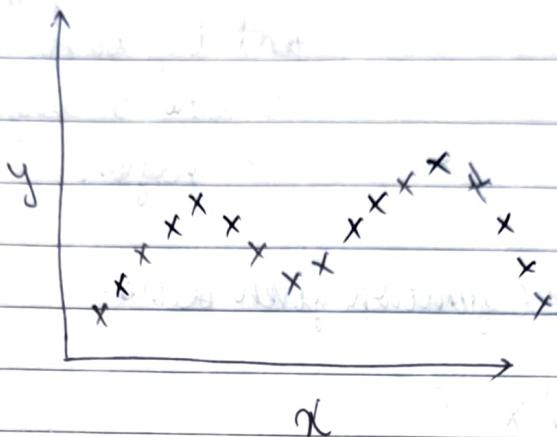


Locally weighted linear regression : →

suppose we have a dataset



we want to find a curve that fits this shape. We will do this using locally weighted L.R.

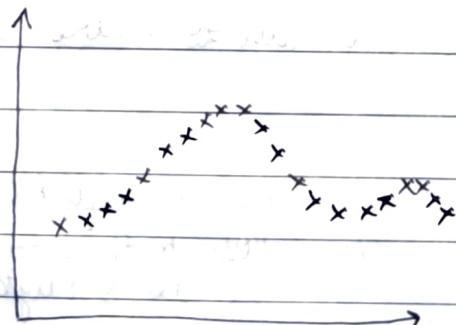
"parametric" learning algorithm → we fit set of parameters (θ_i) to data

→ example linear regression

↳ example linear regression

"Non parametric" learning algorithm → Locally weighted L.R is an example of this

↳ amount of data/parameters you need to keep grows (linearly) with size of data.



"FOR Locally weighted L.R"

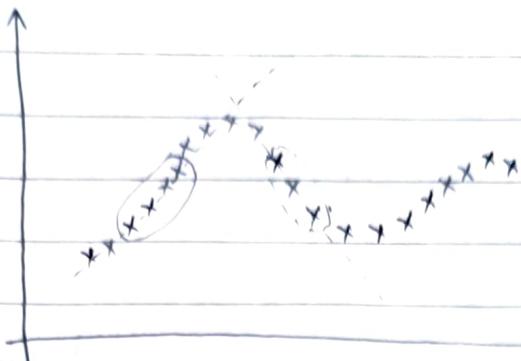
"FOR Linear regression"
to evaluate h at certain x

Fit θ to minimize.

$$\frac{1}{2} \sum_i (y^{(i)} - \theta^T x^{(i)})^2$$

return $\theta^T x$ → gives best st. line.

For locally weighted regression



in locally weighted region regression we select a subregion and try to fit a st. line through that region

Fit θ to minimise cost function given below.

$$\sum_{i=1}^m w^{(i)} (y^{(i)} - \theta^T x^{(i)})^2$$

where $w^{(i)}$ is a "weighting function".

$$w^{(i)} = \exp\left(-\frac{(x^{(i)} - x)^2}{2\tau^2}\right) \quad \text{for each point } (x^{(i)}, y^{(i)})$$

this is a value between 0 and 1.

→ points have more weight if $|x^{(i)} - x|$ is small

" " less " " " " is large

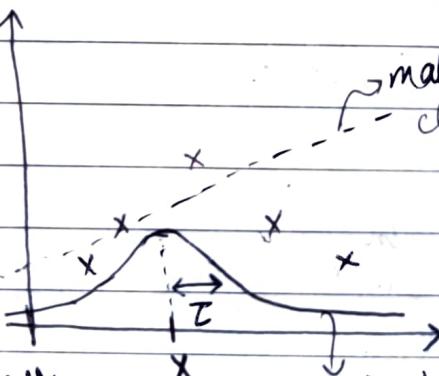
x is the ^{point} at which we have to make a prediction.

affects overfitting and underfitting

↑ lesser width → more local.

↑ to determine the width which determines the weight dist.

↑ τ is the bandwidth parameter.



makes a prediction close to this to fit local straight line

weight dist. function with maxima at X (bell shaped curve)

use locally weighted regression \rightarrow no. of features is low ($n=2$ or 3)
 \rightarrow huge amt. of data.

$\leftarrow \star$

use least squares

Probabilistic interpretation: \rightarrow

Q Why do we use least square in error?

\rightarrow Assume, $y^{(i)} = \theta^T x^{(i)} + \varepsilon^{(i)}$ \hookrightarrow error term to capture unmodelled effects.

$\varepsilon^{(i)} \sim N(0, \sigma^2)$. \hookrightarrow "twiddle" \rightarrow distributed. \hookrightarrow variance σ^2 .

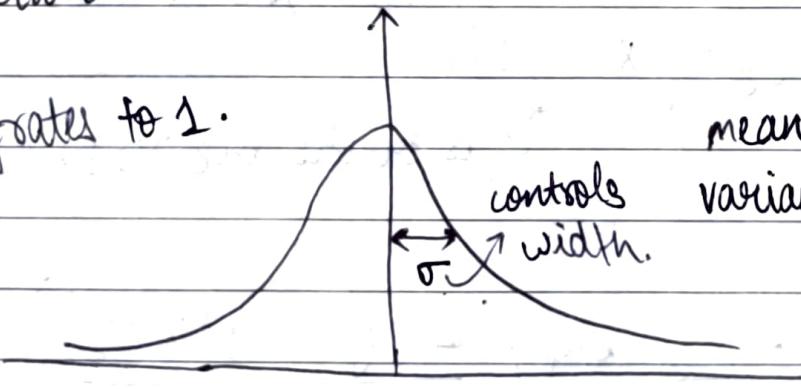
\hookrightarrow mean 0

\hookrightarrow "normal distribution" / "gaussian distribution".

$$P(\varepsilon^{(i)}) = \frac{1}{\sqrt{2\pi} \sigma} \exp\left(-\frac{(\varepsilon^{(i)})^2}{2\sigma^2}\right)$$

prob. density
function.

integrates to 1.



mean zero.

controls variance σ^2 .

\hookrightarrow "Gaussian function".

$\varepsilon^{(i)}$ are IID (Independently and Identically Distributed)

↳ Error term for one input is independent from error input term of another input.

↳ not necessarily true.

This implies : →

$$P(y^{(i)} | x^{(i)}; \theta) = \frac{1}{\sqrt{2\pi} \sigma} \exp\left(-\frac{[y^{(i)} - \theta^T x^{(i)}]^2}{2\sigma^2}\right)$$

↑ "parameterised
as"

i.e,

$$y^{(i)} | x^{(i)}; \theta \sim \mathcal{N}(\theta^T x^{(i)}, \sigma^2).$$

↳ mean ↳ variance.

↳ given x and θ what is the probability of a particular house's price

$L(\theta) = P(\vec{y} | \vec{x}; \theta) \rightarrow$ probability of y given all x 's and θ 's

↓ likelihood of " = $\prod_{i=1}^m p(y^{(i)} | x^{(i)}; \theta)$.

↳ as errors are IID.

$$= \prod_{i=1}^m \frac{1}{\sqrt{2\pi} \sigma} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right)$$

"log likelihood"

↑
 $l(\theta) = \log [L(\theta)]$.

$$l(\theta) = \log \prod_{i=1}^m \frac{1}{\sqrt{2\pi}\sigma} \exp(-\dots).$$

$$l(\theta) = \sum_{i=1}^m \left[\log \left(\frac{1}{\sqrt{2\pi}\sigma} \right) + \log \exp(-\dots) \right]$$

$$l(\theta) = m \log \left(\frac{1}{\sqrt{2\pi}\sigma} \right) + \sum_{i=1}^m -\frac{(y^{(i)} - \theta^\top x^{(i)})^2}{2\sigma^2}.$$

MLE \rightarrow maxm. likelihood estimation.

choose θ to maximise the likelihood (choose θ to maxm. $L(\theta)$)

i.e choose θ that minimises $\frac{1}{2} \sum_{i=1}^m \frac{(y^{(i)} - \theta^\top x^{(i)})^2}{2\sigma^2}$.

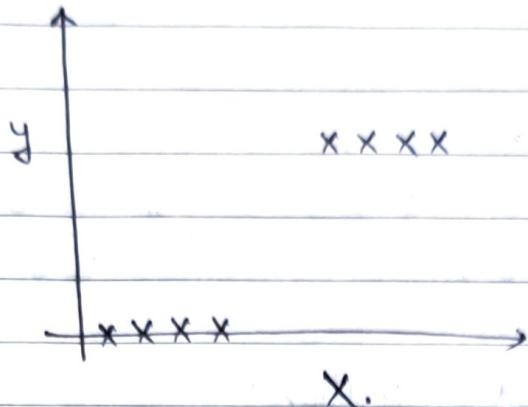
$$= J(\theta).$$

This proof shows finding θ such that $J(\theta)$ is minimum (least square errors) is finding the maximum likelihood estimate for the parameters θ under this set of assumptions that error terms are gaussian and IID.

classification :

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

[binary classification)



Logistic regression : →

want $h_{\theta}(x) \in [0, 1]$

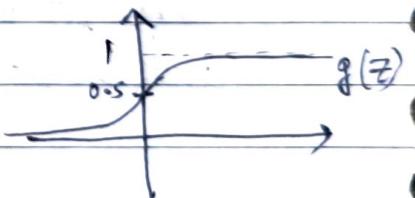
$$h_{\theta}(x) = g(\theta^T x) = \frac{1}{1+e^{-\theta^T x}}$$

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

~~xxxxxx~~

"sigmoid" or "logistic
function")

$$\{ p(y=1 | x; \theta) = h_{\theta}(x). \}$$



$$\{ p(y=0 | x; \theta) = 1 - h_{\theta}(x). \}$$

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

$$\Rightarrow p(y | x; \theta) = h(x)^y (1-h(x))^{1-y}.$$

if $y=1 \rightarrow h(x)$

$y=0 \rightarrow 1-h(x).$

$$L(\theta) = p(\vec{y} | \vec{x}; \theta)$$

$$= \prod_{i=1}^m p(y^{(i)} | x^{(i)}; \theta)$$

$$= \prod_{i=1}^m h_{\theta}(x^{(i)})^{y^{(i)}} \cdot (1 - h_{\theta}(x^{(i)}))^{1-y^{(i)}}$$

we need to find θ that maxim. likelihood of parameters.

~~$$l(\theta) = \sum_{i=1}^m y^{(i)} \log(h_{\theta}(x^{(i)})) + (1-y^{(i)}) \log(1-h_{\theta}(x^{(i)}))$$~~

find θ to maxim. $l(\theta)$.

↓
we will use
gradient ascent

Batch gradient ascent:

~~$$\theta_j := \theta_j + \alpha \frac{\partial}{\partial \theta_j} l(\theta)$$~~

$$\theta_j := \theta_j + \alpha \frac{\partial}{\partial \theta_j} l(\theta)$$

$$\boxed{\theta_j := \theta_j + \alpha \sum_{i=1}^m (y^{(i)} - h_{\theta}(x^{(i)})) \cdot x_j^{(i)}} \quad [\text{proof in lecture notes}]$$

↓ same eqn. as linear

regression but different
 $h_{\theta}(x^{(i)})$

we find parameters such
that $l(\theta)$ is maxim.

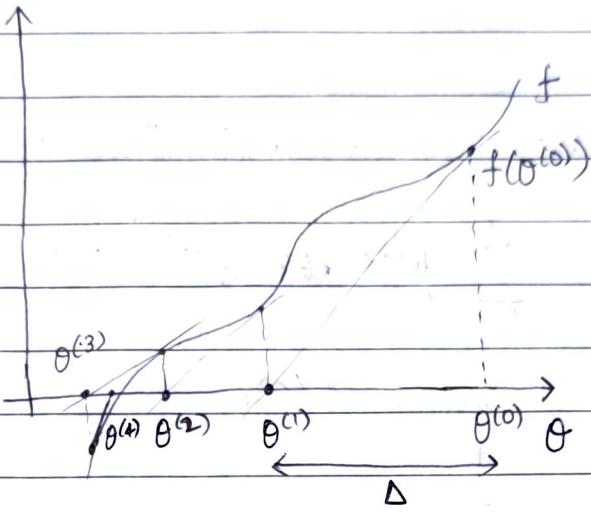
Newtons method

→ faster algorthm

→ might take too long depending on no. of parameters

given a func. $f(\theta)$ find θ s.t. $f(\theta) = 0$

[for our prob. we want to maxim. $\ell(\theta) \Rightarrow \ell'(\theta) = 0$]



at each point create the slope line and use its x-intercept as the new point.

Can overshoot too but returns back.

$$f'(\theta^{(0)}) = \frac{f(\theta^{(0)})}{\Delta}$$

$$\Delta = \frac{f(\theta^{(0)})}{f'(\theta^{(0)})}$$

$$\theta^{(1)} = \theta^0 - \frac{f(\theta^{(0)})}{f'(\theta^{(0)})}$$

$$\Rightarrow \theta^{(t+1)} = \theta^{(t)} - \frac{f(\theta^{(t)})}{f'(\theta^{(t)})}$$

$$\text{let } f(\theta) = \ell'(\theta)$$

$$\theta^{(t+1)} = \theta^{(t)} - \frac{\ell'(\theta^{(t)})}{\ell''(\theta^{(t)})}$$

Newton's method has "quadratic convergence" i.e. if after

one iteration error is 0.01 then error becomes 0.0001 then it becomes 0.00000001 and so on. Hence it takes lesser number of iterations. \rightarrow converges rapidly

when θ is a vector : $\theta \in \mathbb{R}^{n+1}$

$$\theta^{(t+1)} := \theta^{(t)} - H^{-1} \nabla_{\theta} \ell.$$

$H \rightarrow \mathbb{R}^{(n+1) \times (n+1)}$

H is a hessian matrix.

$$H_{ij} = \frac{\partial^2 \ell}{\partial \theta_i \partial \theta_j}$$



"IGNORE"

+ sign as gradient "ASCENT" for logistic regression

otherwise

for grad "descent"
"- " sign

takes large no. of computations to calculate for larger data sets.

\hookrightarrow no. of parameters = 10,000 $\rightarrow 10^{10}$ computations

\hookrightarrow at each step for hessian

\hookrightarrow Hence, gradient descent is better for cases with large number of parameters.

Exponential family:

PDF (prob-density function)

$$p(y, \eta) = b(y) \exp [\eta^T T(y) - a(\eta)]$$

$y \rightarrow$ data

$\eta \rightarrow$ natural parameter

only funct. of y $T(y) \rightarrow$ "sufficient statistic" [$= y$ for this lecture]

$b(y) \rightarrow$ base measure

only funct. of η $a(\eta) \rightarrow$ log-partition

↳ indicates normalising function
of PDF.

$$p(y, \eta) = \frac{b(y) \exp(\eta^T T(y))}{e^{+a(\eta)}}.$$

Bernoulli (Binary PDF)

$\bar{\eta}$ = prob. of an event

$$p(y; \bar{\eta}) = \bar{\eta}^y (1-\bar{\eta})^{1-y}$$

we want to make this into above form

$$= \exp \left(\log (\bar{\eta}^y (1-\bar{\eta})^{1-y}) \right)$$

$$= \exp \left[\log \left(\frac{\bar{\eta}}{1-\bar{\eta}} \right) y + \log(1-\bar{\eta}) \right]$$

$$b(y) = 1.$$

$$\eta = \log \left(\frac{\bar{\eta}}{1-\bar{\eta}} \right) \Rightarrow \bar{\eta} = \frac{1}{1+e^{-\eta}}$$

$$T(y) = y$$

$$a(\eta) = -\log(1-\bar{\eta}) = -\left(1 - \frac{1}{1+e^{-\eta}}\right) = \log(1+e^\eta)$$

Gaussian (with fixed variance)

Assume $\sigma^2 = 1$

$$P(y; \mu) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(y-\mu)^2}{2}\right)$$

↓
canonical
parameters

$$= \frac{1}{\sqrt{2\pi}} e^{-y^2/2} \exp\left(\mu y - \frac{1}{2}\mu^2\right)$$

$b(y)$ η $T(y)$ $a(\eta)$.

$$b(y) = \frac{1}{\sqrt{2\pi}} e^{-y^2/2}.$$

$$T(y) = y$$

$$\eta = \mu$$

$$a(\eta) = \frac{\mu^2}{2} = \frac{\eta^2}{2}$$

Properties of exponential families

- ① MLE w.r.t $\eta \rightarrow$ concave
- (\downarrow
max. likelihood
of exponential
family)

NLL is convex
(negative log)
likelihood

- ② $E[y; \eta] = \frac{\partial}{\partial \eta} a(\eta).$ } easy as no integrals
- ③ $\text{Var}[y; \eta] = \frac{\partial^2}{\partial \eta^2} a(\eta)$ }

GLM (generalised linear models)

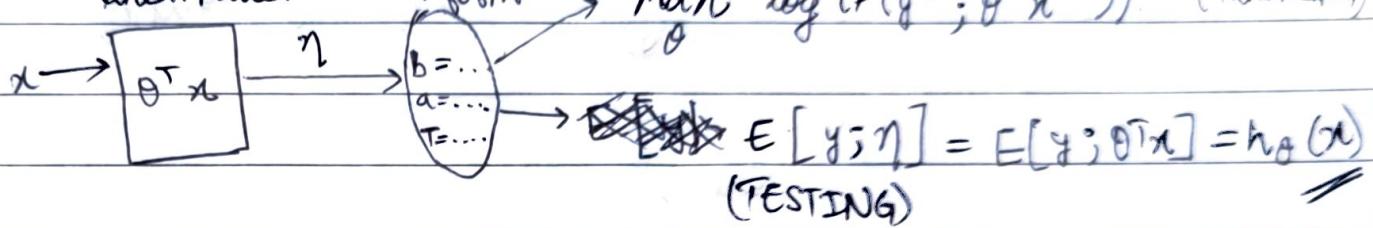
assumptions / design choices for GLM

(i) $y|x; \theta \sim$ member of Exponential family

(ii) $\eta = \theta^T x$

(iii) Test time : output $E[y|x;\theta]$

$\Rightarrow h_\theta(x) = E[y|x;\theta]$ → gradient ascent on log likelihood
linear model exp form $\max_{\theta} \log(P(y^{(i)}; \theta^T x^{(i)}))$ (TRAINING)



θ to train the exponential family model whose mean predicts the ~~value of~~ hypothesis.

