

# Lecture 10 : Pattern Recognition week 47.

## Recap

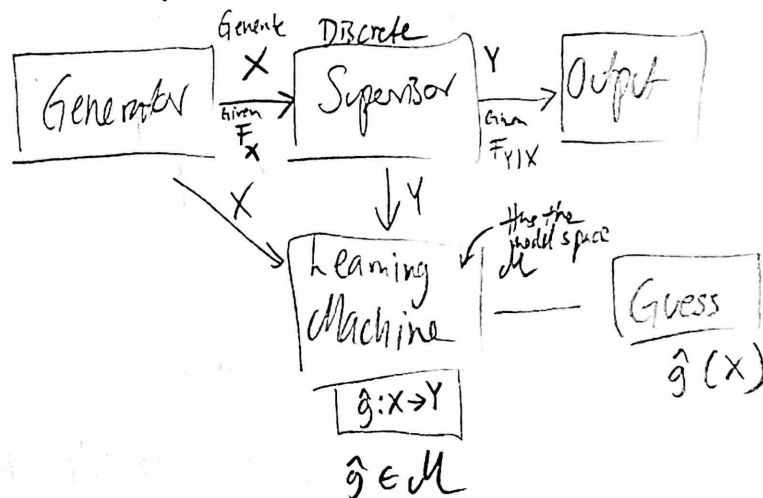
### - Markov chains.

Homogenous Markov chains:

- homogenous in time.
  - dry wet example.
  - Markov can be used on language: what's most likely the next word?
  - Google's Page Rank was based on Markov chains.
- state space ~~X~~

## The pattern recognition problem

Perception later became neural networks



in the pattern recognition problem we are assuming  $F_{Y|X}$  is discrete (a class)

$g_X$  is a decision function / decision rule / class

0-1  
loss

$$L(z, u) = \begin{cases} 0 & \text{if } y = u \\ 1 & \text{if } y \neq u \end{cases}$$

0-1 loss is mistake cost 1, otherwise 0

$z = (x, y)$ ,  $u$  is the guess.

# Pattern Recognition

$$R(\lambda) = \int \underbrace{\ell((x, y), g_\lambda(x))}_{\text{this is the true risk for } g_\lambda \in \mathcal{M}} dF(x, y)$$

this is the true risk for  $g_\lambda \in \mathcal{M}$ .

The goal is to minimize the true risk  $R(\lambda)$

Example

Bayes Classification Rule

Let  $r(x) = E[Y|X]$ , then the Bayes classification rule  $h^*$  is

$$h^*(x) = \begin{cases} 1 & \text{if } r(x) > 1/2 \\ 0 & \text{otherwise} \end{cases}$$

else :

$$Y \in \{0, 1\}$$

$$E[Y|X] = P(Y=1|X) = r(x)$$

$$\text{if } r(x) > 1/2 \Rightarrow h^*(x) = 1$$

$$\text{if } r(x) < 1/2 \Rightarrow h^*(x) = 0$$

• This is what scikit-learn's Logistic Regression does!

Skapad med Tiny Scanner

Function predict

Instance  
&  
label space

$\mathcal{X}$ : instance  
 $\mathcal{Y}$ : label space

$d$  different binary variables  
 $\mathcal{X}$  can be  $\{0, 1\}^d$

$\mathcal{X}$ : data space

$\mathcal{X} = (\mathcal{X}, \mathcal{Y})$  kan skins  
n samples  
sem

Linear  
separator

= a line that separates. Easy in 1D or 2D.  
for higher dimensions we need perceptrons.

The perceptron  
algorithm

- a linear classification rule.

The  
dimension  
trick

A prerequisite.  
Check slides & lecture notes

Perceptron  
algo

Goal is to find:  $(w \cdot x_i) y_i > 0$

$$g_w(x) = \begin{cases} 1 & \text{if } w \cdot x > 0 \\ -1 & \text{if } w \cdot x < 0 \end{cases}$$

$(w \cdot x_i) y_i > 0$  means that you have guessed  
the correctly.

