac{3}{9} \cdot \frac{9}{14} = 5.3 \cdot 10^{-3} \end{aligned}$$

$$\begin{aligned} P(\text{no} | S) &= P(\text{sunny} | \text{no}) P(\text{cool} | \text{no}) P(\text{high} | \text{no}) P(\text{strong} | \text{no}) P(\text{no}) = \\ &= \frac{3}{5} \cdot \frac{1}{5} \cdot \frac{4}{5} \cdot \frac{2}{5} \cdot \frac{5}{14} = 20.6 \cdot 10^{-3} \end{aligned}$$

$$P(\text{no} | S) > P(\text{yes} | S) \rightarrow h_{\text{map}} = (\text{Play Tennis} = \text{no})$$

2D Gaussian



$$\text{Covariance} = \Sigma = \begin{bmatrix} \text{var}(X_1, X_1) & \text{var}(X_1, X_2) \\ \text{var}(X_2, X_1) & \text{var}(X_2, X_2) \end{bmatrix}$$

The largest eigen value of the eigen vector is the direction

↳ If one eigen value is zero

↳ one of the dimensions in the data has completely disappeared

↳ you don't need to store that information at all

↳ data compression

$$\text{var}(X_1, X_2) = \text{var}(X_2, X_1)$$

Don't use distributions, let data speak by itself → Branch of maths

↳ Non-parametric methods

Lecture 8: Neural Networks (Black Box Methods)

Decision Trees \rightarrow very appropriate for medical science because its transparency

Biology

Dendrites \rightarrow Inputs

Axon \rightarrow Outputs

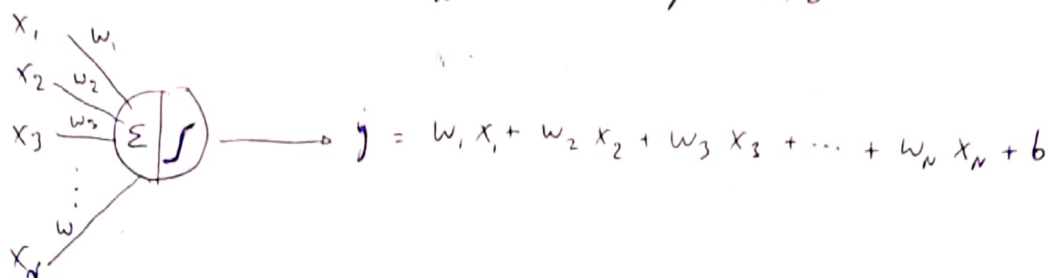
Neuron

Hypothesis: a neuron produces strong output when it is activated by a specific pattern (eg. "I see a bicycle" \rightarrow 1

"no bicycle" \rightarrow 0)

Pattern is specific to certain features ("bicycle features")

linear model: $y = ax + b$

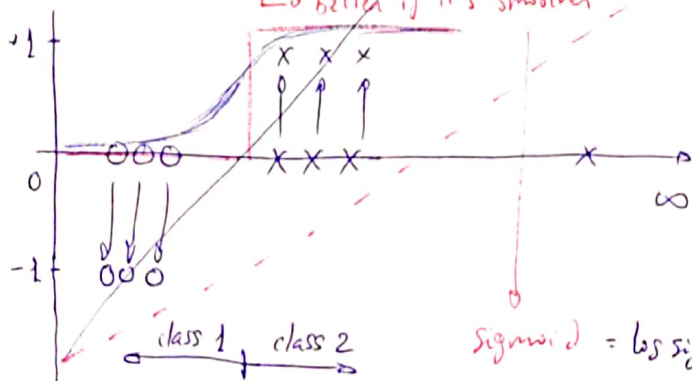


Error: if correctly classified \rightarrow error = 0

else \rightarrow error = 1

Step function?

