

Machine learning.

Statistical learning

$S_n = \{(x_1, y_1), \dots, (x_n, y_n)\}$ training set

$x_i \in X \subseteq \mathbb{R}^d$

$y_i \in Y \rightarrow$ scalar Regression \rightarrow Vectorial Regression
 $Y \subseteq \mathbb{R}^m$

Binary classification: $y = \{-1, 1\}$

Multi-class classification (C)

$Y = \{0, \dots, C-1\} \leftarrow \text{order}$

I want my function to fit & generalize new data

$f: X \rightarrow Y$

$f(x_i) \approx y_i$ fitting

$f(x_{\text{new}}) \approx y_{\text{new}}$ $x_{\text{new}} \in S_n$ generalization

In training set we have not x_{new} but function must work

f must have good approximation.

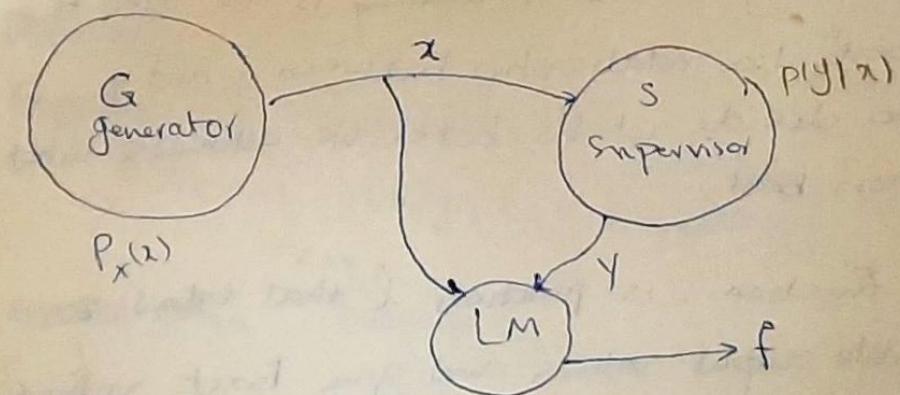
Ass 1: input part of training set generation of probability distribution of X such that the samples

x_i : one i.i.d

x_i : sample according to probability set has been identically and independently sample (i.i.d)

Ass 2: input part of training set

$$P(x, y) = P_x(x)P(y|x)$$



What we looking for f : strongly related to:

$$P(x, y) = P_x(x)P(y|x)$$

generative model try to know probability of

$$P(x, y) = P_x(x)P(y|x)$$

in Supervised learning we want to estimate function not Probability distribution.

Don't forget:

$X \times Y$ is our data space.

for a moment:

function f is the best goal

let's say: Goal of learning is to find the "best" I/O relationship between x and y but to decide what's better we estimate what solution best.

Loss Function: is function l that takes 2 possible output values and give back values $[0, +\infty)$, tells you:

$$l: Y \times Y \rightarrow [0, +\infty)$$

$f(y_i, f(x_i))$ common/popular choice

2. $\rightarrow y_i$
 $\rightarrow f(x_i)$ if $y_i = f(x_i) \rightarrow$ perfect

but we have difference.

$$\|y_i - f(x_i)\|^2 : \text{norm, square loss}$$

if these 2 function are equal prediction and answer would be the same and this is a goal.

but in reality this difference is and this is equal to norm & regression

Definition:

Expected loss (expected risk)

$$E(f) = \int_{X \times Y} p(x, y) l(y, f(x)) dx dy$$

expected value $\underbrace{\text{Data space}}$



$$E_{XY}(l(y, f(x))) \leftarrow \text{for prediction of other}$$

value which is added to the set

Expected value on all data space is add \propto to past & future.

The best I/O function ($f^*: x \rightarrow y$) ✓

$$f^* = \operatorname{argmin}_f \hat{\mathcal{E}}_n(f)$$

$f \in \mathcal{H}$ $f \in \mathcal{Y}$ → Whatever type of f .

I want to find $\min \hat{\mathcal{E}}_n(f)$

argmin of f associated with min function

*** can not be calculated.

Problem:

1, we need probability

If probability change everything
must be changed

2, We don't have access to entire dataspace

What we have a training set with sample
we need function depends on n

if n change the solution changed.

for solve:

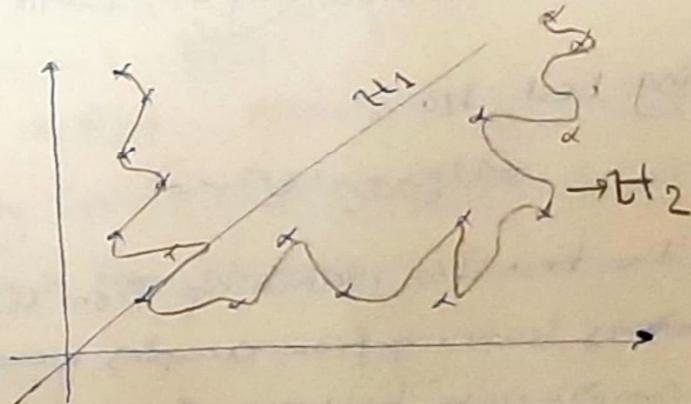
Empirical loss / Risk

$$\hat{\mathcal{E}}_n(f) = \frac{1}{n} \sum_{i=1}^n l(y_i, f(x_i))$$

Limited number of sample so we can compute it

$$\hat{f} = \operatorname{argmin}_f \hat{\mathcal{E}}_n(f) \quad \mathcal{H} \rightarrow \text{Hypothesis space}$$

↓
Best I/O relationship, best function ✓



$$\hat{f}: x \rightarrow y \quad x \in \mathbb{R} \quad y \in \mathbb{R}$$

\mathcal{H} $f(x) = w_0x + b \rightarrow$ class of linear functions
we call \mathcal{H}_1

\mathcal{H}_1 : poor fitting

High stability: \mathcal{H}_1 (Predictions are stable)
→ High stability due to \mathcal{H}_1

\mathcal{H}_2 : very complex function

\mathcal{H}_2 : { High fitting
not stable (poor stability)}

if I change 1 point it changes a lot.

overfitting H_2 "good"

not stability \bar{f}_r $\approx f$

finding best I/O

$E(f^*) \approx E(\hat{f})$

Finding the best I/O relationship from the data
which means learning from examples means to
find a compromise between fitting and
stability generalized.

$E(f^*)$ $E(\hat{f})$
estimated from data
ideal quantity

2 different function:

if I have 2 above function, as I increase
samples, difference of $E(f^*)$, $E(\hat{f})$ is
going to zero.

Consistency:

$$\lim_{n \rightarrow +\infty} E(E(f^*) - E(\hat{f})) = 0$$

number of samples.

- session 3 - 9th Oct

local method Nearest Neighbour [S]

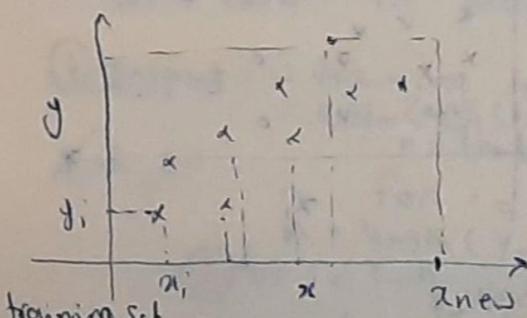
I only consider 1 or few point to take decision
start from assumption: find estimator \hat{f}

$$\hat{f}: X \rightarrow Y$$

if you have 2 in input space should
close in output space.

How explore:

example: talk about function $\hat{f}: X \rightarrow Y$
 $\subseteq \mathbb{R}$ $\subseteq \mathbb{R}$



$$S_n = \{(x_i, y_i)\}_{i=1}^n$$

number of element = 9 in example

lets assume: give a new point (x_{new})
ask estimator what are the act of it
for new point:

$$\hat{f}(x_{\text{new}}) = ?$$

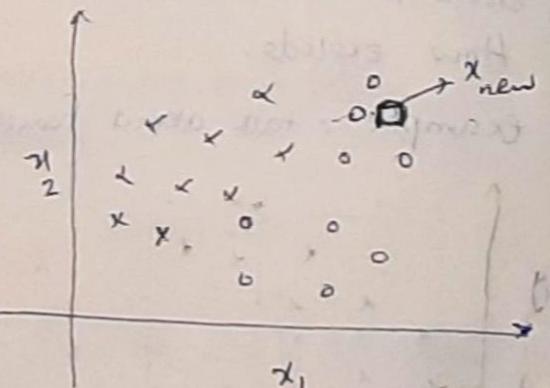
I decide output of x_{new}

what we have Binary classification

2 class x and o

2 coordinate

instead of having x and y



$$\begin{aligned} f: X \rightarrow Y \\ \subseteq \mathbb{R}^2 \\ E\{x, o\} \end{aligned}$$

برای تفکر x_{new} در این حالت زیر کنید که این بین داده ای که نویسید را بازخواهید
نمایش می کنید و این است برای دستور الگوریتم برای
 $\|x_{\text{new}} - x_i\|$ $i=1, \dots, n$ \Rightarrow $(n, 1)$ \rightarrow $(n, 1)$

$\star i^* = \arg \min_{i=1, \dots, n} \|x_{\text{new}} - x_i\|$

$$y_{\text{new}} = \hat{f}(x_{\text{new}}) = y_{i^*}$$

برای این باره از این امر برخوب است دعوه کنیم اما استارتری این برای
سریع

پیشگیرانه

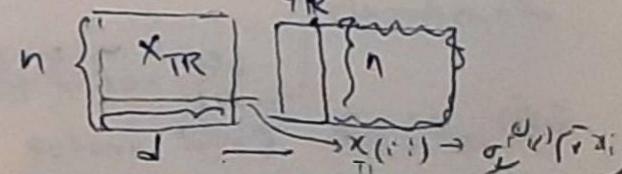
lets turn to code

① INPUT : training set X_{TR} , Y_{TR} and x_{new}
 $\text{DIST} = \text{zeros}(1, n)$
 $n = \text{the length of training set } (x, y)$

$\sqrt{x_{\text{new}}^T x_{\text{new}}}$ i For $i = 1 \dots n$

$n = \text{length}(Y_{\text{TR}})$

② for $i = 1 \dots n$ $x_{\text{TR}}^{(i)} \rightarrow x_{\text{TR}}(i, :)$ $d(x_{\text{new}}, x_{\text{TR}}(i, :)) \rightarrow \text{DIST}(i)$



Nearest neighbours
Input: X_{TR} , Y_{TR} : training set

x_{new}

$\text{Dist} = \text{Zeros}(1, n)$

$n = \text{length}(Y_{TR})$

for $i = 1 \dots n$

$$\hookrightarrow \text{Dist}(i) = d(x_{new}, X_{TR}(i, :))$$

$$(\minVal, \underline{\minIdx}) = \min(\text{DIST})$$

$\Rightarrow Y_{TR} = \text{p}(\underline{\minIdx})$ return $Y_{TR}(\underline{\minIdx})$

$\Rightarrow \min$ *for all classes* \Rightarrow *for all samples*

Computation of complexity of method.

$O(nd)$

\uparrow
all samples of TR \uparrow \rightarrow could be high dimension

K-nearest Neighbours

=

K is a number of closest point to take position

choice of K which are better than others

K: hyper parameter

outside the method.

if K is changed get another distance of method

Input: X_{TR} , Y_{TR} , x_{new} , K \rightarrow is a number

if K is 1: algorithm is nearest neighbour

$n = \text{length}(Y_{TR})$

we need Dist: $\text{Zero}(1, n)$

for $i = 1 \dots n$

$$\text{Dist}(i) = d(x_{new}, X_{TR}(i, :))$$

sorted Dist , dist \uparrow , K is short

① S-Dist = $\text{Sort}(\text{DIST}, \text{'ASCEND'})$

② $\text{dist}_{\text{short}}$ \rightarrow distances \rightarrow short

(S-DIST, S-INV) = $\text{Sort}(\text{DIST}, \text{'ASCEND'})$

③ Regression

$$\text{return } \frac{1}{K} \sum_{j=1}^K Y_{TR}(S_{FDX}(j))$$

④ Binary classification

$$\text{return } \text{Sign} \sum_{j=1}^K Y_{TR}(S_{FDX}(j))$$

	1	2	3	4	
Sort	8	3	5	2	0
	4	3	1	2	0
	0	2	3	5	8

S-DIST S-IDX

interest the first k indexes give me
I have k possible index
if I have regression problem: do average

classification:

sum = k , might have $\{+1, -1\}$ on
boundaries

k must be odd number

Half vote for +1

+1 -1

so k must be odd

$O(n \log n)$ sort

if change k , algorithm to predict output
changes

If try and classify with $k=1$
end up having: (KNN classification)
in python notebook
Implementation: $x \rightarrow$
Session: 14 Oct. 2024

$$(x_i, y_i)_{i=1}^m \rightsquigarrow f: x \rightarrow y$$

$$f(x_{\text{new}}) \sim Y_{\text{new}}$$

$x \in \mathbb{R}^d$ $x = (x^1, \dots, x^d)$: collection of
number
Component of vector

$$Y \subseteq \mathbb{R} \quad Y = \{\pm 1\} \leftarrow \text{if } y \geq 0 \rightarrow +1 \quad \text{if } y < 0 \rightarrow -1$$

$$(x_i, y_i)_{i=1}^m \sim P^m \quad i.i.d.$$

① $P(x, y)$ prob distribution on (x, y)
all data are identical & independent

$$= P_x(x) P(y|x)$$

② $l: \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}_+, l(f(x), y)$

$\min E(f)$
 $f: X \rightarrow Y \downarrow$ expected error $(x, y) \text{ w.r.t.}$
 future data

Given $K \in \mathbb{X}, S_n \sim P^m$
 $S_n \sim P^m$

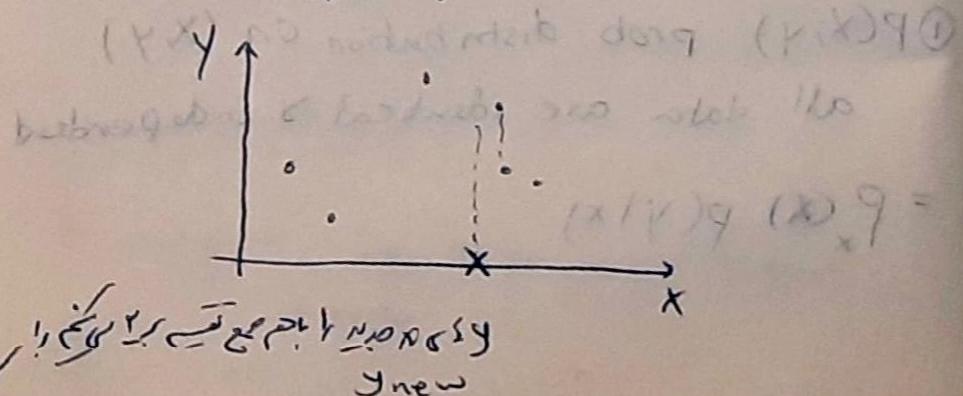
K-NN

$$f(x) = \sum_{i=1}^K y_i / K$$

$I_x \subset \{1 - m\}; d_i = \|x - x_i\|^2; i=1 - m$

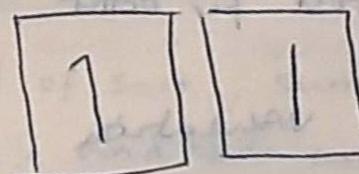
$$(\tilde{I}, \tilde{d}) = \text{Sort}(d_1 \dots d_m)$$

$$I_{x, k} = \tilde{I}(1 - k)$$



① $y = K-\text{NN} (\hat{x}, \hat{y}, x, k)$

close in input, so close output



بینک پیکسل
بینک نکس

هستم که این درست می کند distance

$$(\hat{x}, \hat{y}) = (x_i, y_i)_{i=1}^m$$

$$\hat{y} = (\hat{y}_1, \hat{y}_2, \dots, \hat{y}_m)$$

$$\textcircled{2} \quad \min \|\hat{y}_j - y_{j,k}\|^2 \quad \text{for } j=1 \dots m$$

$$(2) \quad \hat{y} = \|\hat{y}_{j,k} - y_j\|$$

Take Data

half for function

try for find K

an half for test for better

Hold-out cross validation

→ Hold apart of data out

X_{val} → validation set.

$$\hat{Y} = \left(\begin{array}{c} \hat{X}_1 \\ \vdots \\ \hat{X}_T \end{array} \right) \rightarrow \hat{Y}_{\text{validation}}$$

split random

number time split

$$\hat{Y}_{v,k} = K\text{-NN}(\hat{X}_1^q, \hat{Y}_1^q, \hat{X}_v^q, k)$$

Validation set

$$\|\hat{Y}_v^q - \hat{Y}_{v,k}^q\| = \hat{\epsilon}_{v,q}^k$$

$$K=1 \quad \hat{\epsilon}_{v,1}^{(1)} \quad \hat{\epsilon}_{v,1}^{(2)}$$

$$\hat{\epsilon}_{v,2}^{(1)} \quad \hat{\epsilon}_{v,2}^{(2)}$$

$$\hat{K}_1 \quad \hat{K}_2$$

instead of sort, sum the column of
above, and take average

$$\frac{1}{Q} \sum_{K=1}^Q \min \frac{1}{Q} \sum_{q=1}^Q \|\hat{Y}_v^q - \hat{Y}_{v,K}^q\|^2 =$$

$$\frac{1}{Q} \hat{\epsilon}_{v,Q}^k$$

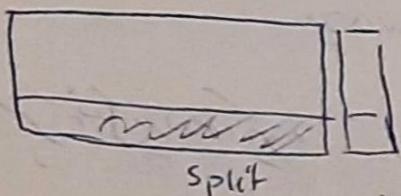
K might be very different.

How choose Q , q

repeat to stabilize it

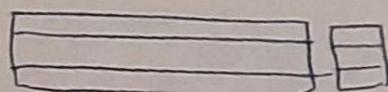
last step caple of variation

First: sequential version



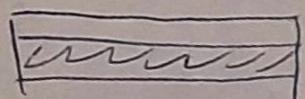
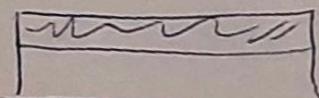
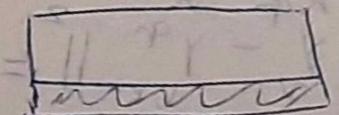
Q fix α , π
fix Π

Sequential split j^*



$Q=3$

Q -Fold



training over all item \rightarrow

$m=12$ $Q=3 \rightarrow$ training = 8 val = 4
Cardinality Cardinality
Large or small cardinality cardinality

$Q=12$ # train = 11 # val = 1

choose
if you second, larger average

splitting is good for a few data

leave-one-out 100

Cross validation.

How choose algorithm:

1 split at random

Q split random

Q split sequentially

may be the leave-one-out

as most large (P, X) tip

where $\rightarrow (P, X) \checkmark$

(P, X) without

15 Oct

$$(\hat{x}_i, \hat{y}_i) \underset{i=1}{\overset{N}{\rightsquigarrow}} \hat{f}(x) = \sum_{i \in I_x} \frac{y_i}{K}$$

$\hat{x}, x, K \rightarrow I_{x,K} \subset \{1 - n\}$ function

for $i=1 \dots n$

$$d_i = \|x - x_i\| \quad (d = d_1 \dots d_n)$$

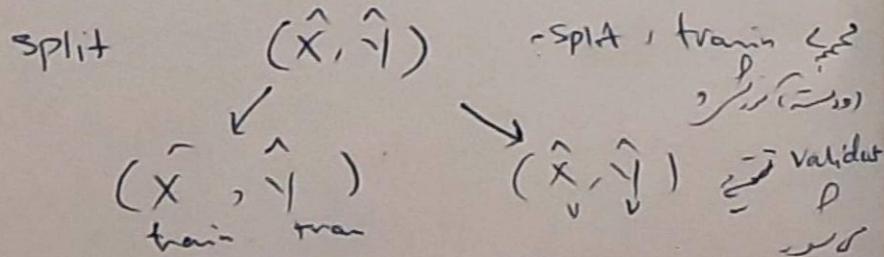
$$(\hat{f}, \hat{d}) = \text{sort } (d_1 \dots d_n)$$

↓ index of closes
closes get sorted number

$$I_x = \tilde{I}(1-K)$$

$$k(\hat{x}, \hat{y}, x, K) = y$$

How do u choose K :



Hold cross validation

$$\hat{K} = \text{HO-CV}(\hat{x}, \hat{y})$$

~~SPLIT~~ (\hat{x}, \hat{y})

$$(\hat{x}_v, \hat{y}_v), (\hat{x}_v, \hat{y}_v) = \text{SPLIT}(\hat{x}, \hat{y})$$

for $K=1 \dots n_{\text{max}}$ for v (j)
 $\text{sort } \{\hat{d}\}_{v=1}^{n_{\text{max}}}$

$$KNN(\hat{x}_v, \hat{y}_v, \hat{x}_v, K_v) = \hat{y}_{v,K}$$

$$\epsilon_K = \|\hat{y}_{v,K} - \hat{x}_v\|^2$$

$$(\tilde{I}(1), \epsilon(1)) \leftarrow \text{the first one}$$

$$= \text{Sort } (\epsilon)$$

$$\hat{K} = \tilde{I}(1)$$

values sorted
first is important

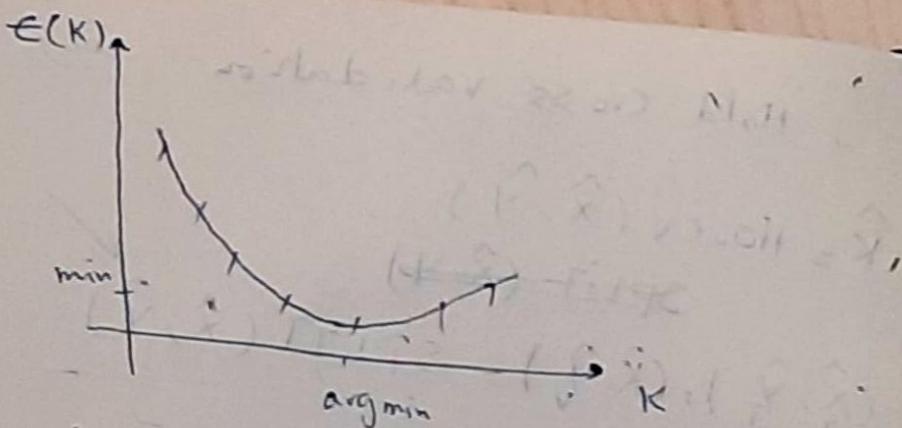
for $K=1 \dots n$

$$K-NN(\hat{x}_v, \hat{y}_v, \hat{x}_v, K_v) = \hat{y}_{v,K}$$

$$\epsilon(K) = \|\hat{y}_{v,K} - \hat{x}_v\|^2$$

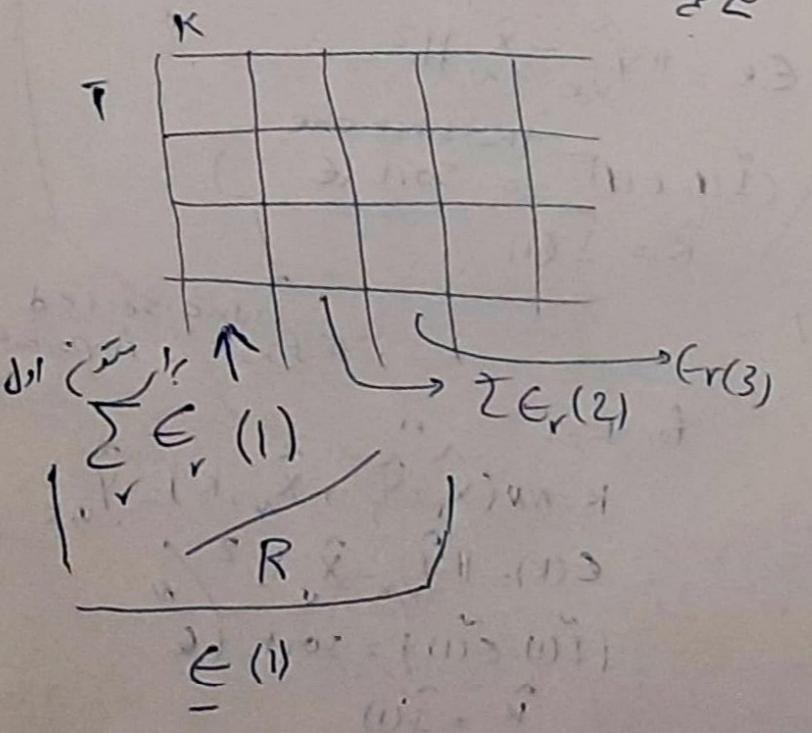
$$(\tilde{I}(1), \epsilon(1)) = \text{sort } (\epsilon)$$

$$\hat{K} = \tilde{I}(1)$$



\leftarrow round to vector for \hat{Y}

$$\epsilon_r(K) = \|\hat{Y}_{v,k} - \hat{X}_v\|_n^2 \rightarrow \hat{Z}_v$$



not big dataset, so you must repeat after time
split.