GLM Training

for each GLM learning update rule is the same,

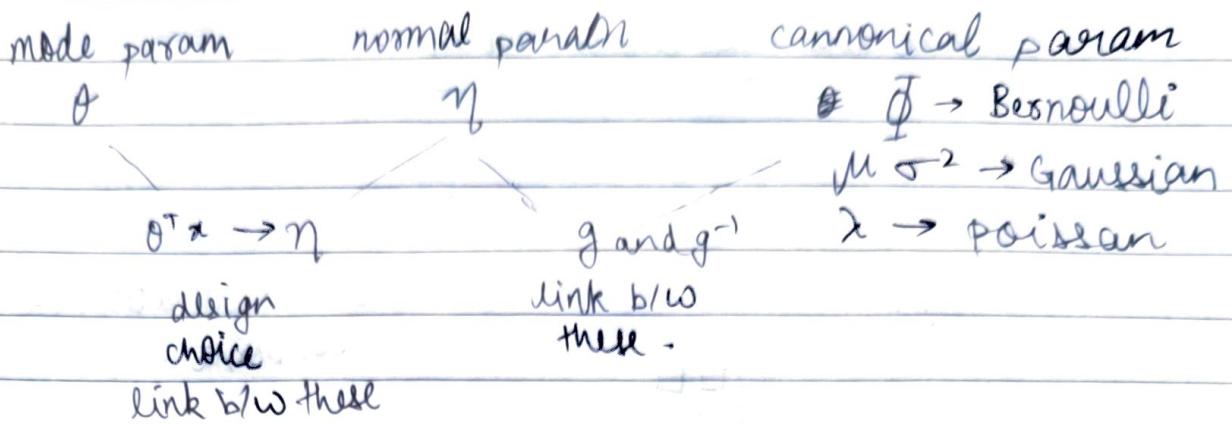
$$\theta_j := \theta_j + \alpha (y^{(i)} - h_\theta(x^{(i)})) x_j^{(i)}$$

Terminology $\eta \rightarrow$ natural parameter

$$\mu \quad E[y; \eta] = g(\eta) \rightarrow \text{canonical response function}$$

$$\eta = g^{-1}(\mu) \rightarrow \text{"link function"}$$

3 parameterisations for GLM



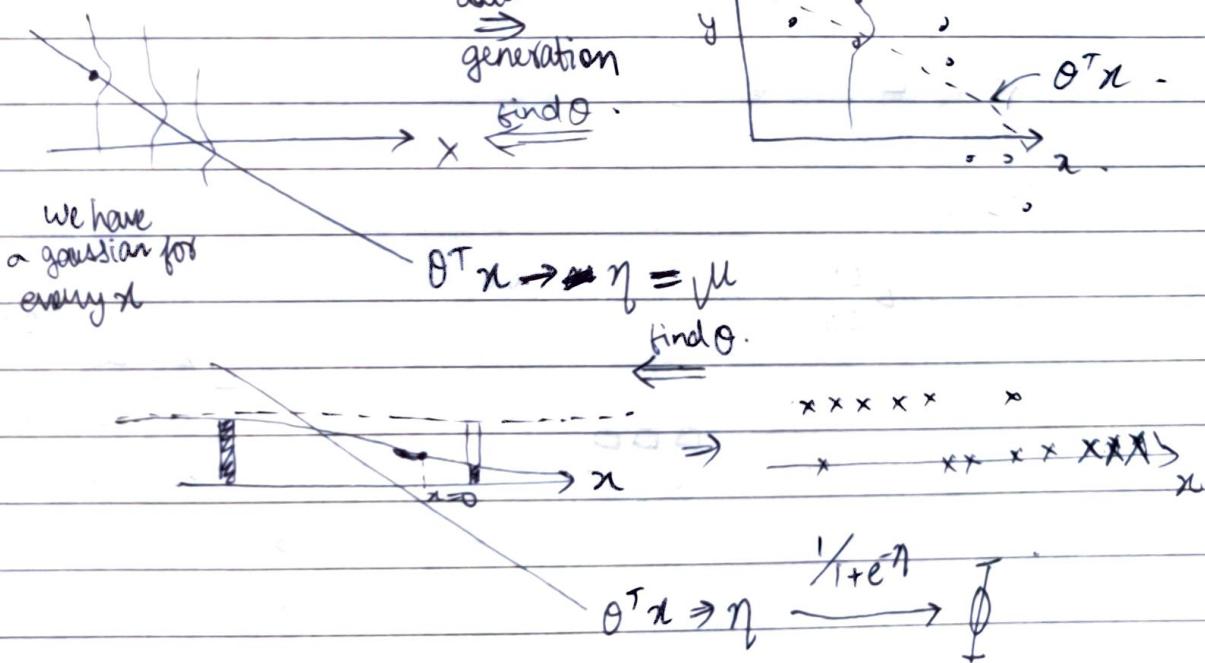
Logistic regression \rightarrow Bernoulli

$$h_\theta(x) = E[y|x; \theta] = \bar{\phi} = \frac{1}{1+e^{-\eta}} = \frac{1}{1+e^{\theta^T x}}$$

We pick GLM distribution depending on the task at hand.

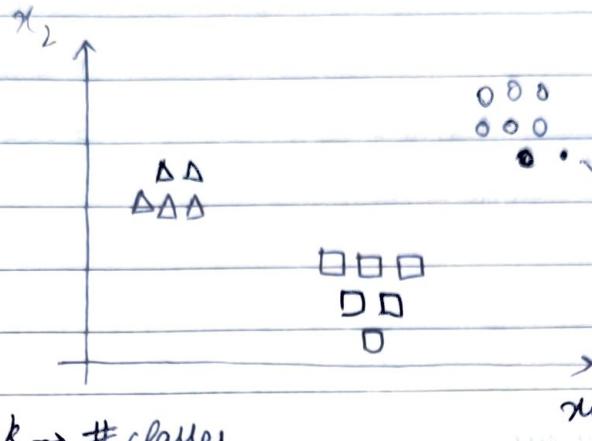
Assumptions

Regulation



softmax regression

is non GLM but we get similar eqns.
now softmax does cross entropy min.



• predict classification of new point.

$k \rightarrow \# \text{ classes}$
 $\mathbf{x}^{(i)} \in \mathbb{R}^n$

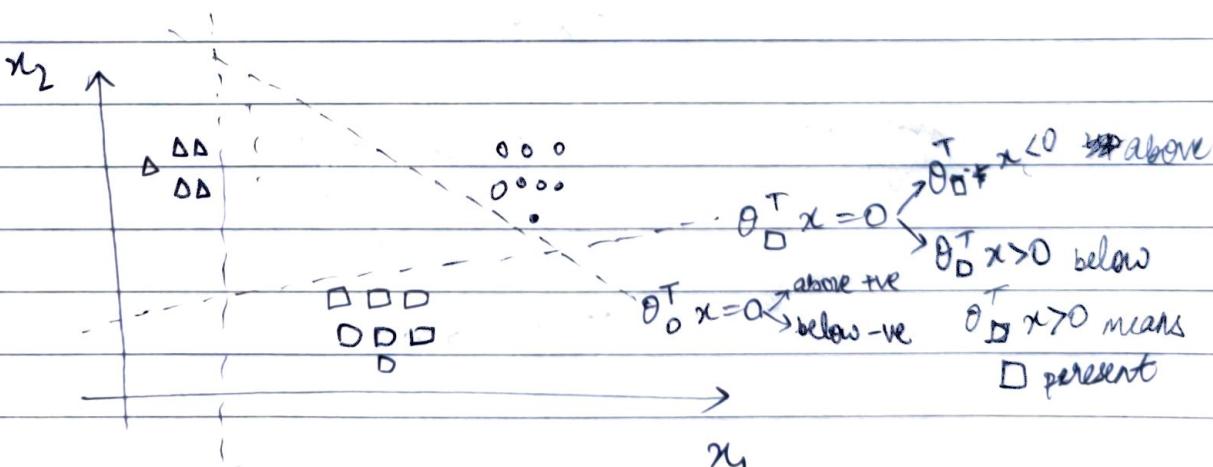
label $y = [\xi_0, \xi_1]^k$ eg. $[0, 0, 1, 0]$

for $k=4$

tells state for each classification

$\theta_{\text{class}} \in \mathbb{R}^n$ k such class $\in \{\Delta, \circ, \square, \dots\}$

$$\theta_i = \begin{bmatrix} \theta_{i1} \\ \theta_{i2} \\ \vdots \end{bmatrix} \quad \theta_i \text{ is all class parameters for } i$$



$\theta_\Delta^T x = 0$
 $\theta_\Delta^T x > 0 \Rightarrow \Delta \text{ left}$
 $\theta_\Delta^T x < 0 \Rightarrow \Delta \text{ right}$
 $\theta_\Delta^T x > 0 \Rightarrow \Delta \text{ present}$

given x

logit
space

let's say

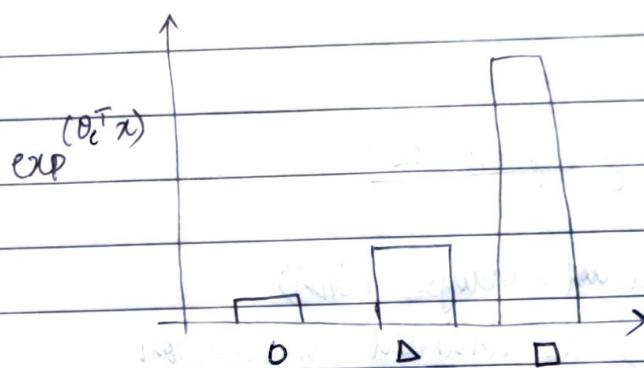
class here is \square

we need to
find a distribution

s.t. for $0 \& \Delta$ the
value is ~~less than~~ less than \square

$\downarrow \exp^{(\theta_i^T x)}$

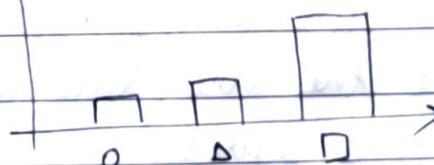
first exponentiate over logit
func.



\downarrow normalise this \Rightarrow divide all by their sum

$$\frac{e^{\theta_i^T x}}{\sum_{i=1,2,3} e^{\theta_i^T x}}$$

$\hat{p}(y) \rightarrow$ is function that gives
us prob. dist. over classes

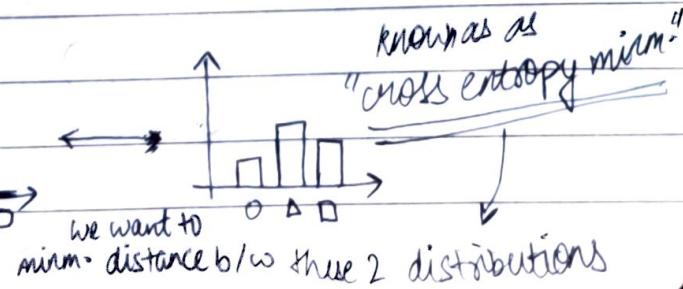
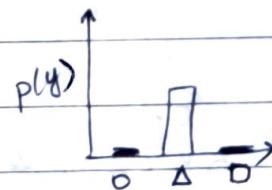


sum of heights

= 1

then y will look like

$p(y)$ for $\Delta \rightarrow$



we want to
min. distance b/w these 2 distributions

$$\text{Cross Ent}(p, \hat{p}) = \sum_{y \in \{0, 1, D\}} p(y) \log \hat{p}(y)$$

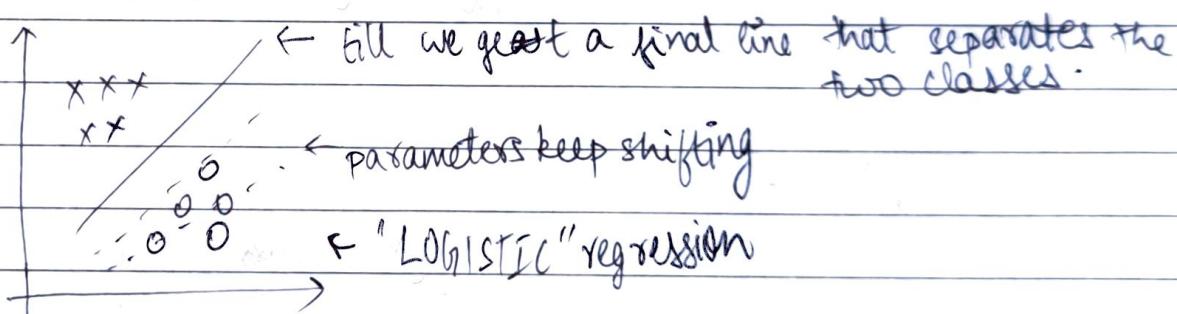
$$\text{for let's say } \Delta \rightarrow = -\log(\hat{p}(y_\Delta))$$

$$\text{Cross Ent}(p, \hat{p}) = -\log \left(\frac{e^{\theta_\Delta^T x}}{\sum_{c \in \{0, 1, D\}} e^{\theta_c^T x}} \right)$$

do gradient descent w.r.t. parameters.

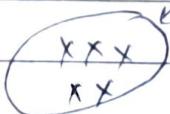
Generative learning algorithms

- Gaussian discriminant analysis (GDA)
- Generative & discriminant comparison
- Naive Bayes

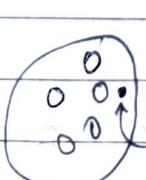


Generative learning algorithm → rather than trying to get the line, this algorithm looks at all the classes separately

all malignant live in this area



all benign live in this



predicts for new datapoint

Discriminative:

learn $p(y|x)$

(or learn $h_0(x) = \begin{cases} 0 \\ 1 \end{cases}$ directly)

Generative LA:

learns $p(x|y)$ → given a class → what would its
features ↑ ∈ class features look like

$p(y)$ → prob. of a class.

↳ a patient walks in → what is
prob. that he has
cancer malignant

acc. to bayes rule:

$$P(y=1|x) = \frac{P(x|y=1)p(y=1)}{P(x)}$$

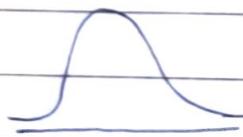
$$P(x) = P(x|y=1)p(y=1) + P(x|y=0)p(y=0)$$

Gaussian discriminative analysis. (GDA)

$x \in \mathbb{R}^n$ (drop $x_0 = 1$ convention)

assume $p(x|y)$ is Gaussian

multivariable



single var.
gaussian.

$$z \sim N(\vec{\mu}, \Sigma) \quad z \in \mathbb{R}^n$$

expected value \downarrow mean $\vec{\mu} \in \mathbb{R}^n$
 $E[z] = \mu$

$$\Sigma \in \mathbb{R}^{n \times n}$$

$$\text{cov}(z) = E[(z - \mu)(z - \mu)^T]$$

$$\text{cov}(z) = E[zz^T] - (Ez)(Ez)^T$$

multivariable gaussian distribution

$$p(x) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu)\right)$$

(Refer lecture notes.)

GDA model

$$p(x|y=0) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2} (x-\mu_0)^T \Sigma^{-1} (x-\mu_0)\right)$$

↓
same Σ

$$p(x|y=1) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2} (x-\mu_1)^T \Sigma^{-1} (x-\mu_1)\right)$$

parameters for both eqns. are $\mu_0, \mu_1, \Sigma, \Phi$

$\Phi \in \mathbb{R}^{n \times n}$ $\mu_0 \in \mathbb{R}^n$ $\mu_1 \in \mathbb{R}^n$ $\Sigma \in \mathbb{R}^{n \times n}$

$$p(y) = \Phi^y (1-\Phi)^{1-y} \quad \leftarrow \text{random bernoulli variable.}$$

suppose we have a training set: $\{(x^{(i)}, y^{(i)})\}_{i=1}^m$

Joint likelihood:

$$\begin{aligned} L(\Phi, \mu_0, \mu_1, \Sigma) &= \prod_{i=1}^m p(x^{(i)}, y^{(i)} | \Phi, \mu_0, \mu_1, \Sigma) \\ &= \prod_{i=1}^m p(x^{(i)} | y^{(i)}, \mu_0, \mu_1, \Sigma) p(y^{(i)}) \end{aligned}$$

we maxmize the above function.

for generative algo we do not find the parameters.

↳ we try to know $p(x|y)$.

for discriminative algos

conditional likelihood

$$L(\theta) = \prod_{i=1}^m p(y^{(i)} | x^{(i)}; \theta)$$

we try to find θ
s.t. cond. likelihood
is maximized.

Maximum Likelihood estimation

$$\max_{\theta, \mu_0, \mu_1, \Sigma} l(\theta, \mu_0, \mu_1, \Sigma) = \log L(\dots).$$

take derivative ; equate it to 0 ; solve for parameters

final answer $\Rightarrow \hat{\theta} = \frac{\sum_{i=1}^m y^{(i)}}{m}$ "indicator"

$$= \frac{\sum_{i=1}^m 1 \{y^{(i)} = 1\}}{m}$$

$$\mu_0 = \frac{\sum_{i=1}^m 1 \{y^{(i)} = 0\} x^{(i)}}{\sum_{i=1}^m 1 \{y^{(i)} = 0\}}$$

$1(\text{true}) = 1$
 $1(\text{false}) = 0$.

$$\mu_1 = \frac{\sum_{i=1}^m 1 \{y^{(i)} = 1\} x^{(i)}}{\sum_{i=1}^m 1 \{y^{(i)} = 1\}}$$

← sum of feature vectors
for all ex. with $y=1$

$$\sum_{i=1}^m 1 \{y^{(i)} = 1\}$$

← no. of examples
with $y = 1$.

$\bar{\theta}$
mean of parameters
of $y=1$ values

$$\Sigma = \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \mu_{y(i)}) (x^{(i)} - \mu_{y(i)})^T$$

↳ trying to fit ellipses around the mean for a class.

Prediction: $\min_z (z - s)^2 = 0.$ $\min_z \rightarrow \min_{\text{over } z}$
 $\arg \min_z (z - s)^2 = s //$

$$\arg \max_y p(y|x) = \arg \max_y \frac{p(x|y) p(y)}{p(x)}.$$

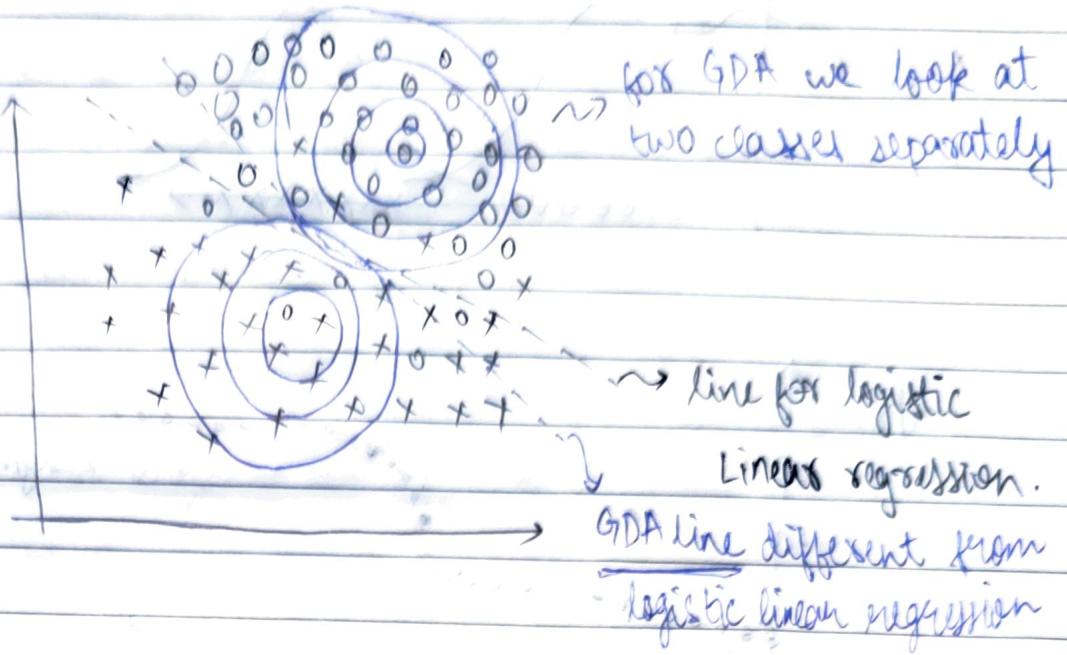
$\min_z f(z) \rightarrow$ returns smallest attainable value

$\arg \min_z f(z) \rightarrow$ returns the argument for the smallest value.

$$\arg \max_y p(y|x) = \arg \max_y \frac{p(x|y) p(y)}{p(x)}$$

" ↳ constant as
 $\arg \max_y p(x|y) p(y)$ independent of y
 as we only want the argument.