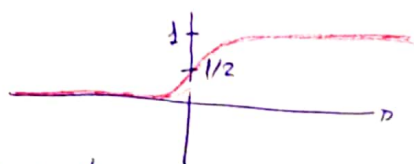
Lo better if it's smoother



Problem with linear model for classification

Lo error must NOT depend on how far on the correct/wrong side you are

Lo \uparrow a \uparrow the quicker the jump is



$$\text{logsig}(x) = \frac{1}{1 + e^{-x}} \rightarrow g = \text{logsig}\left(\sum_{i=1}^N w_i x_i + b_i\right)$$

Homework 6:

Define the formula for the number of free parameters in D-dimensional Gaussian distribution

$$\mu_{D \times 1} = D, \Sigma_{D \times D} = \binom{D}{2} = \frac{D!}{(D-2)! 2!} = \frac{D(D-1)(D-2)!}{(D-2)! 2!} = \frac{D(D-1)}{2} \rightarrow \mu + \Sigma = D + \frac{D(D-1)}{2}$$

$$MSE = \frac{1}{N} \sum_i (y_i - \hat{y}_i)^2 = \frac{0.1 + 0.1 + 0.1}{3} = \frac{1}{3} \sum_i (y_i - \log \text{sig}(\hat{y}_i))^2$$

$$\text{Error} = \sum_{i=1}^N (y_i - \log \text{sig}(\sum_{j=1}^D w_j x_j - b))^2 = MSE$$

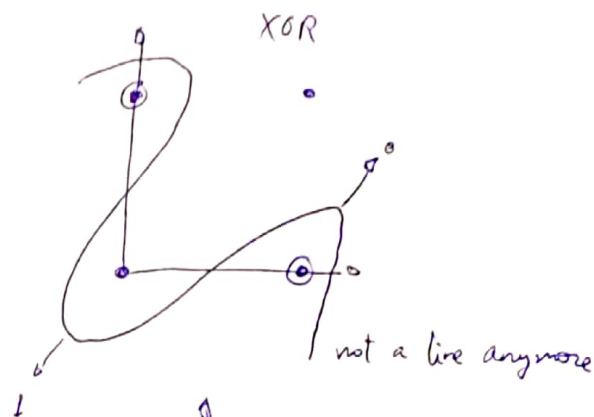
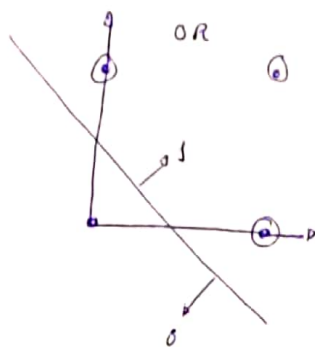
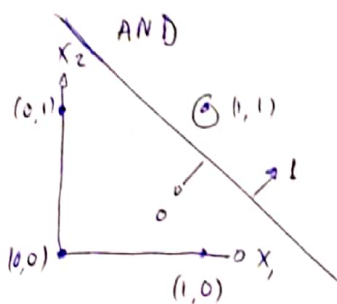
Initialize $w_i = \text{rand}(0, 1)$, $b = \text{rand}(0, 1)$

Then go towards the negative gradient

$$w_i^{k+1} = w_i^k - \alpha \frac{\partial MSE}{\partial w_i} \quad b^{k+1} = b^k - \alpha \frac{\partial MSE}{\partial b}$$

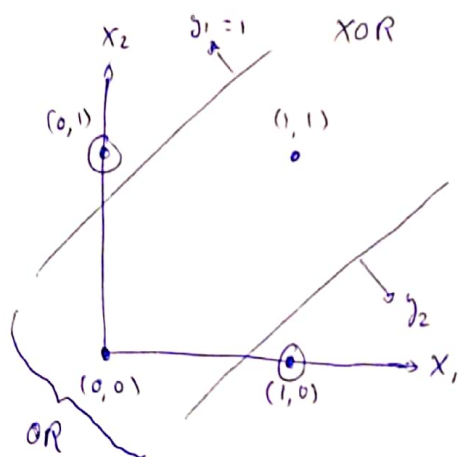
Gradient descent

Example What can be classified?