$$10^{10} \rightarrow \frac{999999+1}{999998+2} \rightarrow k_1 \quad k_2$$

best of prediction

$$\mathcal{D}_n = (X_i, Y_i)_{i=1}^n \sim P^n$$

$$\min_{f(x,y)} E(f) = \mathbb{E}_{(x,y) \sim P} [(Y - f(x))^2] =$$

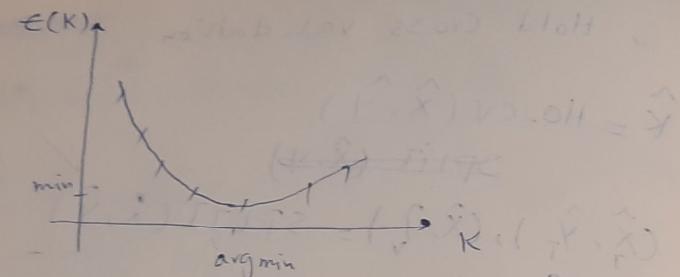
$$\int (Y - f(x))^2 P(x, y) dx dy$$

↓
best of best

can't compute

$$Y_i = f_x(x_i) + \delta_i \quad x_i \sim P_x \quad \delta_i \sim N(0, \sigma^2)$$

$$Y_i \sim N(f_x(x_i), \sigma^2)$$



$$\epsilon_r(k) = \|\hat{Y}_k - \hat{X}\|_n^2$$

K	1	2	3	4	5	6	7	8	9	10
$\epsilon_r(k)$	13	13	13	13	13	13	13	13	13	13
$\sum \epsilon_r(k)$	130	130	130	130	130	130	130	130	130	130
R	130	130	130	130	130	130	130	130	130	130
$\epsilon_r(1)$	13	13	13	13	13	13	13	13	13	13
$\epsilon_r(2)$	13	13	13	13	13	13	13	13	13	13
$\epsilon_r(3)$	13	13	13	13	13	13	13	13	13	13
$\epsilon_r(4)$	13	13	13	13	13	13	13	13	13	13
$\epsilon_r(5)$	13	13	13	13	13	13	13	13	13	13
$\epsilon_r(6)$	13	13	13	13	13	13	13	13	13	13
$\epsilon_r(7)$	13	13	13	13	13	13	13	13	13	13
$\epsilon_r(8)$	13	13	13	13	13	13	13	13	13	13
$\epsilon_r(9)$	13	13	13	13	13	13	13	13	13	13
$\epsilon_r(10)$	13	13	13	13	13	13	13	13	13	13

$$\sum \epsilon_r(k) = 130$$

$$\epsilon_r(1) = 13$$

not big dataset, so you must repeat after time
split.

$$999999 + 1 \rightarrow k_1$$

$$999998 + 2 \rightarrow k_2$$

best of prediction

$$\sum_{i=1}^n (x_i, y_i)^n \sim P^n$$

$$\min E(f) = \mathbb{E}_{(x,y) \in P} [(y - f(x))^2]$$

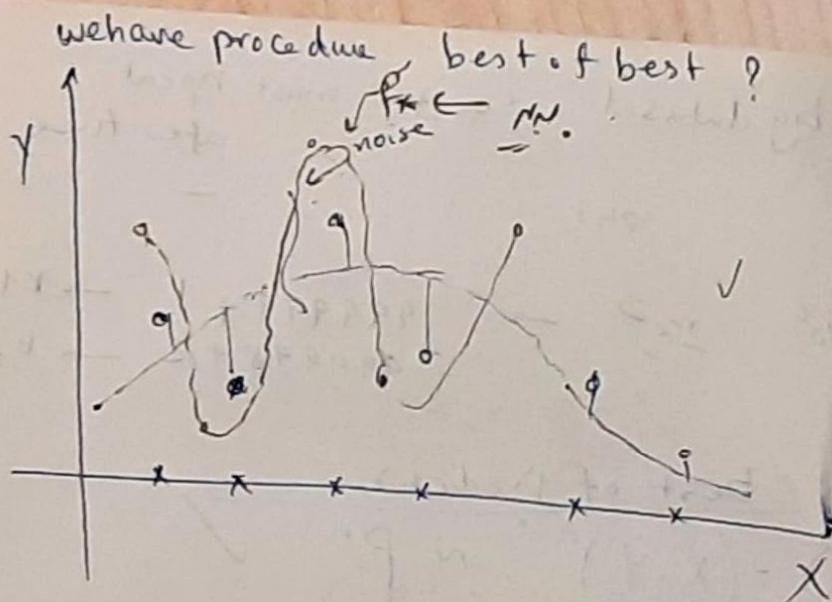
$$(y - f(x))^2 P(x, y) dx dy$$

best of best

can't compute

$$y_i = f_*(x_i) + \delta_i \quad x_i \sim P_x \quad \delta_i \sim N(0, \sigma^2)$$

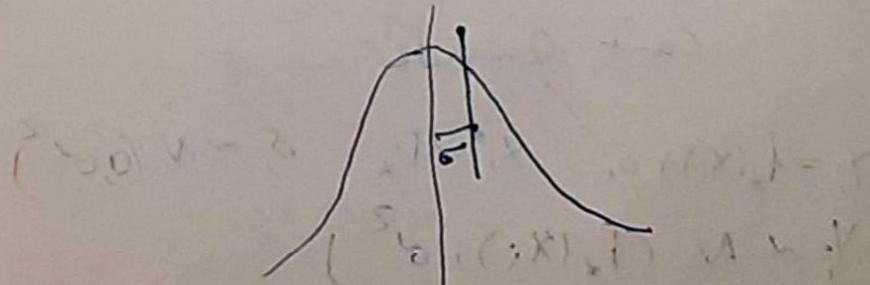
$$y_i \sim N(f_*(x_i), \sigma^2)$$



best: good prediction;

$$K = \arg \min_{K=1 \dots n} \epsilon(\hat{f}_K)$$

Can't compute because we don't have infinite data.



$$Y_i \sim N(f_x(x_i), \sigma^2) \quad \text{noise level}$$

فقط برصاص این طور میتوانیم این داده را در نظر گیری کرد

$$f_K(x) = \sum_{i \in I_{x,K}} f_x(x_i) \quad K \sim f_*(x)$$

و کمpute x

و define \check{f}

و gentle function

can get very big change

و ave fun

$$\frac{\Delta K}{n}$$

$$|\bar{f}_K(x) - f_*(x)|^2 = \frac{(\Delta K)^2}{n}$$

maximum change of function

پس از اینجا اگر ΔK بزرگ شود

$$\frac{|f_k(x) - f_k(x')|}{|x-x'|} \leq L_*$$

I want to compare my algorithm with noise or not

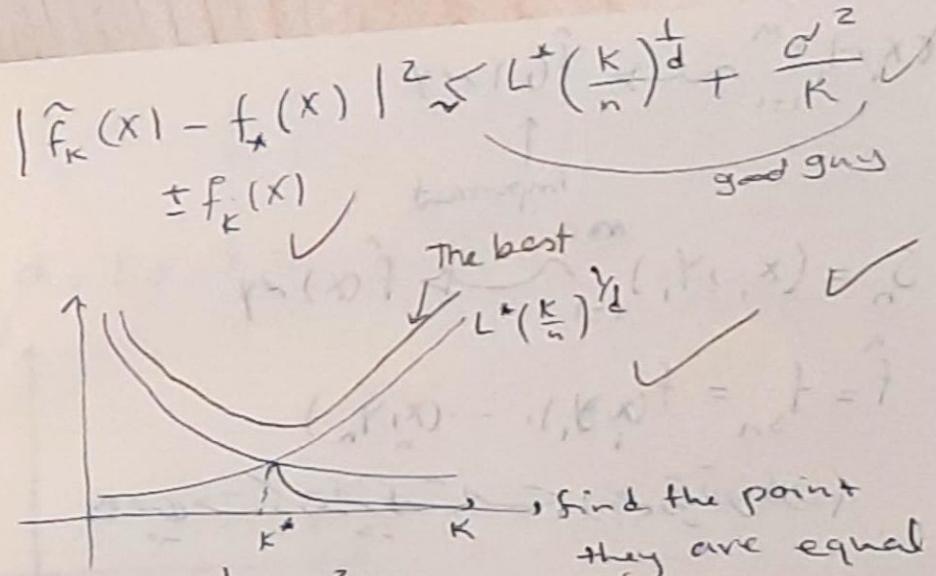
$$\hat{f}_k(x) = \sum_{\substack{i \in I \\ x_i \in K}} \frac{y_i}{K} \quad \begin{matrix} \rightarrow \text{with noise} \\ \text{according to IV} \end{matrix}$$

$$f_k(x) = \sum_{\substack{i \in I \\ x_i \in K}} \frac{f_k(x_i)}{K} \quad \begin{matrix} \text{without noise} \\ \text{according to IV} \end{matrix}$$

$$|\hat{f}_k(x) - f_k(x)|^2 \leq ?$$

$$= \sum_{i \in I_{J \times K}} \frac{\delta_i}{K} \rightarrow \text{variance}$$

$$\mathbb{E} |\hat{f}_k(x) - f_k(x)|^2 = \frac{1}{K^2} \mathbb{E} \left| \sum_{i \in I_{x, \epsilon}} \delta_i \right|^2 = \frac{\sigma^2}{K}$$



$$L^* \left(\frac{K}{n} \right)^d = \frac{\sigma^2}{K}$$

$$K_* = \left(\frac{\sigma^2 n^d}{d} \right)^{\frac{1}{d+1}}$$

\hat{K} is good way to find good guy.

when function change?

for $f_k(x)$

K LIST	1	1	1	1	1	1	1	1
Dist	0	1	2	3	4	5	6	7
59	58	16	-	-	(*)	(*)	(*)	8

Best key

$$(x_i, y_i)_{i=1}^n \rightsquigarrow \hat{f}(x) \approx y$$

↑
Important

$$S_n = (x_i, y_i)_{i=1}^n \rightsquigarrow \hat{f}(x) \approx y$$

$$\hat{f} = \hat{f}_{S_n} = f(x_1, y_1), \dots, (x_n, y_n)$$

دستورات داده شده برای تابع پیش‌بینی

$$S_n = (x_i, y_i)_{i=1}^n \sim P^n \rightarrow \hat{f}(x) \approx y, x \in \mathbb{R}^d, y \in \mathbb{R}$$

solve

$$\min_{\hat{f}} L(\hat{f}) \quad L(\hat{f}) = E_{(x,y) \sim P} [\ell(\hat{f}(x), y)]$$

$$f: x \rightarrow \mathbb{R}$$

$$\text{given } S_n \sim P^n$$

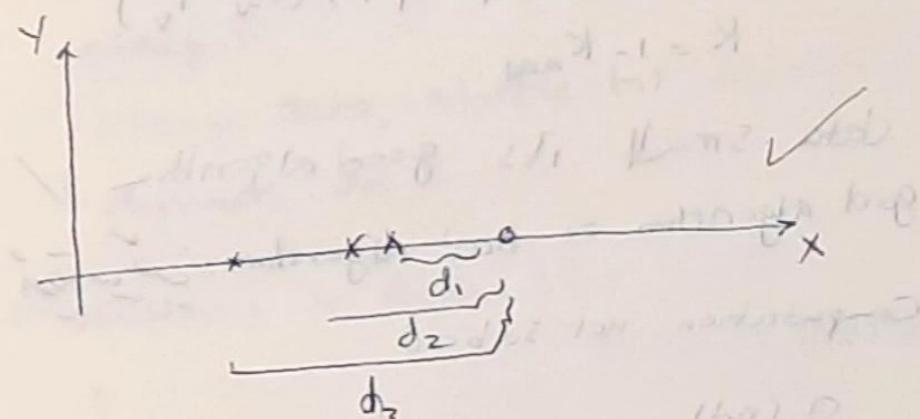
$$\begin{bmatrix} ? \\ P^n \end{bmatrix}$$



$$\hat{f}_K(x) = \frac{1}{K} \sum_{i=1}^K y_i$$

برابر

$$d_i = \|x - x_i\|, d = (d_1, \dots, d_n)$$

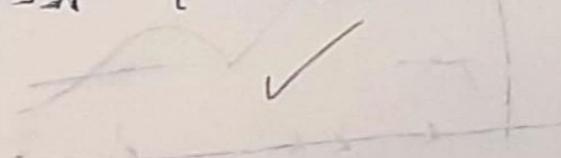


$$(\tilde{I}, \tilde{d}) = \text{Sort}$$

$$\tilde{x} = \tilde{I}(1, \dots, K)$$

اگر لایبل از y بتواند بر K بخشید

$$[\text{Sign } f(x) =]_n = \begin{cases} 1 & f(x) > 0 \\ -1 & f(x) \leq 0 \end{cases}$$



finding K :

$$K = \text{Ho-CV}(\hat{x}, \hat{y}) \quad \checkmark$$

$$\text{split}(\hat{x}, \hat{y}) \rightsquigarrow (\hat{x}_T, \hat{y}_T), (\hat{x}_V, \hat{y}_V)$$

$$K = 1 - K_{\text{max}}$$

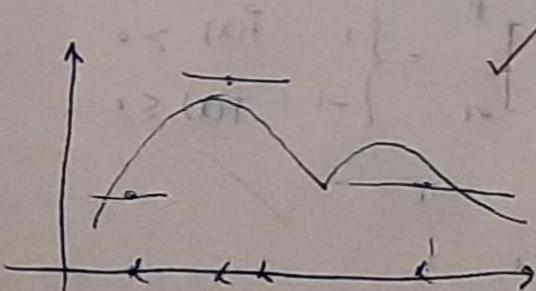
data small its good algorithm
good algorithm or bad algorithm, ~~it's up to you~~

computation not so bad.

$$\mathcal{O}:(nd)$$

big data,

nearest neighbour is flexible if you
have enough data.



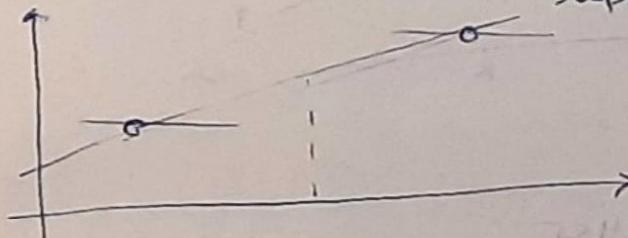
KNN can approximate many functions
provided enough data.
nervous! - (KNN is very nervous!)

change K

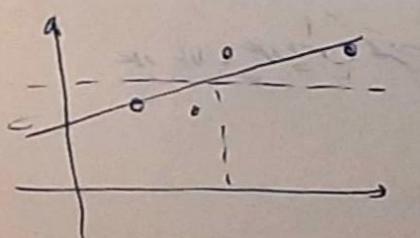
change data, a little bit always
depends on K

sensitive to data! (for small K)

suppose give you
this pic
(2)



The only way to start to ~~be~~ is
small - \rightarrow KNN on \sqrt{p} , p is the size



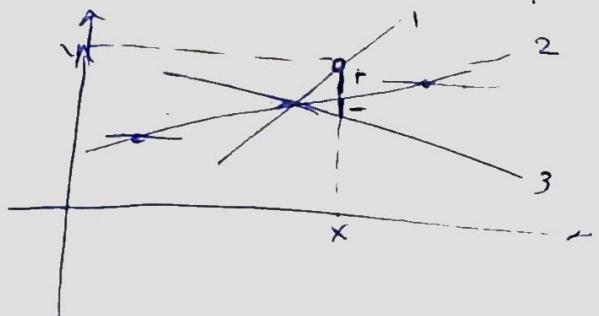
? Small data
↳ line
big data
K-NN

لعنود ون ایجاد فریم

to find line : 2 point

finding line is global , not good
not bad

if I assume linear model this is
global , am present ... ?



2 is better

$$wx + b \rightarrow$$

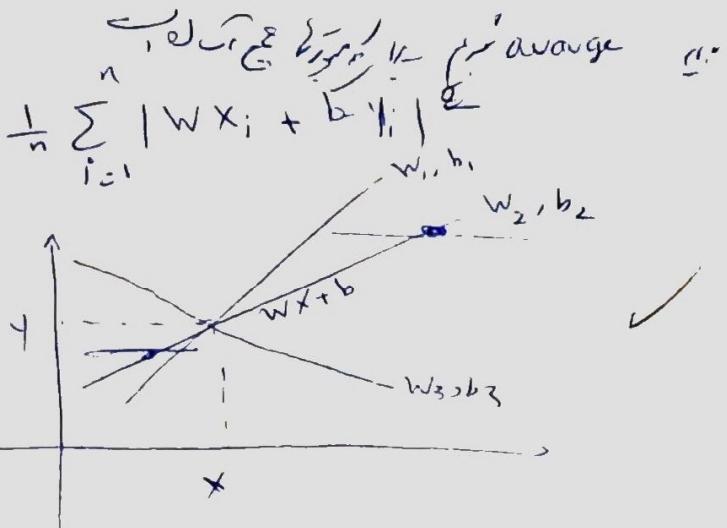
$$bx$$



$$wx_i + b - y_i$$

عده داشتیم که می خواهیم این را کم کنیم

$$\left| wx_i + b - y_i \right|^2$$



I judge line if close to my data
or not

فیض

$$\frac{1}{n} \sum_{i=1}^n (wx_i + b - y_i)^2$$

How to you find good w

delete b:

$$\frac{1}{n} \sum_{i=1}^n (wx_i - y_i)^2$$

lets write function and take derivative.

function of w



Loss of data

$F(w)$, instructive name is $L(w)$

$\hat{L}(w)$ is agreed name for loss function of data $\Rightarrow \frac{1}{n} \sum_{i=1}^n (wx_i - y_i)^2 = \hat{L}(w)$

$$\frac{d\hat{L}(w)}{dw} = \frac{1}{n} \sum_{i=1}^n 2x_i (wx_i - y_i) = 0 \Rightarrow$$

$$0 = \frac{2}{n} \sum_{i=1}^n x_i (wx_i - y_i) \Rightarrow \left(\sum_{i=1}^n x_i^2 \right) w =$$

$$\frac{2}{n} \sum_{i=1}^n x_i y_i$$

$$\left(\sum_{i=1}^n x_i^2 \right) w = \sum_{i=1}^n x_i y_i$$

$$\hat{w} = LS(\hat{x}, \hat{y})$$

سچیز: تحریک ایندیکی

$$\frac{1}{n} \sum_{i=1}^n (w^T x_i - y_i)^2 = \hat{L}(w)$$

scatter product, turn w to w^T to make scatter

derivative of many dimension
its vector first of all

built vector, frequency gradient

I:

$$\nabla \hat{L}(w) = \left\{ \frac{\partial}{\partial w_i} \hat{L}(w) = \frac{2}{n} \sum_{i=1}^n 2(w^T x_i - y_i) x_i \right\}$$

$$w^T x = \sum_{i=1}^n w_i x_i$$

$$\frac{\partial}{\partial w_d} \hat{L}(w) = \frac{2}{n} \sum_{i=1}^n 2(w^T x_i - y_i) x_i$$

\Rightarrow



$$\nabla \hat{L}(\omega) = \left\{ \begin{array}{l} \frac{\partial}{\partial w_i} \hat{L}(\omega) = \frac{1}{n} \sum_{i=1}^n 2x_i^T (\omega^T x_i - y_i)^2 \\ \frac{\partial}{\partial w^d} \hat{L}(\omega) = \frac{1}{n} \sum_{i=1}^n 2x^d (w^T x_i - y_i) \end{array} \right.$$

$$\nabla \hat{L}(\omega) = \left\{ \begin{array}{l} \frac{\partial}{\partial w_i} \hat{L}(\omega) = \frac{1}{n} \sum_{i=1}^n 2x_i^T (\omega^T x_i - y_i)^2 \\ \frac{\partial}{\partial w^d} \hat{L}(\omega) = \frac{1}{n} \sum_{i=1}^n 2x^d (w^T x_i - y_i) \end{array} \right.$$

$$\nabla w^T x = \begin{cases} x \\ x^d \end{cases} = x$$

$$\nabla \hat{L}(\omega) = \frac{2}{n} \sum_{i=1}^n x_i (w^T x_i - y_i)$$

$$F(x) = F(x = w^T x) = \sum_{j=1}^n w^j x^j$$

← Transpose

$$\nabla F(\omega) = \left(\frac{\partial F(\omega)}{\partial \omega}, \dots, \frac{\partial F(\omega)}{\partial \omega^d} \right)$$

$$\Rightarrow \nabla w^T x = x$$

which means:

$$\frac{2}{n} \sum_{i=1}^n x_i (w^T x_i - y_i) = 0$$

$$\sum_{i=1}^n x_i (w^T x_i) = \sum_{i=1}^n x_i x_i^T w \quad * \\ \text{and } Z = \text{cov}(x_1, \dots, x_m, w)$$

$$(\dots) \sum_{i=1}^n x_i x_i^T w = \sum_{i=1}^n x_i y_i$$

Compute vector 2 ways: (seudocode has 2 ways)

$$*: \underbrace{x_i^T x_i}_{\substack{1 \times d \\ d \times 1}} \quad , \quad \underbrace{x_i x_i^T}_{\substack{d \times 1 \\ d \times d}}$$

⇒

for $i=1 \dots n$

$$\hat{C}_i = \hat{C}_{i-1} + x_{i-1} x_{i-1}^T$$

$$S_i = S_{i-1} + x_{i-1} (x_{i-1}^T \omega)$$

Code

Linearity of these operations

$$\Rightarrow (\dots) \left(\sum_{i=1}^n x_i x_i^T \right) \omega = \sum_{i=1}^n x_i y_i = h$$

$$\hat{C} \hat{\omega} = \hat{h}$$

taking gradient,

of linear function -

$$\hat{C} \omega = \hat{y}$$

$$d\hat{C} d\omega \quad d\hat{y}$$

Linear system

$$\hat{\omega} = \hat{C}^{-1} \hat{y}$$

it should be invertable.

\Rightarrow you have code.

$$\text{task} \left\{ \begin{array}{l} \hat{\omega} = \text{LS}(\hat{x}, \hat{y}) \\ \text{compute } \hat{C}, \hat{y} \\ \hat{\omega} = \hat{C}^{-1} \hat{y} \end{array} \right.$$

2 minutes.