Discriminative vs generative learning algo : →



Compare GDA to logistic regression: →

for fixed $\Phi, \mu_0, \mu_1, \Sigma$, let's

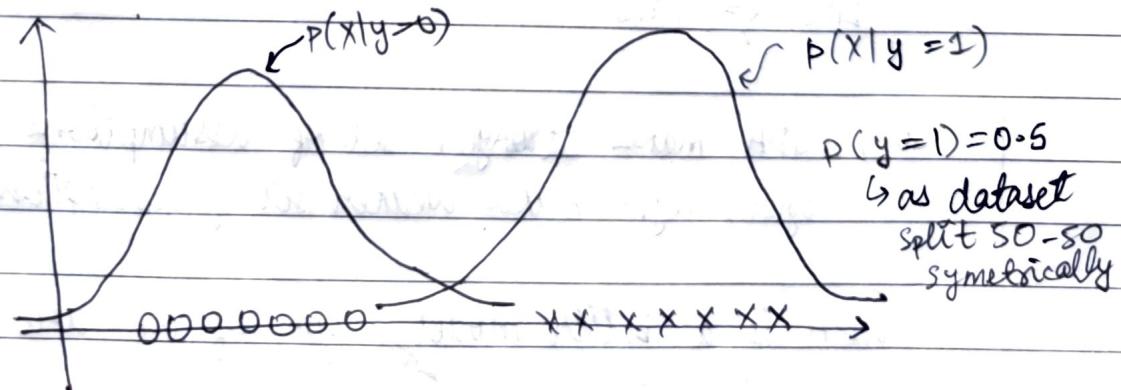
plot $p(y=1|X; \Phi, \mu_0, \mu_1, \Sigma)$ as a function of X .
→ Bernoulli prop.

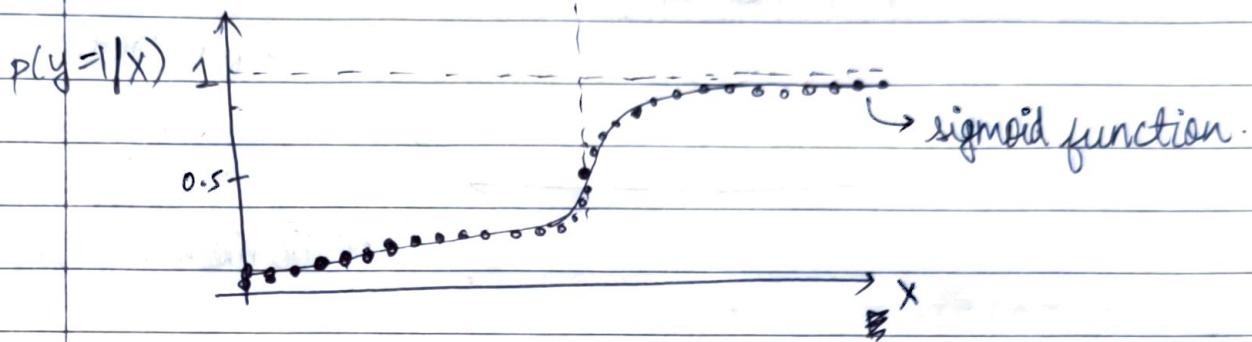
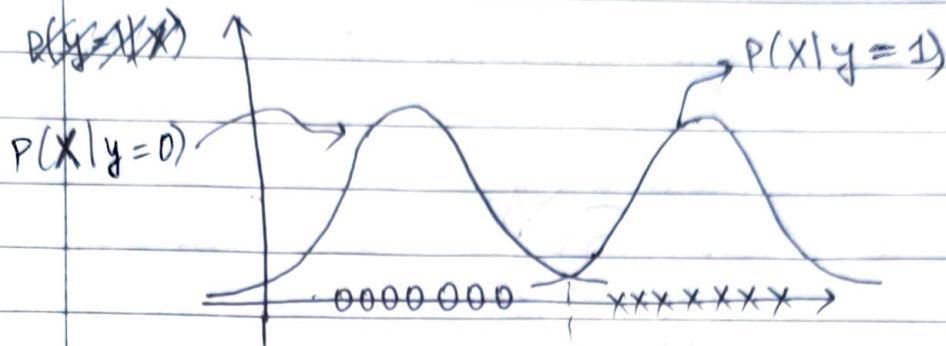
$$= p(x|y=1; \mu_1) / [p(x; \Phi, \mu_0, \mu_1, \Sigma)]$$

$$p(x; \Phi, \mu_0, \mu_1, \Sigma)$$

case where we only have 1 feature

— OOOOOOOO — X X X X X X X — x





when is GDA superior or inferior to Logistic regression:

$$\begin{aligned}
 & \text{GDA} \\
 & x|y=0 \sim N(\mu_0, \Sigma) \\
 & x|y=1 \sim N(\mu_1, \Sigma) \\
 & y \sim \text{Ber}(\theta) \\
 & \text{LHS} \\
 & \text{implies:} \\
 & \text{RHS} \\
 & \text{Logistic} \\
 & \text{RHS}
 \end{aligned}$$

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

LHS \Rightarrow RHS

LHS $\not\Rightarrow$ RHS

so, GDA makes stronger set of assumptions
logistic reg. makes weaker set of assumptions

GDA is a better model ~~for~~ if we are given LHS condition.

GDA is also computationally efficient.

lets say;

$$\left. \begin{array}{l} x|y=1 \sim \text{poisson}(\lambda_1) \\ x|y=0 \sim \text{poisson}(\lambda_0) \\ y \sim \text{bernoulli}(p) \end{array} \right\} \Rightarrow p(y=1|x) \text{ is logistic}$$

LHS \Rightarrow RHS

logistic regression is a rough fit if we don't know if our data is poison or gaussian.

Naive Bayes Learning algorithm :-

This is our first intro to NLP.
we have been given email text and we want
to classify it as spam or not spam.

What will be the feature vector?

pick the 10000 most occurring words from your email-

↳ binary feature vector

$$x \in \{0, 1\}^n$$

$x_i = 1$ { word i appears in email }

in naive bayes we finally want;

$p(y|x)$ so we'll need to model $p(x|y)$ and $p(y)$

$2^{10,000}$ possible values of $X \rightarrow$ not possible X

Assume x_i 's are conditionally independent given y

$$p(x_1, x_2, \dots, x_{10000}|y) = p(x_1|y) \cdot p(x_2|x_1, y) \cdot p(x_3|x_1, x_2, y) \cdots p(x_{10000}|\dots)$$

nothing assumed yet
this is chain rule of prob.

we here assume that

$$\text{# LHS} = p(x_1|y) \cdot p(x_2|y) \cdot p(x_3|y) \cdots p(x_{10000}|y).$$

as we assume that they are independent
of each other

$$= \prod_{i=1}^n p(x_i|y)$$

parameters: if a spam email \rightarrow what is chance of j appearing

$$\hat{\Phi}_{j|y=1} = p(x_j=1 | y=1)$$

\sim if not a spam \rightarrow what is chance of j appearing.

$$\hat{\Phi}_{j|y=0} = p(x_j=1 | y=0)$$

$\hat{\Phi}_y = p(y=1) \sim$ prob. that next email you receive is a spam email

Joint likelihood:

$$L(\hat{\Phi}_y, \hat{\Phi}_{j|y}) = \prod_{i=1}^m p(x^{(i)}, y^{(i)}; \hat{\Phi}_y, \hat{\Phi}_{j|y})$$

maxm. likelihood estimates (MLE):

$$\hat{\Phi}_y = \frac{\sum_{i=1}^m 1 \{ y^{(i)} = 1 \}}{m}$$

$$\hat{\Phi}_{j|y=1} = \frac{\sum_{i=1}^m 1 \{ x_j^{(i)} = 1, y^{(i)} = 1 \}}{\sum_{i=1}^m 1 \{ y^{(i)} = 1 \}}$$

\sim find all spam containing j th word

$$\hat{\Phi}_{j|y=0} = \frac{\sum_{i=1}^m 1 \{ x_j^{(i)} = 1, y^{(i)} = 0 \}}{\sum_{i=1}^m 1 \{ y^{(i)} = 0 \}}$$

$$\hat{\Phi}_{j|y=0} = \frac{\sum_{i=1}^m 1 \{ x_j^{(i)} = 1, y^{(i)} = 0 \}}{\sum_{i=1}^m 1 \{ y^{(i)} = 0 \}}$$

proof in lecture notes : 2

at prediction time

$$p(y=1|x) = \frac{p(x|y=1) p(y=1)}{p(x|y=0)p(y=0) + p(x|y=1)p(y=1)}$$

$$P(x|y=0)p(y=0) + p(x|y=1)p(y=1)$$

This algorithm almost works but raises a very important problem.

↳ what happens when we get an email containing a new word or a less occurred word.

↳ What do we take the prob. as

↳ eg: $\rightarrow x_{6017} = "NIPS"$ which we don't get in mails
that often

$$\underset{6017|y=1}{\rightarrow} p(x_{6017}=1|y=1) = \frac{0}{\#\{y=1\}}$$

bad idea to assume it is zero if we haven't seen it

$$\underset{6017|y=0}{\rightarrow} p(x_{6017}=1|y=0) = \frac{0}{\#\{y=0\}}$$

so $\cancel{p(x|y=1)} = \boxed{\frac{1}{m} \sum_{i=1}^m p(x_i|y)}$ becomes zero so

gives an error.

so, what do we take our new probability to be?

we use Laplace's smoothing

Laplace smoothing

Let's say a football team play 5 games against A, B, C, D, E.

	win / lost (X)	0 → lost	1 → won
A	0		
B	0		
C	0		
D	0		
E	2		

what is the prob of them winning E.

$$P(X=1) = \frac{\#\text{"1"s}}{\#\text{"1"s} + \#\text{"0"s}}$$

$$= \frac{0}{4} = 0 \quad X \text{ not a good estimate.}$$

what Laplace smoothing does?

↪ imagine # "1"s was # "1"s + 1

and # "0"s was # "0"s + 1

optimal
estimate

more generally

$$X \in \{1, \dots, k\}$$

estimate

$$P(X=j) = \frac{\sum_{i=1}^M 1_{\{X^{(i)}=j\}} + 1}{M+k}$$

we applied
Laplace smoothing

$$\hat{P}_{j|y=0} = \frac{\sum_{i=1}^m \mathbb{1}\{x_j^{(i)}=1, y^{(i)}=0\}}{m} + 1$$

j can be 0 or 1.

$$\sum_{i=1}^m \mathbb{1}\{y^{(i)}=0\} + 2$$

Laplace smoothing

New representation of Naive Bayes →

Multivariate Bernoulli event model

Old:	index
0	a ↠ 1
0	aardvark ↠ 2
:	:
1	buy ↠ 800
:	
1	drugs ↠ 1600
:	
i	now ↠ 6200

New: Multinomial event model representation

$$X = \begin{bmatrix} 1600 \\ 800 \\ 1600 \\ 6200 \end{bmatrix} \in \mathbb{R}^{n_i}$$

$n_i \rightarrow$ no. of words in our email

$$x_j \in \{1, 2, \dots, 10000\}$$

n_i is length of our email

Drugs: buy drugs now

Above is the conventional representation

We lost the # time a word occurs

drugs is stored

as 1 only

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

$$P(X, y) = P(X|y) \cdot P(y)$$

$$= \prod_{j=1}^n P(x_j|y) \cdot P(y)$$

$$\hat{P}_y = P(y=1)$$

$$\hat{P}_{k|y=0} = P(x_j=k|y=0)$$

Chance of word jth being k if $y=0$

$$\hat{\Phi}_{k|y=1} = p(x_j=k | y=1)$$

MLE:

$$\hat{\Phi}_{k|y=1} = \sum_{i=0}^m 1 \{y^{(i)}=0\} \sum_{j=1}^{n_i} 1 \{x_j^{(i)}=k\} + 1$$

$$\sum_{i=1}^m 1 \{y^{(i)}=0\} = 0.5 \cdot n_i + 10,000$$

no. of outcomes
possible

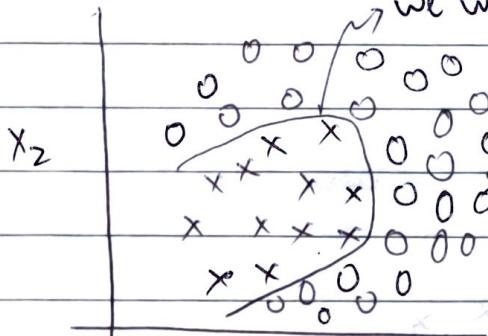
look at all words in all non-spam emails

↳ what fraction of them have x_j occurring in them.

↳ also takes repetition into account in a mail



Support vector machine



we want to get a non linear decision boundary

can't use basic logistic regression

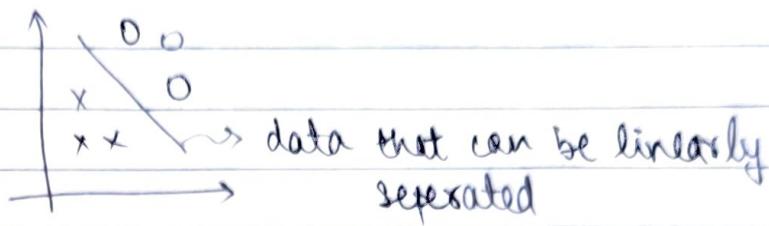
feature vector
given

$$f(x) = \begin{bmatrix} x_1 \\ x_2 \\ x_1^2 \\ x_2^2 \\ x_1 x_2 \end{bmatrix}$$

we map our original feature vector to a higher dimensional feature vector

then we apply logistic regression for this new vector to get non linear decision boundary.

Optimal margin classifier (separable case)

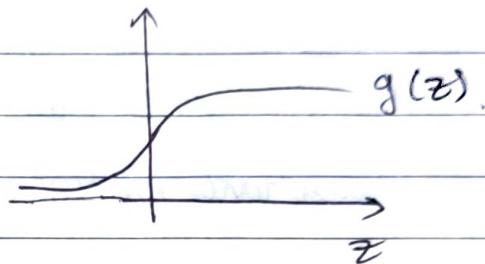


functional margin of a classifier → how well/confidently we classify an example

$$h_{\theta}(x) = g(\theta^T x)$$

predict "1" $\theta^T x > 0$

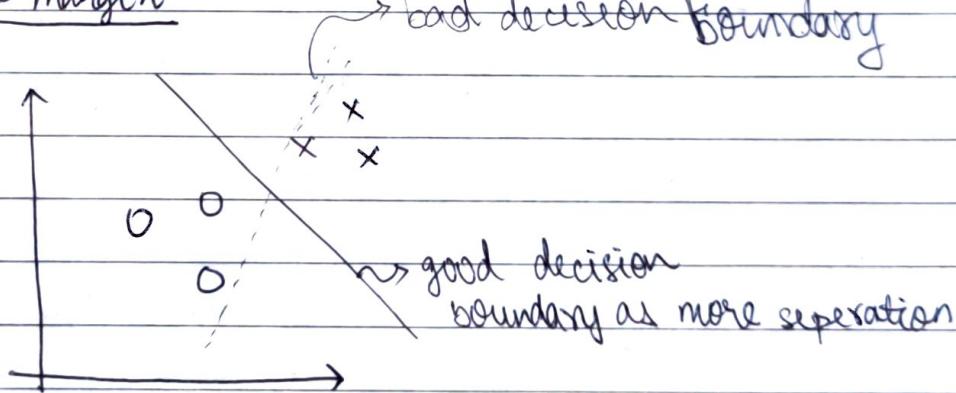
"0" otherwise



If $y^{(i)} = 1$ we want $\theta^T x \gg 0$

If $y^{(i)} = 0$ we want $\theta^T x \ll 0$

Geometric margin



Notations for SVM

labels $y \in \{-1, +1\}$

n outputs $\{-1, +1\}$.

$g(z) = 1$ if $z \geq 0$

-1 otherwise

previously for \mathbb{R}

$$h_{\theta}(x) = g(\theta^T x)$$

$\mathbb{R}^{n+1}, x_0 = 1$

$$h_{w,b}(x) = g(w^T x + b) \quad \begin{matrix} \downarrow \mathbb{R}^n \\ \downarrow \mathbb{R} \end{matrix} \quad \begin{matrix} \text{drop } x_0 = 1 \\ \text{convention} \end{matrix}$$

Function margin of (w, b) wrt $(x^{(i)}, y^{(i)})$

read as "gamma head" $\hat{y}^{(i)} = y^{(i)} (w^T x^{(i)} + b)$ \rightarrow for one training sample

if $y^{(i)} = 1$ we want $w^T x^{(i)} + b > 0$

if $y^{(i)} = -1$ we want $w^T x^{(i)} + b < 0$

\therefore we always want $\hat{y}^{(i)} \gg 0$

if $\hat{y}^{(i)} > 0$, means $h(x^{(i)}) = y^{(i)}$

functional margin wrt training set

$$\hat{\gamma} = \min_i \hat{y}^{(i)} \quad i = 1, 2, \dots, m$$

for training set. \rightarrow pick minimum from all training cases.

we want $\|w\| = 1 \rightarrow$ to normalise the length of functional margin

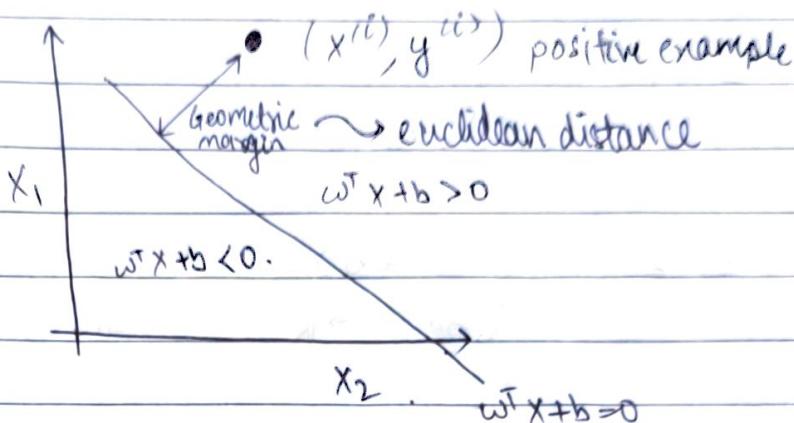
$$(w, b) \rightarrow \left(\frac{w}{\|w\|}, \frac{b}{\|w\|} \right)$$

because suppose if we multiply all param by 2

functional margin doubled tho

\leftarrow decision boundary doesn't change

Geometric margin:



Geometric margin of hyperplane (w, b) wrt $(x^{(i)}, y^{(i)})$

$$\check{y}^{(i)} = \frac{(w^T x^{(i)} + b) \cdot y^{(i)}}{\|w\|}$$

no cap

$$y^{(i)} = \frac{\hat{y}^{(i)}}{\|w\|}$$

$\hat{y} \rightarrow \text{functional}$
 $y \rightarrow \text{geometric}$

Geometric margin wrt training set.

$$y = \min_i y^{(i)}$$

Optimal margin classifier → choose w, b to maximise y

$$\max_{w, b} y$$

$$\text{s.t. } \frac{y^{(i)}(w^T x^{(i)} + b)}{\|w\|} \geq y \quad i=1, 2, \dots, m$$

optimisation prob.

$$\left. \begin{array}{l} \min_{w, b} \|w\|^2 \\ \text{s.t. } y^{(i)}(w^T x^{(i)} + b) \geq 1 \end{array} \right\}$$

LHS can be rewritten as ↓
"convex" optimisation problem
↪ solvable easily.

derivation: →

$$y^{(i)} = \frac{y^{(i)}(\mathbf{w}^\top \mathbf{x}^{(i)} + b)}{\|\mathbf{w}\|}$$

$$\|\mathbf{w}\|$$

~~$$\min(\gamma^{(i)}) = \gamma^*$$~~

so, we want to find a value γ such that it is always less than $\gamma^{(i)}$ and we want to maximise this value

$$\max \gamma \text{ s.t.}$$

$$\frac{y^{(i)}(\mathbf{w}^\top \mathbf{x}^{(i)} + b)}{\|\mathbf{w}\|} \geq \gamma \quad \text{for all } i$$

example: → geometric margins are ≥ 1 γ will be 2 here

$$\geq 2$$

$$\geq 5$$

↓
max value ~~set~~ for γ such
that $\gamma \leq \gamma^{(i)}$

We can scale \mathbf{w} by any factor and decision boundary remains the same

$$\text{let } \gamma' \in \|\mathbf{w}\| = \frac{1}{\gamma}$$

we want

~~$$\max_{\gamma, \mathbf{w}, b} \gamma$$~~

$$y^{(i)}(\mathbf{w}^\top \mathbf{x}^{(i)} + b) \geq \gamma \quad \text{for all } i$$

$$\max \gamma$$

$$\gamma, \mathbf{w}, b$$

$$y^{(i)}(\mathbf{w}^\top \mathbf{x}^{(i)} + b) \geq 1 \Rightarrow \min_{\gamma, \mathbf{w}, b} \|\mathbf{w}\| \quad y^{(i)}(\mathbf{w}^\top \mathbf{x}^{(i)} + b) \geq 1$$

$$\min_{\mathbf{w}, b} \|\mathbf{w}\|^2$$

$$\text{s.t. } y^{(i)}(\mathbf{w}^\top \mathbf{x}^{(i)} + b) \geq 1$$

— X —

Optimal margin classifier only works when $\mathbf{x}^{(i)} \in \mathbb{R}^{100}$

↳ finite

what if our dataset has 100 trillion or infinite features

↳ use SVM

assume that,

$$\mathbf{w} = \sum_{i=1}^m \alpha_i \mathbf{x}^{(i)} y^{(i)}$$

→ \mathbf{w} can be represented
as a linear combination

of the features

this assumption turns
out to be the optimum value → proof in lecture notes

$$w = \sum_{i=1}^m \alpha_i y^{(i)} x^{(i)}$$

$$\begin{array}{ll} \min_{w,b} & \frac{1}{2} \|w\|^2 \\ \text{s.t.} & y^{(i)}(w^T x^{(i)} + b) \geq 1 \quad i=1, \dots, m \end{array}$$

$$\min \frac{1}{2} \left(\sum_{i=1}^m \alpha_i y^{(i)} x^{(i)} \right)^T \left(\sum_{i=1}^m \alpha_i y^{(i)} x^{(i)} \right)$$

$$= \min \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y^{(i)} y^{(j)} \underbrace{x^{(i)} \cdot x^{(j)}}_{\langle x^{(i)}, x^{(j)} \rangle} \xrightarrow{\text{"inner product of vectors"}}$$

s.t.

$$y^{(i)} \left(\sum_j \alpha_j y^{(j)} \langle x^{(i)}, x^{(j)} \rangle + b \right) \geq 1 \rightarrow \text{putting val. of } w$$

can be written as

$$\max \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j y^{(i)} y^{(j)} \alpha_i \alpha_j \langle x^{(i)}, x^{(j)} \rangle$$

s.t. $\alpha_i \geq 0$

$$\sum_i y^{(i)} \alpha_i = 0$$

"dual optimisation prob."