perceptron
a model of a single neuron
cannot solve this

↓
Multilayer perceptron model → MLP → invention of NN

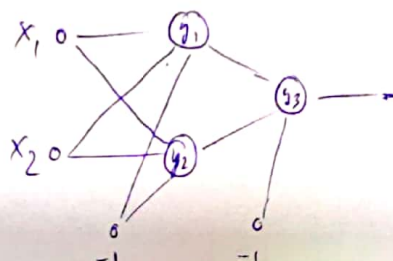


$$y = \log \text{sig}(w_1 x_1 + w_2 x_2 - b)$$

$$y_1 = \log \text{sig}(-2x_1 + 2x_2 - 1)$$

$$y_2 = \log \text{sig}(2x_1 - 2x_2 - 1)$$

$$y_3 = \log \text{sig}(2y_1 + 2y_2 - 1)$$



Lecture 7 - Reinforcement Learning \rightarrow how robot can learn

It started in the robotics field

learning for embodied agents

\hookrightarrow somebody who is active in real world

Optimizing wents in the future

Layers of behaviours in robots

\hookrightarrow Lowest ones make sure it does not destroys itself nor others

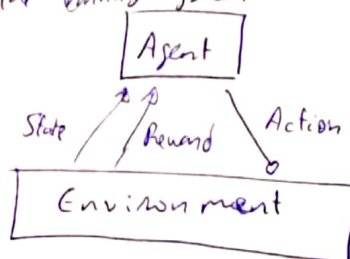
Intelligence \rightarrow Multilayered perception-action patterns

Goal: as much reward as possible and ASAP

$$r_0 + \gamma r_1 + \gamma^2 r_2 \quad 0 \leq \gamma < 1$$

\hookrightarrow changes the power of the learning system

$$S_0 \xrightarrow[a_0]{r_0} S_1 \rightarrow \dots$$



State transition can be or not known

S	a_0	a_1	a_2	...
S_0	S_1			
S_1				
S_2				
\vdots				

Always unknown

R	a_0	a_1	a_2	...
S_0	r_0			
S_1				
S_2				
\vdots				

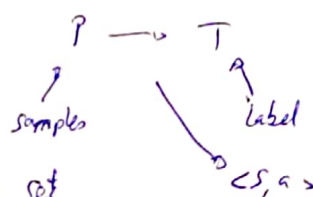
\rightarrow take the maximum action
take the action whose reward is maximum

$$S_{t+1} = \mathcal{S}(S_t, a_t) \quad \hookrightarrow \text{random states}$$

$$r_t = r(S_t, a_t)$$

Problem is to find the policy π

So far our training data \rightarrow

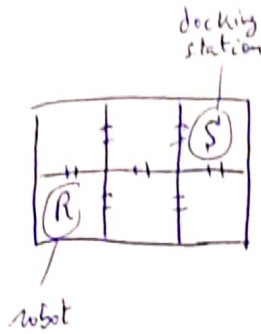


In robotics we don't have that training set

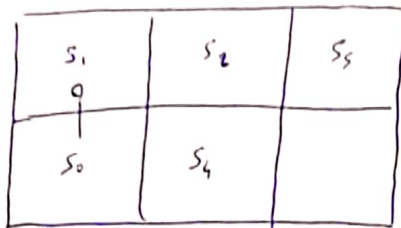
If you start a sequence of actions maybe at some point you will get the reward

Evaluation function $V^\pi(s)$ for some policy π

$$\pi^*(s) = \arg \max_a [r(s,a) + \gamma V^\pi(s(s,a))]$$



$a = \text{up, down, left, right}$
 $a_0 \quad a_1 \quad a_2 \quad a_3$



Prior knowledge makes quicker the learning progress

If there are huge amount of states and actions, finding the proper way can take ages in the context of robotics (non-simulated world)

S	a_0	a_1	a_2	a_3	...
s_0	s_1	n/a	n/a	n/a	
s_1		s_0			
s_2				s_5	
s_3					
s_4					
s_5					
R	a_0	a_1	a_2	a_3	...
s_0	0	0	0	0	
s_1	0	0	0	0	
s_2	0	0	0	100	
s_3					
s_4	s_0				
s_5					

$\gamma = 0.9$

Lecture 10: Unsupervised learning

So far we've been solving:

Find f so that

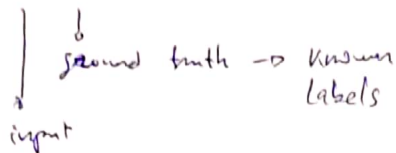
$$f: X \rightarrow Y$$

• Reinforcement learning

$$f: S \rightarrow A$$

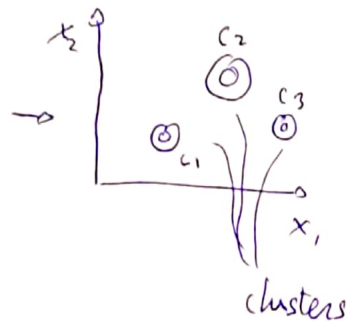
• Supervised learning

Given $\langle \bar{x}_1, y_1 \rangle, \langle \bar{x}_2, y_2 \rangle, \dots, \langle \bar{x}_N, y_N \rangle$ training set



• How about if we only have $\bar{x}_1, \bar{x}_2, \dots, \bar{x}_N$?