Probability $\exists \tilde{f}, \tilde{\omega}$

$L(f)$ is too \rightarrow

$$\hat{L}(f) = \frac{1}{n} \sum_{i=1}^n l(f(x_i), y_i)$$

Solve

$$\min_{f: x \rightarrow \mathbb{R}} L(f) \quad L(f) = \mathbb{E}_{(x, y) \sim P} [\ell(f(x), y)]$$

EMPIRICAL RISK MINIMIZATION

$$\min_{\omega \in \mathbb{R}^d} \hat{L}(f_w) \quad \hat{L}(f) = \frac{1}{n} \sum_{i=1}^n \ell(f(x_i), y_i)$$

$$\ell(a, b) = (a - b)^2 \quad \text{Least - SQUARES}$$



LEAST SQUARES - ERM

$$\min_{f: x \rightarrow \mathbb{R}} L(f) \quad L(f) = \mathbb{E}_{(x, y) \sim P} [\ell(f(x), y)]$$

Given $(x, y_1)^n \sim P^n$

$$\begin{cases} \min_{\omega \in \mathbb{R}^d} \hat{L}(f_\omega) \\ f_\omega(x) = \omega^T x \\ \hat{L}(f) = \frac{1}{n} \sum_{i=1}^n \ell(f(x_i), y_i) \\ \ell(a, b) = (a - b)^2 \\ \ell(f(x), y) = (f(x) - y)^2 \\ y \in \{\pm 1\} = (y_f(x) - 1)^2 \end{cases}$$

$$f(x) = w^T x \quad f: x \rightarrow \{\pm 1\}$$

even if data look like,



$$\begin{aligned} \text{Sign}(f(x)) &= 1 \\ \text{Sign}(f(x)) &= -1 \end{aligned}$$

$$\text{Sqr} \quad \textcircled{1} \quad l(f(x), y) = (f(x) - y)^2$$

$$\textcircled{2} \quad y \in \{\pm 1\} = (y f(x) - 1)^2$$

$$\text{Solve } \rightarrow n \times 4 / y \rightarrow \frac{(y f(x) - 1)^2}{y} = 1$$

~~linear system defined by matrix~~

$$*\nabla \hat{L}(w) = 0 \Rightarrow C w = h$$

$$\text{TSU } \hat{C} = \frac{1}{n} \sum_{i=1}^n x_i x_i^T = \frac{1}{n} \hat{X}^T \hat{X}$$

contain all input as rows
dim n & columns each simple
and coordinate vectors

$$\hat{h} = \frac{1}{n} \sum_{i=1}^n x_i y_i = \frac{\hat{X}^T \hat{Y}}{n}$$

Vector number

Covariance matrix of data

second moment of ...

$$\frac{1}{n} \sum_{i=1}^n x_i = \cdot$$

* I'm looking for w

$$\hat{w} = \hat{C}^{-1} \hat{h}$$

$$\nabla \hat{L}(w) = 0 \Rightarrow \hat{C} w = \hat{h}$$

$$\frac{1}{n} \| \hat{X} w - \hat{Y} \|_2^2$$

$$\hat{X} w = \hat{Y}$$

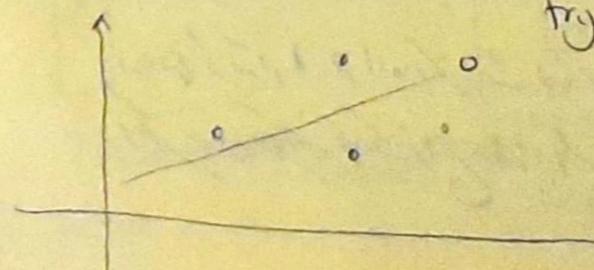
if w satisfy of this exist \rightarrow

$$\hat{L}(w) = 0$$

$$\hat{X} w = \hat{Y} \Leftrightarrow x_i^T w = y_i$$

$$y_n^T w = y_n$$

try to find line through node



do you think you find w that satisfy $\textcircled{**}$: generally no !

According to the graph:

$$n=5, d=1$$

$$n > d$$

number of w is less than the number of point we called

under parameterized

$n > d \Rightarrow$ over determined

$n < d$ under " / over "

$\frac{d-n}{n}$
Solution we have

$n < d$ ~~interpol~~ Can plot

friction (with $d > n$)
- linear regression

one thing you can do

\hat{c} is defined this way:

$$\hat{c} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T = \frac{1}{n} \hat{\mathbf{X}}^T \hat{\mathbf{X}}$$

if matrix is full rank its invertible
can invert it happy.

What is the 2 different problem
 $n > d$ not two differ

$$\begin{bmatrix} \mathbf{x}^T \\ \mathbf{x} \end{bmatrix} = \begin{bmatrix} \hat{c} \\ \mathbf{x} \end{bmatrix} \quad \begin{matrix} \text{dxd} \\ \text{nxd} \end{matrix} \quad \begin{matrix} \text{1 Matrix} \\ \text{sque is dependent} \end{matrix}$$

$n < d$

$$\begin{bmatrix} \mathbf{x}^T \\ \mathbf{x} \end{bmatrix} = \begin{bmatrix} \hat{c} \\ \mathbf{x} \end{bmatrix} \quad \begin{matrix} \text{dxd} \\ \text{nxd} \end{matrix} \quad \begin{matrix} \text{2} \\ \text{no square is independent} \end{matrix}$$

I don't have square number ↑ here

What is Rank of * , ** huge matrix

for * : Rank is d

for ** : Rank is n

how to fix this:

change the algorithm:

$$\min_{w \in \mathbb{R}^d} \hat{L}(w) + \lambda \|w\|^2$$

↓ square norm
↓ SVD

$$\hat{w} = w - \underbrace{\sum_{i=1}^n \sigma_i u_i v_i^T}_{\text{right singular vectors}}$$

$F(\theta)$

$$\min F(\theta) = F(\theta_*)$$

in our case

$$\min_w \hat{L}(w) = \hat{L}(\hat{w})$$

↓ right singular
↓ left singular

what is going on here?

the idea is minimize error

$$\min_w \hat{L}(w) + \lambda \|w\|^2, \hat{L}(w) = \frac{1}{n} \sum_{i=1}^n (y_i - x_i^T w)^2$$

$\hat{L}(w)$

$$\nabla \hat{L}_\lambda(w) = 0 \Rightarrow \hat{C}w = \hat{h}$$

$$\nabla \hat{L}(w) = +2 \hat{X}^T (\hat{X}w - \hat{Y}) = +2(\hat{C}w - \hat{h})$$

$$\nabla \|w\|^2 = 2w$$

λ is constant

remember that s is a norm

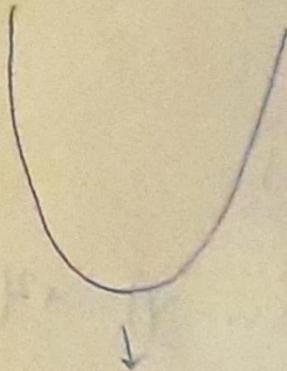
$$\|w\|^2 = \sum_{j=1}^d (w_j)^2 \Rightarrow \text{from (J)}$$

$$\nabla \hat{L}_\lambda(w) = 2(\hat{C}w - \hat{h}) + \lambda 2w = 2(\hat{C} + \lambda I)w - 2\hat{h} = 0$$

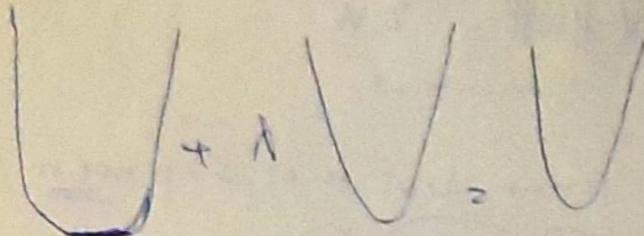
$$(\hat{C} + \lambda I)w = \hat{h}$$

$$\hat{w}_\lambda = (\hat{C} + \lambda I)^{-1} \hat{h}$$

if you want to draw in one dimension
what kind of function is this:



many direction
without minimized



first observation
larger λ
smaller λ

addition

Second ~~and~~ observation

$$\hat{w} = (\hat{C} + \hat{\lambda}I)^{-1} \hat{h} \quad \text{II}$$

$$\hat{C} \rightarrow (\hat{C} + \lambda I)$$

$$\hat{C} = \begin{pmatrix} b_1 & & \\ & \ddots & \\ & & b_d \end{pmatrix}, \quad \hat{C} \rightarrow (\hat{C} + \lambda I)^{-1}$$

$$\hookrightarrow \begin{pmatrix} \frac{1}{b_1} & & \\ & \frac{1}{b_2} & \\ & & \frac{1}{b_d} \end{pmatrix}$$

if rank $< d$

- λ achieves of value,

$$\text{if } \hat{C} \rightarrow (\hat{C} + \lambda I)^{-1}$$

$$\begin{pmatrix} \frac{1}{b_1+\lambda} & & \\ & \frac{1}{b_2+\lambda} & \\ & & \ddots & \frac{1}{b_d+\lambda} \end{pmatrix}$$

$$\delta \gg \lambda \quad \frac{1}{\delta+\lambda} \approx \frac{1}{\delta}$$

$$\delta \ll \lambda \quad \frac{1}{\delta+\lambda} \approx \frac{1}{\lambda}$$

Small value I curve

In special matrix, you can see
what we did,
Uniqueness & invertibility?

its better λ is big or small

2 question: how to PEAK λ

\hat{C} is not diagonal

never complex
~~real~~

$$\begin{pmatrix} \sigma_1 & & & \\ & \sigma_2 & & \\ & & \sigma_3 & \\ & & & \sigma_d \end{pmatrix} V^T = V \begin{pmatrix} \frac{1}{\sigma_1 + \lambda} & & & 0 \\ 0 & \frac{1}{\sigma_2 + \lambda} & & \\ & & \frac{1}{\sigma_3 + \lambda} & \\ & & & \frac{1}{\sigma_d + \lambda} \end{pmatrix}$$

how to solve problem in practice

- solve linear system II

built C

built h

$$\hat{w}_\lambda = \text{RLS}(\hat{x}, \hat{y}, \lambda)$$

$$\hat{C} = \hat{x}^T \hat{x}$$

$$\hat{h} = \frac{\hat{x}^T \hat{y}}{m}$$

$$\{ [D, v] = \text{eig}(\hat{C}) \}$$

↓ procedure in numpy

which return D

$$B = V(D + \lambda I)^{-1} V^T$$

$$\hat{w}_\lambda = B \hat{h}$$

instead of this

$$w = \hat{h} \sqrt{\hat{C} + \lambda I}$$

↑
not inverse

Logistic Regression & SVM

ERM Empirical Risk Minimization

$$\min \mathcal{E}(f) \quad \mathcal{E}(f) = \mathbb{E}_{(x,y) \sim P} [l(f(x), y)]$$

$f: x \rightarrow \mathbb{R}$

Given $(x_i, y_i)_{i=1}^n \sim P^n$

$$f_w(x) = w^T x$$

$w \in \mathbb{R}^d$

$$\hat{\mathcal{E}}(f) = \frac{1}{n} \sum_{i=1}^n l(f(x_i), y_i)$$

$$\min_{w \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n l(w^T x_i, y_i)$$

$\hat{\mathcal{E}}(w)$

$$f(x) = w^T x + b$$

w

$$(w, b) \leftarrow w^T x \sim (x, 1)$$

$$l(w^T x, y) = (w^T x - y)^2$$

$$\min_{w \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n l(w^T x_i, y_i) + \lambda \|w\|^2$$

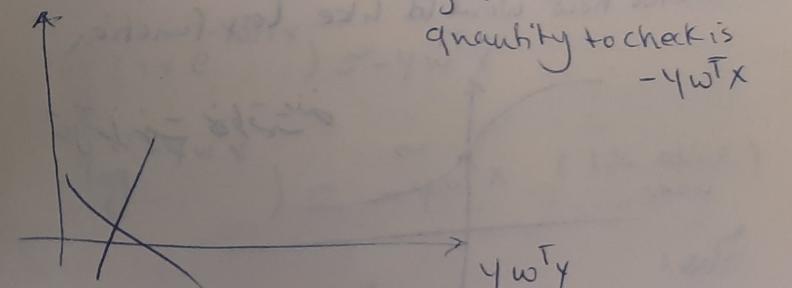
$$\hat{\mathcal{E}}(w)$$

$w \in \mathbb{R}^d$, and we want to plot $y - w^T x$

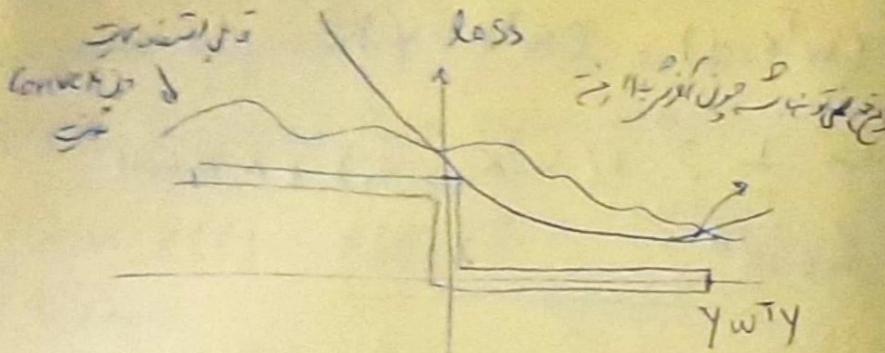
logistic / cross entropy loss

How come up with sth like this:

figure out how to plot quantity to check is $-y w^T x$



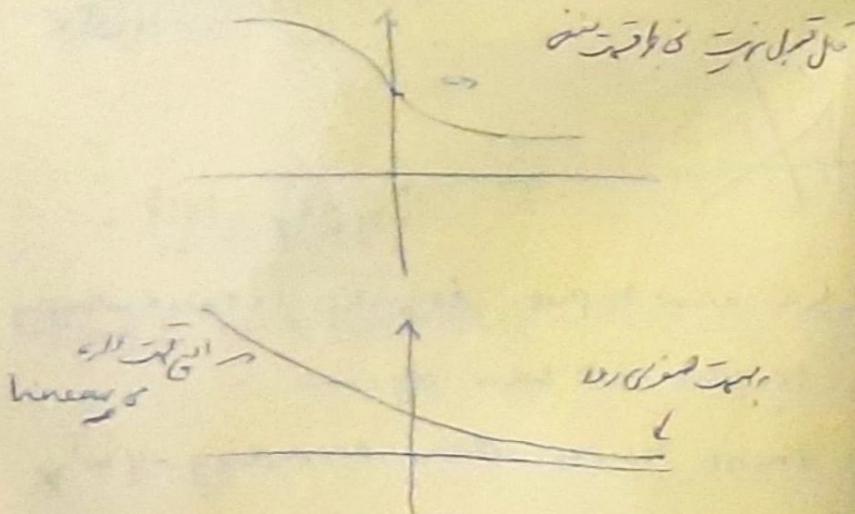
We want to plot logistic (cross entropy) to plot we saw log — think how to draw according $-y w^T x$



$$D y w^T x \leq 0$$

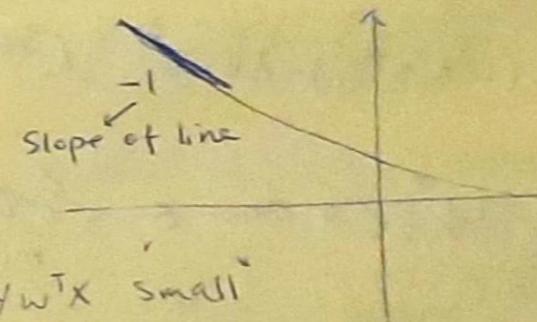
2 problems: Smoothness, Convexity
flat & smile ↗

Second how would like loss function



when $y w^T x$ is very small

$$\log \left(1 + \frac{1}{e^{y w^T x}} \right) \underset{\text{close to } +\infty}{\approx}$$



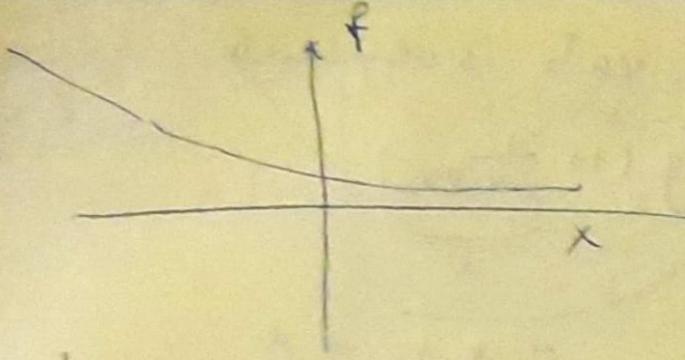
$$\log \left(1 + e^{-y w^T x} \right) \approx -y w^T x$$

$$\log \left(e^{-y w^T x} \right) \approx -y w^T x \quad (\text{like above})$$

$$f(x) = \log(1 + e^{-x})$$

$$\log \left(e^{-x} \right) \approx \log(e^{-x})$$

$$\log \left(1 + e^{-x} \right) \approx \log 1 \approx 0$$



$\log(1+e^{-x})$

$$\min_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{w}^T \mathbf{x}_i, y_i) + \lambda \|\mathbf{w}\|^2$$

gradient descent \leftarrow min

$$\hat{\nabla}_{\mathbf{w}} \hat{\mathcal{E}}(\mathbf{w}) =$$

$$\nabla_{\mathbf{w}} \|\mathbf{w}\|^2 = 2\mathbf{w}$$

$$\nabla_{\mathbf{w}} \hat{\mathcal{E}}(\mathbf{w}) = \nabla_{\mathbf{w}} + \frac{1}{n} \sum_{i=1}^n \log(1+e^{-y_i \mathbf{w}^T \mathbf{x}_i})$$

\leftarrow gradient of $\mathcal{E}(\mathbf{w})$

$$= \frac{1}{n} \sum_{i=1}^n \nabla_{\mathbf{w}} \log(1+e^{-y_i \mathbf{w}^T \mathbf{x}_i})$$

$$= \frac{1}{n} \sum_{i=1}^n \frac{1}{1+e^{-y_i \mathbf{w}^T \mathbf{x}_i}}$$

$-e^{-\mathbf{x}}$

$$= \frac{1}{n} \sum_{i=1}^n \frac{e^{-y_i \mathbf{w}^T \mathbf{x}_i}}{1+e^{-y_i \mathbf{w}^T \mathbf{x}_i}} (y_i \mathbf{x}_i)$$

$$= \frac{1}{n} \sum_{i=1}^n -y_i \mathbf{x}_i / 1+e^{y_i \mathbf{w}^T \mathbf{x}_i}$$

\Rightarrow we descended the gradient.

$$\nabla_{\mathbf{w}} \hat{\mathcal{E}}(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \frac{-y_i \mathbf{x}_i}{1+e^{y_i \mathbf{w}^T \mathbf{x}_i}}$$

$$\nabla_{\mathbf{w}} \hat{\mathcal{E}}(\mathbf{w}) + \lambda \|\mathbf{w}\|^2 = \frac{1}{n} \sum_{i=1}^n \frac{-y_i \mathbf{x}_i}{1+e^{y_i \mathbf{w}^T \mathbf{x}_i}} + 2\lambda \mathbf{w} = 0$$

question is what do we do?

Find \mathbf{w} ?

linear system, a system of linear equation

not linear, non linear in \mathbf{w}

No closed form solution!

optimal condition

Invert square lot to get linear system
linear system is bucket

Gradient Methods

(first order method)

use

first derivative not second derivative,
so call it first order:

$$w_{t+1} = w_t - \gamma_t \nabla_w \hat{\epsilon}(w_t) \quad \gamma_t = \text{constant}$$

$$w_0 = 0$$

Logistic_train($\hat{\epsilon}_\lambda(w)$)

compute
 $w_0 = 0, \gamma_t = (?)$

for $t=1 \rightarrow T_{MAX}$

$$g(t) = \nabla \hat{\epsilon}_\lambda(w_t); \text{ if } g(t) = 0$$

$$w_{t+1} = w_t - \gamma_t g(t)$$

We want w to find the main algorithm
minimize.

$$\nabla_w (\hat{\epsilon}(w) + \lambda \|w\|^2) = \frac{1}{n} \sum_{i=1}^n \frac{-y_i x_i}{1 + e^{-y_i w^T x_i}} + 2\lambda w_0$$

$$\min_{w \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \hat{\epsilon}_\lambda(w^T x_i, y_i) + \lambda \|w\|^2$$

$$\hat{\epsilon}(w)$$

$$F(w)$$

$$F: \mathbb{R} \rightarrow \mathbb{R}$$

$$d=1$$

there is function

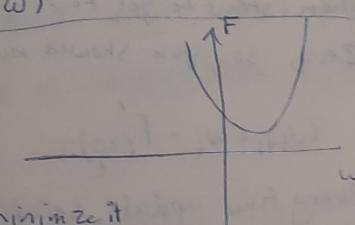
& I want to minimize it

we want to come up to
algorithm.

invent iteration?
iteration?

Go from w_0 to w_1

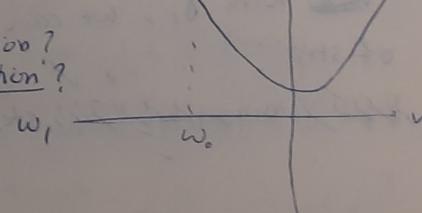
$w_0 \rightarrow w_1$



take a point near w_0

$$\min F(w_1) \leq F(w_0)$$

check if $F(w_1) \leq$



جدا خوبی نداریم w_t , γw_t است خوبی داشتیم
 $F(w)$ نیز خوبی داشتیم $\gamma F(w)$

این را می‌دانیم که w_t را بروز رسانی کنیم

$$w_1 = w_t - \gamma F(w_t)$$

when start to get to min, gradient is zero, so you should not be worry.

$$w_{t+1} = w_t - \gamma F'(w_t)$$

Every time update solution, get steps

if arrive to min, ~~then~~ how to control

~~with~~ with γ , we can control size of step., ✓ ✓

just 1 min or 1000 step ~~or big~~

چون γ کوچک

when function change,

first derivative change a lot,

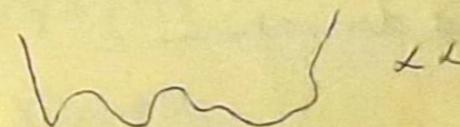
second derivative & big you want step to be small

newton comes up to this equation:

$$w_{t+1} = w_t - \frac{F'(w_t)}{F''(w_t)}$$

~~so~~ γ is not constant here.

local search is good & good for ∇ ~~to~~



Newton method ~~or~~ in $d=1$ / (NM)
if dimension is bigger than 1,

NM $d>1$

$$w_{t+1} = w_t - \nabla F(w_t)$$

Hessian H

$$w_{t+1} = w_t - H(w_t)^{-1} \nabla F(w_t)$$

inverse Hessian at w_t !

Hessian $H(w)$
ATM

in practice, to break it enough

Gradient descent

$$w_{t+1} = w_t - \gamma_t \nabla F(w_t)$$

① $\gamma_t \rightarrow 0$

② if you know maximum change in
maximum second derivative

③ $\gamma_t \approx$ Decaying step size

$$\|\nabla F(w) - \nabla F(w_t)\| \leq c \|w - w_t\|$$

max, max change in 1st derivative

← get second derivative

if you take square

$$\gamma = \frac{1}{2c}$$

$$F(w_t) = \min_{w \in \mathbb{R}^d} F(w) \leq \begin{cases} \frac{1}{t} & \lambda = 0 \\ e^{-t} & \lambda > 0 \end{cases}$$

A large better ??