- ① solve for α_i and b
- ② making prediction from this

$$h(x) = g(w^T x + b) = g\left(\left(\sum_i \alpha_i y^{(i)} \langle x^{(i)}, x \rangle + b\right)\right)$$

we basically need the ability to compute these inner products quickly \rightarrow kernel helps with this.

Kernel trick:

$\langle x, z \rangle$

- 1) write algo in terms of $\langle x^{(i)}, x^{(j)} \rangle$
- 2) let there be mapping from $x \rightarrow \Phi(x)$.

- 3) find way to compute

$$K(x, z) = \Phi(x)^T \Phi(z)$$

- 4) Replace $\langle x, z \rangle$ in algorithm with $K(x, z)$

ex. let $x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_1x_2 \\ x_3x_4 \\ x_1x_2x_3 \end{bmatrix}$

Examples of above trick: →

$$\text{let } x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

$$x \rightarrow \Phi(x) = \begin{bmatrix} x_1x_1 \\ x_1x_2 \\ x_1x_3 \\ x_2x_1 \\ x_2x_2 \\ x_2x_3 \\ x_3x_1 \\ x_3x_2 \\ x_3x_3 \end{bmatrix}$$

$$z \rightarrow \Phi(z) = \begin{bmatrix} z_1z_1 \\ z_1z_2 \\ z_1z_3 \\ z_2z_1 \\ z_2z_2 \\ z_2z_3 \\ z_3z_1 \\ z_3z_2 \\ z_3z_3 \end{bmatrix}$$

$$z = \begin{bmatrix} z_1 \\ z_2 \\ z_3 \end{bmatrix}$$

$$\begin{bmatrix} z_1z_1 \\ z_1z_2 \\ z_1z_3 \\ z_2z_1 \\ z_2z_2 \\ z_2z_3 \\ z_3z_1 \\ z_3z_2 \\ z_3z_3 \end{bmatrix}$$

$$x \in \mathbb{R}^n$$

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

$$\Phi(x), \Phi(z) \in \mathbb{R}^{n^2}$$

we need n^2 time to compute $\Phi(x)$ or $\Phi(x)^T \Phi(z)$

$$K(x, z) = \Phi(x)^T \Phi(z) = (x^T z)^2$$

$O(n)$

$O(n)$

$O(n)$.

Proof: →

$$\left(\sum_{i=1}^n x_i z_i \right) \left(\sum_{j=1}^n x_j z_j \right)$$

$$= \sum_i \sum_j x_i z_i (x_j z_j) = \sum_{i=1}^n \sum_{j=1}^n (x_i x_j)(z_i z_j) = \Phi(x)^T \Phi(z)$$

suppose our data was

$$\Phi(x) = \begin{bmatrix} x_1 x_1 \\ x_1 x_2 \\ x_1 x_3 \\ \vdots \\ x_3 x_3 \\ \sqrt{2c} x_1 \\ \sqrt{2c} x_2 \\ \sqrt{2c} x_3 \\ c \end{bmatrix}$$

$$\Phi(z) = \begin{bmatrix} z_1 z_1 \\ z_1 z_2 \\ z_1 z_3 \\ \vdots \\ z_3 z_3 \\ \sqrt{2c} z_1 \\ \sqrt{2c} z_2 \\ \sqrt{2c} z_3 \\ c \end{bmatrix}$$

$$K(x, z) = \Phi(x)^T \Phi(z) = (x^T z + c)^2$$

$c \in \mathbb{R}$

$$\text{if } K(x, z) = (x^T z + c)^d$$

$\Phi(x)$ has all features of ~~order~~ order $\leq d$.

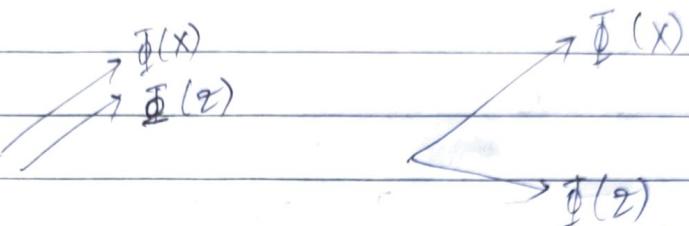
$d=5 \rightarrow x_1 x_2 x_5 x_6 x_7 \checkmark \quad \text{order} \leq 5$

$x_1 x_3 x_9 \checkmark \quad \text{order} \leq 5$

Support vector machine \Rightarrow optimal margin classifier + kernel trick

how to make kernels?

If x, z are "similar" $K(x, z) = \Phi^T(x)\Phi(z)$ is "large"
 " " " " dissimilar" $K(x, z)$ is "small"



$K(x, z) = \exp\left(-\frac{\|x-z\|^2}{2\sigma^2}\right)$ is a ^{kernel} func. that satisfies above both criterias
 gaussian kernel

We can use the above ^{kernel} function only when there exists Φ s.t.

$$K(x, z) = \Phi(x)^T \Phi(z)$$

suppose,

$\{x^{(1)}, \dots, x^{(d)}\}$ be d inputs.

$K \in \mathbb{R}^{d \times d} \rightarrow$ "kernel matrix"

$K_{ij} = K(x^{(i)}, x^{(j)}) \rightarrow$ kernel function applied to this

K should be "positive semi definite"

\hookrightarrow necessary and sufficient condition.

given any vector z ,

$$\begin{aligned} z^T K z &= \sum_i \sum_j z_i K_{ij} z_j = \sum_i \sum_j z_i \Phi(x^{(i)})^T \Phi(x^{(j)}) z_j \\ &= \sum_i \sum_j \sum_k z_i (\Phi(x^{(i)}))_k (\Phi(x^{(j)}))_k z_j \\ &= \sum_k \sum_i \sum_j z_i \Phi(x^{(i)})_k z_j \Phi(x^{(j)})_k \\ &= \sum_k \left(\sum_i z_i \Phi(x^{(i)})_k \right)^2 \geq 0 \end{aligned}$$

$\therefore K \geq 0 \rightarrow$ positive semi definite

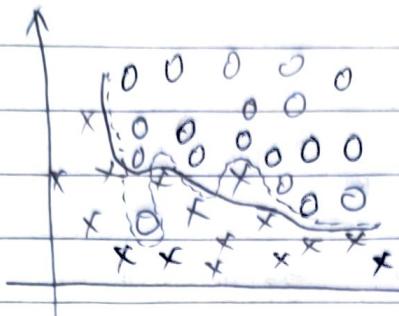
Mercer's Theorem: $\Rightarrow K$ is valid kernel func. i.e. $\exists \Phi$
s.t. $K(x, z) = \Phi(x)^T \Phi(z)$
and for any d inputs $\{x^{(1)}, \dots, x^{(d)}\}$
if and only if $K \geq 0$

Linear kernel $\rightarrow K(x, z) = x^T z$.
 $\Phi(x) = x$.

Gaussian kernel $\rightarrow \Phi(x) \in \mathbb{R}^\infty$ meaning can be of ∞ order



Sometimes we have noisy data that we don't have to try too hard to separate



— → our boundary/decision because data is noisy

--- → what

$$\min \frac{1}{2} \|w\|^2$$

↳ optimal margin classifier

We use "L₁ norm soft margin sum" for this

instead of $\min \frac{1}{2} \|w\|^2$

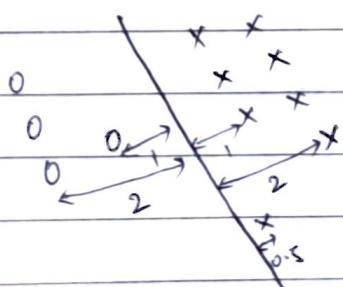
$$\text{s.t. } y^{(i)} (\underbrace{w^T x^{(i)} + b}_{\text{functional}}) \geq 1$$

for L₁ norm soft margin sum we take

$$\min \frac{1}{2} \|w\|^2 + C \sum_{i=1}^m \xi_i \quad \rightarrow \text{"psi"}$$

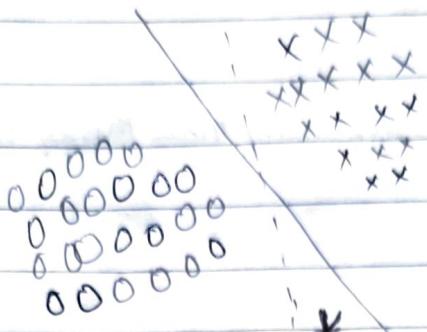
$$\text{s.t. } y^{(i)} (w^T x^{(i)} + b) \geq 1 - \xi_i \quad i=1, \dots, m$$

$$\xi_i \geq 0$$



according to L₁ norm

let's us get away with functional margins less than 1.



outlier caused a huge change in decision boundary

L_1 norm helps with this and helps keep decision boundary close to - instead of --

Our new eqn. for L_1 norm becomes

$$\max \sum_{i=1}^m \alpha_i - \sum_{i=1}^m \sum_{j=1}^m y^{(i)} y^{(j)} \alpha_i \alpha_j \langle x^{(i)}, x^{(j)} \rangle$$

$$\text{s.t. } \sum_{i=1}^m y^{(i)} \alpha_i = 0$$

$$\text{s.t. } \boxed{0 \leq \alpha_i \leq C} \quad \rightarrow \text{only change to prev eqn.}$$

Same eqn.

Change in condn.

protein sequence classifier using SVM

A, B, C, ..., Z \rightarrow amino acids

B A S T A I B A S T A U ...

$\Phi(x) = ?$

$x \mapsto y$

we make all sequences using 4 amino acids

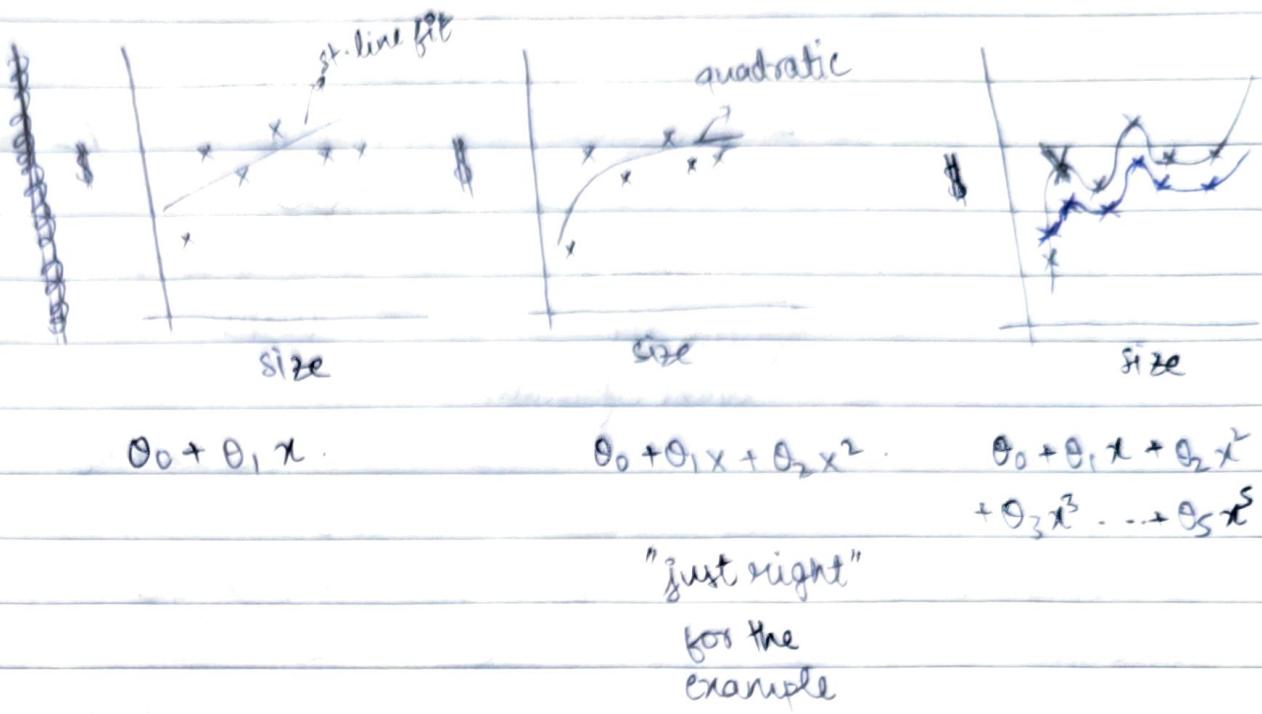
feature is how many times this seq occurs

AAAAA	0
AAAAB	0
AAAC	0
AAAZ	0
AABA	0
AABB	0
AAAB	0
ABAA	0
ABAB	0
ABBA	0
ABBB	0
ABAZ	0
ABBB	0
BAJT	2
TSIA	01
zzzz	0

BAJT, SIAI, BAJT, AU

$\Phi(X)^T \Phi(Z) = K(X, Z) \rightarrow$ using a dynamic programming algorithm

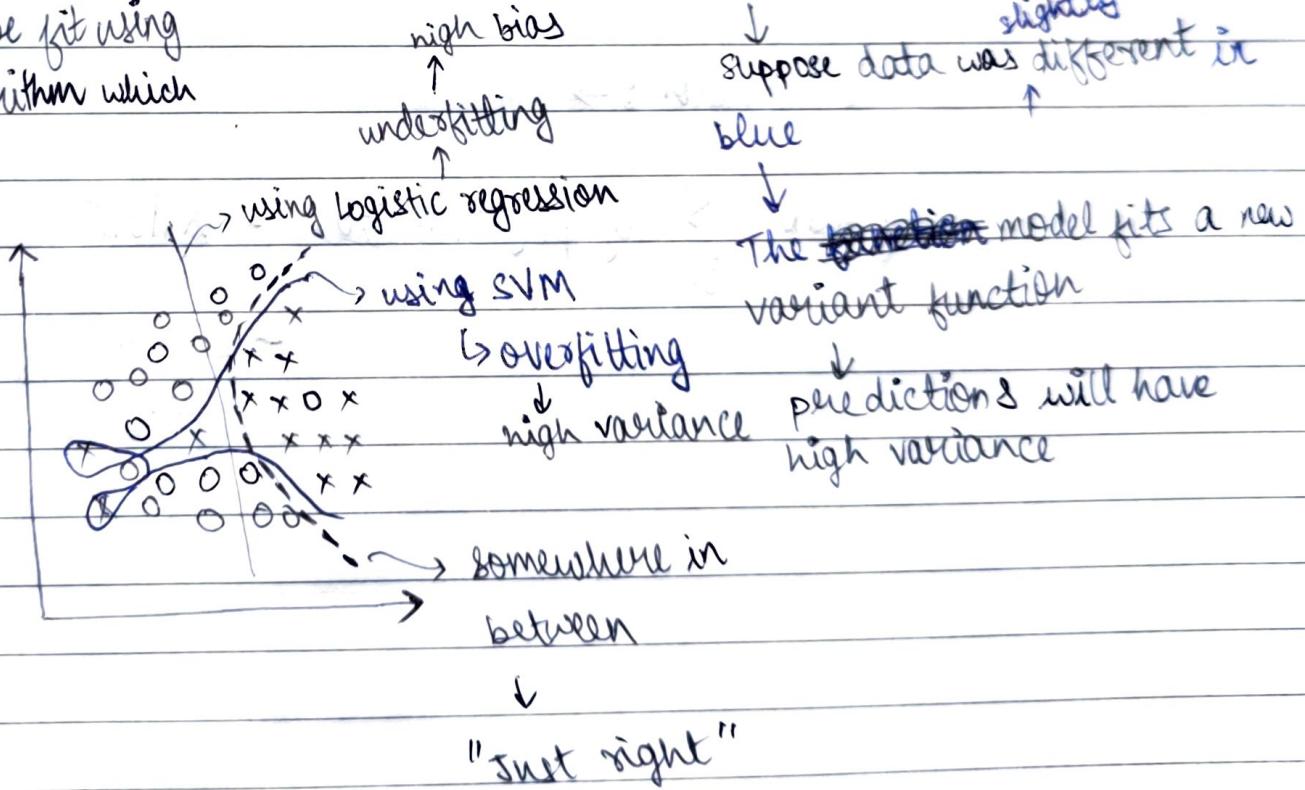
bias / variance in ML models



doesn't capture
the trend
"Underfits" the data
"High-bias"

This learning
algo had a strong
preconception that
data could be fit using
linear algorithm which
isn't true