Example: Data $\begin{pmatrix} x_1^i \\ x_2^i \end{pmatrix}$ is points in 2D \rightarrow Euclidean data



A phone company hires you and wants to know what are the features and types of clients who drop out their service

x_1 : age

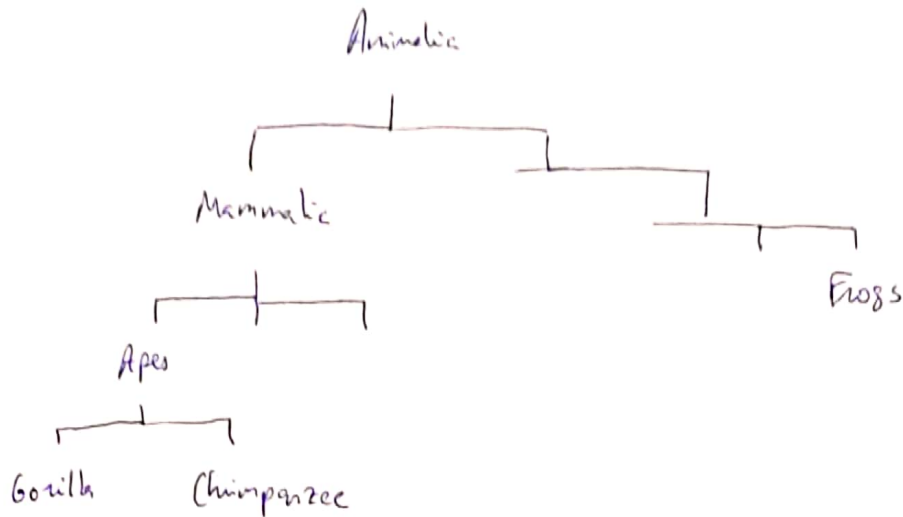
x_2 : income

c_1 : $\bar{\mu}_{c_1} = (15, 200)^T$ kids with smallish amount of money

c_3 : $\bar{\mu}_{c_3} = (78, 1500)^T$ elderly with OK income

c_2 : $\bar{\mu}_{c_2} = (30, 2000)^T$ taking variance (which is large) into account, this is various types

Example → Hierarchical Data



Domain transfer → If I know how to classify objects, can I use this information for unknown data?

↳ Semi-supervised learning

Clustering approaches → Hierarchical clustering

- Agglomerative (bottom-up)
 - Divisive (top-down)
- ↳ Every sample is its own cluster

We need:

1. Distance function

2. Linkage criteria

- Mean distance
- Distance to mean
- Min distance
- Max distance
- Median

$$L1: \sum |a_i - b_i|$$

$$L2 \text{ (euclidean)}: \sqrt{\sum (a_i - b_i)^2}$$

$$L_\infty: \max \{a_i - b_i\}$$

$$\text{Cosine}: \frac{a \cdot b}{\|a\| \|b\|}$$

→ best norm for nearest neighborhood classifier for sparse data

Median → Median Square error adequate for sparse data

Mixture model clustering - Clustering approach

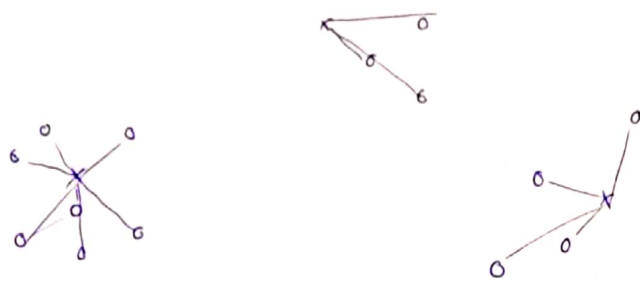
Points are samples from a probability distribution

$$P(\bar{x}) = \sum_{i=1}^K \pi_i \underbrace{\text{prob}(\bar{x}; \theta_i)}_{\substack{\text{basic distribution} = \text{Gaussian}, \\ \text{a priori}}} \quad \theta_i = \bar{\mu}_i, \Sigma_i$$