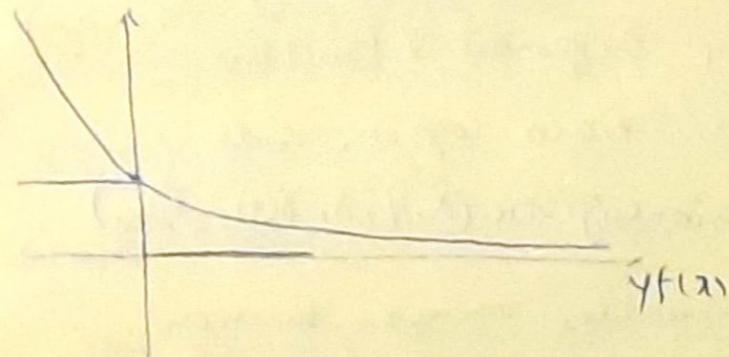
that

$$\|\nabla E_\lambda(w) - \nabla E_\lambda(\underline{w})\| \leq c \|w - \underline{w}\| ?$$

Logistic Regression \rightsquigarrow SVM (Support vector

Machines) loss function

$$y \in \{\pm 1\} \quad l(y, f(x))$$



ERM

$$\min_{w \in \mathbb{R}^d} \frac{\sum_{i=1}^n l(y_i, w^T x_i)}{n} + \underbrace{\lambda \|w\|^2}_{E_\lambda(w)}$$
$$\log(1 + e^{-y f(x)})$$

$$\nabla \hat{E}_\lambda(w) = \frac{1}{n} \sum_{i=1}^n \frac{y_i x_i}{1 + e^{-y_i w^T x_i}} + 2\lambda w = \phi$$

$w_0 = 0 \rightarrow w_t \rightarrow \infty$

$$w_0 = 0, w_{t+1} = w_t - \gamma_t \nabla \hat{E}_\lambda(w_t) \rightarrow T_{\max} \text{ and}$$

have d number against d number
cost of this and intime

time complexity

how many arguments & function

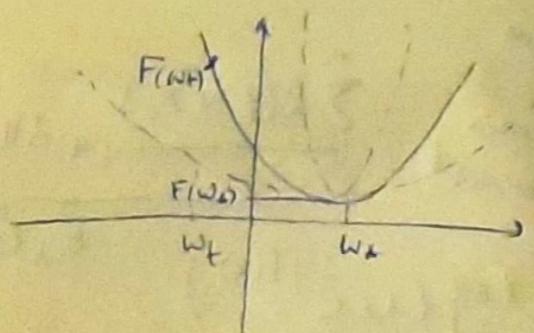
get w , train logistic code

$$w = \text{Train-Logistic}(X, y, \lambda, \gamma(t), T_{\max})$$

think generalize, minimize function

$$F: \mathbb{R}^d \rightarrow \mathbb{R}$$

$$\min_{w \in \mathbb{R}^d} F(w)$$



w_t , compute gradient

$$\nabla F(w_t) \leq \gamma$$

↑
below γ

$$F(w_x) = \min_{w \in \mathbb{R}^d} F(w)$$

$$\left\{ F(w_t) - \min_{w \in \mathbb{R}^d} F(w) \leq \gamma \dots \right.$$

$$\left\{ \|w_t - w_x\| \leq \gamma \right.$$

$$F(w_t) - F(w_{t+1}) \leq \gamma$$

You can not compute

$$\|w_t - w_x\|$$

look at $\{\}$. You start from w_t & get

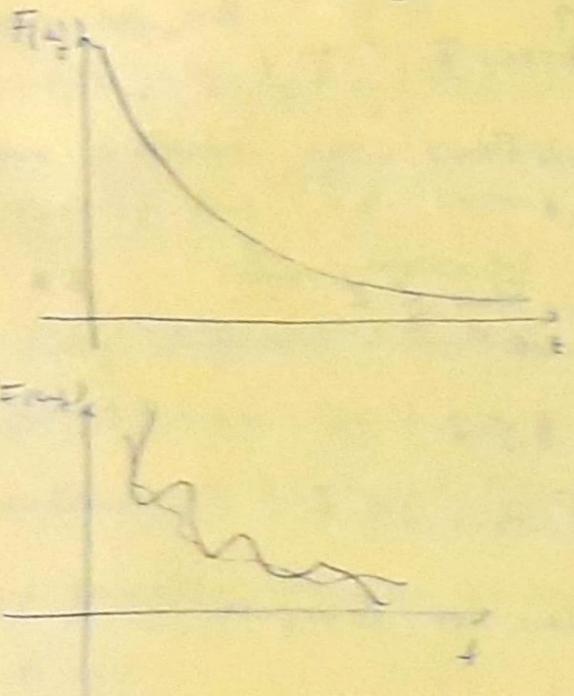
$$F(w_t),$$

according to diagram,

gradient just the direction of $F(w)$

circles

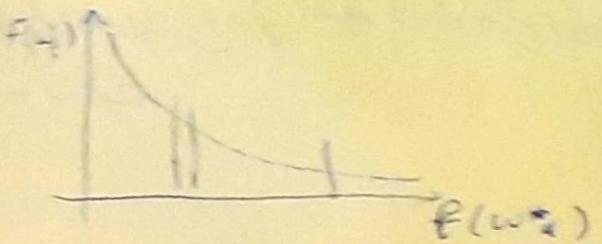
مارشین لرننگ
جیوچیست فورمیٹ اے



so we can have

$$w = \text{Train_logistic}(x, y, \lambda, \Delta t, T, T_{\max})$$

↓
threshold



step size

$$w_{t+1} = w_t - \gamma_t \nabla F(w_t)$$

$$\gamma_t \text{ if } d=1 \Rightarrow \gamma_t = \frac{1}{F''(w_t)}$$

$F''(w_t)$ because لکھا?
how much derivative changed

$$dH, \gamma_t = H(w_t)^{-1} \quad O(d^3)!!$$

→ no $\rightarrow \delta q$ with every iteration

How do break it?
smaller step?

$$\frac{|F(w) - F(v)|}{|w-v|} < c$$

in any dimension

$$\frac{|\nabla F(w) - \nabla F(v)|}{|w-v|} \leq c$$

if I take this $\|\nabla \hat{E}(w) - \nabla \hat{E}(v)\| \leq$

$$\left\| \frac{1}{n} \sum_{i=1}^n \frac{-y_i x_i}{1 + e^{y_i w^\top x_i}} - \frac{1}{n} e^{w^\top x_i} \right\| \leq 2A \|w-v\|$$

Suppose big data set 10,000,000 point

dimension: ~~1000~~¹⁰⁰⁰
 $O(n^d)$

how can we break this.

$$w_{t+1} = w_t - \gamma \nabla \hat{E}_\lambda(w_t)$$

if have 1000000 point in 1000 dimension
how much memory do you need?

when we are working in hurry

dataset big, not time to have
iteration

if you don't have memory how you want
implement?

Split it, I can't compute ∇ ---
compute this for a part of data

for b point

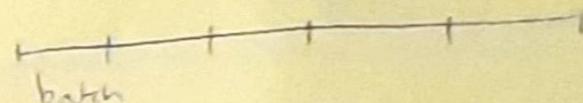
$$\nabla \hat{E}_\lambda(w_t) = \frac{1}{b} \sum_{i=1}^b \underbrace{\frac{-y_i x_i}{1 + e^{y_i w_t^T x_i}}}_{g(i)} + 2\lambda w_t + 0$$

compute a solution, exactly

must comput 1000

≈ 1000

≈ 1000



do & update solution.

per batch

say: $w_t \rightarrow \nabla \hat{E}_\lambda(w_t) = \frac{1}{b} \sum_{i=1}^b \frac{-y_i x_i}{1 + e^{y_i w_t^T x_i}} + 2\lambda w_t$

& for w_2 : per batch

$$\nabla \hat{E}_\lambda(w_2) = \frac{1}{b} \sum_{i=b+1}^{2b} \frac{-y_i x_i}{1 + e^{y_i w_2^T x_i}} + 2\lambda w_2 = 0$$

so for t (per $\frac{1}{b}$) batch = t

$$\nabla \hat{E}_\lambda(w_t) = \frac{1}{b} \sum_{i=t-b+1}^{tb} \frac{-y_i x_i}{1 + e^{y_i w_t^T x_i}} + 2\lambda w_t = 0$$

compute $\nabla \hat{E}_\lambda(w_t)$
= $\frac{1}{b} \sum_{i=t-b+1}^{tb} \frac{-y_i x_i}{1 + e^{y_i w_t^T x_i}}$

Shift

Z

X

C

V

B

N

M

<

>

?

Shift

1

2

3

End

Down

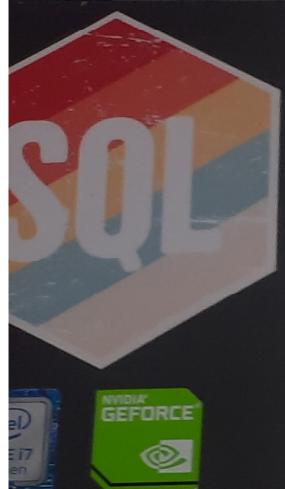
PgUp

PgDn

Enter

Ctrl

Fn



Suppose big dataset. 10,000,000 points
 dimension: $10^{10} \times 10^6 \times 10^3$
 $\approx 10^{19}$ points
 $\approx n^d$

How can we break this?

$$w_{t+1} = w_t - \gamma \nabla_{\lambda} \hat{E}(w_t)$$

If have 1000,000 point in 1000 dimension,
 how much memory do you need?

When we are working in hurry

dataset big, not time to have
 iteration

If you don't have memory how you would
 implement?

Split it, I can't compute ∇E ---

compute this for a part of data

for b point

$$\nabla_{\lambda} \hat{E}(w_t) = \frac{1}{b} \sum_{i=1}^b \frac{-y_i x_i}{1 + e^{y_i w_t^T x_i}} + 2\lambda w_t + 0$$

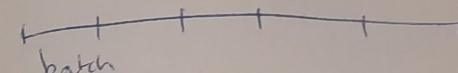
$\underbrace{\quad}_{g(b)}$

compute a solution, exactly

must compute 1000

≈ 1000

≈ 1000



do & update solution,

$$w_t \rightarrow \nabla_{\lambda} \hat{E}(w_t) = \frac{1}{b} \sum_{i=1}^b \frac{-y_i x_i}{1 + e^{y_i w_t^T x_i}} + 2\lambda w_t$$

& for w_2 :

$$\nabla_{\lambda} \hat{E}(w_2) = \frac{1}{b} \sum_{i=b+1}^{2b} \frac{-y_i x_i}{1 + e^{y_i w_2^T x_i}} + 2\lambda w_2 = 0$$

so for t (part of batch = t)

$$\nabla_{\lambda} \hat{E}(w_t) = \frac{1}{b} \sum_{i=t-b+1}^t \frac{-y_i x_i}{1 + e^{y_i w_t^T x_i}} + 2\lambda w_t = 0$$

$\underbrace{\quad}_{= 1/b}$

$T_{\max}(z)$
 ↓
 threshold
 γ_t
 b

how to choose staff

to compute \star , \star must have all data

$\|w_t - w_{t+1}\|$ is cheap but \star, \star are
so cheap but it is easy to stop in chunks
of data.

$$\frac{n}{b} = \text{epoch}$$

talk about iteration

number of receive data (epoch)

I am not changing the?
main issue when descend.
how inverted Newton



What goes on in just doing SGD(-)

$$\gamma_t = \frac{\text{const}}{\sqrt{t}} \leftarrow \text{how ever / very important}$$

mini batches

$$\|\nabla F(w_t)\|^2 \leq C^2 ?$$

I split data,
← mini randomize

stochastic ?

how can change order

Permutation
Version

randomly shuffle data,
when get to the end, shuffle
again, sampling with replacement,
sampling all of them & put it back

Very ^{so} Catastrophic Version

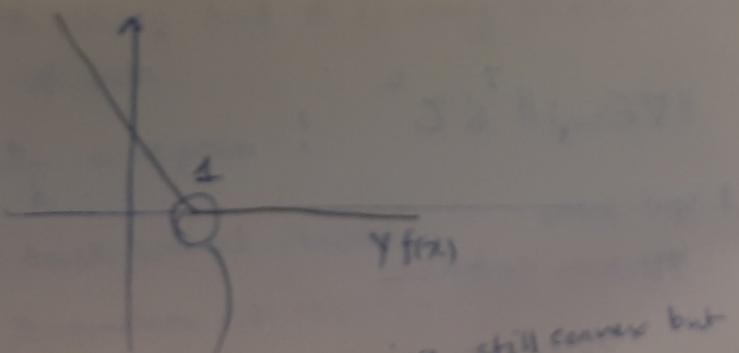
Select 2 points of bigger y with back, 1 point of lower, if not repeated, after a epoch, ...

what you should get:

SVMi.

is another loss function

Vietor like



what are we losing, still convex but
not differentiable

Can not do the gradients

if you do not have derivatives but left
derivatives, you have angle, $\frac{d}{dx} \theta = \frac{1}{\cos \theta}$
every thing set in left derivative

$$|1 - y w^T x|_+ = \max \{0, 1 - y w^T x\} *$$

zero & take number

when $y w^T x$ is $> 1 \rightarrow$ zero

$y w^T x < 1 \rightarrow$ this is line with
slope -1

this $*$ is expression of loss func
just know

$$\min \sum_{i=1}^n |1 - y_i w^T x_i| + \lambda \|w\|^2$$

You can not take gradient anymore

$$\begin{aligned} \frac{\partial E(w)}{\partial w_j} &= \pm \sum_{i=1}^n -y_i x_i^j \text{ s.t. } \lambda \neq 0 \\ \text{exp (left derivative) } &\stackrel{?}{=} \end{aligned}$$

$$w_0 = 0, w_2 = w_1 - \gamma \nabla_{w_1} \hat{J}(w_t) = \gamma x_i^T y_i - 1$$

gradient descent,

ML Algorithmic Toolbox

locality KNN
close Input \rightarrow same output
you average the y_i :

$$\hat{f}(x) = \sum_{i \in I(x)} \frac{y_i}{K}$$

there are a few computation:
what you do to calculate $I(x)$

- {1- sort
- {2- find K

the main problem: local
good \rightarrow you can adapt
simple function (line)

We look at linear models
what are the algorithm you know

RLS,

What is RLS:
just tell: best line to compute
least square error.

forward. Take a line

check if good or not
if good
sum them up

& try to search to find min

$$\min_{w \in \mathbb{R}^d} \sum_{i=1}^n (w^T x_i - y_i)^2 + \lambda \|w\|^2 \quad \lambda \geq 0$$

do you need to λ equal to zero?

if λ is small you have data

if λ is big: how decide it is?

check & & try different λ & compute

solution & check error & and check error

on training set. the
test

how to compute:

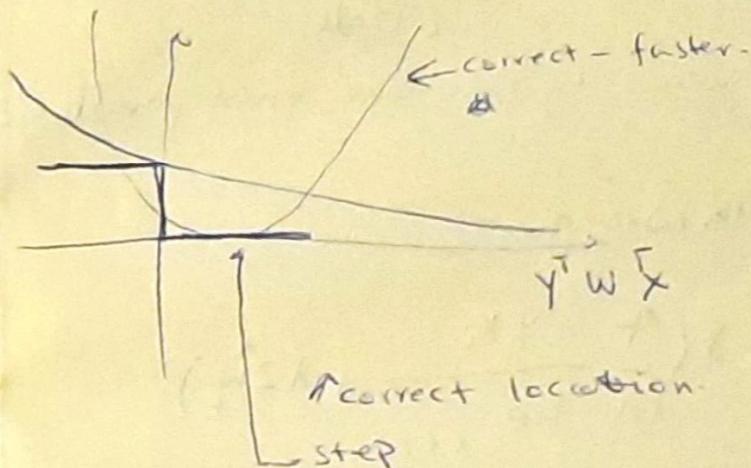
gradient.

to find min & set equal to zero

& find a linear equation

$$\hat{w} = (\hat{x}^T \hat{x} + \lambda I)^{-1} \hat{x}^T \hat{y}$$

RLS : logistic Regression
 $-y_i w^T x_i$
minimize $\log(1+e^{-y_i w^T x_i})$



problem: neither convex & nor differentiable

differentiable

mostly
in regression & classification

$$\sum_{i=1}^n (w^T x_i - y_i)^2 + \lambda \|w\|^2 \quad \lambda \geq 0$$

in classification:

$$\sum_{i=1}^n \log(1+e^{-y_i w^T x_i}) + \lambda \|w\|^2 \quad \lambda \geq 0$$

$$\hat{w}_{t+1} = \hat{w}_t - \gamma x^T (\hat{x} w_t - \hat{y})$$

$\mathcal{O}(mdt)$

↑
size
↓
replace by

If you want to calculate must parallel it.

X sth wrong.

$$\hat{w}_{t+1} = \hat{w}_t - \gamma \left(\sum_{i=1}^n \frac{-y_i x_i}{\|x\|^2} + \lambda 2 \hat{w}_t \right)$$

$$\text{ERM. } \min_{w \in \mathbb{R}^d} \sum_{i=1}^n \ell(y_i, w^T x) + \lambda \|w\|^2$$

SGD

$$\hat{w}_{t+1} = \hat{w}_t - \gamma \nabla \ell(y_i, \hat{w}_t^T x_i) + 2\lambda \hat{w}_t$$

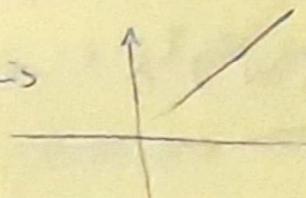
$$f(x) = w^T x$$

function which are linear

there are 2 ways:

what non linear means

instead of this



I have wave

linearity: sum, & multiply by number

sin, abs, cos, square

we have $w \neq x$, we have to prove linearity for $x \neq w$

- non linear:

sin, abs,
cos, square

make the formula not linear

so we have many option to do non linear.

$$f(x) = \vec{w}^T \vec{x}$$

() $\in \mathbb{R}$

$$1, f(x) = \vec{w}^T (\vec{x})^2 \leftarrow \text{respect to } x \text{ non linear}$$

$$2, f(x) = (\vec{w}^T \vec{x})^2 \leftarrow \text{both non linear}$$

$$3, f(x) = (\vec{w}^2)^T \vec{x} \rightarrow \vec{w} \text{ non linear}$$

1

non-linear

You have to code,

instead of using line, I want to
use polynomial & ... and now
no additional code.

we have $x, -x \in \mathbb{R}$

$$f(x) = \vec{w}^T \vec{x}$$

I ask to do $f(x) = ax^2 + bx + c$
invent a way to recycle the code you

have to make this

$$\text{RLS}(\hat{x}, \hat{y}, \lambda) = \hat{\vec{w}}_x$$

use the same code fit the

$$f(x) = ax^2 + bx + c$$

use RLS code & find a, b, c

$$f(x) = \vec{w}^T \vec{z} \quad \hat{\vec{w}}_x = \text{RLS}(\hat{x}, \hat{y}, \lambda)$$

idea is:

$$f(x) = ax + b$$

$$\vec{w}(a, b), \vec{z} = (x, 1)$$

$$\text{new } \vec{w} = (a, b, c) \quad \vec{z} = (x^2, x, 1)$$

fit polynomial in some sense is
line in \mathbb{Z} , linear combination
of monomials, it's a line up
to change values,
 $\vec{w}(a, b, c), \vec{z} = (x^2, x, 1)$

$$\hat{z} = \text{feature} - \text{map}(\hat{x})$$

for any input, create corresponding \hat{z}

How to write polynomial degree Three
number of parameters in model ↑

$$f(x) = ax^3 + bx^2 + cx + d$$

$$w = (a, b, c, d) \quad z = (x^3, x^2, x, 1) \in \mathbb{R}^4$$

$$\hat{z} = \text{Feature} - \text{map}(\hat{x}) \quad w = (s)^T$$

we want to see how powerful these

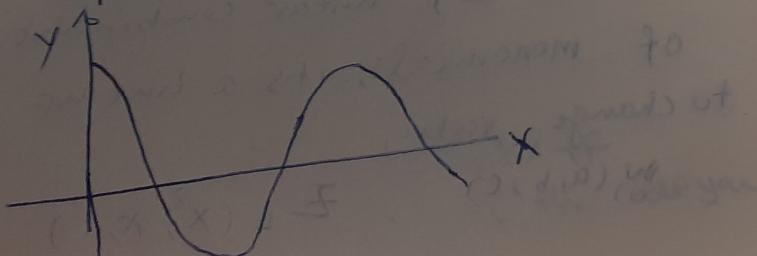
trick is!

this allowed us to go from polynomial

all to a linear form.

even in 1 dimension, z changes.

suppose that we have this function



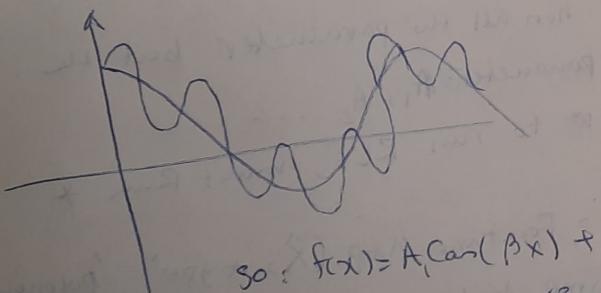
if you need polynomial, which is sin.

$$f(x) = \cos(x)$$

$A \cos(x)$
or we can have amplitude or frequency

$$f(x) = A \cos(Bx)$$

we have sin which modulated by faster sin (or cosi) - S



$$\text{so: } f(x) = A_1 \cos(B_1 x) + A_2 \cos(B_2 x)$$

This function, can we write it with $f(z) = w^T z$, so:

$$w = (A, B, A_2, B_2)$$

$$Z =$$

$$(A)x_2 + (B)t$$

this trick does not work for B
we cannot find β , if we want
to do this:

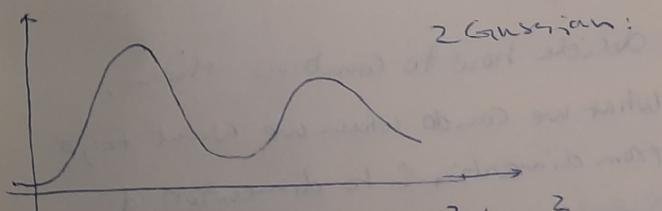
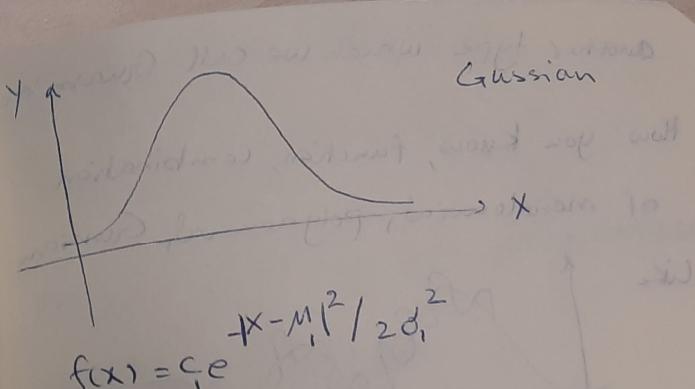
$$w = (A, A_2)$$

$$Z = (\cos(\beta_1 x), \cos(\beta_2 t))$$

non all the parameter but the
parameters A_1, A_2, \dots

to run RLS, must run *

\hat{Z} = Feature-MRP (\hat{x} , 'type', 'polynomial')
as same before

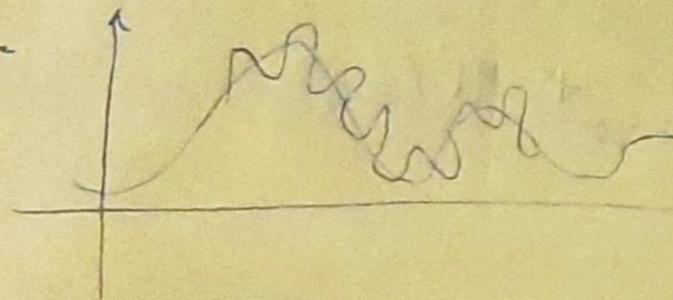


the parameter to find $c_1, c_2, M_1, M_2, \sigma_1^2, \sigma_2^2$
Can I fit (x, y) ?
we can find w , no because only c_1, c_2
are linear, instead of pick 2, pick
many, put over place & combine it.