"Overfits" the data
"High variance"
we get accurate
values for the 6 values



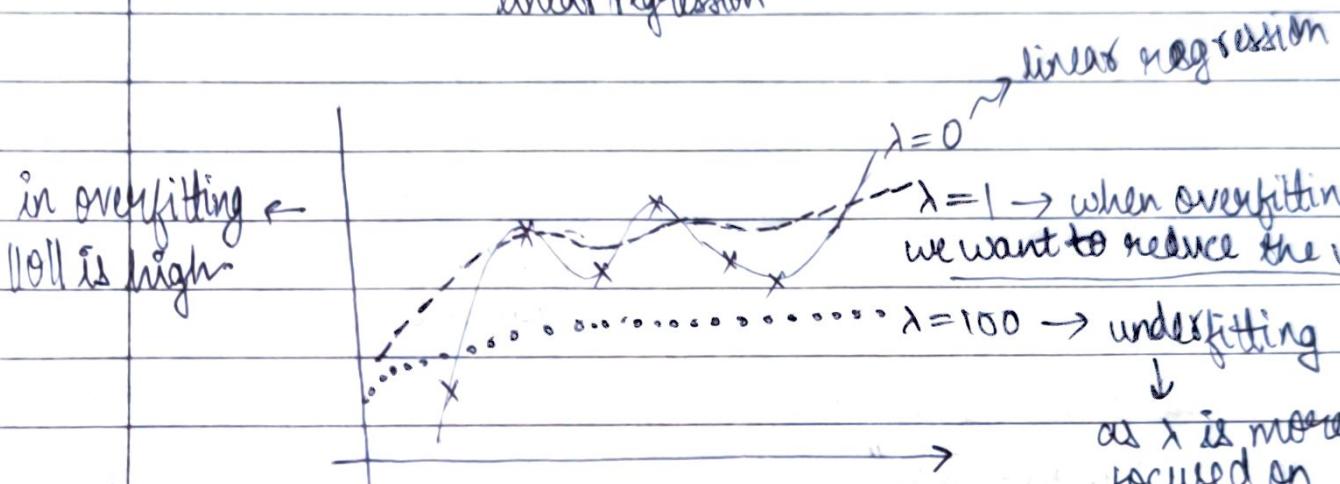
Regularisation

$$\min_{\theta} \left(\frac{1}{2} \sum_{i=1}^m \|y^{(i)} - \theta^T x^{(i)}\|^2 \right) + \left(\frac{\lambda}{2} \|\theta\|^2 \right)$$

optimisation
Objective for
linear regression

regularisation
term

' $\frac{1}{2}$ sometimes taken
to help with
derivation'



for logistic regression:-

$$\arg \max_{\theta} \sum_{i=1}^n \log p(y^{(i)} | x^{(i)}; \theta) - \lambda \|\theta\|^2$$

'-' sign as max instead of min

example: →

Text classification

$$m = 100$$

$$n = 10,000$$

} we generally use logistic regression when m is a similar order to n

$$X = \begin{bmatrix} 1 & \text{aardvark} \\ 0 & \vdots \\ \vdots & \vdots \\ 0 & \text{diddle} \\ \vdots & \vdots \\ - & n \times 1 \end{bmatrix}$$

we can use logistic regression for this classification problem

but it overfits the data //

However logistic regression with regularization gives a good fit for our data

another way to think of regularization

$$S = \{x^{(i)}, y^{(i)}\}_{i=1}^m \quad S \rightarrow \text{some training set}$$

we need to find θ given training set

$$p(\theta|S) = \frac{p(S|\theta) \cdot p(\theta)}{p(S)}$$

$$\arg \max_{\theta} p(\theta|S) = \arg \max_{\theta} p(S|\theta) \cdot p(\theta).$$

we get the regularization technique we discussed

$$= \arg \max_{\theta} \left(\prod_{i=1}^m p(y^{(i)} | x^{(i)}, \theta) \right) p(\theta).$$

→ GLM model

solving for this using complex methods

$$p(\theta) = \frac{1}{\sqrt{2\pi^T |Z^2 I|^{1/2}}} \exp\left(-\frac{1}{2}\theta^T (Z^2 I)\right)$$

$p(\theta)$: $\theta \sim N(0, Z^2 I) \rightarrow$ "assume" $p(\theta)$ is gaussian.

schools of statistics



Frequentist

Bayesian

we determine θ that makes our data as likely as possible

i.e. maximise $p(s|\theta)$.

"Maximum Likelihood Estimate"

θ is unknown

but we have a preconception about θ and how it is generated in the real world

This preconception is captured by a probability distribution

[ex]: →

$p(\theta)$

s.t. $\theta \sim N(0, \sigma^2 I)$

$\leftarrow p(\theta) \rightarrow$ prior

distribution

is a reasonable practice mostly.

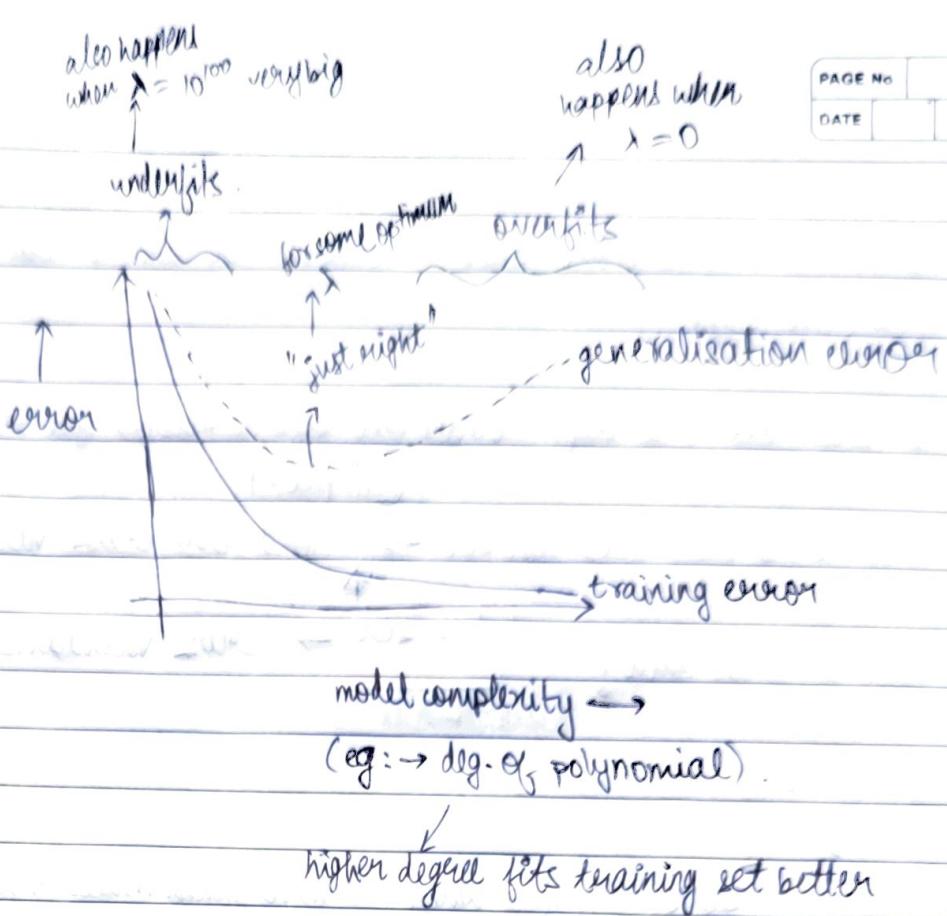
goal is to find theta that is ~~posterior~~ most likely upon seeing the data.

find θ s.t. $\max_{\theta} p(\theta|s)$

known as maximum a posteriori

(MAP)

estimation



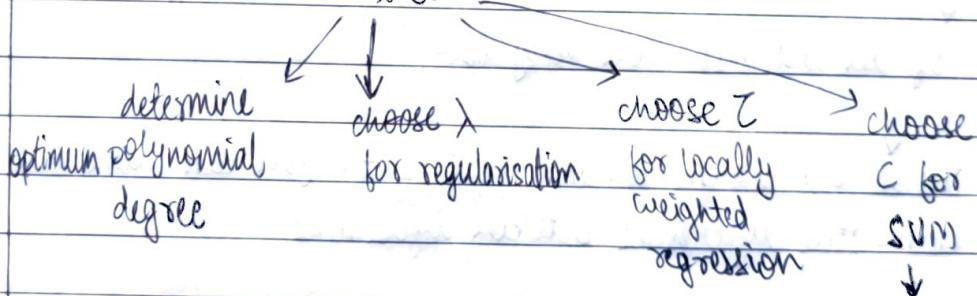
how to determine "just right" λ ?

→ take your dataset and split it into $\begin{cases} \text{train} \\ \text{dev} \\ \text{test} \end{cases}$ } 3 sets

eg: → 10,000 examples.

and we are trying to do model selection

i.e.



so we split S as

S

$\rightarrow S_{\text{train}}, S_{\text{dev}}, S_{\text{test}}$

$$\min \frac{1}{2} \|w\|^2 - C \sum |y_i - h(x_i)|$$

$S \rightarrow \text{dataset}$

Train each model on S_{train} i (options on what we want to optimise)

get some hypothesis $h_i \rightarrow$ hypothesis for i th model

degree of polynomial

→ measure error on Sdev. Pick model with lowest error on Sdev

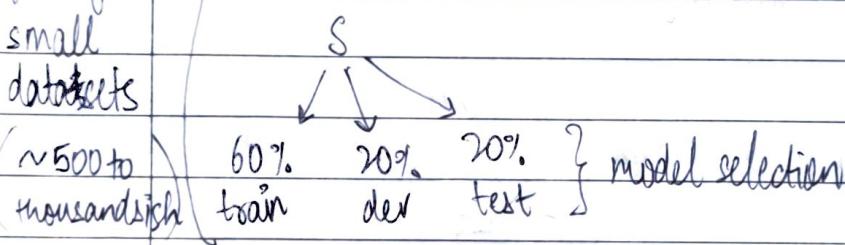
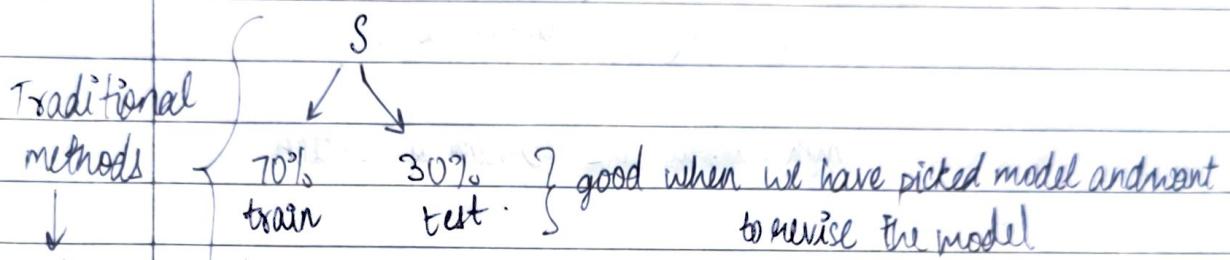
→ Evaluate the algorithm on set (Stest) for report.

(simple)

what we did with dev & test set → "Hold-out cross validation"

how to split data?

dev set → "cross validation set"



if S is very large ~ 10 million or so



% of dev and test gets very less



choose big enough dev and test s.t. you can make out meaningful differences between algorithms.

eg:

90% } can be
95% } distinguished
by say

1000 examples in dev and test

If this is meaningful
90-10% } → to our project we
90% } need to increase
cannot be distinguished
by
1000 examples in dev & test

~~Suppose we have a small dataset of 100 inputs~~

70 Train
30 Test

not a good way
of splitting as we lost
critical data by splitting
like this.

So we use the k-fold cross validation: →

we split the dataset into $k=5$ here in this example
 $x^{(1)}, y^{(1)}$ but $k=10$ is the norm

for each group we train on $k-1$ inputs and
test on the k th input then take the average the
 $x^{(100)}, y^{(100)}$ error estimates

For $i=1, \dots, k$
 \sum

Train (fit param) on $k-1$ pieces
test on remaining 1 piece }

Average

if $m=20 \rightarrow$ very small dataset

"leave one out cross validation"

we do our training on $m-1$ inputs and test on one
input and we do this m times leaving each input
exactly once and averaging over m .

"Feature selection" → find a small subset of features that are most
useful → "takes judgement"

feature selection \rightarrow start with $F = \emptyset$

repeat:

1) Try adding each feature i to F and see which single feature addition makes improves dev set performance

2) Add that feature to F

this procedure is known as
"forward search"

\downarrow
we start empty and keep adding features



Learning theory

Assumptions:

1) \exists data distribution D

$$(x, y) \sim D$$

i.e. all examples / inputs we get in our training set is from a data distribution.

2) Independent samples

$x^{(1)}$	$y^{(1)}$
$x^{(2)}$	$y^{(2)}$
$x^{(3)}$	$y^{(3)}$
:	
$x^{(m)}$	$y^{(m)}$

$$\sim D$$

learning algo

deterministic function

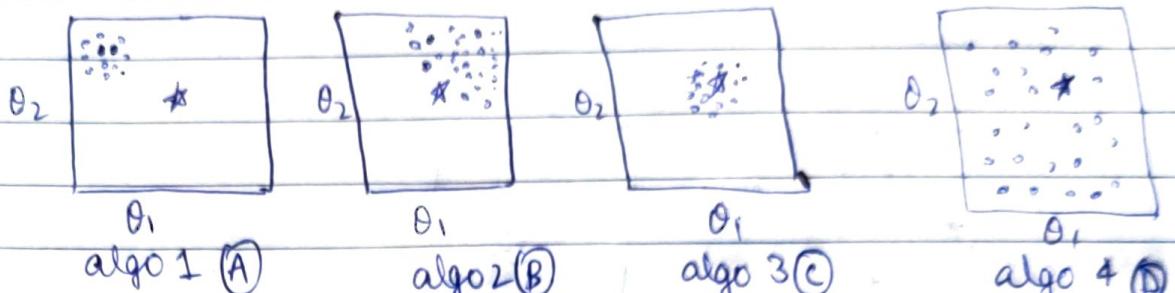
$S \rightarrow$ Random variable

$\hat{\theta} \rightarrow$ estimate
 $\theta^* \rightarrow$ true

in this model $\exists \theta^* \text{ on } h^*$ which is "true param"

\hat{h} or $\hat{\theta} \rightarrow$ "sampling distribution"
Random variable

let's say we have 4 different learning algs. with θ_1 , θ_2 as features
we make the following parameter view



$*$ → true θ

we run above algos for multiple samples to obtain different $\hat{\theta}$

(C) and (D) → low bias → "bias" means if it is centred around true parameter

(A) and (C) → low variance → "variance" measures how dispersed the parameters are

\$ as $m \rightarrow \infty \rightarrow \text{var}[\hat{\theta}] \rightarrow 0$ as for sufficiently large no. of inputs data will be more rate at which $\text{var}[\hat{\theta}]$ decreases clustered
ie called

"statistical efficiency"

so on.
less clustered
more clustered

then algo is

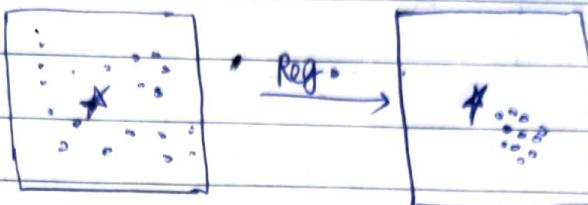
* if $\hat{\theta} \rightarrow \theta^*$ as $m \rightarrow \infty$ Consistent algorithm

* if $E[\hat{\theta}] = \theta^*$ regardless of $m \rightarrow$ the algorithm is called unbiased algo.

reducing variance

i) $m \rightarrow \infty$ (i.e get more data)

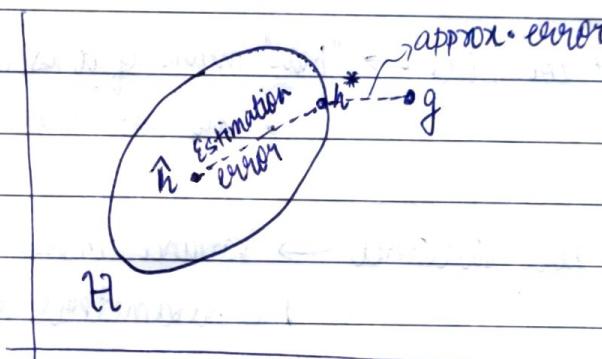
ii) regularisation



low bias

high variance

small bias \rightarrow slightly increased
low variance



space of hypothesis

H is the set of hypothesis

g^* is best possible hypothesis

h^* is best hypo in class H

\hat{h} is random variable learnt from finite data

$E(h)$: Risk / Generalisation error = $E_{(x,y) \sim D} [1\{h(x) \neq y\}]$

expected number

of "non-matches" from data } \rightarrow infinite

data dist.-wed
to create samples

\leftarrow distribution

steps

$\hat{\mathcal{E}}_S(h)$: Empirical risk

$$= \frac{1}{m} \sum_{i=1}^m \mathbb{1}_{\{h(x^{(i)}) \neq y^{(i)}\}}$$

$\mathcal{E}(g)$: Bayes Error / Irreducible error

$\mathcal{E}(h^*) - \mathcal{E}(g) = \text{Approximation Error} \rightarrow \text{i.e. price to pay for limiting to a class } \mathcal{H}$

$\mathcal{E}(\hat{h}) - \mathcal{E}(h^*) = \text{Estimation error}$

$\mathcal{E}(\hat{h}) = \text{Estimation error} + \text{approx. error} + \text{Irreducible error}$

because

we are working
on limited data

depends on

choice
class

cannot reduce
no matter
what

Estimation error

Estimation
~~error~~ +
Var
err

Estimation + Approx. error + Irreducible

$\mathcal{E}(\hat{h}) = \text{Var} +$

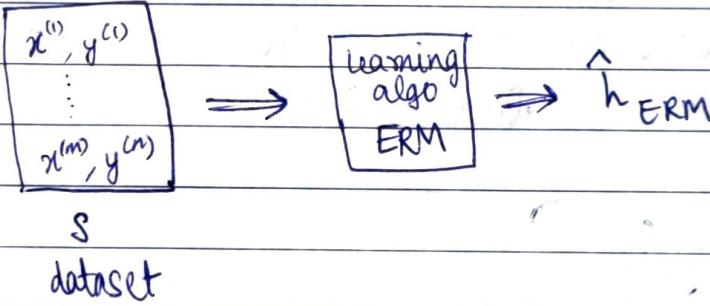
bias + Irreducible

Reducing high bias

① make H bigger



Empirical risk minimisation (ERM) → is a learning algo



$$\hat{h}_{ERM} = \arg \min_{h \in H} \frac{1}{m} \sum_{i=1}^m \mathbb{1}\{h(x^{(i)}) \neq y^{(i)}\}$$

Uniform convergence

random variable empirical error

① consider a hypothesis with training error, $\hat{\epsilon}(h)$ then what can we say about its generalisation error $\epsilon(h)$? $\hat{\epsilon}(h) \text{ vs } \epsilon(h)$

② $\epsilon(\hat{h})$ vs $\epsilon(h^*)$ → how does generalisation error of our "learnt hypothesis" compare to best possible generalisation error of our class.