E.g. unsupervised Gaussian mixture model

↳ needs a lot of data

Sum-of-squares method



Assume k clusters ($k=3$)
Find such μ_1, \dots, μ_k that
distances from points in
each cluster are minimized

This is brute force \rightarrow every datum is compared to every cluster

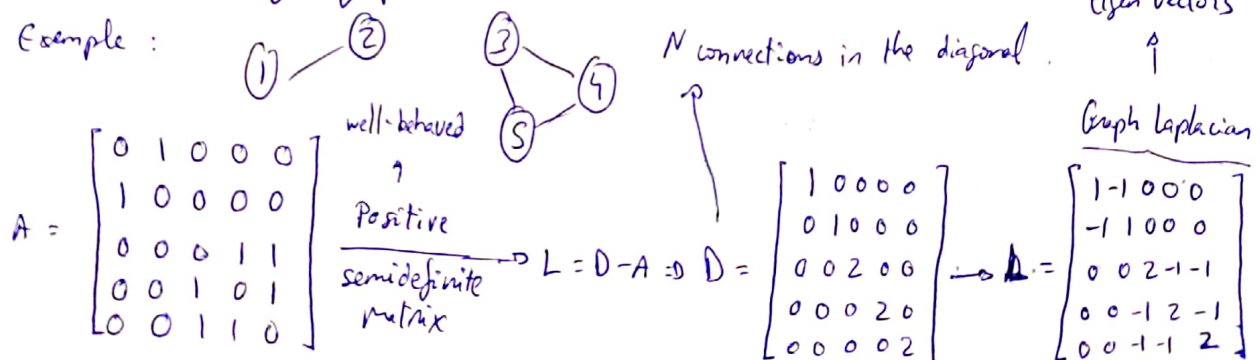
K-means (maximum likelihood solution)

1. randomly pick k points $\bar{\mu}_1, \dots, \bar{\mu}_k$
2. assign each point to its closest cluster center
3. calculate the mean of each cluster and update $\bar{\mu}_1, \dots, \bar{\mu}_k$
4. Repeat 2 and 3 until converges
↳ assignments do not change

Spectral clustering

Our data is of graph/network type

Example:



Many properties:

If any eigenvalues are zero, there is this amount of unconnected graphs and their eigenvectors reveal their nodes

$\lambda_1 = \lambda_2 = 0 \rightarrow$ the corresponding eigenvectors of those must be

① and ② and ③, ④ and ⑤
 ① ②

Instead of finding clusters, can we somehow visualize high dimensional data?

Self-organizing map (SOM)

1. Repeat M times
2. Select a sample \bar{x}
3. Find the closest SOM cell \bar{m}_j for \bar{x}_i (\bar{m}_{BMU})
4. Make \bar{m}_{BMU} ^{best matching unit} more similar to \bar{x}
5. Make neighbours of \bar{m}_{BMU} more similar
6. End repeat
7. Return \bar{m}_j

$$\frac{0.5 \quad 6}{2 \quad 3}$$

$$\bar{m}_{BMU} = \underset{\bar{m}_j}{\operatorname{argmin}} \|\bar{x}_i - \bar{m}_j\|^2$$

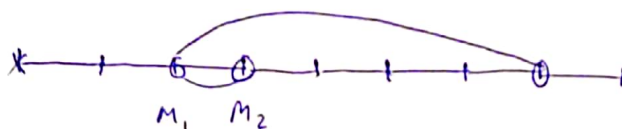
$$\bar{m}_j^{(l+1)} = \bar{m}_j^{(l)} + \underset{0.5}{\alpha} \underset{6}{(\bar{x}_i - \bar{m}_j^{(l-1)})} \underset{8}{}$$