Another type which we call Gaussian

How you know, function, combination
of monomial, polynomial, Gaussian

like



decide how to combine them,
what we can do when we want to go
from dimension 1 to dimension d:

$$f(x) = ax + bx + c \rightarrow 2 \text{ dimension}$$

$$x = (x_1, x_2)$$

$$f(x) = a(x_1)^2 + b(x_2)^2 + c x_1 x_2$$

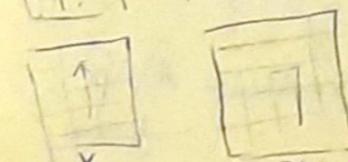
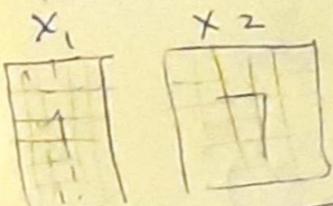
* same trick just the same,

$$f(x) = A_1 C_1 e^{-(\beta_1^T x)} + A_2 C_2 e^{-(\beta_2^T x)}$$

$$f(x) = c e^{-\|x - \mu_1\|^2 / 2\sigma_1^2}$$
$$f(x) = c e^{-\|x - \mu_2\|^2 / 2\sigma_2^2}$$

think always > big over parameterize.

You have 2 images:



Can we use least squares to solve problem

take regression & run it

pixels are not very good, when you
treat as a vector ? ? ?

take windows & take differences,
because there is difference that
change

$$\frac{\begin{array}{|c|c|}\hline x & x \\ \hline \end{array}}{z^1} - \frac{\begin{array}{|c|c|}\hline x & x \\ \hline \end{array}}{z^2}$$

so, we built Z , and bunch of numbers you obtain by manipulating with

$$\boxed{f} \rightarrow \begin{array}{c|c} \text{grid} & - \\ \hline \text{grid} & \left(\begin{array}{c} Z \\ Z \\ Z \\ \vdots \end{array} \right) = Z \end{array}$$

$f(x) = (x)^T$
Gaussian
polynomial
mixed of two,

notice that:

$$\hat{Z} = \text{Feature-Map}(\hat{x}, 'typ', 'p')$$

$\Phi : x \rightarrow Z \rightsquigarrow \text{feature Map}$

$$\Phi(x) = \begin{pmatrix} f_1(x) \\ \vdots \\ f_p(x) \end{pmatrix}$$

called features

$$\int_j X \rightarrow \mathbb{R} \quad j=1 \dots p$$

this is just abstraction, it's just a function

give me three & fours of f

f most could be sin, cos, polynomial, Gaussian, mixed of two,

every thing we discuss:

Whenever you see X ,

Z or $\Phi(x)$ are the same

$$X \rightarrow Z = \Phi(x)$$

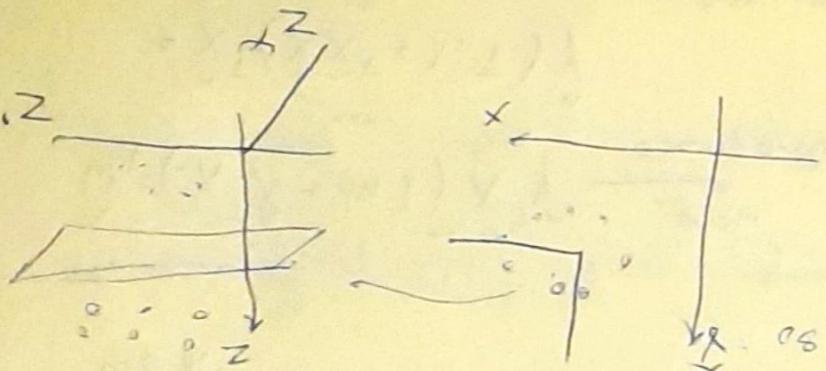
~~mag~~

Feature MAPS \rightarrow KERNELS

$$\min_{W \in \mathbb{R}^d} \sum_{i=1}^n l(x_i^T w, y_i) + \lambda \|w\|^2$$

$$\Phi(x)^T w$$

You can use parameterization to make nonlinear not depend on w



If you can use linear for classify,
then cannot because think of a plate



$\phi(x)^T f - \phi(x')^T f = \phi(x) - \phi(x') = z$ also same class
how do you know?

new idea: The normal choice of (nr) linear
vector structure

also be with else, in which case
dual specific: ϕ
you can use many

$\phi(\mathbf{x})$: closure of graph

each person is graph. If we can
what are some numbers classify people
and ages, can we define the model?

If x is graph, the set of x 's

$\{x\}$, has to be vector

~~anywhere you can~~

$\phi(x) = (x)^T f - \phi(x')^T f = \phi(x) - \phi(x')$
what ever given by fit to feature

Sigh, abs, exponential,

what are except of ϕ
the boundary is a line

$\phi_{\perp}(\mathbf{x})$

We will start w/ \mathbf{x}

but it not linear

this model is parametrize linearly

$\Phi(x) \in \mathbb{R}^p$

$$W^T B^T = (B^T W)^T$$

if Φ is nonlinear,

how to know Φ or linear.

I not know

main point : cross validation,
try on " " to get best

purpose of today

reparameterize problem

and infinit:

$$\min_{\substack{i=1 \\ w \in \mathbb{R}^p}} \sum l(\Phi(x_i^T w, y_i) + \lambda \|w\|^2)$$

remember:

$$\hat{w}_\lambda = (\hat{X}^T \hat{X} + \lambda n I)^{-1} \hat{X}^T \hat{y} \xrightarrow{\text{swap order}} \hat{w}_\lambda$$

$$= \hat{X}^T (\underbrace{\hat{X} \hat{X}^T}_{n \times n} + \lambda n I)^{-1} \hat{y}$$

give name $C \in \mathbb{R}^{d \times d}$

SVD

$$\hat{X} = UDV^T$$

$$\hat{X}^T = VD^T U^T$$

~~skip and prove~~

finish the collection

$$\hat{w}_\lambda = \hat{X}^T C = \sum_{i=1}^n x_i c_i \quad C = (\hat{X} \hat{X}^T + \lambda n I)^{-1} \hat{y}$$

$$\hat{f}_\lambda(x) = x^T \hat{w}_\lambda$$

Combine & & &

$$\hat{f}(x) = x^T \hat{w}_\lambda = \sum_{i=1}^n x^T x_i c_i$$

~~we want to replace x with $\Phi(y)$ in $\hat{f}(y)$~~

$$\hat{\Phi}_{n \times p}$$

$$\hat{w}_\lambda =$$

tediously: rewrite just put $\hat{x} \rightarrow \hat{\Phi}$,
 $x \rightarrow \hat{y}$

Just:

$$\hat{w}_\lambda = (\hat{\Phi}^T \hat{\Phi}_{n \times p} + \lambda n I)^{-1} \hat{\Phi}^T \hat{y}$$

$$(\cdots) = \sum_{i=1}^n \phi(x_i) c_i, \quad c = (\hat{\Phi}^T \hat{\Phi}_{n \times p} + \lambda n I)^{-1} \hat{y}$$

$$\hat{f}_1(\lambda) = \sum_{i=1}^n \phi(x_i)^T \phi(x_i) c_i$$

so just change notation.

observation: if \hat{w}_λ , but $f(w)$ just look at f just means look at prediction of w , what about c , how compute what is the entry of $\hat{q} \hat{q}^T$ matrix

$$(\hat{\Phi}^T \hat{\Phi})_{ij} = \phi(x_i)^T \phi(x_j)$$

inner product between

$$\hat{x} = \begin{pmatrix} x \\ \vdots \\ \hat{x} \end{pmatrix}, \quad \hat{\Phi} = \begin{pmatrix} \phi(x_1) \\ \vdots \\ \phi(x_n) \end{pmatrix}$$

nxd 2vectors

$$\phi(x)^T \phi(x') = \sum_{j=1}^p f_j(x) f_j(x') = K(x, x')$$

give a name

inner product of any x ,

question can we imagine to compute this operation when $j \rightarrow +\infty$

answer if do, doesn't pay attention?

Geometric Series

$$\sum_{j=0}^{\infty} a^j = \frac{1}{1-a} \quad a < 0 < 1$$

Exponential Series

$$\sum_{j=0}^{\infty} \frac{a^j}{j!} = e^a$$

we want to use these 2 to invent, write down the K .

44

$$X = \begin{pmatrix} 0, 1 \end{pmatrix}$$

not only one dimensional

$$\text{invent } \phi_j$$

to do the sum

$$\phi(x)^T \phi(x') = \sum_{j=0}^{\infty} f_j(x) f_j(x') = K(x, x')$$

$$\text{consider } f_j(x) = x^j \quad j = 0, \dots, \infty$$

$$\sum_{j=0}^{\infty} x^j x^j = \sum_{j=0}^{\infty} (x)_j^j = \frac{1}{1-x} \sum_{j=0}^{\infty} w_j x^j$$

$$f(x) = \sum_{j=0}^{\infty} w_j, \quad f_j(x) =$$

if example if $X(0, \infty)$

$$f_j(x) = \left(\frac{x}{c}\right)^j \quad j = 0, 1, \dots, \infty$$

$$\sum_{j=0}^{\infty} \left(\frac{x}{c}\right)^j \left(\frac{x}{c}\right)^j = \sum_{j=0}^{\infty} \left(\frac{x}{c}\right)^j = \frac{1}{1-\frac{x}{c}} = \sum_{j=0}^{\infty} w_j x^j$$

Kernel is easy
but look weird

OK for Exponential series:

$$K(x, x') = e^{-\frac{(x-x')^2}{2\sigma^2}}$$

lets develop square:

$$K(x, x') = e^{-\frac{(x-x')^2}{2\sigma^2}} = e^{\frac{-x^2}{2\sigma^2}} e^{\frac{-x'^2}{2\sigma^2}} e^{\frac{2xx'}{2\sigma^2}}$$

$$e^{\frac{-x^2}{2\sigma^2}} e^{\frac{-x'^2}{2\sigma^2}} \sum_{j=0}^{\infty} \underbrace{\frac{1}{j!} \left(\frac{2xx'}{\sigma^2}\right)^j}_{\text{we have a and need to massage to find } \phi_j}$$

how can

$$= e^{\frac{-x^2}{2\sigma^2}} e^{\frac{-x'^2}{2\sigma^2}} e^{\frac{2xx'}{\sigma^2}} = e^{\frac{-x^2}{2\sigma^2}} e^{\frac{-x'^2}{2\sigma^2}} \sum_{j=0}^{\infty} \frac{1}{j!} \left(\frac{2xx'}{\sigma^2}\right)^j$$

$$= \sum_{j=0}^{\infty} \left(\frac{\sqrt{2} x^j e^{\frac{-x^2}{2\sigma^2}}}{\sqrt{2} \sqrt{j!}} \right) \left(\frac{\sqrt{2} x'^j e^{\frac{-x'^2}{2\sigma^2}}}{\sqrt{2} \sqrt{j!}} \right)$$

$$= \sum_{j=0}^{\infty} f_j(x) f_j(x')$$

Select the relevant expression to
computable code: if I give you

ϕ is Φ

$$f(x) = \sum_{i=1}^n \phi(x)^T \phi(x_i) c_i$$

$$y = C = (\hat{\Phi}\hat{\Phi}^T + n\lambda I)^{-1} \hat{Y} \rightarrow$$

$$(\hat{\Phi}\hat{\Phi}^T)_{ij} = \phi(x_i)^T \phi(x_j)$$

\Rightarrow

$$y = \hat{f}(x) = \sum_{i=1}^n K(x, x_i) c_i$$

$$C = (\hat{K} + \lambda n I)^{-1} \hat{Y}$$

$$\hat{K}_{ij} = K(x_i, x_j)$$

After turn these 3 expression to
pseudo code.

$C = \text{KRLS1}(\hat{X}, \hat{Y}, \lambda, 'Kernel')$
call it $Ker-type$

another function to call

$\hat{K} = \text{Kernel}(\hat{X}, \hat{X}, 'Ker-type')$

↓
call it K

$\Rightarrow C = \text{KRLS1}(\hat{X}, \hat{Y}, \lambda, 'Ker-type')$

$\hat{K} = \text{Kerl}(\hat{X}, \hat{X}, \dots, \dots)$

$$C = \frac{Y}{(\hat{K} + n\lambda I)}$$

~~for loop~~
~~for loop~~
 $K = \text{Ker}(\hat{X}, \hat{X}, 'Ker-type')$

$= \text{KRLS-test}(C, X, \hat{X}, 'Ker-type')$

give new X → then inside can do 3

You have to compute expression.

for $i=1 \dots n$

$K(i) = \text{Kernel}(X, X_i, 'type')$

$y_{temp} = y_{temp} + K(i)C(i)$

Neural networks

KNN is first algorithm - nearest neighbor
local algorithm.

depend on K, we might choose
smooth function, & wish to have
linear function.

No constraints on the family of function
for estimator \hat{f} .

If you have lots of data, its
its local method,

↳ number of data

$\lambda_1 \approx 7, 18 \text{ nov}$

↳ no han of dimension ~~not~~ increase
number of data
If dimension is big, distance become very
large.

or Network
Regularization methods - global

looking for estimator, introduce RLS

We want to minimize empirical Risk:

$$\hat{f} = \arg \min_{\hat{f}} \left(\sum_{i=1}^n \ell(y_i, \hat{f}(x_i)) + \lambda \|\hat{f}\|^2 \right)$$

regularization term

\mathcal{F} : family of function,

this has been general.

↳ RLS (RLS with linear function)
↳ we made specific least square,
square distance loss

first choice:

$$\hat{f} = \arg \min_{\hat{f}} \left(\sum_{i=1}^n (y_i - f(x_i))^2 + \lambda \|\hat{f}\|^2 \right)$$

H : is class of linear function \Rightarrow

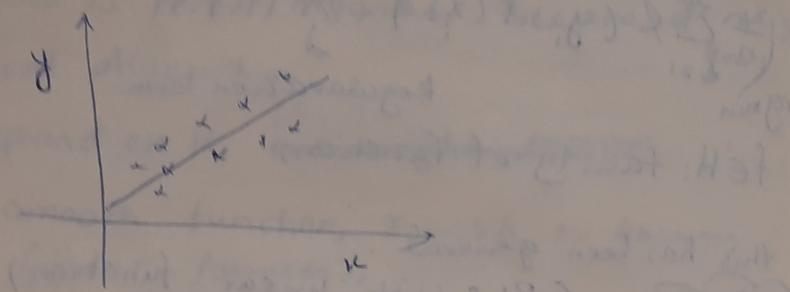
$$\hat{f} = \arg \min_{\hat{f}} \left(\sum_{i=1}^n (y_i - \hat{f}(x_i))^2 + \lambda \|\hat{f}\|^2 \right)$$

$\hat{f} \in H$

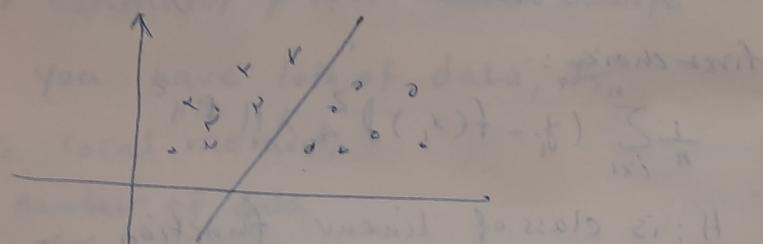
D: Dimension of x

LR logistic function

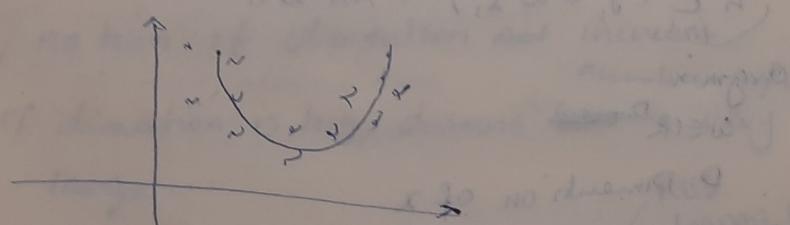
for regression



in binary classification



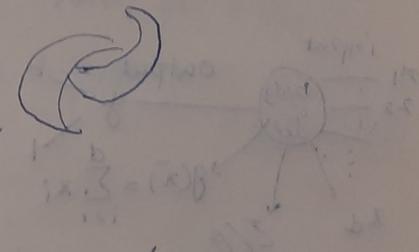
Non-linear



feature map

if problem is not linear so increase dimension d, I can map point to higher dimensional space that one each point is linear.

$$\Phi: \mathbb{R}^d \rightarrow \mathbb{R}^P$$



$$P \gg d$$

replace x with $\Phi(x)$

$$\hat{\omega} = \operatorname{argmin}_{\omega \in \mathbb{R}^P} \frac{1}{n} \sum_{i=1}^n (y_i - \omega^T \Phi(x_i))^2 + \lambda \|\omega\|^2$$

Kernel

If K is Function:

$$K(x_i, x_j) = \Phi(x_i)^T \Phi(x_j)$$

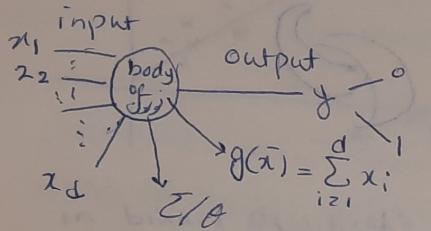
(27:15)

You can describe solution ω as wanted combination with term, ~~temp~~

Neural Networks:

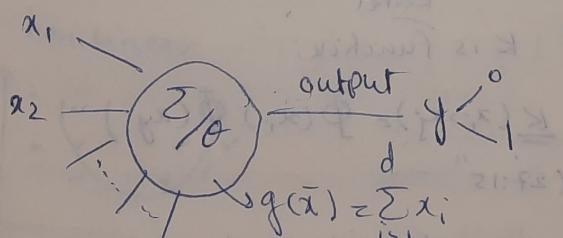
first model (1943) (single layer perceptron)

main element: artificial neuron;



this artificial neuron accept vector $\bar{x} \in \mathbb{R}^d$

$$P(\bar{x}) = \begin{cases} 1 & \text{if } g(\bar{x}) > 0 \\ 0 & \text{otherwise} \end{cases}$$



all input have same importance.

one connect

Complexity of model:

~~connected with AND, OR, XOR, ...~~

Consider simple situation:

A	B	AND	Sum
0	0	0	0
0	1	0	1
1	0	0	1
1	1	1	2

try to use 2 input, on plane:



linear

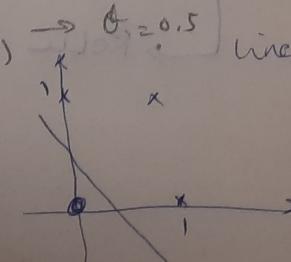
We would like to find separator like *

$$\theta = 1.5$$

OR

A	B	OR
0	0	0
0	1	1
1	0	1
1	1	1

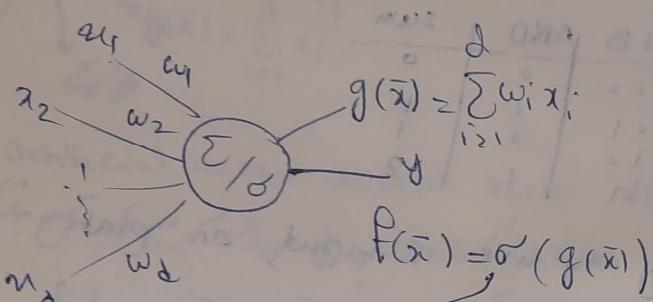
$\theta = 0.5$ linear



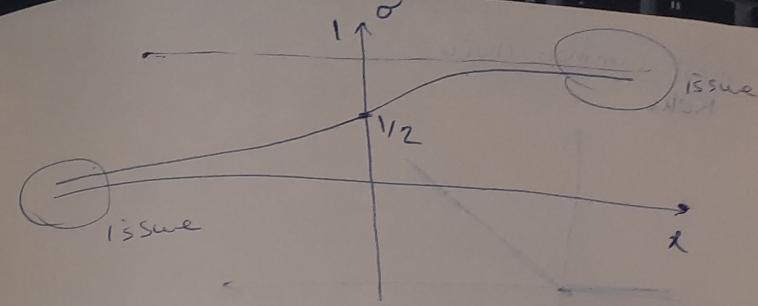
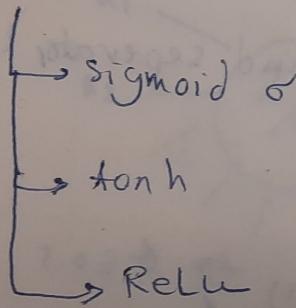
A	B	XOR
0	0	0
0	1	1
1	0	1
1	1	0

non-linear

Put weight



non linear activation function.



$$\sigma(x) = \frac{1}{1+e^{-x}}$$

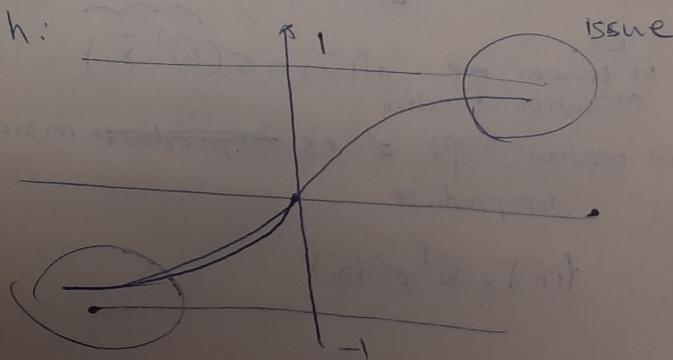
for any x , output is in Range $(0, 1)$
easy to understand

for $x < 0$

$x > 0$

we all need to optimise the weight
need of computing gradient

\tanh :



most common choice

ReLU



$$\text{ReLU}(x) = \max(0, x)$$

4

$$f(\bar{x}) = \sigma(\bar{w}^T \bar{x})$$

We can use feature maps & kernels

In case of feature map:
Feature map:

$$f(\bar{x}) = \bar{w}^T \sigma(\bar{x})$$



non linear model

linear model

non
linear \approx activation function

$$f(\bar{x}) = \sigma(\bar{w}^T \bar{x})$$

in feature maps σ is ~~dot product~~ inside
dot product

$$f(\bar{x}) = \bar{w}^T \sigma(\bar{x})$$

in linear model, extra term, which is Bias
call it w_0 in our model

$$f(x) = w_0 + \underline{\underline{q}}$$

one thing: $f(x) = w_0 + w_0$

re write our function: $f(\bar{x}) = \sigma(\bar{w}^T [1, \bar{x}])$
 $\bar{w} = [w_0, w_1, \dots, w_d]$

in our model we have ~~extra~~

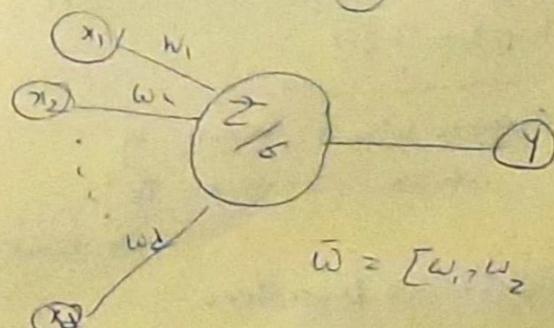
take data & standardize ~~a~~ with respect
of min

trick: solution goes through origin
you have it.