Tools to solve these questions:

- 1) Union bound (we have k events)
 A_1, A_2, \dots, A_k (need not be independent)

$$P(A_1 \cup A_2 \cup \dots \cup A_k) \leq P(A_1) + P(A_2) + \dots + P(A_k).$$

2) Hoeffding's Inequality

Let $z_1, z_2, \dots, z_m \sim \text{Bern}(\bar{\Phi})$ we have parameters from bernoulli dist.

$$\hat{\Phi} = \frac{1}{m} \sum_{i=1}^m z_i$$

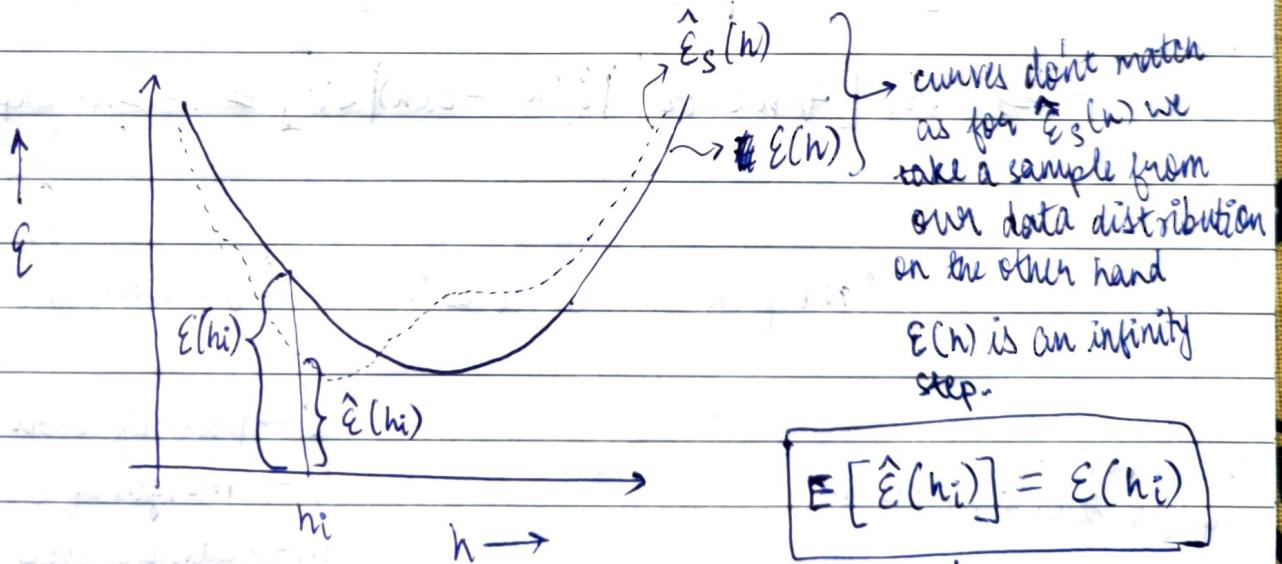
$\hat{\Phi}$ is average of these.
↳ "learnt hypothesis"

let $r > 0$ [margin]

$$Pr[\underbrace{|\hat{\Phi} - \bar{\Phi}|}_{\text{error}} > r] \leq 2 \exp(-2r^2 m)$$

↳ prob. that $|\hat{\Phi} - \bar{\Phi}| > r$ is RHS
proof in lecture notes

↳ prob. that $|\hat{\Phi} - \bar{\Phi}| > r$ is RHS



$$E[\hat{\epsilon}(h_i)] = \epsilon(h_i)$$

as $\hat{\epsilon}$ is over some m random inputs and when we keep doing this throughout our dataset, expected value is $\epsilon(h_i)$

on applying Hoeffding's Ineq.

$$\Pr \left[|\hat{\varepsilon}(h_i) - \varepsilon(h_i)| > r \right] \leq 2 \exp(-2r^2 m)$$

as we increase m , the dotted curve ($\hat{\varepsilon}(h_i)$) gets closer to solid curve ($\varepsilon(h_i)$).

we got above bound for a fixed \mathcal{H}
we want bound for all \mathcal{H}

Finite hypothesis class \mathcal{H} :

$$|\mathcal{H}| = K \rightarrow \text{cardinality}$$

$$\Pr \left[\exists h \in \mathcal{H} \mid \hat{\varepsilon}_s(h) - \varepsilon(h) > r \right]$$

prob that atleast

one hypothesis from class \mathcal{H} $\leftarrow \Pr \left[\exists h \in \mathcal{H} \mid |\hat{\varepsilon}_s(h) - \varepsilon(h)| > r \right] \leq K \cdot 2 \exp(-2r^2 m)$

$$\Rightarrow \Pr \left[\forall h \in \mathcal{H} \mid |\hat{\varepsilon}_s(h) - \varepsilon(h)| < r \right] \geq 1 - \underbrace{K \cdot 2 \exp(-2r^2 m)}_{\delta}$$

"proof in lecture notes".

$$\boxed{\delta = 2K \cdot \exp(-2r^2 m)}$$

$\delta \rightarrow$ prob. of error

$r \rightarrow$ margin of error

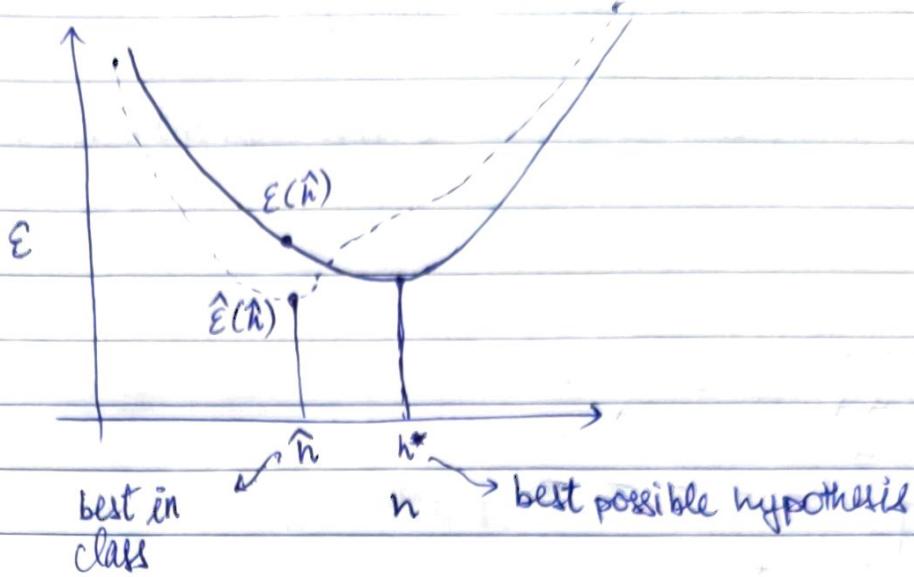
$m \rightarrow$ sample size

\uparrow fix $r, \delta > 0$

"sample complexity result"

$$\boxed{m \geq \frac{1}{2r^2} \log \left(\frac{2K}{\delta} \right)}$$

margin of error between empirical and generalisation error will be less than r if $m \geq$ RHS



relation b/w $\hat{E}(\hat{h})$ and $E(h^*)$

$$E(\hat{h}) \leq \hat{E}(\hat{h}) + \gamma$$

$$E(\hat{h}) \leq \hat{E}(h^*) + \gamma \quad \text{as } \hat{E}(\hat{h}) \text{ is minima of } \hat{E} \text{ func.}$$

$$E(\hat{h}) \leq E(h^*) + 2\gamma \quad \text{as } |\hat{E}(h) - E(h)| \leq \gamma \text{ is true throughout } h$$

with probability $1 - \delta$ and training size m

$$E(\hat{h}) \leq E(h^*) + 2 \sqrt{\frac{1}{2m} + \log \frac{2K}{\delta}}$$

VC dimension \rightarrow tries to assign size to ∞ size class

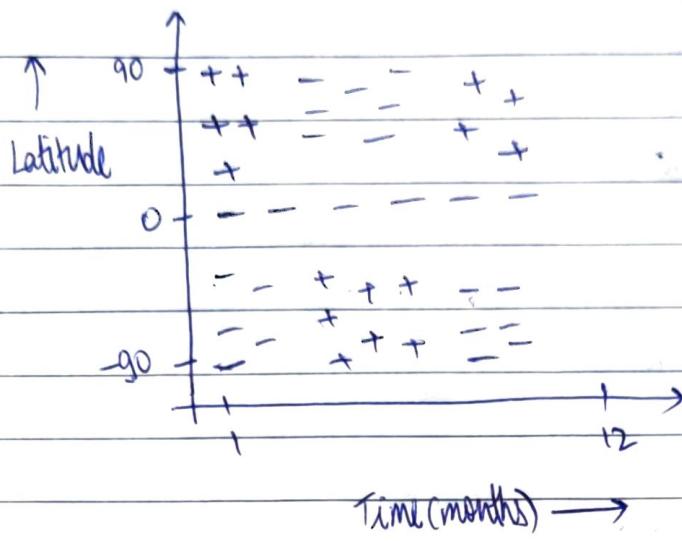
$VC(\mathcal{H}) \rightarrow$ ~~means~~ can be taken as size of ∞ size class

$$E(\hat{h}) \leq E(h^*) + O\left(\sqrt{\frac{VC(\mathcal{H}) \log\left(\frac{m}{VC(\mathcal{H})}\right)}{m}} + \frac{1}{m} \log\left(\frac{1}{\delta}\right)\right)$$

\hookrightarrow for ∞ ~~size~~ size class \mathcal{H} .

Decision trees \rightarrow non linear algo

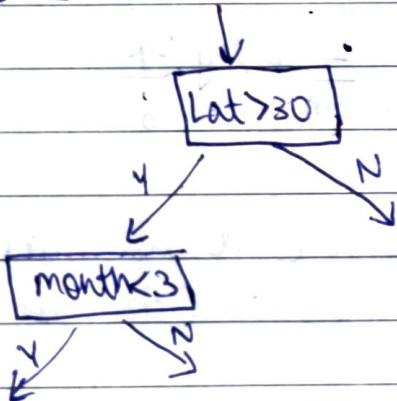
suppose we like skiing, and want to know when and where is it possible



We might be able to use SVM to split this data but decision trees give us a more natural way of handling this problem.

Decision trees \rightarrow Greedy, topdown, recursive partitioning

Basically we ask certain questions to separate out the data into sets. Then we repeat the steps for the sets obtained.



more formally

given a region R_p
we are trying to find a
split function S_p s.t.
 j^{th} feature of $X \sim R_1$

$$S_p(j, t) = \begin{cases} \{x | x_j < t, x \in R_p\}, \\ \{x | x_j \geq t, x \in R_p\} \end{cases}$$

j is feature number threshold R_1
 R_2

How to choose these splits?

→ Define loss on a region $\rightarrow L(R)$

→ Obv example \rightarrow miss classification loss

assume given C classes, define \hat{P}_c to be proportion of examples in R that are of class c

$$L_{\text{missclass}} = L(R) = 1 - \max_c \hat{P}_c \rightarrow \text{reasoning:} \rightarrow$$

we have C classes in our region R

we want to pick a split that minimises the loss as much.

we find the class which has the maximum proportion and we say the rest are missclassified.

refers
notes
as to
why?

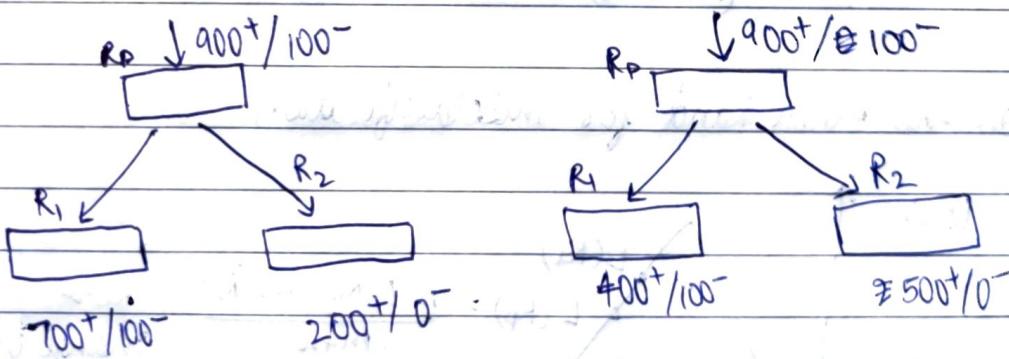
$$\leftarrow \max_{\text{first}} \text{loss}_{\text{parent}} = L(R_p) - (L(R_1) + L(R_2))$$

in
parent
children
loss

we pick the parameters

to min loss

missclassification loss has issues (here we are talking about absolute numbers rather than fraction)



$$L(R_1) + L(R_2) = 100 + 0 = 100$$

$$L(R_1) + L(R_2) = 100 + 0 = 100$$

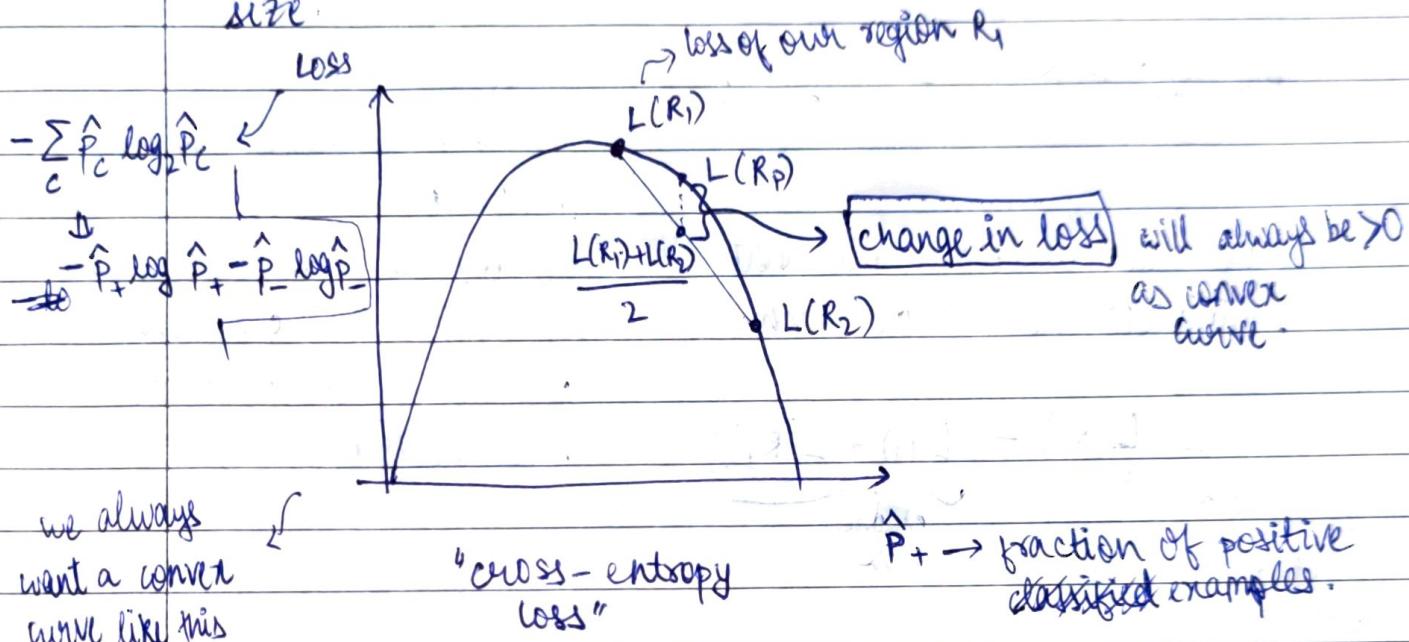
$$L(R_p) = 100$$

both have same loss but right one is better \rightarrow missclass-loss isn't sensitive enough //

we define a cross entropy loss

$$L_{\text{cross}} = - \sum_c \hat{P}_c \log_2 \hat{P}_c$$

suppose we split a region into two classes of equal size



we always want a convex curve like this

"cross-entropy loss"

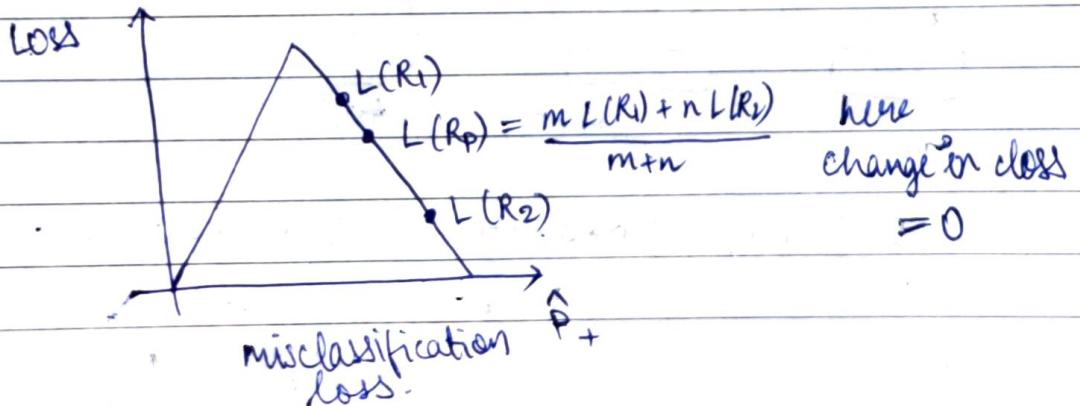
\hat{P}_+ → fraction of positive classified examples.

suppose the classed went of same size → then we'd opt
suppose R_1 has size m and R_2 has size n

$$\frac{m L(R_1) + n L(R_2)}{(m+n)}$$

would be our new point

on the other hand for misclassification

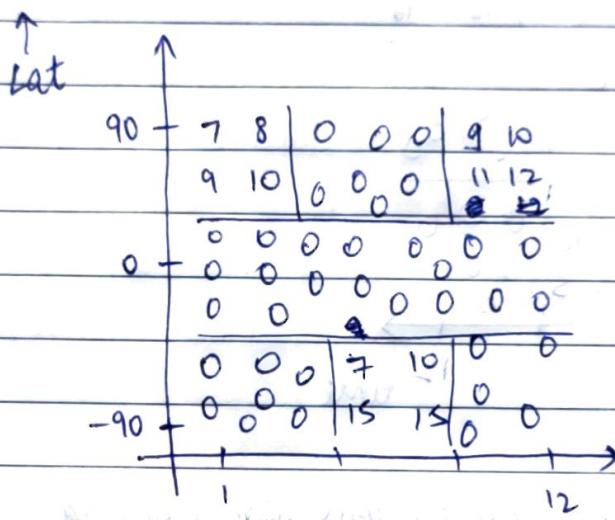


$$\text{Gini-loss} = - \sum_c \hat{P}_c (1 - \hat{P}_c) \rightarrow \text{also has a curve similar to cross-entropy loss.}$$

Regression tree

Instead of classification we try regression

e.g.: → we predict amount of snowfall rather than whether we can ski or not.



Time (months) →

84

$$\text{polardist} \quad \hat{y}_m = \sum_{i \in R_m} (y_i - \hat{y}_m)^2$$

we have a region R_m for which we want to predict \hat{y}_m

given R_M ,

$$\hat{y}_m = \frac{\sum_{i \in R_m} y_i}{|R_m|} \rightarrow \text{summing values in region and take the average}$$

The loss we take in this case is

$$L_{\text{squared}} = \frac{\sum_{i \in R_m} (y_i - \hat{y}_m)^2}{|R_m|}$$

Categorical variables

if we have given example as

North	8 6 0 0 12
	7 14 0 0 14 13 16
	0 0 0 0 0 0
Equator	0 0 0 0 0 0
	0 0 9 10 0 0 0
	0 0 15 15 6 0 0
South	0 12 time

for these kinds of problem we don't have threshold

\downarrow
 lat $\in \mathbb{N}^3$ we ask questions is latitude is northern hemisphere

if we have q categories we deal with 2^q possible split.

If we keep running a decision tree \rightarrow we end up with each point having its own class \rightarrow overfitting
 so, we need to reduce the variance for this decision by regularization.

Regularization of decision Tree

heuristic meaning \rightarrow proceeding to solution by trial and error.

We regularize decision trees by using a number of heuristics

- i) min leaf size
- ii) max depth
- iii) max number of nodes
- iv) minimum decrease in loss \rightarrow generally not a good idea

before split $\rightarrow L(R_p)$

after split $\rightarrow L(R_1) + L(R_2)$

but sometimes we need to ask multiple questions to get to a good split

we do not pick a split if decrease is big

- v) grow full tree then ~~prune out nodes~~

pruning done using validation set, evaluate ~~using~~ misclassification error on validation set for each example/leaf we remove.

Runtime

n examples

f features

d depth

Test time: $O(d)$ typically $d < \log_2 n$

basically has $O(d)$ parents

Training time: Each point is a parent of $O(d)$ nodes.

A ~~feature~~ point has a cost $O(f)$ at each node.

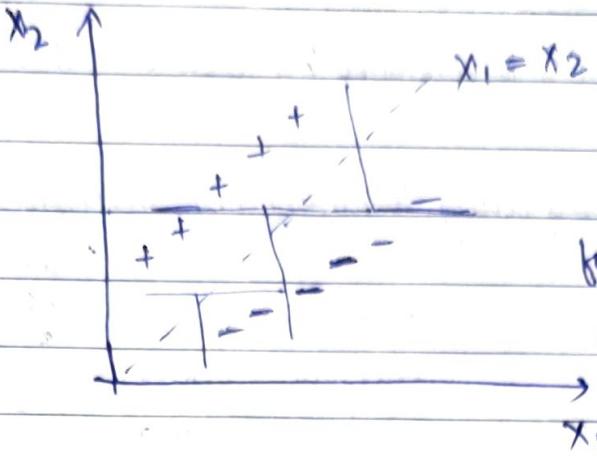
\hookrightarrow consider f splits

we have n points

after a split
a point belongs to the left or the right of it

Total cost = $O(nfd)$

No additive structure



easy for logistic regression

for decision trees we ~~need~~ need
a lot of splits to approximate
the ~~one~~ split.