example

$$\alpha = 0.5$$



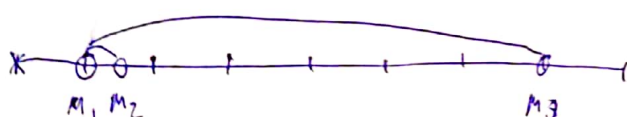
Input $x_1 = 6$



Input $x_2 = 0$

Step $\Delta m_1 = \alpha (x - m_1) = 0.5 (0 - 2) = -1$

Step $\Delta m_2 = 0.5 (0 - 3) = -1.5$



Lecture 13: Learning sets of rules

P	Q	$P \rightarrow Q$
T	T	T
T	F	F
F	T	T
F	F	T

Think: "If rains then drive carefully"

equivalent form: $P \rightarrow Q \iff \neg P \vee Q$

\rightarrow not OR

Learning first-order rules

search-based, but the search-space explodes so no feasible solutions

ProLog \rightarrow mother(x, y)
father(x, y)

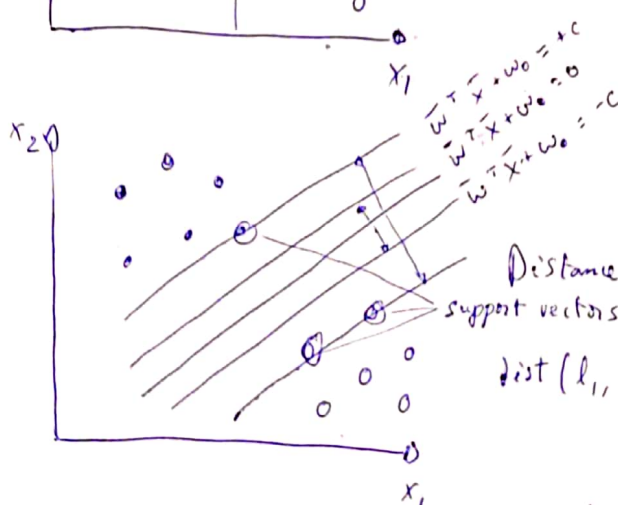
Lecture 9: Support Vector Machines



$$MSE(H_2) = MSE(H_3) = 0$$

For MLP they are equally good solutions

For all correctly classified points: $y_i(\bar{w}^T \bar{x}_i + w_0) > 0$



Goal: Maximize the distance of the two parallel lines (of the decision boundary $\bar{w}^T \bar{x} + w_0 = 0$)

\hookrightarrow push c as far as possible

Distance of two parallel lines:

$$\text{dist}(l_1, l_2) = \frac{|c_2 - c_1|}{\sqrt{a^2 + b^2}} = \frac{2c}{\sqrt{\bar{w}^T \bar{w}}} = \frac{2}{\sqrt{\bar{w}^T \bar{w}}}$$

$$ax_1 + bx_2 + c = 0 \rightarrow \frac{a}{c}x_1 + \frac{b}{c}x_2 + 1 = 0$$

Problem: maximize $\frac{2}{\sqrt{\bar{w}^T \bar{w}}}$ $\bar{w} = \begin{pmatrix} w_0 \\ w_1 \\ w_2 \end{pmatrix}$

\hookrightarrow minimize $\sqrt{\bar{w}^T \bar{w}} \rightarrow \min_{\bar{w}} \bar{w}^T \bar{w} \rightarrow$ is 0

subj. to:

$$y_1(\bar{w}^T \bar{x}_1 + w_0) \geq 1 - \epsilon_1$$

$$y_2(\bar{w}^T \bar{x}_2 + w_0) \geq 1 - \epsilon_2$$

$$y_m(\bar{w}^T \bar{x}_m + w_0) \geq 1 - \epsilon_m$$

Homework 7: Give weights w_1 , w_2 and w_b so that a perceptron model establishes: x_1 or x_2 , x_1 and x_2 , $\neg x_1$

$$x_1 \text{ and } x_2 : w_1 = w_2 = 0.5, w_b = -0.5$$

$$x_1 \text{ or } x_2 : w_1 = w_2 = 0.5, w_b = 0.5$$

$$\text{not } x_1 : w_1 = -0.5, w_2 = 0, w_b = 0.5$$

$$y = w_1 x_1 + w_2 x_2 + w_b$$

$$\text{ssn}(y) = \begin{cases} 1, & y > 0 \\ -1, & y < 0 \end{cases}$$

Optimization

1. Linear optimization \rightarrow e.g. simplex
 $\max w_1 x_1 + w_2 x_2$ ($-c \sum E_i$) \rightarrow a restrictive is with points in the
 sub; $x_1 + x_2 \leq B$ $\xrightarrow{\text{SVM addition}}$ penalty for misclassifications wrong side
 $x_1 \geq 0, x_2 \geq 0$