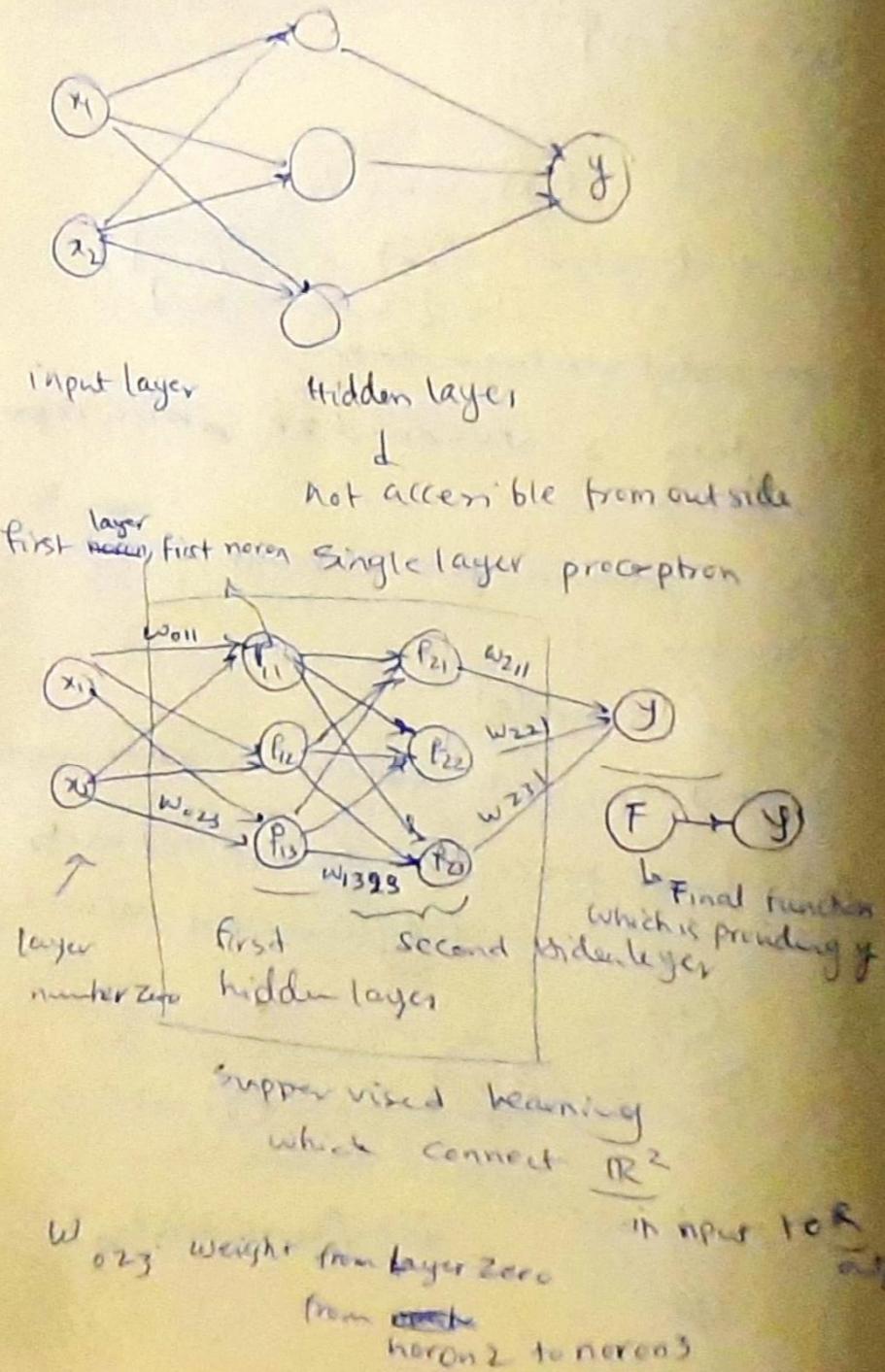
we have a model

& start to stack neurons together

Multi-layer perceptron (dense network)
fully connected network



$$\bar{w} = [w_1, w_2, \dots, w_d]$$



each neuron which is hidden is like the first neuron of the lesson.

$$f_{11}(\bar{x}) = \sigma\left(\sum_{i=1}^2 w_{0i1} x_i\right)$$

$$\bar{x} = [x_1 \ x_2]$$

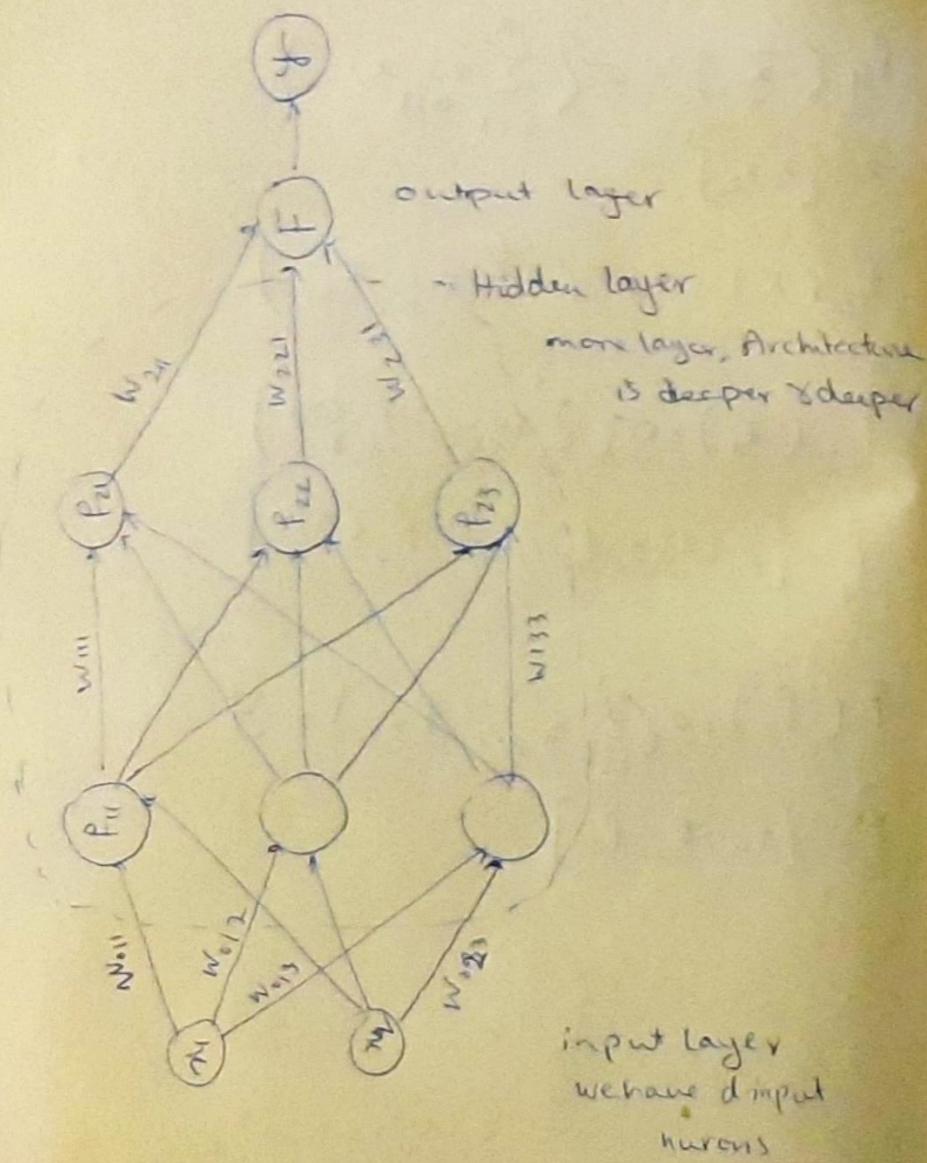
$$f_{12}(\bar{x}) = \sigma\left(\sum_{i=1}^2 w_{0i2} x_i\right)$$

$$f_{23}(\bar{p}_1) = \sigma\left(\sum_{i=1}^3 w_{1i3} \cdot f_{1i}(\bar{x})\right)$$

$$\bar{p}_1 = [f_{11}(\bar{x}), f_{12}(\bar{x}), f_{13}(\bar{x})]$$

$$F(\bar{p}_2) = \sigma\left(\sum_{i=1}^3 w_{2i1} \cdot p_{2i}\right) = y$$

$$\bar{p}_2 = [f_{21}, f_{22}, f_{23}]$$



$$\hat{F}: \mathbb{R}^d \rightarrow \mathbb{R}$$

$x \in \mathbb{R}^d$ → in this example
in general \mathbb{R}^n

$$y \in \mathbb{R}$$

We need to train.

initialize network with random weight

$F_{\bar{w}_0}$: initialization of weight

$$F_{\bar{w}_0}(x_i) = \hat{y}_i \quad \forall i=1, \dots, n$$

~~we can do the same without~~

we must compute loss. loss would one of the loss we had before.

it means that, estimate empirical Risk

Forward Propagation

$$\frac{1}{n} \sum_{i=1}^n (y_i - F_{\bar{w}}(x_i))^2$$

$$\bar{w} = \underset{w \in \mathbb{R}^m}{\operatorname{argmin}}$$

In Regularized Network

$$\bar{w} = \underset{w \in \mathbb{R}^m}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n (y_i - F_w(x_i))^2 + \lambda \|w\|^2$$

do not activation function

only parameter according to weight

With more neurons or more layers number of parameter to learn is explode.

if number of parameter m is super high with respect to complexity of problem

& data set is poor, you end up with overfitting.

training is complicate, a good idea to balance the architecture

$$\hat{W} = \arg \min_{W \in \mathbb{R}^m} \frac{1}{n} \sum_{i=1}^n (y_i - F_w(x_i))^2$$

Solve :

gradient = 0
respect to w

but it's hard, may be we have a local min locally.

there is advantage which is called chain rule of derivation to compute gradient

$$F(x) = f(g(h(x))), \frac{dF}{dx} = \frac{df}{dg} \cdot \frac{dg}{dh} \cdot \frac{dh}{dx}$$

$$\frac{\partial f(w)}{\partial w_k}$$

$k = 1 \dots m$

in training we have we have data Y with Forward propagation

& try & minimize each problem with $\frac{\partial f(w)}{\partial w_k}$

with partial derivative, is it possible to say write down explicit solution : $\hat{W} = \dots$
not because of complexity,

w_0

$$w_t = w_{t-1} - \gamma \nabla J(w_{t-1})$$

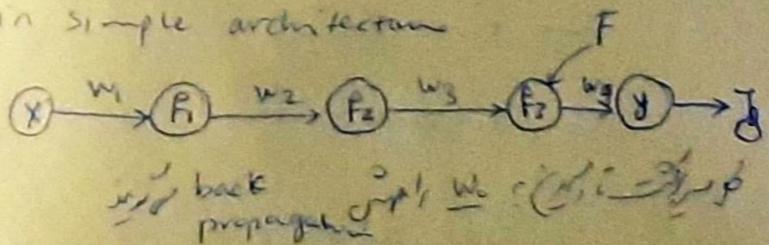
fixed as logistic regression
multiplant the gradient of function

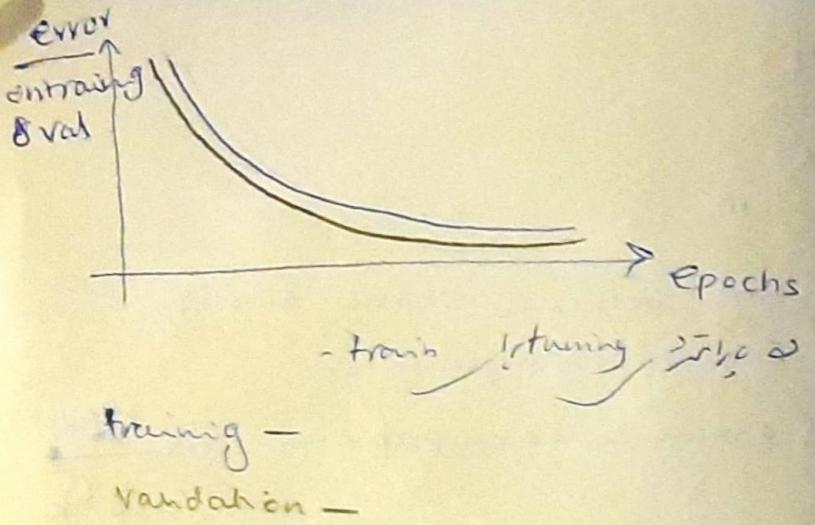
Back propagation

refer to process, Compute the root

& back

in simple architecture





Gradient
descent
entire training set

In logistic regression,
use entire training

Compute gradient

&

Conclude an epoch

Stochastic GD

basis of single
sample of training
& compute grad
only according to
that point

(the ~~path~~
is unstable
but close
to min)

for each epoch we need 1 epoch

a usual choice is not Gradient or
stochastic

even in very large dataset is ??

there is intermediate solution
not entire training & not 1
subset of training set



mini batch GD

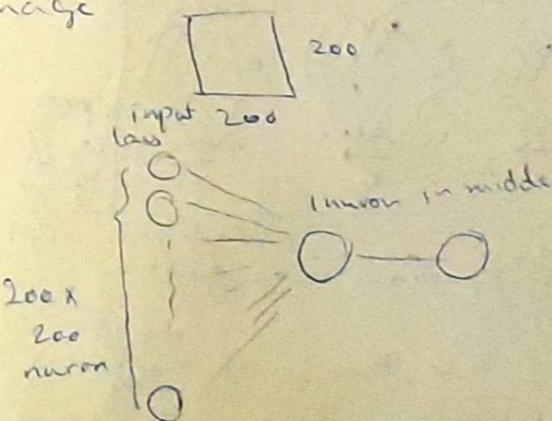
Enough subset to cover entire
training set

Introducing Convolutional Neural Network

input data are characterized by topology.

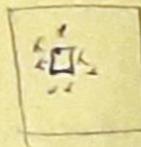
We have:

image



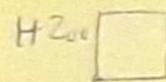
2 observation

when you have image Y consider what happens
in 1 of the pixels.



21.47: 20 Nov
image

change to visualize.



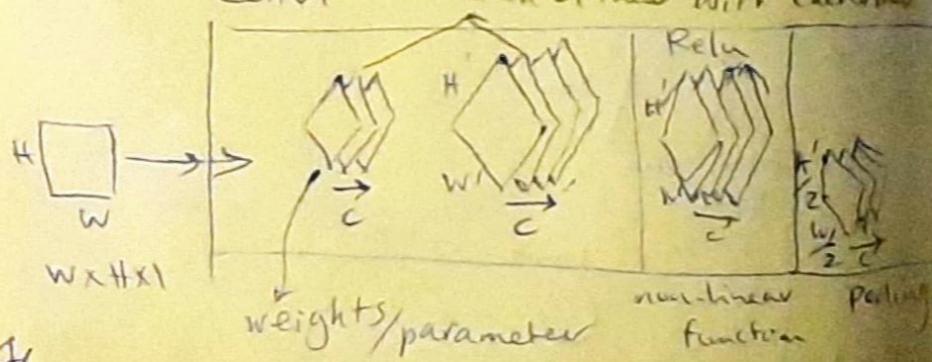
200
W

would turn to many input of
neuron
this image provide layer which
is called convolutional layer

convolutional layer has different steps

we use filter by size $t \times K$, result of
applying it's another visual representation

$W' \times H'$ (depend on size of filter -
Convl each of them with each other)



spare intervals

if I take this set of neurons & apply
to different position of image.

Second: parameter sharing,

apply very safe parameter to
image.

input neuron to some part of the image
each neuron in input have own set of
weight, ... ? everything is independent

B, How many parameter would you have
~~so poor~~ in convolution layer, have more
filter (with 1 its poor)

⇒ in representation have different representation
How many parameters?
 $K \times K \times C$

~~training & initialize, refine task~~

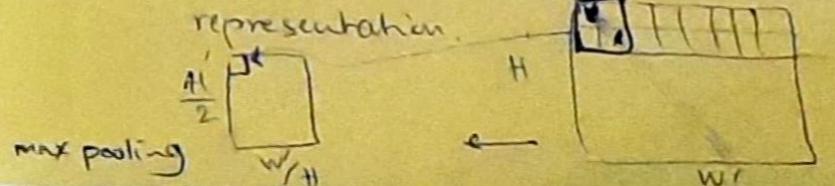
apply second second step: non linear activation function

ReLU to each element of $W'H'$

It doesn't change the size, remain the same

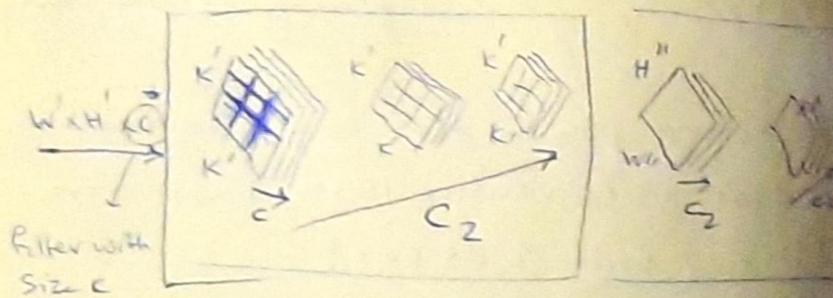
Pooling: avoid if you had 2×2 patch.

doesn't change the strength of
representation.

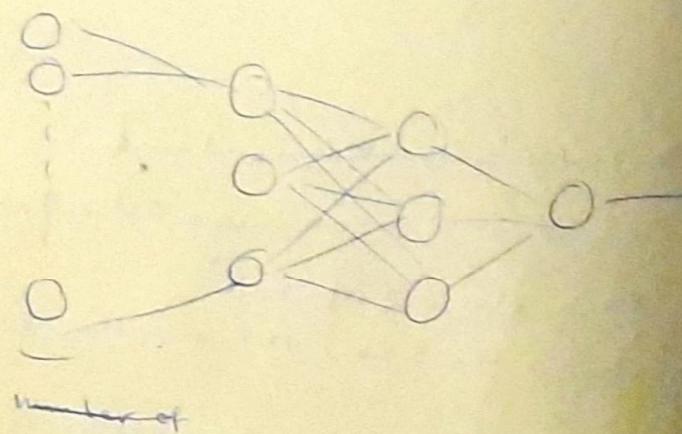


only parameters have to learn are
the ones of convolution. the
convolution layer play role in layer

Conv₂



In Conv₂ same as Conv₁



Title: Machine Learning 25 Nov

$$P(x, y) \quad f: x \rightarrow y \quad E[l(y, f(x))] = \\ \text{ERM} \quad \sum_{i=1}^n l(y_i, f(x_i)) + \lambda R(\theta) \\ \text{Regularization} \quad \|\theta\|^2$$

1 FREQUENTISTS

2 DISCRIMINATIVE : take decision in $P(x, y)$

3 GENERATIVE / PROBABILISTIC

4 BAYESIAN

Probability :

$P(x, y)$: parametrized, function would be parameterized by θ .

Typically,

$$P(x, y) = P_\theta(x, y)$$

we are going to

MLE

Maximum Likelihood

MAP

BAYESIAN
LEARNING

How we can estimate mean in distribution

$$IR \Rightarrow z \sim p(x, y)$$

$$R \Rightarrow z \sim P(z), M = E[z] = \frac{1}{n} \sum_{i=1}^n z_i \quad (z_1, \dots, z_n)$$

why it is good: the estimator is unbiased

$$E[\hat{M}] = M$$

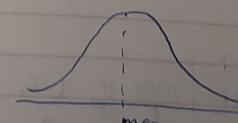
$$P(|\hat{M} - M| > \epsilon) \rightarrow 0 \quad n \rightarrow \infty$$

I want to find P_θ

$$P(z) = P_{\theta}(z) \quad z_1, \dots, z_n \sim P_{\theta}$$

$$\text{example of probability: } P_\theta(z) = \frac{e^{-|z-\mu|^2/2\sigma^2}}{\sqrt{2\pi}\sigma}$$

distributions (mean & variance)



could be linear

Gaussian distribution.

How determine models:

$$P(z) = P_{\theta}(z) \quad z_1, \dots, z_n \sim P_{\theta}$$

$$\text{which } \theta \text{ generated my } \theta \rightarrow \theta \rightarrow z_1^\theta, z_n^\theta \sim P_{\theta} \quad \left[\begin{array}{l} \theta_1 \\ \theta_2 \\ \vdots \end{array} \right]$$

what is probability of data under this distribution

θ_x : right one;

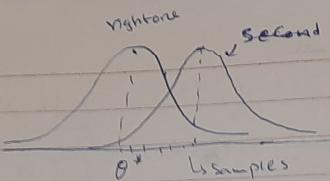
$$P_{\theta_x} = P(z_1, \dots, z_n) = \prod_{i=1}^n P_{\theta_x}(z_i) \geq P_{\theta_1}(z_1, \dots, z_n) = \prod_{i=1}^n P_{\theta_1}(z_i)$$

in general $\rightarrow \prod_{i=1}^n P_{\theta_i}(z_i) \quad (\theta_1, \dots, \theta_n)$

if data are gaussian $-z_i = \mu^T z_i / \sigma^2, \theta_i = \mu^T z_i$

$$P_{\theta_x}(z_1, \dots, z_n) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} e^{-|z_i - \mu^T z_i / \sigma^2|^2 / 2\sigma^2}$$

Title.