Recap

↓
Advantage

- Easy to explain
- can see exactly what is happening
- Interpretable
- can deal with cat. vars
- Categorical vars
- Fast

↓
disadvantage

- High variance
- Bad at additive
- low predictive accuracy.

Ensembling

Take X_i 's which are random variables (RV)

that are independently identically distributed (iid)

$$\text{Var} = \left(\frac{1}{n}\right)^2 \text{Var}(\sum_i X_i) = \left(\frac{1}{n}\right) (\sigma^2)$$

as they are independent.

$$\text{Var}(X_i) = \sigma^2$$

$$\text{Var}(\bar{X}) = \text{Var}\left(\frac{1}{n} \sum_i X_i\right) = \frac{\sigma^2}{n}$$

↓

variance of one of variables is $\sigma^2 \Rightarrow$ var. of mean of many variables = $\frac{\sigma^2}{n}$

variance of all variables will be σ^2 and they will have same mean $\frac{\sigma^2}{n}$
as they are identically distributed.

Drop the independence assumption.

now x_i 's are identically distributed

x_i 's correlated by ρ

$$\text{var}(\bar{x}) = \rho \sigma^2 + \frac{(1-\rho)}{n} \sigma^2$$

Ways to ensemble

- 1) use different algorithms
- 2) different training sets
- 3) Bagging (Random forests)
- 4) Boosting (Adaboost, XG boost)

} commonly used

Bagging \rightarrow Bootstrap aggregation.

Have a true population P

Training set $S \sim P$ (S is sampled from P)

for bootstrapping do: Assume $P = S$ (population is training set)

Bootstrap samples $Z \sim S$

we can separate several
train our models on bootstrap samples and then average
them.

bootstrap samples z_1, \dots, z_M

Train model G_m on z_1, z_2, \dots, z_M

$$G(\bar{x}) = \frac{1}{M} \sum_{m=1}^M G_m(x)$$

bias variance analysis on above result

$$\text{Var}(\bar{x}) = p\sigma^2 + \frac{(1-p)\sigma^2}{M}$$

bootstrapping is reducing down p

more M is \rightarrow less variance as it reduces p

the variance is reduced but the bias is increased.

as bootstrap samples have lesser data and the models are thereby less complex. (Random subsampling)

Decision trees & bagging :-

D.T are high variance and low bias

→ Ideal fit for bagging.

Random forests

At each split \rightarrow consider only a fraction of your total features.

this decreases p .

will cause all models to be correlated at first split → suppose we have a very good predictor which gives very good performance on its own.

helps decorrelate models; regardless of bootstrap sample we use model with same predictor as first split

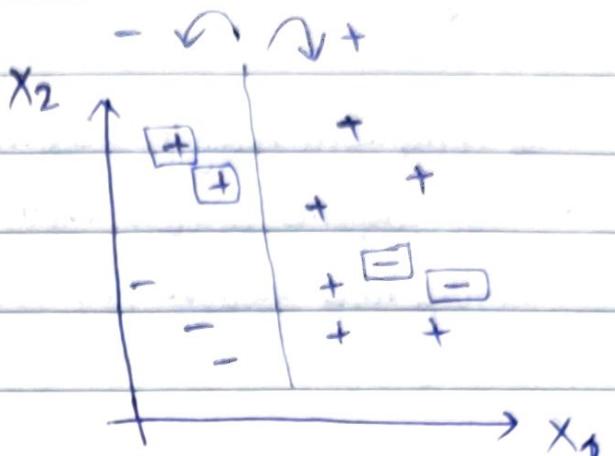
Boosting

Decrease bias

more Additive

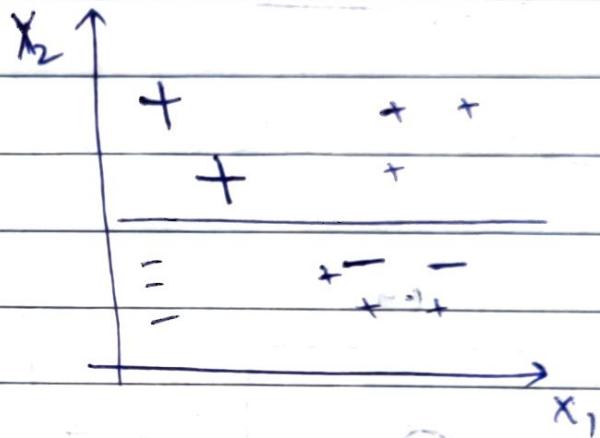
loop

First we ask a question
and come up with a decision
boundary



we highlight the mistakes and
increase their weight.

next model trained on
the dataset with previous
mistakes having high weight.



Determine for classifier G_m a weight α_m

for adaboost

$$\alpha_m = \log \left(\frac{1 - \epsilon_m}{\epsilon_m} \right)$$

$$G(X) = \sum_m \alpha_m G_m$$

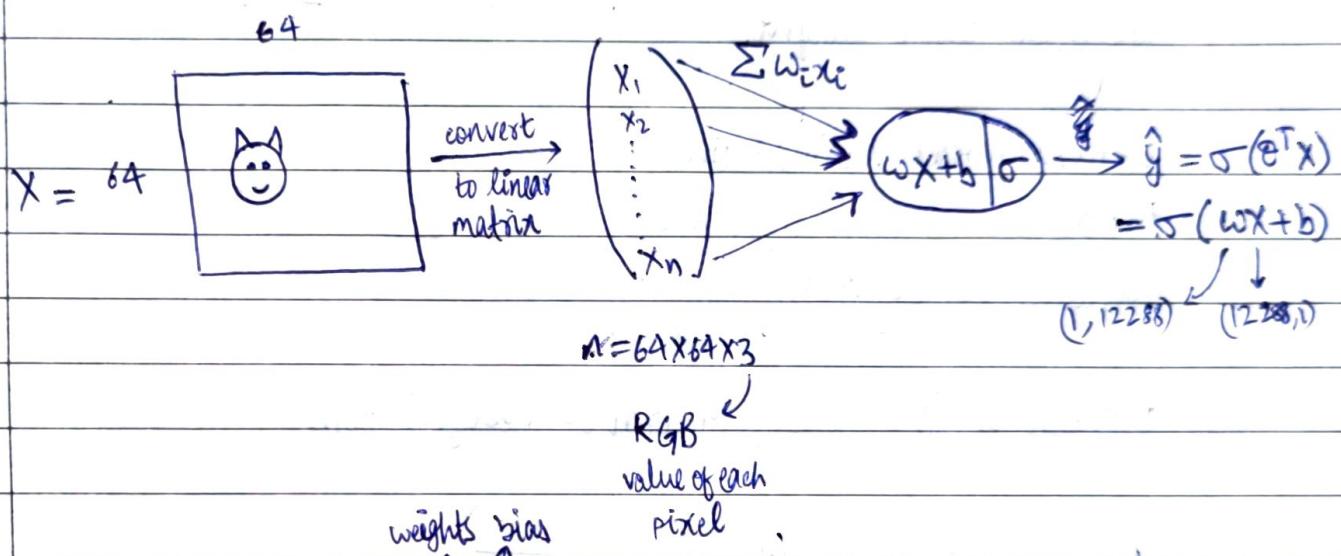
Each G_m trained on reweighted training set.

Deep learning :-

- is computationally expensive
- algo performs better with more data
- algorithms which help reduce the computational power needed

Logistic regression

goal → Find cats in images $\left\{ \begin{array}{l} 1 \rightarrow \text{presence} \\ 0 \rightarrow \text{absence} \end{array} \right.$



i) Initialise w, b

minimise the loss function

$$L = -[\hat{y} \log \hat{y} + (1 - \hat{y}) \log (1 - \hat{y})]$$

ii) Find optimal w, b

$$\left\{ \begin{array}{l} w = w - \alpha \frac{\partial L}{\partial w} \\ b = b - \alpha \frac{\partial L}{\partial b} \end{array} \right.$$

iii) use $\hat{y} = \sigma(w+b)$ to predict

neuron = linear + activation

$y = \sigma(\sum w_i x_i)$ in this case sigmoid function

model = architecture + parameters

the neurons
we have
 w, b here

goal 2 → find cat / lion / iguana in images → we use 3 neurons here

will the loss
function given
earlier work
for this goal?

NO

as prev 1 has

only two
outcomes
and can't
match the
labeling

dataset will
look like

$$\begin{array}{|c|} \hline \text{cat} \\ \hline \end{array} \Rightarrow \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}$$

cat
lion
iguana

we train the network as earlier
and get all our weights and biases

suppose we had an input like

$$\begin{array}{|c|} \hline \text{cat} \\ \hline \text{lion} \\ \hline \end{array} \Rightarrow \begin{pmatrix} 1 \\ 1 \\ 0 \end{pmatrix}$$

will our neural network be robust for inputs as above?

represent layers → where neurons don't communicate with each other

$$y_1 = a_1^{[1]} = \sigma(w_1^{[1]} x + b_1^{[1]})$$
$$z_1^{[1]}$$

$$y_2 = a_2^{[1]} = \sigma(w_2^{[1]} x + b_2^{[1]})$$
$$z_2^{[1]}$$

$$y_3 = a_3^{[1]} = \sigma(w_3^{[1]} x + b_3^{[1]})$$
$$z_3^{[1]}$$

here output is a 3×1 vector. → gives exact same derivative

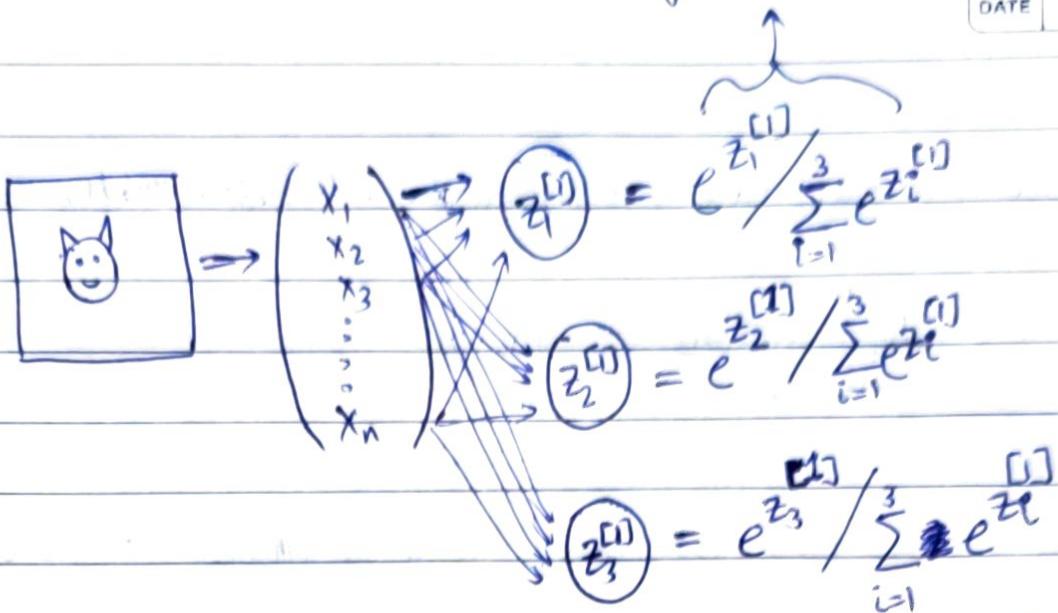
$$L_{3N} = -\sum_{k=1}^3 [y_k \log \hat{y}_k + (-y_k) \log (1 - \hat{y}_k)]$$

yes, as the neurons don't communicate amongst themselves.

goal 3 → + constraint
unique animal on an image
(atmost one animal)

softmax formula

PAGE NO.	
DATE	



now these prob. are dependent
on each other as they sum upto one.

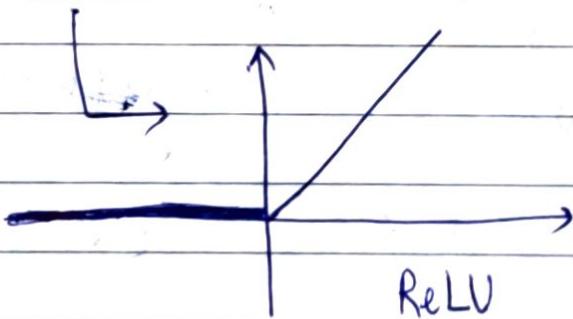
This model is called a softmax multiclass network.

Loss function here is cross entropy loss

$$L_{CE} = - \sum_{k=1}^3 y_k \log \hat{y}_k$$

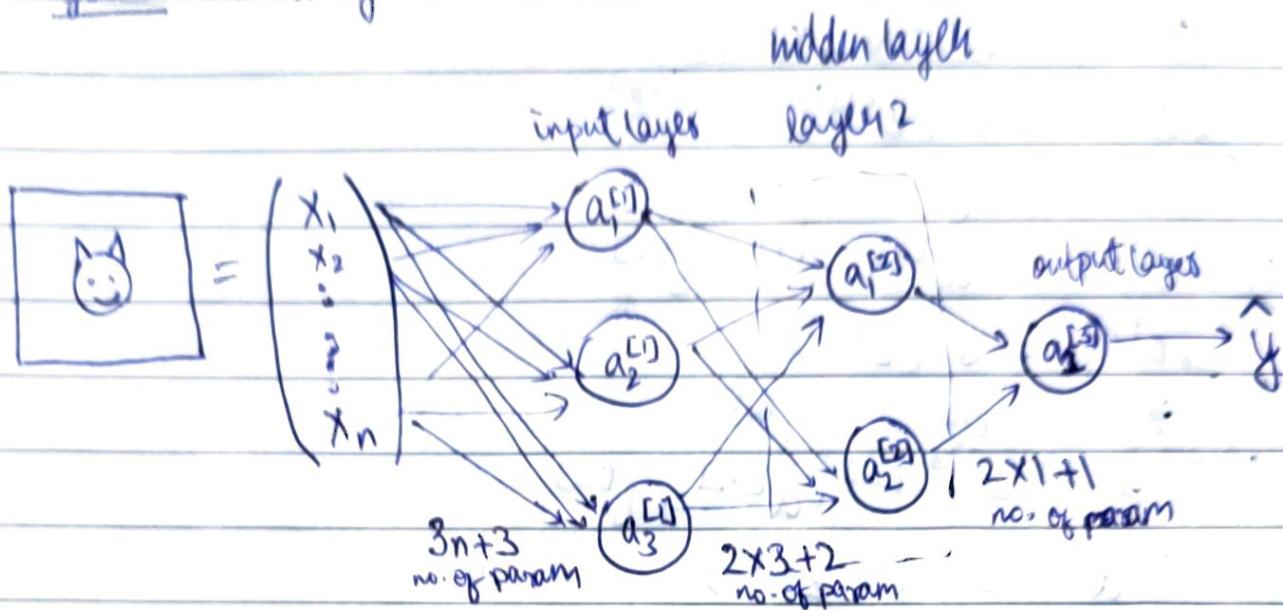
ReLU

ReLU function \rightarrow Rectified Linear Unit



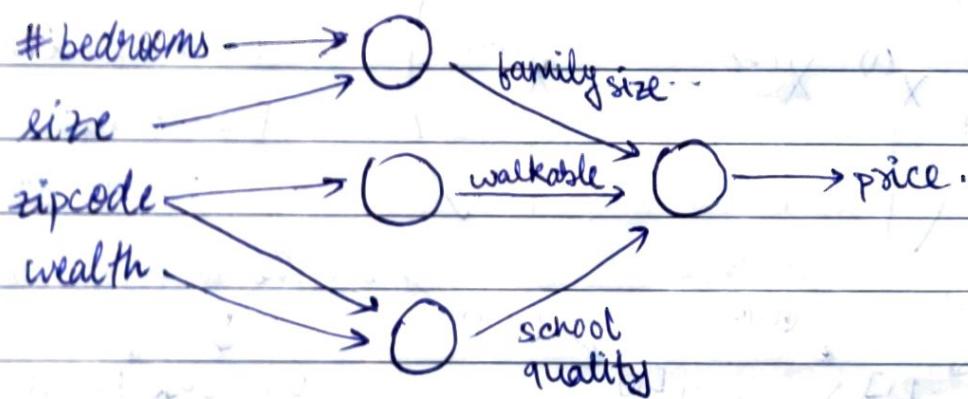
Neural Networks → also called end to end learning
also known as blackbox models

goal → image has cat or no cat.



last layer should have as many neurons as the no. of classifications

house price prediction



[how humans would construct a neural network]

"we let the next layer figure out which neuron is an important feature"

In practice we fully connect two layers

propagation equation

$$\begin{aligned}
 & (3,1) \rightarrow z^{[1]} = w^{[1]} X + b^{[1]} \quad (n,1) \\
 & (3,1) \rightarrow a^{[1]} = \sigma(z^{[1]}) \quad (3,1) \\
 & (2,1) \rightarrow z^{[2]} = w^{[2]} \cdot a^{[1]} + b^{[2]} \quad (3,1) \\
 & (2,1) \rightarrow a^{[2]} = \sigma(z^{[2]}) \quad (2,1) \\
 & (1,1) \rightarrow z^{[3]} = w^{[3]} \cdot a^{[2]} + b^{[3]} \quad (2,1) \\
 & (1,1) \rightarrow a^{[3]} = \sigma(z^{[3]}) \quad (1,1)
 \end{aligned}$$

What happens for an input batch of m examples

$$X = \begin{pmatrix} | & | & | \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ | & | & | \end{pmatrix}$$

$m \rightarrow$ inputs

(\cdot) \rightarrow represent id of input

$[] \rightarrow$ represents layer number

$$z^{[1]} = w^{[1]} X + b^{[1]} \quad (3,1) \quad \text{but the dimensions don't add up}$$

\downarrow we use "BROADCASTING"

$$(3,m) = \begin{bmatrix} | & | & | \\ z^{[1]}(1) & \dots & z^{[1]}(m) \\ | & | & | \end{bmatrix}$$

$$\tilde{b}^{[1]} = \begin{pmatrix} | & | & | \\ b^{[1]} & b^{[1]} & \dots & b^{[1]} \\ | & | & | \end{pmatrix} \quad m \text{ times}$$

$$(3,m)$$

optimising $w^{[1]}, w^{[2]}, w^{[3]}, b^{[1]}, b^{[2]}, b^{[3]}$

define loss / cost function

example ↴ multiple examples in the batch.

$$J(\hat{y}, y) = \frac{1}{m} \sum_{i=1}^m L^{(i)} \quad \hat{y} = a^{[3]}$$

$$\text{with } L^{(i)} = -[y^{(i)} \log \hat{y}^{(i)} + (1-y^{(i)}) \log (1-\hat{y}^{(i)})]$$

Backward propagation

$$\forall l = 1, \dots, 3 \quad \begin{cases} w^{[l]} = w^{[l]} - \alpha \frac{\partial J}{\partial w^{[l]}} \\ b^{[l]} = b^{[l]} - \alpha \frac{\partial J}{\partial b^{[l]}} \end{cases}$$

we start this process from $w^{[3]}$

↳ as it has a direct impact on the output

↓ easier to see how $w^{[3]}$ changes impact output than how $w^{[1]}$ changes affect output

$$\frac{\partial J}{\partial w^{[3]}} = \underbrace{\frac{\partial J}{\partial a^{[3]}} \cdot \frac{\partial a^{[3]}}{\partial z^{[3]}} \cdot \frac{\partial z^{[3]}}{\partial w^{[3]}}}_{\text{Chain rule}}$$

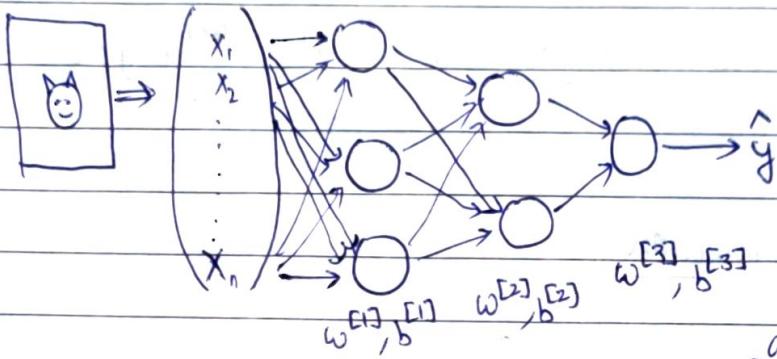
$$\frac{\partial J}{\partial w^{[2]}} = \underbrace{\frac{\partial J}{\partial z^{[3]}}}_{\text{Circular}} \cdot \frac{\partial z^{[3]}}{\partial a^{[2]}} \cdot \frac{\partial a^{[2]}}{\partial z^{[2]}} \cdot \frac{\partial z^{[2]}}{\partial w^{[2]}}$$

$$\frac{\partial J}{\partial w^{[1]}} = \underbrace{\frac{\partial J}{\partial z^{[2]}}}_{\text{Circular}} \cdot \frac{\partial z^{[2]}}{\partial a^{[1]}} \cdot \frac{\partial a^{[1]}}{\partial z^{[1]}} \cdot \frac{\partial z^{[1]}}{\partial w^{[1]}}$$

cost function : $J(\hat{y}, y) = \frac{1}{m} \sum_{i=1}^m L^{(i)}(\hat{y}, y)$

for ith data
with $L^{(i)} = -[y^{(i)} \log(\hat{y}^{(i)}) + (1-y^{(i)}) \log(1-\hat{y}^{(i)})]$

assume our neural network



and update each as $w^{[l]}$ as

we need to optimise these w s and b s $w^{[l]} := w^{[l]} - \alpha \frac{\partial J}{\partial w^{[l]}}$

we start by optimising outermost param first.