2. Unconstrained optimization

$$\max/\min f(x)$$

e.g. method gradient descent

3. Constrained problems

$$\min f(x)$$

$$\text{subject to } h_i(x) = 0$$

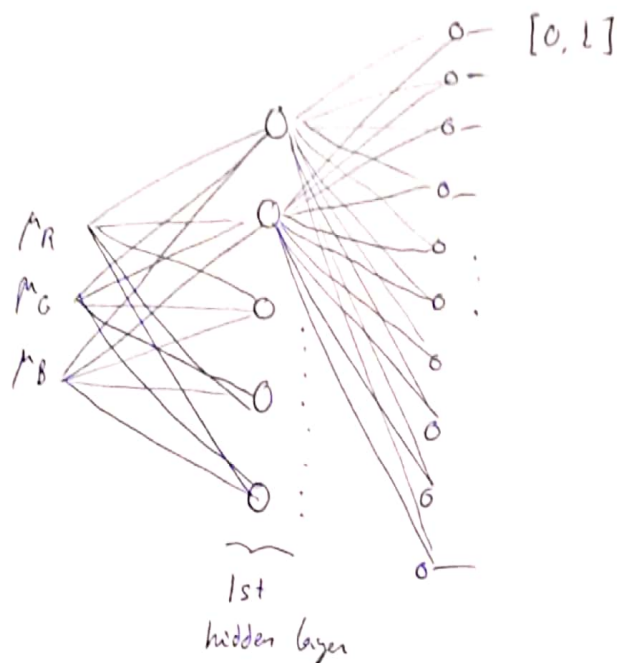
$$g_i(x) \leq 0 \quad \Rightarrow \text{SVM}$$

4. Discrete optimization

e.g. travelling salesman, ...

$\left. \begin{array}{l} \text{samples_train} \\ \text{labels_train} \end{array} \right\} \text{ learn a } \overbrace{\text{Multi-layer perceptron network}}^{\text{MLP}}$

samples_test \hookrightarrow need to define # neurons of hidden layer
 labels_test



When the network is trained, labels_train need to be reshaped

GUI \rightarrow see patternet in wiki matlab

\gg nn_start

\gg est_labels_test = sim(net, samples_test);

$\text{vec2ind}([0 \ 0 \ 1 \ 0 \ 0]) \rightarrow 3$, see also ind2vec

$\text{patternet}([10 \ 10])$

$\underbrace{\hspace{2cm}}$
 2 hidden layers
 with 10 neurons each

take a subset of 10% of the training set while training
 \hookrightarrow when everything looks OK \rightarrow the whole set

	1	2	3	4	(5)	6
1	0	4	13	24	12	8
2		0	10	22	11	10
(3)			0	7	(3)	9
4				0	6	18
5					0	8.5
6						0

	1	2	3-5	4	6
1	0	(4)	12	24	8
2		0	10	22	10
3-5			0	6	8.5
4				0	18
6					0

	1-2	3-5	4	6
1-2	0	10	22	8
3-5		0	(6)	8.5
4			0	18
6				0

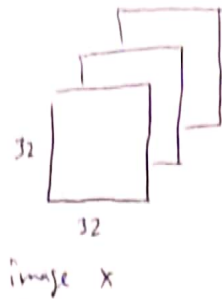
	1-2	3-4-5	6
1-2	0	10	(8)
3-4-5		0	8.5
6			0

	1-2-6	3-4-5
1-2-6	0	(8)
3-4-5		0

	1-2-3-4-5-6
1-2-3-4-5-6	0

Decision Tree

a and b



too big $\rightarrow f = \text{afar} = 10 \text{ features}$
7
= 3072
dimension

Lo mean
of every channel

$$f_{-1} = \text{mean}(x(:, :, 1))$$

$$f_{-3} = \dots \quad \text{Lo 1-end}$$

$$f_{-b} = \dots$$

$$f = [f_{-1}, f_{-3}, f_{-b}]$$

Bayesian classifier

$$c = \text{afar} - 10 - \text{bayer} - \text{classify}(f, \mu, \sigma, p)$$

sample

$$[\mu, \sigma, p] = \text{afar} - 10 - \text{bayer} - \text{learn}(f_{\text{-train}}, \text{label}_{\text{-train}})$$