~~P~~ if compute probability but I can not find better than the right mass

$$P_{\theta}^M = P_{\theta}(z - z_n) = \prod_{i=1}^n P_{\theta}(z_i), \quad \max_{\theta} P_{\theta}(z_1, \dots, z_n) \Leftrightarrow \max_{\theta} \prod_{i=1}^n P_{\theta}(z_i)$$

Likely hood of θ

Maximum likelihood estimation (MLE)

We want to chase probability

$$\hat{L}(\theta) = P_{\theta}(z_1 - z_n) = \prod_{i=1}^n P_{\theta}(z_i) = \prod_{i=1}^n \frac{1}{(\sqrt{2\pi})} e^{-|z_i - \mu|^2 / 2}$$

lets change the expression with log.

log likely hood:

$$\log \hat{L}(\theta) = \sum_{i=1}^n \left[\log \left(\frac{1}{\sqrt{2\pi}} + \log e^{-|z_i - \mu|^2 / 2} \right) \right] = \sum_{i=1}^n \log P_{\theta}(z_i)$$

Product is like sum

$$\log \hat{L}(\theta) = \sum_{i=1}^n \left(\log \left(\frac{1}{\sqrt{2\pi}} \right) + (-|z_i - \theta|^2 / 2) \right)$$

we call it Δ

to maximize this expression must maximize for II part

$$\min_{\theta} \sum_{i=1}^n \frac{(z_i - \theta)^2}{2}$$

$$\sum_{i=1}^n \frac{2(z_i - \theta)}{2} = 0, \quad \sum_{i=1}^n \theta = \sum_{i=1}^n z_i$$

$$\theta_n = \frac{1}{n} \sum_{i=1}^n z_i$$

~~so we get~~

$$P(X, Y)$$

$$P_{\theta}(X, Y) = P_{\theta}(Y|X) P_{\theta}(X)$$

Regression

make a model:

$$y_i = w^T x_i + \epsilon_i, \quad \epsilon_i \sim N(0, \sigma^2)$$

noise is gaussian.

$$y_i \sim P_{\theta}(y_i | x_i) = N(w^T x_i, \sigma^2)$$

$$\left[P_{\theta}(y_i | x_i) = N(\theta^T x_i, \sigma^2) = \frac{1}{c} e^{-|y_i - \theta^T x_i|^2 / 2\sigma^2} \right]$$

x is fixed $\leftarrow -|y_i - \theta^T x_i|^2 / 2\sigma^2$

$$\hat{L}(\theta) = \prod_{i=1}^n \frac{1}{c} e^{-|y_i - \theta^T x_i|^2 / 2\sigma^2}$$

now take log:

$$\log \hat{L}(\theta) = \sum_{i=1}^n \left[\text{constant} - \frac{(y_i - \theta^T x_i)^2}{2\sigma^2} \right]$$

$$\max \hat{L}(\theta) \Leftrightarrow \min_{\theta \in \mathbb{R}^d} \sum_{i=1}^n \frac{(y_i - \theta^T x_i)^2}{2\sigma^2}$$

Min least square 

Title: Classification $y \in \{-1, 1\}$

$$P_{\theta_*}(y|x) = \sigma(y\theta^T x), \quad P(1|x) = \frac{1}{1+e^{-\theta^T x}} = \frac{1}{1+e^{-\theta^T x}}$$

$$\sigma(z) = \frac{1}{1+e^z}$$

$$P_{\theta}(-1|x) = \frac{1}{1+e^{\theta^T x}} \times \frac{e^{-\theta^T x}}{e^{-\theta^T x}} = e^{-\theta^T x}$$

Logistic model:

$$P_{\theta}(-1|x) = \frac{e^{-\theta^T x}}{1+e^{-\theta^T x}}$$

$\hat{Y}_{\theta} \in \{-1, 1\}$

$$P_{\theta}(y_i|x_i) \Leftrightarrow P_{\theta_*}(y|x) = \frac{1}{1+e^{-\theta^T x}}$$

$$\hat{L}(\theta) = \prod_{i=1}^n \frac{1}{1+e^{-y_i \theta^T x_i}}$$

$$\log \hat{L}(\theta) = \sum_{i=1}^n \log \left(\frac{1}{1+e^{-y_i \theta^T x_i}} \right)$$

$$\max_{\theta \in \mathbb{R}^d} \hat{L}(\theta) \Leftrightarrow \min_{\theta \in \mathbb{R}^d} \sum_{i=1}^n \log \left(\frac{1}{1+e^{-y_i \theta^T x_i}} \right)$$

logistic regression when
 $y \in \{0, 1\}$

$$P(y|x) = \frac{1}{1+e^{-\theta^T x}}$$

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n z_i$$

$$\hat{w} = \arg \min_{w \in \mathbb{R}^d} \sum_{i=1}^n (y_i - w^T x_i)^2$$

$$\hat{w} = \arg \min_{w \in \mathbb{R}^d} \sum_{i=1}^n \log \left(\frac{1}{1+e^{-y_i w^T x_i}} \right)$$

MLE

z (random value)
 $\text{distribution random value}$

$$\text{IR}^d \ni z \sim P(z) = P_{\theta_*}(z)$$

Date

$$(z_1, \dots, z_n) \sim (P_{\theta_*})^n \text{ (iid)}$$

$$\hat{P} \sim P \Leftrightarrow \hat{\theta} \approx \theta_*$$

$$\hat{L}(\theta) = \prod_{i=1}^n P_{\theta}(z_i)$$

$$\log \hat{L}(\theta) = \sum_{i=1}^n \log P_{\theta}(z_i)$$

because log is easy to write 'sum of product'

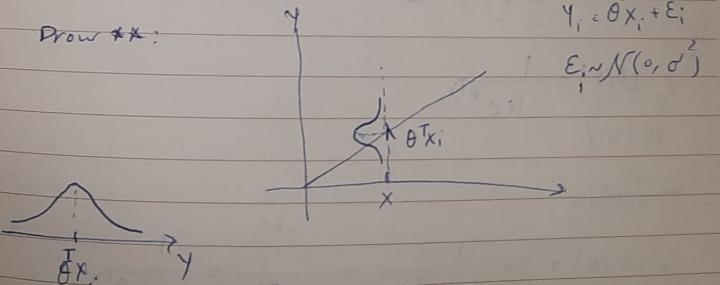
maximize likelihood, $\max \log \hat{L}(\theta)$

$$\# P_{\theta}(z) = \frac{1}{C} e^{-\|z-\theta\|^2/2\sigma^2} \quad (\sigma^2 \text{ known}), \quad \log \hat{L}(\theta) \propto \sum_{i=1}^n (z_i - \theta)^2$$

$$\hat{w} = \arg \min_{w \in \mathbb{R}^d} \sum_{i=1}^n (y_i - w^T x_i)^2 \quad P_{\theta}(1|x) = \frac{1}{C} e^{-\|x - \theta\|^2/2\sigma^2}$$

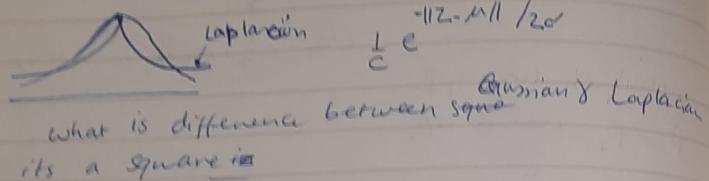
$$\log \hat{L}(\theta) \propto - \sum_{i=1}^n (y_i - \theta^T x_i)^2$$

Draw #:



Title:

instead of Gaussian, goes Laplacian. (SVM)



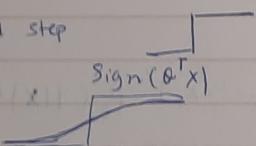
? { use absolute value when ...
use when
Laplacian $\propto -A \log$

$$P_\theta(Y|X) = \frac{1}{1 + e^{-Y_i \theta^T X_i}}$$

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

$$P(1|X) + P(-1|X) = 1$$

idea: shape: smoothed step



$$\log \hat{L}(\theta) = - \sum_{i=1}^n \log (1 + e^{Y_i \theta^T X_i})$$

GF

$$P_\theta(1|X) = \frac{1}{1 + e^{-\theta^T X_i}}$$

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

$$\log \hat{L}(\theta) =$$

$$P_\theta(0|X) = 1 - \frac{1}{1 + e^{-\theta^T X}} = \frac{e^{-\theta^T X}}{1 + e^{-\theta^T X}} = \frac{1}{1 + e^{\theta^T X}}$$

$$P_\theta(Y|X) = \left(\frac{1}{1 + e^{-\theta^T X}} \right)^Y \left(\frac{1}{1 + e^{\theta^T X}} \right)^{1-Y}$$

$$\log \hat{L}(\theta) = - \sum_{i=1}^n \left[Y_i \log (1 + e^{\theta^T X_i}) + (1 - Y_i) \log (1 + e^{-\theta^T X_i}) \right]$$

negative of crossentropy.

$$\hat{w} = \underset{w \in \mathbb{R}^d}{\operatorname{argmin}} \sum_{i=1}^n (y_i - w^T x_i)^2 + \lambda \|w\|^2$$

$$\hat{w} = \underset{w \in \mathbb{R}^d}{\operatorname{argmin}} \sum_{i=1}^n \log (1 + e^{-y_i w^T x_i}) + \lambda \|w\|^2$$

BAYES:

Z, U

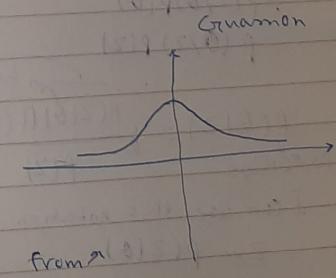
(2 random variable instead of X, Y)

$$P(Z, U) \Rightarrow P(U|Z) = \frac{P(Z, U)}{P(Z)}$$

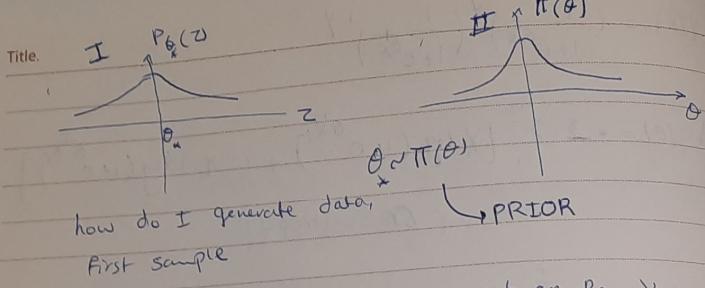
$$P(Z|U) = \frac{P(Z, U)}{P(U)}$$

$$P(Z, U) = P(Z|U)P(U) \\ = P(U|Z)P(Z)$$

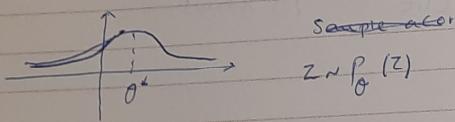
$$P(U|Z) = \frac{P(Z|U)P(U)}{P(Z)}$$



Data comes from true model from



Pick 1 guy (sample) from π , put on P_θ & get distribution. (generate it according to I)



(data Generation)

PRIOR

Can you learn it, No, Choice it.

$$p(z, \theta) = p(\theta, z) = \frac{P(z, \theta)}{P(z)}$$

$$= p(z|\theta)p(\theta)$$

$$= p(\theta|z)p(z)$$

$$p(\theta|z) = \frac{p(z|\theta)\pi(\theta)}{p(z)}$$

fixed likelyhood
prior

I can use this notation, play role of likelyhood ($p(z|\theta)$)

$$z \sim p(z|\theta)$$

for many z :

$$p(\theta|z_1, \dots, z_n) = \frac{p(z_1, \dots, z_n|\theta)\pi(\theta)}{p(z_1, \dots, z_n)}$$

Date.

$z \in \mathbb{R}^d$

$\mathbb{R}^d \ni \theta \sim P(z|\theta), \theta \sim \pi(\theta)$

$z_1, \dots, z_n \sim (P(z|\theta))^{n \text{ i.i.d}}$

derive new principle
idea: try to make prediction, goal: posterior proportion

$p(\theta|z_1, \dots, z_n) \propto \int \pi(\theta) \cdot P(z|\theta)$

exactly probability

$\prod_{i=1}^n P(z_i|\theta) \pi(\theta)$

REGRESSION

$y_i \sim N(\theta^T x_i, \sigma^2)$, $P(y_i|\theta) = \frac{1}{c} e^{-\frac{(y_i - \theta^T x_i)^2}{2\sigma^2}}$

$\Leftrightarrow (y_i = \theta^T x_i + \epsilon_i, \epsilon_i \sim N(0, \sigma^2))$

$p(\theta|z_1, \dots, z_n) \propto \int \pi(\theta) \cdot P(z|\theta)$

$= \frac{1}{c^n} e^{-\frac{\sum_{i=1}^n (y_i - \theta^T x_i)^2}{2\sigma^2}}$

how pick for θ . freedom, simple choice is
 $\|\theta\|^2$ Guassian (standard one)

take this and plug it here:

$$p(\theta|z_1, \dots, z_n) \propto \int \pi(\theta) \cdot P(z|\theta)$$

$$= \frac{1}{c^n} e^{-\frac{\sum_{i=1}^n (y_i - \theta^T x_i)^2}{2\sigma^2} \cdot \frac{1}{c} e^{-\frac{\|\theta\|^2}{2\sigma^2}}}$$

$$= \frac{1}{c^n c'} e^{-\left(\sum_{i=1}^n \frac{(y_i - \theta^T x_i)^2}{2\sigma^2} + \frac{\|\theta\|^2}{2\sigma^2}\right)}$$

least square + regularization

Title

this is distribution, generate data,
try to find good estimator θ

MAP: MAXIMUM POSTERIOR

$$\text{So: } -\left(\sum_{i=1}^n (y_i - \theta^T x_i)^2 + \|\theta\|^2 \right)$$

$$\max_{\theta} = \frac{1}{c^n C'} e^{-\left(\sum_{i=1}^n (y_i - \theta^T x_i)^2 + \|\theta\|^2 \right)}$$

can be for logistic, ---

$$z_i = z_m \sim (P_z)^n$$

$$P_z = P_{\theta_z} \quad \theta_z \in \mathbb{R}^p$$

$$\text{MLE } \hat{\theta}(\theta) = \prod_{i=1}^n P_{\theta}(z_i)$$

$$z = y | x$$

$$y \sim P_{\theta}(y|x)$$

REGRESSION

$$P_{\theta}(y|x) = N(\theta^T x)$$

$$P_{\theta}(y|x) = N(\theta^T x, \sigma^2), \frac{1}{c} e^{-\sum_{i=1}^n (y_i - \theta^T x_i)^2 / 2\sigma^2}$$

Classification

$$y = (\pm 1) \quad P_{\theta}(y|x) = \frac{1}{1 + e^{-\theta^T x}} \cdot \frac{1}{c} e^{-\sum_{i=1}^n \log(1 + e^{x_i \theta_i})}$$

$$P_{\theta}(\theta|z) = \frac{1}{C'} e^{-\frac{\sum (y_i - \theta^T x_i)^2 + \|\theta\|^2}{2\sigma^2}}$$

noise fix the amount of regularization
you could do variance formula.

there is another technique for cross validation

MARGINALISED LIKELIHOOD, (in place of CV)
cross validation

we don't go to it.

story:

$$(x_i, y_i) \stackrel{i.i.d. \text{ iid}}{\sim} f(x_{\text{new}}) \sim y_{\text{new}}$$

X-AI (explainable AI)
interpretable
accountability

SPARSITY: bunch of zero.

the target depends on few building blocks

$$\text{linear models} \quad f(x) = w^T x = \sum_{j=1}^d w_j x_j$$

take vector consider j from 1 to d (all w_j zero norm)

$$\|w\| = \#\{j \in \{1-d\} \mid w_j \neq 0\} = \sum_{j=1}^d \mathbb{1}_{\{w_j \neq 0\}}$$

sparse vector
How many

$$(\text{sum of steps}) - \#\{x_j \leq 1 \mid j=1-d\}$$

Delija

Title.

$$\begin{cases} -1 \leq x_j \leq 1 & j=1-d \\ x_j \rightarrow E[x_j] = 0 \\ E|x_j|^2 = 1 \end{cases}$$

SCALING MATTERS!

algorithm knows for sparsity
n for w that contains zero

$$\min_{w \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n (y_i - w^T x_i)^2 + \lambda \|w\|_1$$

to calculate: take derivative equal to zero
the second part is sum of steps

$$\begin{aligned} f(x) &= w^T x + b \\ &= \bar{w}^T \bar{x} \\ &= (\bar{w}, b) \\ &\bar{x} = (x, 1) \end{aligned}$$

Best subset selection (BSS)

Date.

Forturing the data until they confess

$$J = \{1, \dots, d\}$$

\hat{X} , \hat{x}_j j-th column of \hat{X}

nxd

$$\hat{Y}_{n \times 1}$$

BSS :

$$\hat{J} = \text{BSS}(\hat{X}, \hat{Y})$$

$$\hat{w} = \underset{\text{least square}}{\text{LS}} (\hat{X}, \hat{Y})$$

$\hat{w}_i \quad i=1, \dots, 2 \rightarrow$ best error on training
 $\hat{w}_i \quad i=1, \dots, d \rightarrow$ best possible error
 $n < d \rightarrow$

built model, check on error $\approx 30\%$ of data validation.

we could implement it, but you need subset function

$$\{\hat{x}_1, \dots, \hat{x}_d\} = \text{subset}(\hat{X}, J)$$

$$\{\hat{x}^1, \dots, \hat{x}^d\} = \text{subset}(\hat{X}, J)$$

solution: $\|\hat{y} - \sum_{i=1}^d \hat{w}_i \hat{x}_i\| \Rightarrow$

Title.

$$\text{error}_j = \min_{\text{VIR}} \| \hat{Y} - V \hat{X}^j \|_2^2 \quad j=1 \dots d$$

just order them, and best simple 1

initialize $j_0 = \{\emptyset\}$ = empty set.

$$j_* = \text{sort}(\text{err}_j), \text{ update } j = j_0 \cup \{j_*\}$$

\Rightarrow

$$\hat{j} = \text{BSS}(\hat{X}, \hat{Y})$$

$$\{\hat{x}^1 \dots \hat{x}^d\} = \text{subset}(\hat{X}), j_0 = \emptyset$$

$$\text{error}_j = \min_{\text{VIR}} \| \hat{Y} - V \hat{X}^j \|_2^2 \quad j=1 \dots d$$

$$j_* = \text{sort}(\text{err}_j), j_1 = j_0 \cup \{j_*\}$$

$$\text{MP}(\hat{X}, \hat{Y})$$

MATCHING PURSUIT (GREEDY METHOD)

select 1 variable, add & update solution

$$(\hat{W}, \hat{r}) = \text{MP}(\hat{X}, \hat{Y})$$

Selected variable

Date.

$$\hat{j} = \emptyset \quad \hat{r} = \hat{Y}$$

$$\hat{w} = 0 \quad v_x^j = \arg \min_{\text{VIR}} \| \hat{Y} - V \hat{X}^j \|_2^2 \quad \text{Find best coefficient}$$

$$\text{err}_j = \min_{\text{VIR}} \| \hat{r} - V \hat{X}^j \|_2^2 \quad \text{check error for each variable}$$

$$j_* = \text{sort}(\text{err}_j) \rightarrow \text{best column}$$

I have to update j

$$\hat{j} = \hat{j} \cup \{j_*\}, \hat{r} = \hat{Y} - V \hat{X}^j, \hat{w} = \hat{w} + v_x^j e_j$$

add column to select set update

this is one iteration
 $T = T+1$ first

update in right place

need to coefficient update

$$\hat{w}_2 = \begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \\ \vdots \\ i \end{pmatrix} + v_x^j \begin{pmatrix} 0 \\ 0 \\ 0 \\ 1 \\ \vdots \\ 0 \end{pmatrix} \hat{j}_k$$

orthogonal matching pursuit (recompute everything all the time)

$$\hat{w} = \hat{w} + v_x^j * e_j$$

$$\hat{e} = \arg \min_{\text{VIR}^M} \| \hat{X}^j a - \hat{r} \|_2^2$$

$$m = |\hat{j}|$$

instead of compute separate for every variable
use all the selected one in 1

w, put right number in right position

Title. $\hat{r} = \hat{y} - \hat{x}_j \cdot \hat{\theta}$

it cost more - solve better

regularization parameter: $F \gg d$, decide when to stop, trade off fitting data and sparsity

$$\min_{w \in \mathbb{R}^d} \sum_{i=1}^n (y_i - w^T x_i)^2 + \lambda \|w\|_1$$

$$\|w\|_1 = \sum_{j=1}^d |w_j|$$

LASSO

$$\textcircled{1} \quad w_{t+1} = \underset{\lambda, \eta}{\text{arg min}} \left(w_t - \eta \nabla \|w_t - \hat{y}\|^2 \right)$$

Gradient descent to pick some loss

soft thresholding of vector

entry small kill it

big: leave alone, decrease abit

Small

λ : tells you how much threshold &

1: small & big ??

more sparsity

(5)

iterative soft thresholding

$$\textcircled{2} \quad \delta_{\lambda \eta}(w) = \begin{cases} w - \lambda \eta & w^i > \lambda \eta \\ 0 & \|w^i\| < \lambda \eta \\ w + \lambda \eta & w^i < \lambda \eta \end{cases}$$

Review Sparse method

Date.

$$f_0(x) = w^T x \quad x \in \mathbb{R}^d$$

$$= \sum_{i=1}^d w_i x_i$$

$x = x$ position x_i = area x_3 ... features

x_4 = How many cats #cats the ~~feature~~ feature doesn't not related to the question

$$z_1 - z_n \sim P^n$$

$$\hat{L}(w) = \sum_{i=1}^n (y_i - w^T x_i)^2 + \lambda \|w\|$$

could be large small → pick 1, $w \rightarrow$
means

large →
could be small

Consider all feature

With λ you can have trade off between error & ??

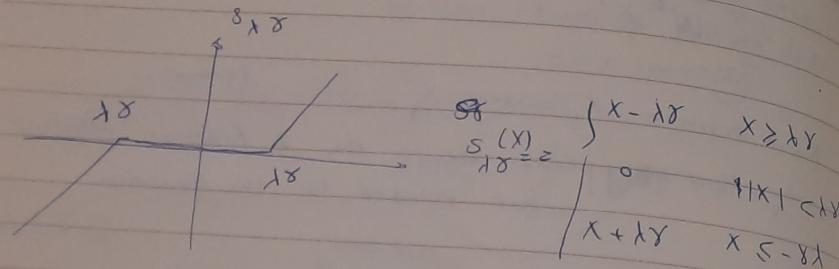
$$\hat{L}(w) = \sum_{i=1}^n (y_i - w^T x_i)^2 + \lambda \sum_{j=1}^d |w_j|$$

LASSO

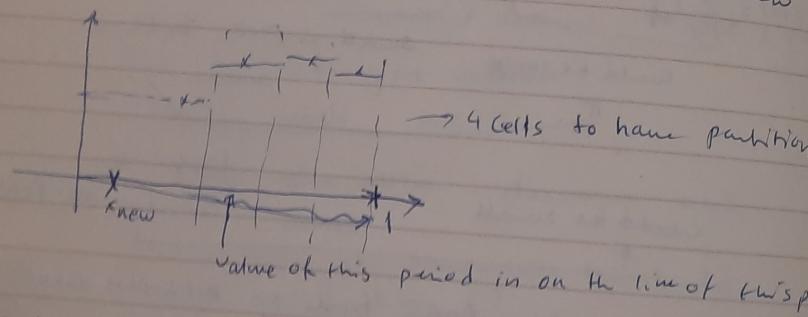
Title: ISTA: iterative soft thresholding Algorithm

$$w_0 = 0$$

$$w_i = \text{soft}_{\lambda\gamma} (w_{i-1} - \frac{2\gamma}{m} X^T (Y - w_{i-1}^T X))$$



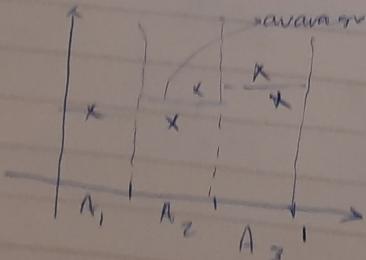
You only care about neighbours close to him
local estimate algorithm



Def (Partition)

$$A = \left\{ A_j \cup A_i \mid j \neq i ; A_j \cap A_i = \emptyset \text{ if } j \neq i \right\}$$

Given



after partition.

if we consider our estimator:

Piece wise constant function:

in partition we have different data

piecewise constant function

$$\min_c \sum_{x_i \in A_j} (y_i - c_j)^2$$

Consider data just in cell

$$\min_c \sum_{x_i \in A_j} (y_i - c_j)^2 \Rightarrow c_j = \frac{1}{n_{A_j}} \sum_{x_i \in A_j} y_i$$

how many points you have
in A_j (many points)
→ indicator function here

$$f(x) = \sum_{j=1}^J c_j \mathbb{1}_{x \in A_j}(x)$$

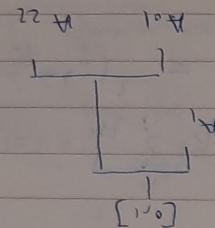
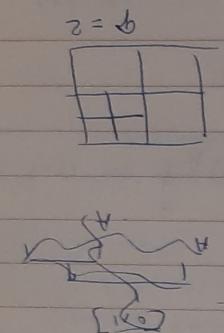
~~$x \in A_2$~~ $x \in A_2 \rightarrow f(x) \in C$
 $J \rightarrow$ how many partition you have

$$f(x) = \sum_{j \geq 1} (w_j x) \mathbb{1}_{x \in A_j}(x)$$

Piece wise linear function (or any function)

if partition is not good, find different partition?
you can design partition?