if  $(w = x_i) y_i > 0 \quad \forall i = 1, \dots, n$

$$\frac{1}{n} \sum_{i=1}^n L((x_i, y_i), g_w(x_i)) = 0.$$

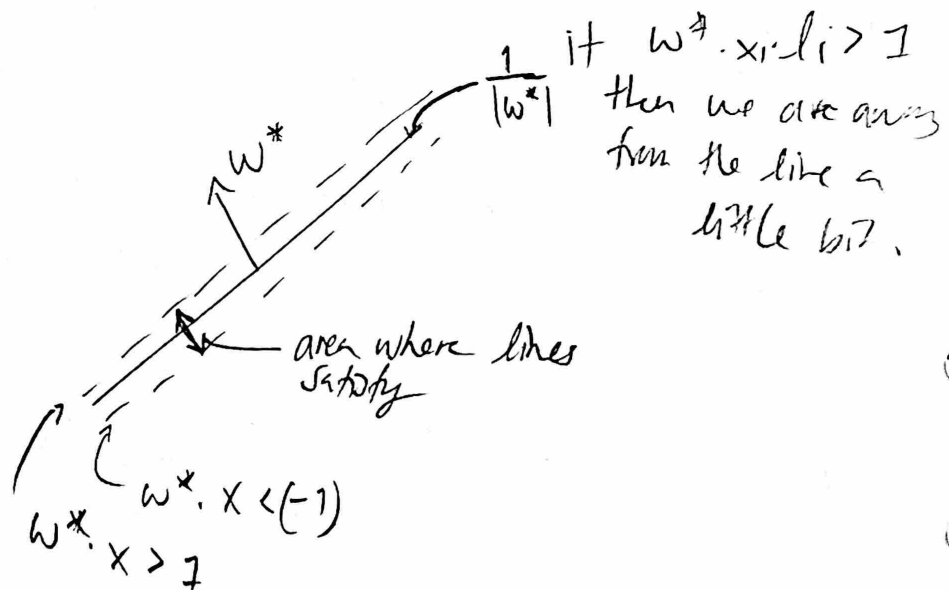
the algo

if you misclassify something  
that should've been 1,  
then  $w := w + x_i y_i$

if you misclassify something that  
should've been (-1) then  
 $w := w - x_i y_i$ .

Feedback  
End. ass.  
Problems 17-18  
Chapter reader  
Computer exercises.  
40p. 20 godkand.  
extra peng. = 14p.

$w^*$

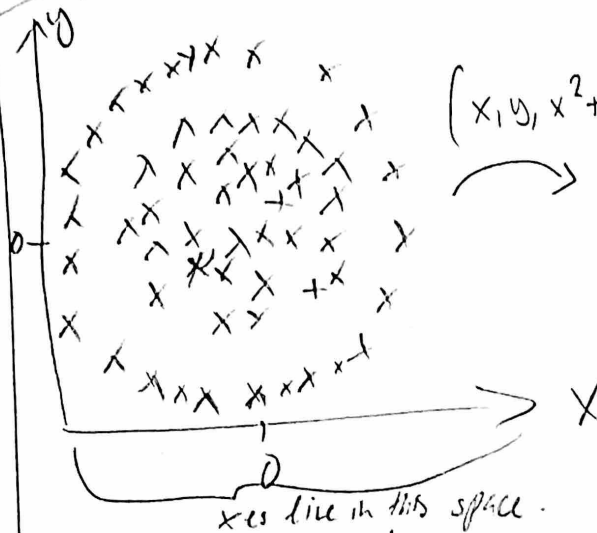


the closer  $w^* \cdot x < -1$  or  $w^* \cdot x > 1$   
are to the line, there are fewer lines  
that satisfy.

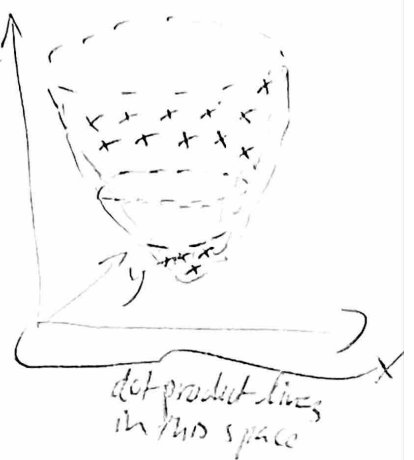
Convergence is guaranteed if the data is separated?

Kernelization

→ check slides.



$$(x, y, x^2 + y^2)$$



you have to have as many parameters as you have dimensions.

$l_i = y_i$  i has votes.

A Kernelized perceptron algorithm

$$K_{ij} = \Phi(x_i) \cdot \Phi(x_j) = \underbrace{\Phi(x) \cdot \Phi(y)}_{\text{dot product}}$$

$$1. c = 0$$

where i has slides

2. While there exists an  $i$  such that

$$(Kc)_i y_i \leq 0 \quad \text{update} \quad c_i := c_i + y_i$$

Kernel functions

we only need to know the kernel to use it, we do not need to know  $\Phi$ .

$c_i$  holds the number of times we have added or subtracted  $x_i$ .