10 x 3 10 x 1 50000 x 3 50000 x 1

Lo how many ones divided by total # labels
two
three
var(x)
function mean(x)

$$\mu = \text{mean}(x)$$

$$\begin{bmatrix} x_1^{(1)} & x_2^{(1)} \\ x_1^{(2)} & x_2^{(2)} \\ \vdots & \vdots \\ x_1^{(n)} & x_2^{(n)} \end{bmatrix}$$

for class = 1... 10
for ex = 1... length(p)
posteriori(class) = normpdf(f, mu(class, 1), sigma(class, 1)) + normpdf(...)

Naive assumption \rightarrow 3 single variable gaussian function

Correlated assumption \rightarrow 1

$$\text{var} \rightarrow \text{cov}(x)$$

$$\text{sign} \rightarrow \Sigma$$

now instead of normpdf(...) \rightarrow mvnpdf(), cov()

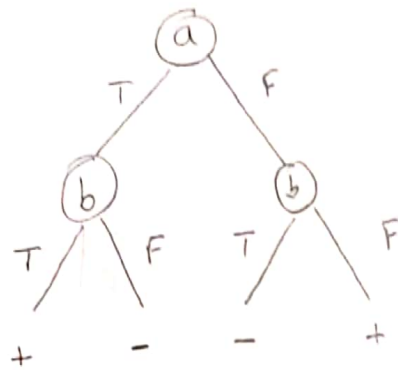
Lo matrix

Divide the image 4 features \rightarrow training images \rightarrow if it works \rightarrow subdivide \rightarrow if you keep subdividing \rightarrow original image

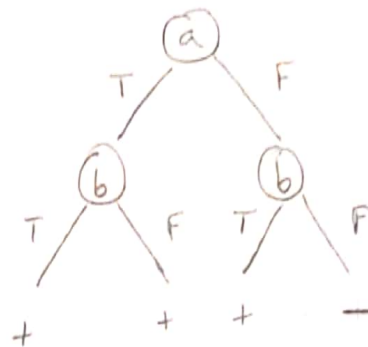
$$\Sigma_{3 \times 3} = \begin{bmatrix} \sigma_{RR}^2 & \sigma_{RG}^2 & \sigma_{RB}^2 \\ \sigma_{GR}^2 & \sigma_{GG}^2 & \sigma_{GB}^2 \\ \sigma_{BR}^2 & \sigma_{BG}^2 & \sigma_{BB}^2 \end{bmatrix}$$

Decision Tree

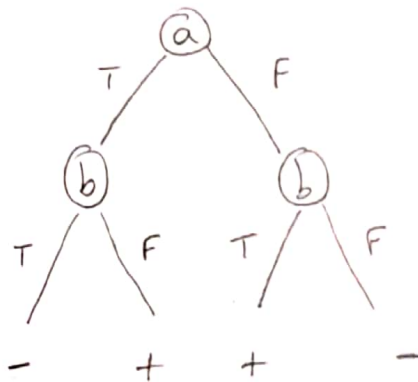
• a and b



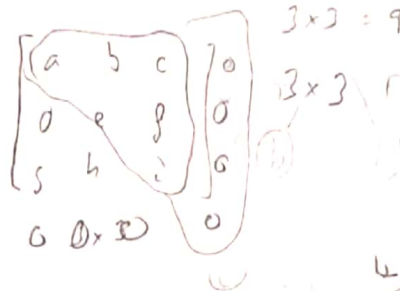
• a or b



• a xor b



• (a and b) or (c and not d and e)



formula # of free parameters in D-dimensional gaussian distn

$\vec{\mu}_{D \times 1} \rightarrow D$ par: $D + \frac{D^2 + D}{2}$

$\Sigma_{D \times D} = \begin{bmatrix} \sigma & \sigma \\ \sigma & \sigma \end{bmatrix}$ $D + \frac{D^2 + 1}{2}$

$D + \frac{D^2}{2} \rightarrow D + \frac{D^2 + 1}{2}$

$U = (D^2/2 - D)$

for every 3 parameters

↳ you need at least 3 examples