$\frac{2}{q^2}$ \leftarrow Cut point
Cut point $\in \mathbb{Z}$



$$A_2 = \{A_{21}, A_{22}\}$$

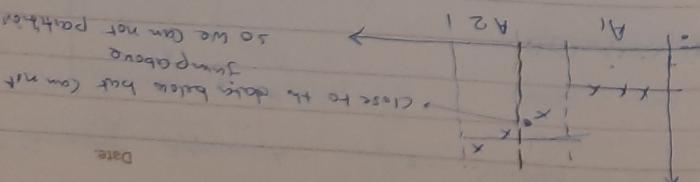
do more path then

$$\frac{1}{q} A_2 = \frac{1}{q} \sum_{i=1}^{q^2} (y_i - c_i)^2 < \text{threshold}$$

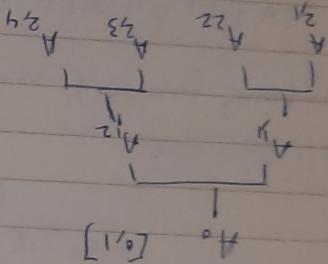
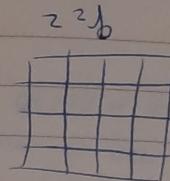
If we want partition A_1 ,
it is if we are in Cutform, because if we partition

This case cannot partition
piece wise function

If you get partition, compute less



ADAPTIVE PARTITION



A we have deep $\frac{q}{2}$

$$A_3 = \{A_{31}, A_{32}\}$$

$$A_2 = \{A_{21}, A_{22}\}$$

$$A_1 = \{A_{11}, A_{12}\}$$

$$A = \{A_{01}\}$$

$$A_{11} = \{A_{111}, A_{112}\}$$

$$A_{111} = \{A_{1111}, A_{1112}\}$$

$$A_{1111} = \{A_{11111}, A_{11112}\}$$

$$A_{11111} = \{A_{111111}, A_{111112}\}$$

$$A_{111111} = \{A_{1111111}, A_{1111112}\}$$

$$A_{1111111} = \{A_{11111111}, A_{11111112}\}$$

$$A_{11111111} = \{A_{111111111}, A_{111111112}\}$$

Title partition tree

$A_{111111111} = \{A_{1111111111}, A_{1111111112}\}$

AC \neq future exist q...Get

such that $A_2 \cup C_1$

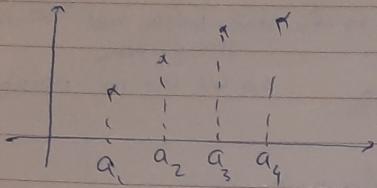
$A_{1111111111} = \{A_{11111111111}, A_{11111111112}\}$

$A_{11111111111} = \{A_{111111111111}, A_{111111111112}\}$

$A_{111111111111} = \{A_{1111111111111}, A_{1111111111112}\}$

$A_{1111111111111} = \{A_{11111111111111}, A_{11111111111112}\}$

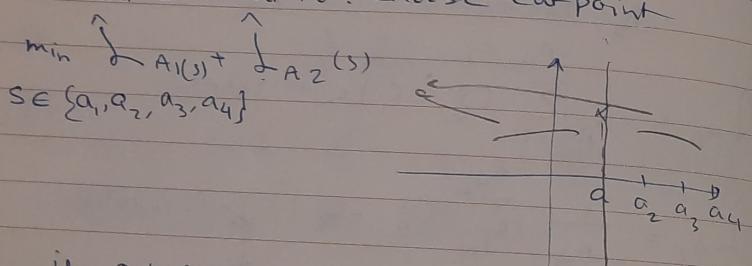
Title: Greedy approach: I give samples. I don't want to dispaly, what I want to do.



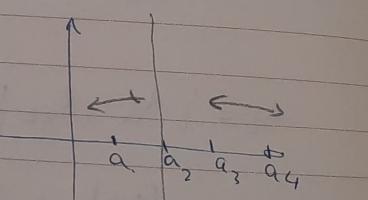
$$A(S) = \{x \in X \mid x \leq s\}$$

$$A_1(S) = \{x \in X \mid x \leq s\}, A_2(S) = \{x \in X \mid x > s\}$$

What I want to: choose cut point



if cut from a2

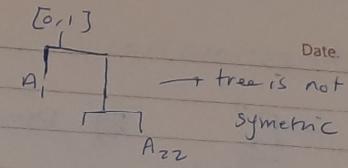


$$\text{if } s = a_1 \Rightarrow \hat{\int}_{A_1}(a_1) + \hat{\int}_{A_2}(a_1) = \frac{1}{n_{A_1}} \sum_{i \in I_1} (y_i - c_1)^2 + \frac{1}{n_{A_2}} \sum_{i \in I_2} (y_i - c_2)^2 = L_1$$

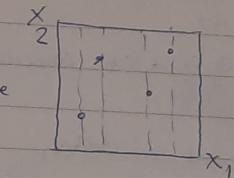
average value

Cut a2
Compare $L_1, L_2, L_3, L_4 \rightarrow$ which is the smallest
Cut a3
Cut a4

What tree took like



in dim = 2 \rightarrow



If you have more dimension,

so minimizer not only depend on s, it depends on j

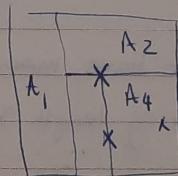
$$\text{so } (s^*, j^*) \min_{S, j} \hat{\int}_{A_1(S)} + \hat{\int}_{A_2(S)}$$



not good, loss is not good if we cut in this form

need compare loss &
choose smallest,
you can cut vertically
or horizontally

maybe:



each cut make error
smallest

Title. DENSITY ESTIMATION:

$$\text{I) } (X_i)_{i=1}^n \sim P_*^n$$

unknown

Generative AI

chatgpt & kind of generative AI

$$P_*^n$$

Find \hat{P} based on data,

Cats I observed

$$\text{II, } (\hat{X}_i)_{i=1}^N \sim \hat{P}^N$$

find approximation \hat{P}

SAMPLING part

$$\# X = N \underset{\downarrow}{\underbrace{P_*(1) \dots P_*(N)}} \text{, } \sum P_*(i) \leq 1; \sum_{i=1}^N P_*(i) = 1$$

all the possible cats in universe

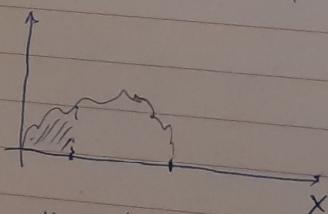
$$X \subseteq \mathbb{R}^d$$

$A \subset X$
chunk of X

$$P(A) = \int_A P_*(x) dx \quad P_*(x) \geq 0$$

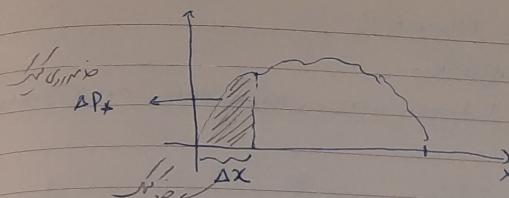
$$\int_X P_*(x) dx = 1$$

probability, in each point don't have meaning must take interval,
prob measuring length & value

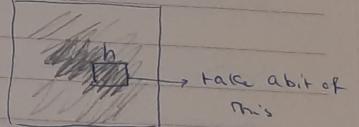


Normalize amount of stuff (all the parts) how much
stuff is in (///) according to all the part

Density distribution part



$$\frac{\Delta P_*}{\Delta x}$$



$$\Delta x \sim h^d$$

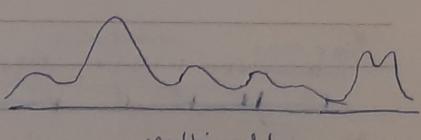
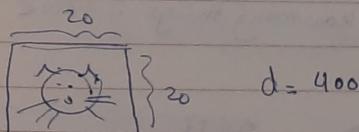
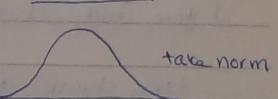
with out creme, we have just cake and its Uniform. Normalize total amount.

$$\text{Ex. } U([0, 1]^d)$$

$$- \|x - M\|^2 / 2\sigma^2$$

lets draw in 1 dimension:

unimodel



multimodel

~~main~~ main goal maximum likelihood

main reason that we haven't got label,
an example of unsupervised learning.

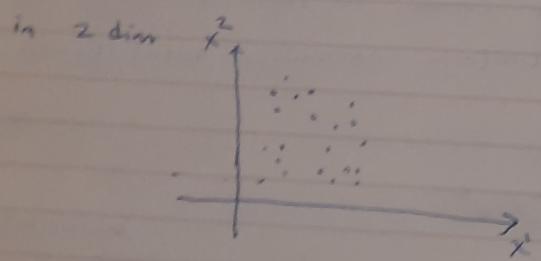
lets invent way to estimate P_x

how could be invent algorithm.

lets draw it, in 1 dimension, we have data on x



in 2 dim



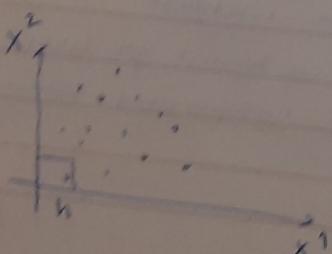
I give data, 1 or 2 dim,

take chunk, go in check how many things it have



take another interval: shift the first

(is it good)



$$\hat{P}_h(x) = \sum_{j=1}^{N_h} \frac{n_j(h)}{h^n} \prod_{i \in B_j(h)} (x_i)$$

number of point inside jth cell

$$n_j(h) = \sum_{i \in B_j(h)} 1$$

$$B_i(h) \cap B_j(h) \neq \emptyset \quad \text{if } i \neq j \text{ does not overlap.}$$

Histograms

$$X = \bigcup_{j=1}^{N_h} B_j(h)$$

$$B_1(h) \cap B_2(h) \dots B_j(h) \neq \emptyset \quad \text{if } i \neq j \text{ in } \bigcup_{j=1}^{N_h}$$

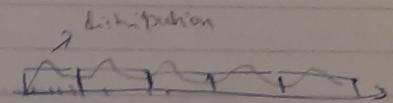
→ Binning $B_j(h)$ bins

now select h : cross-validation.

there is a good h .

if $h \gg$, point whole \rightarrow Uniform distribution

in 1 dimension:



$\frac{1}{n}$ for each data: have a lot of spike.

Concern: partitioning & selection bandwidth (h)

dimension

Title.

3

10

$\gg 100$

in dimension of 2, 3 it's OK

but think of d dimension must think different

$d=10$: high for mathematic

very high dimension, completely of maybe neural network

~~softmax~~ beat voting

split domain, value? everybody vote.

Kernel Density estimator (KDE)

$\hat{P}_h(x)$

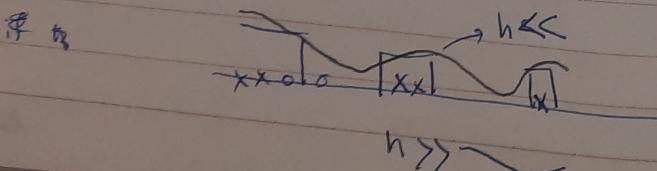
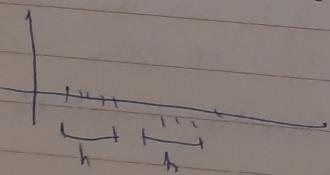
if I ask you how much is density is here based on data,

$$\hat{P}(x) = \sum_{i=1}^n \frac{1}{h} \mathbb{I}_{\|x - x_i\| \leq h}$$

change name.
Count anything $\leq \epsilon$

divide by n

It's like window, move

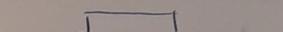


$$\hat{P}_h(x) = \sum_{i=1}^n \frac{\mathbb{I}(x - x_i \leq h)}{h} \rightarrow$$

$$\hat{P}_h(x) = \frac{1}{n h^d} \sum_{i=1}^n K\left(\frac{\|x - x_i\|}{h}\right)$$

$K_h(\|x - x_i\|)$

shapes:



$$K(s) = \mathbb{I}(s)$$

$$K(s) = \mathbb{I}_{|s| \leq 1}$$

What are other shapes:



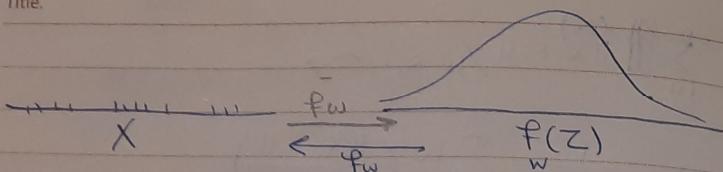
in all of this there are h which may stretch it.

every time give a number, compute distance to any,
draw shape & that's all.

Generative approach:

Generative approach:

Title.



Can I transfer Gaussian into sth.

function of random variable

if random variable of Z is Gaussian \rightarrow
 $\text{ran } \sim \sim \sim X \text{ in linear}$

$$P_X = ? P_Z \quad Z \sim N(0, I)$$

Given density (P_Z) go to another density P_X
 f_z is invertable.

so I can go back. & transform Gaussian into
sth else.

measure distance between this two ($f_w(z), X$)

Why we call it Generative:

Suppose I have f_w how to generate X
sample Z , get f_w and get X

this idea, Generative networks,

Normalizing Flow,

f_w is invertable so we have \rightarrow formula
if Z is Gaussian, (Normal distribution)
so we can have f_w

Generative idea.

normalizing flow

Date.

also called 'anomaly detection'.
anomaly have unlabeled data.

Anomaly Detection:

$$\hat{P}(x) \quad \hat{P}(z)$$

try to make algorithm

1; one of them is far away

2; if sparsely number, or high
find anomaly dimension

How to find anomaly

$\hat{P}(x) \leftarrow$ put x , check $\hat{P}(x) \leq \tau$ then say
Anomaly
if big not anomaly

$$a(x) = \prod_{\hat{P}(x_i) \leq \tau} (X_i)$$

$\rightarrow \rightarrow$ big anomaly

How to choose τ .

Clustering: (grouping) group stuff

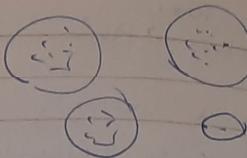
How could grouping, lot of algorithm

Find regions (distance small, Density high)

$$X_{\tau} = \{x \in X \mid \hat{P}(x) > \tau\}$$

Title: selected region

$$I, X_C = \{x \in \hat{P}(x) > C\}$$



II, find connected regions

DBSCAN

Unsupervised learning:

$$\{x_i \rightarrow x_n\} \text{ Data, inputs, no output here}$$

$$\{x_i \rightarrow x_n\} \sim P_*^n \text{ Data come from Distribution.}$$

$$\{x_i \rightarrow x_n\} \sim P_*^n$$

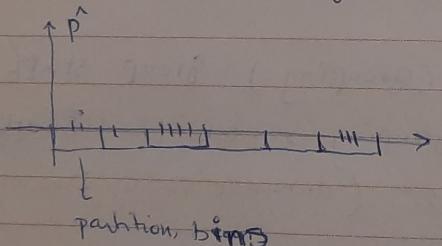
We would like to find \hat{P} which is good approximation of P_*^n

$$\Rightarrow \{x_i \rightarrow x_n\} \sim P_*^n$$

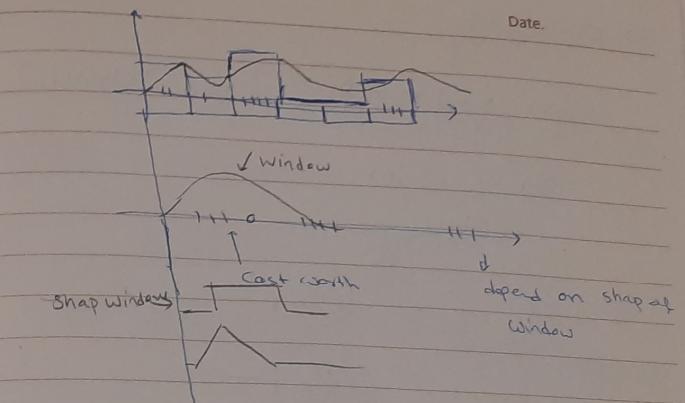
$\hat{P} \approx$

Idea:

① Partition Based \rightarrow histograms



Date:



for 2: go inside each cell & count

call Cell j

$$M_j = \sum_{i=1}^n \prod_{x_i \in B_j} \frac{1}{h^d} \rightarrow \text{Density of each cell}$$

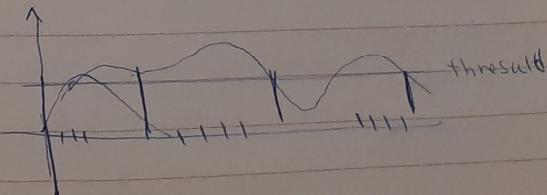
$\frac{M_j}{n / h^d} = P_j$

for

$$I \quad \hat{P}(x) = \sum_{j=1}^{\# \text{Bins}} \prod_{x \in B_j} \hat{P}_j$$

for

$$III \quad \hat{P}(x) = \sum_{i=1}^n \frac{K(x - x_i \| / h)}{h^d}$$

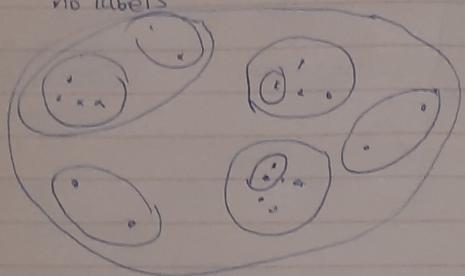


Title. Clustering

Given $x_1 \rightarrow x_n \in P^n$ find "clusters"

Set of points

no labels

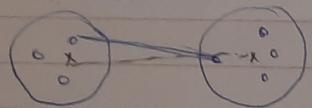


We can have different cluster, but do not when do statistic, its in sum kind its in some kind complicated.

give algorithm to group the points.

~~ستEPS لـ CLUSTERING~~ -> ~~الخطوات~~

Hierarchical clustering



We call it: AGGLOMERATIVE CLUSTERING

between cluster: linkage

Linkage {
AVGAG
Single
Complete

DENDROGRAM

Date.

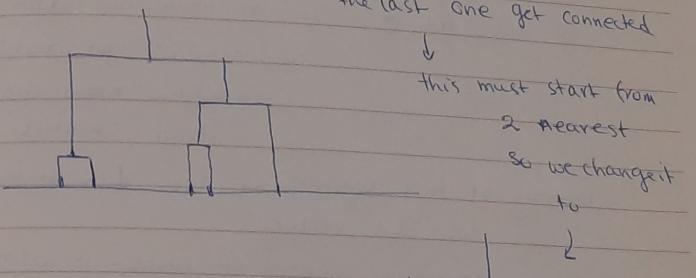
Five point in 1 dimension

II II I

two points get connected

the second 2 point get connected

the last one get connected



this must start from
2 nearest
so we change it
to

calculate distance in binary string

$$d = \sum_{i=1}^n \|x_i^K - x_j^K\| = d_H(x_i, x_j) \quad x_i, x_j \in \{0,1\}^K$$

second algorithm:

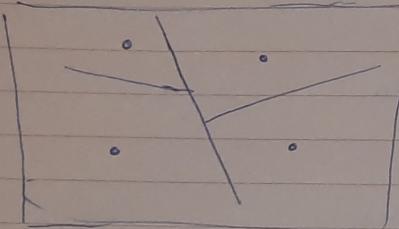
K-means

Statistical Algorithm ↗ idea
in computation call Lloyd's Algorithm ↗
or Kmeans++

Algorithm

Deljia

Title: big challenge: how could partition:
optimize possible partition:



We are in cube
and there are
partition
we can define capital
of each partition

if take all \mathbf{x} : if get mean can define partition
and also if get partition
can define mean and go to
find Capital

beauty of this optimizing, lets do gradients
and we need vectors!

so partition, centre of partition is bunch of
vectors and optimize centre of masses

Call V_j : regions

$j=1 \rightarrow K$ of them we need name for centers call
them c_j

which of this point is in clusters:

$$V_j = \{x \in X \mid \|x - c_j\| \leq \|x - c_i\| \text{ if } i \neq j\}$$

if I get you 2 partition, get different error
lets put some centers & fine means. we can say
pick 2,

8

Partition for X
Date.

for triangle

0 0
0 0

put number
I want function $L(c_1, \dots, c_K)$:

take each point have to decide in which
partition is this. So fix center, I :

$$\|x_i - c_j\|^2$$

among all centers, which is closest. so

$$L(c_1, \dots, c_K) = \sum_{i=1}^n \min_{j=1 \rightarrow K} \|x_i - c_j\|^2$$

this is one way to do it.

this is very similar to loss function

$$\sum_{i=1}^n \min_{j=1 \rightarrow K} \|x_i - c_j\|^2$$

$$e(x_i, c_j)$$

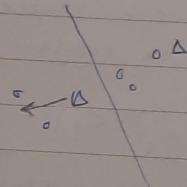
Title: what to do, minimize, Gradient & equal to zero,
in infinit data, Partition are different.

the most intuitive thing:
let's write down:

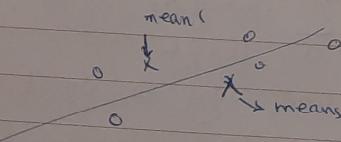
$$\min_{c_1 \dots c_k} L(c_1 \dots c_k)$$

we have 2 minimization, one inside & one outside
initialize the centers, compute the error
and update means.

find center of mass and keep.



what you see get to right place and it's easy



assign the means, reassign point →

update means

change the
partition



init $c_i^0 = c_k^0$

assign $(x_i - x_n) \mapsto c_i^0 \quad c_i^0 \mapsto \hat{v}_i^0 - \hat{v}_k^0$

update means $(\hat{v}_i^0 - \hat{v}_k^0) \mapsto c_i^0 \mapsto c_k^0$

Date:

so generalize:

θ init $c_i^0 = c_k^0$

assign $(x_i - x_n) \mapsto c_i^{t-1} \quad c_k^{t-1} \mapsto \hat{v}_i^t - \hat{v}_k^t$

update means $(\hat{v}_i^t - \hat{v}_k^t) \mapsto c_i^t = c_k^t$

$$c_j^t = \sum_{x \in \hat{v}_j^t} \frac{1}{\#\hat{v}_j^t}$$

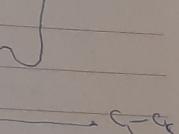
what can you say about:

use

in situation

local minimal

$$\sum_{i=1}^n (c_i - c_k)^2$$



int & its important for initialization

idea: choose initialization:

bad n put all centers

together is bad idea,

choose point from data,

and second is the furthers from it

third guy should be the furthers from 2
points

it's no this algorithm is robust, by randomization
1 one pick at random, second one is pick of random.

Title.

- pick second one is farther away.
- if do average, shows that robust algorithm is easy but robust

I do:

$$c_i^* = x_j \sim U[1 - \alpha]$$

$$d_i = \|x_i - c_i^*\| \quad i \neq j$$

Probability distribution:

so:

$$d_i = \|x_i - c_i^*\| \quad i \neq j \rightarrow p_i = \frac{d_i}{\sum_{i \neq j} d_i}$$

Sample the next mean

$$c_2^* \sim [p_1 \rightarrow p_{n-1}]$$

again compute.

K-means it's initialization itself.

1, parameter tuning

2, K-means

Suppose that
image,

make image
compress

take pixel, clusters

10 Clusters



For example the ears is in 1 cluster. Date.
the ones is black.

Very good in images.

its not good examples of classifying - for example cats & dogs
its good in 1 simple image for clustering