↳ as outermost have direct effect on cost.

we need $\frac{\partial J}{\partial w^{[3]}}$

we find $\frac{\partial L}{\partial w^{[3]}}$ and we can add this for m examples

$$\frac{\partial L}{\partial w^{[3]}} = - \left[y^{(i)} \frac{\partial}{\partial w^{[3]}} (\log(\sigma(w^{[3]} a^{[2]} + b^{[3]}))) \right]$$

here $a^{[2]} = (w^{[3]} \cdot a^{[2]} + b^{[3]})$ $+ (1-y^{(i)}) \frac{\partial}{\partial w^{[3]}} (\log(1-\sigma(w^{[3]} a^{[2]} + b^{[3]})))$

$$\frac{\partial L}{\partial w^{[3]}} = - \left[y^{(i)} \frac{1}{a^{[3]}} \cdot a^{[3]} \cdot (1-a^{[3]}) \cdot (a^{[2]})^T + (1-y^{(i)}) \frac{1 \cdot a^{[2]} \cdot (1-a^{[3]}) \cdot a^{[2]}}{(1-a^{[3]})} \right]$$

$$\frac{\partial L}{\partial w^{[3]}} = - \left[y^{(i)} (1-a^{[2]}) a^{[2]}^T + (1-y^{(i)}) a^{[3]} a^{[2]}^T \right]$$

$$= - \left[y^{(i)} a^{[2]}^T - a^{[3]} a^{[2]}^T \right] = \boxed{-(y^{(i)} - a^{[3]}) a^{[2]}^T}$$

$$\boxed{\frac{\partial J}{\partial w^{[3]}} = -\frac{1}{m} \sum_{i=1}^m (y^{(i)} - a^{[3]}) a^{[2]}^T}$$

we update
 $w^{[3]}$ like this

$$\frac{\partial L}{\partial w^{[2]}} = \underbrace{\frac{\partial L}{\partial a^{[3]}}}_{\text{we did this}} \cdot \underbrace{\frac{\partial a^{[3]}}{\partial z^{[3]}}}_{\text{as we want the error to backpropagate}} \cdot \underbrace{\frac{\partial z^{[3]}}{\partial a^{[2]}}}_{\text{we did this}} \cdot \underbrace{\frac{\partial a^{[2]}}{\partial z^{[2]}}}_{\text{we did this}} \cdot \underbrace{\frac{\partial z^{[2]}}{\partial w^{[2]}}}_{\text{we did this}}$$

$$(a^{[3]} - y) \cdot a^{[2]}^T \cdot a^{[2]} (1-a^{[2]}) \cdot w^{[3]}^T$$

$$\frac{\partial L}{\partial w^{[2]}} = \underbrace{(a^{[3]} - y)}_{(1,1)} \cdot \underbrace{w^{[3]}^T}_{(2,1)} \cdot \underbrace{a^{[2]}^T}_{(2,1)} \cdot \underbrace{(1-a^{[2]})}_{(2,1)} \cdot \underbrace{a^{[1]}^T}_{(1,3)}$$

$$\frac{\partial L}{\partial w^{[2]}} = \underbrace{w^{[3]}^T}_{(2,1)} \cdot \underbrace{a^{[2]}^T}_{(2,1)} \cdot \underbrace{(1-a^{[2]})}_{(2,1)} \cdot \underbrace{(a^{[3]} - y)}_{(1,1)} \cdot \underbrace{a^{[1]}^T}_{(1,3)}$$

Element wise multiplication.

"REFER NOTES FOR RIGOROUS PROOF."

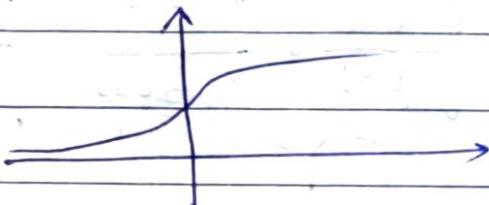
$$\frac{\partial J}{\partial w^{[2]}} = \frac{1}{m} \sum_{i=1}^n \frac{\partial J}{\partial w^{[2]}}$$

Improving your NN (Neural network)

I Experiment with activation functions.

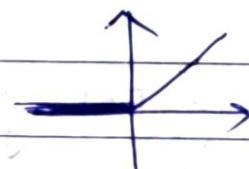
sigmoid ($\sigma(z)$)

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



$$\sigma'(z) = \sigma(z)(1 - \sigma(z))$$

$$\underline{\text{ReLU}(z)} = \begin{cases} 0 & \text{if } z \leq 0 \\ z & \text{if } z > 0 \end{cases}$$

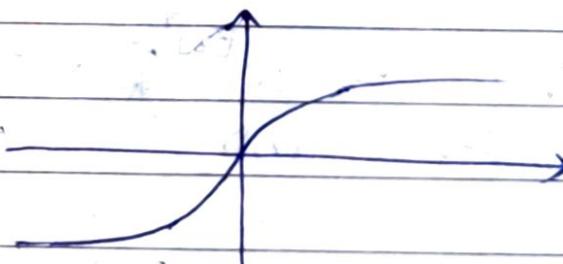


$$\text{ReLU}'(z) = 1_{\{z > 0\}}$$

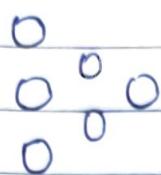
~Indicator function

$$\underline{\tanh(z)} = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

$$\tanh'(z) = 1 - (\tanh(z))^2$$



why do we need activation functions?



let activation = Id func. $z \mapsto z$.

$$\hat{y} = a^{[3]} = z^{[3]} = w^{[3]} a^{[2]} + b^{[3]}$$

$$\hat{y} = w^{[3]}(w^{[2]} a^{[1]} + b^{[2]}) + b^{[3]}$$

$$\hat{y} = w^{[3]} w^{[2]} (w^{[1]} x + b^{[1]}) + w^{[3]} b^{[2]} + b^{[3]}$$

$$\hat{y} = w x + b.$$



happens if we

regardless of
depth NN becomes
linear.

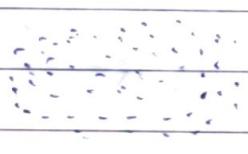
have no activation
function

② Use initialisation or normalisation

Normalizing your input

$$\text{assume data} \rightarrow x = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$$

$x_2 \uparrow$



$$x = x - \mu$$



$$x_i = \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix} = \mu = \frac{1}{n} \sum_{i=1}^n x^{(i)}$$

for fig 2 loss function



for grad. descent the slope
bounces around

$$x = \frac{x}{\sigma}$$

$$\sigma = \sqrt{\frac{1}{n} \sum_{i=1}^n (x^{(i)})^2}$$

for fig 3 loss function.



hence we converted the
data

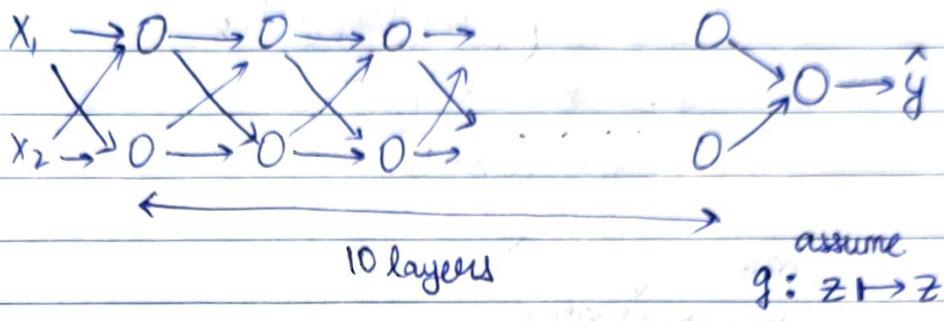
points consistently inward
for grad. descent



fig 3

Vanishing / Exploding gradients.

assume a NN



$$\hat{y} = w^{[L]} \cdot a^{[L-1]} \neq w^{[L]} \cdot w^{[L-2]} \cdot a^{[L-2]}$$

$$\hat{y} = w^{[L]} \cdot w^{[L-1]} \cdot \dots \cdot w^{[1]} \cdot x$$

$\underbrace{\quad\quad\quad}_{(2,2) \text{ matrices}}$

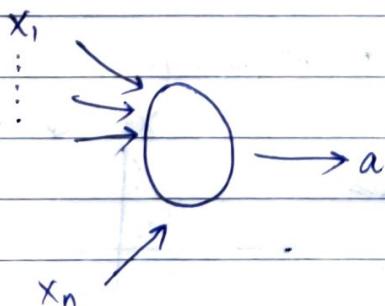
assume that $w^{[L]}$ is slightly larger than Identity

matrix $\rightarrow \begin{pmatrix} 1.5 & 0 \\ 0 & 1.5 \end{pmatrix}$

all mult. become
"explodes" for large values of \downarrow

$$\begin{pmatrix} 1.5^L & 0 \\ 0 & 1.5^L \end{pmatrix}$$

Example with 1 neuron :-



$$a = \sigma(z)$$

$$z = w_1 x_1 + \dots + w_n x_n$$

large $n \rightarrow$ small w_i as we want around 1 because of exploding gradients

initialise $w_i \sim \frac{1}{n}$

some initialisation schemes are :- (practical obs)

$$w_i = \text{np.random.randn(shape)} * \text{np.sqrt}\left(\frac{1}{n^{[L-1]}}\right)$$

for sigmoid

$L \rightarrow$ layers we are at

$n \rightarrow$ no. of neurons in a layer

for ReLU \rightarrow same above eqn. with

$$\frac{2}{n^{[L-1]}} \text{ instead of } \frac{1}{n^{[L-1]}}$$

Xavier initialisation

$$w^{[L]} \sim \sqrt{\frac{1}{n^{[L-1]}}} \text{ for tanh}$$

He Initialisation:

$$w^{[L]} \sim \sqrt{\frac{2}{n^{[L]} + n^{[L-1]}}}$$

less computationally expensive

computing J for each minibatch
only takes 1000 steps

$$J = \frac{1}{1000} \sum_{i=1}^{1000} L^{(i)}$$

algo

for iteration $t = 1, \dots$

select batch $(x^{[t]}, y^{[t]})$

update $w^{[L]}, b^{[L]}$ by
backprop and such

$$X = (x^{(1)}, \dots, x^{(m)})$$

$$Y = (y^{(1)}, \dots, y^{(m)})$$

mini-batch

$$X = (x^{[t]}, \dots, x^{[t]})$$

each minibatch
is of 1000 examples

$$Y = (y^{[t]}, \dots, y^{[t]})$$

Gradient descent + Momentum Algorithm

assume we have an ^{spread out} extended loss func.

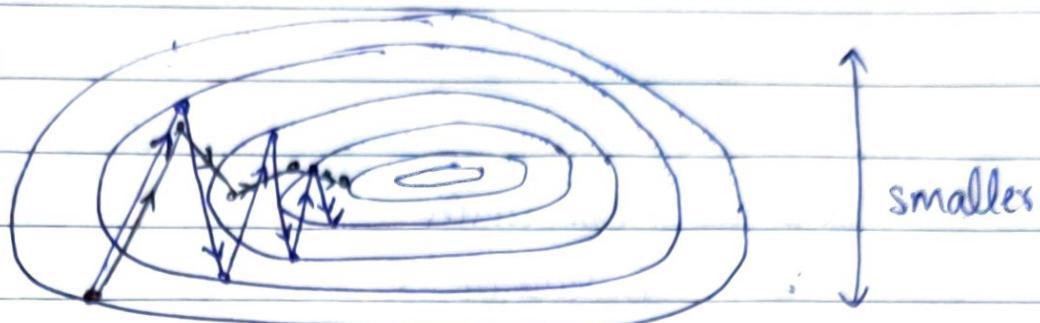
"INTUITION"

- → normal G-D

- → using momentum

looks at how
we moved in the
past by taking
average movement

this G-D bounces all over the place



larger

smaller

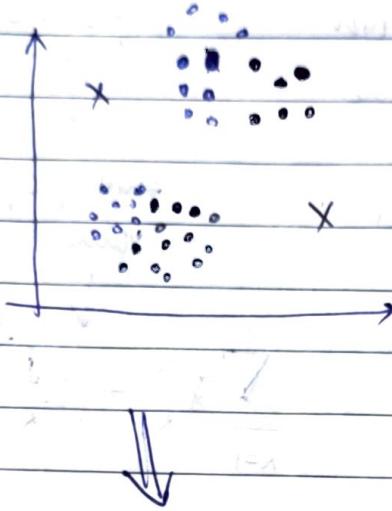
↓
average we want more movement in the horizontal dirn. than
movement is little in vertical dirn. as that'd allow us to get optimum
along y axis value earlier
and more along
x-~~axis~~ axis.

$$w = w - \alpha \frac{\partial J}{\partial w} \rightarrow \text{normal gradient descent.}$$

$$\left. \begin{aligned} v &= \beta v + (1-\beta) \frac{\partial J}{\partial w} \\ w &= w - \alpha v \end{aligned} \right\} \text{gradient descent with momentum}$$

Unsupervised Learning

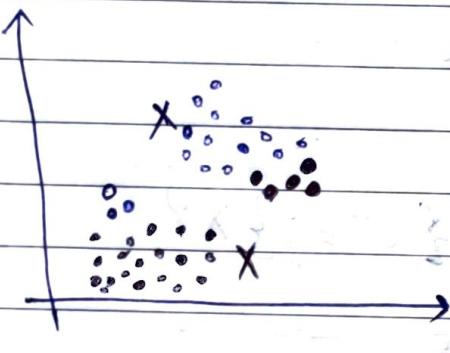
K-means clustering



• → data point

X → centroid

we pick 2 centroids randomly here
and we colour all data points depending
on which centroid it is closer to.



centroid is updated as new centre of
coloured values.

new black centroid is centre of
all black points.

and we keep on doing this
process.

Mathematical notation for K-means clustering

Data : $\{x^{(1)}, \dots, x^{(m)}\}$.

1) Initialize cluster centroid,

$\mu_1, \mu_2, \dots, \mu_k \in \mathbb{R}^n$ (for low dimensions)

how to initialise → randomly pick assign them value.

OR

pick k examples from dataset and make
them as centroids. (High dimensions like)

2) Repeat until convergence

a) Set $c^{(i)} := \underset{j}{\operatorname{argmin}} \|x^{(i)} - \mu_j\|_2^2$

$\|\cdot\|_2$ subscript denotes L₂ norm

number here for each $x^{(i)}$ in our dataset we assign it a color $c^{(i)}$ depending on the index number of centroid closest to it.
like μ_1 is closer or μ_2 is closer.

sum of squares of values of the vector

$$\sum_{k=1}^n (x_k^{(i)} - \mu_k^{c(i)})^2$$

b) For $j = 1, \dots, k$

$$\mu_j := \frac{\sum_{i=1}^m \mathbb{1}\{c^{(i)} = j\} x^{(i)}}{\sum_{i=1}^m \mathbb{1}\{c^{(i)} = j\}}$$

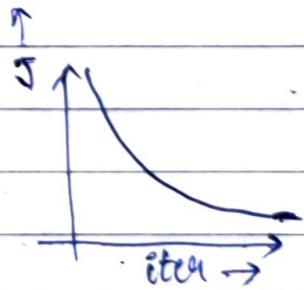
"updates the centroid" by relocating it as the centre of all points coloured that colour.

* cost function \rightarrow

$$J(c, \mu) = \sum_{i=1}^m \|x^{(i)} - \mu_{c^{(i)}}\|^2$$

"assigned colors"

"centroid"

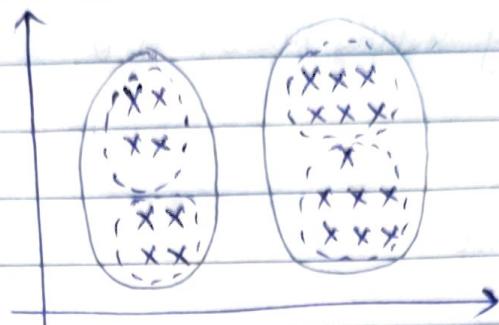


proof in lecture notes.

$\rightarrow J$ converges as it goes down every iteration

how to decide how many clusters to pick?

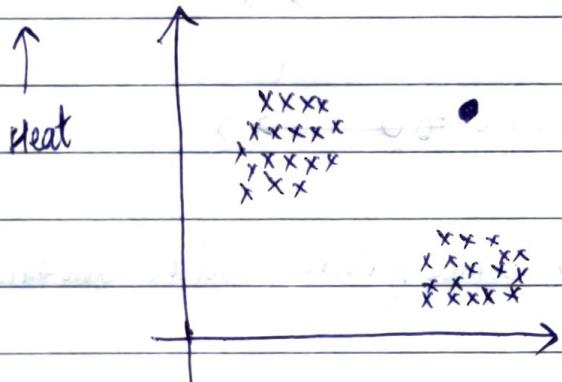
↳ you can pick it by hand.



someone might see 2 segments
clusters while someone may
see 4.

Density estimation

suppose we have an aircraft engine manufacturer.



$x \rightarrow$ good engines

• \rightarrow we want to test this
engine against the ones we
know are good

vibration \rightarrow

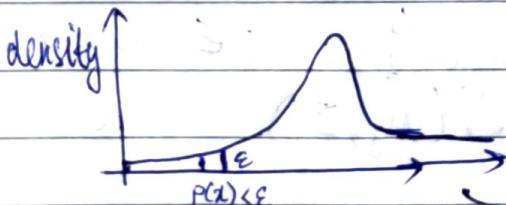
probability
we make a density distribution curve

example:

just consider vibration

we find $p(x)$ according to that
model.

if $p(x) < \epsilon$ \rightarrow anomaly



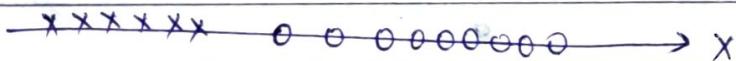
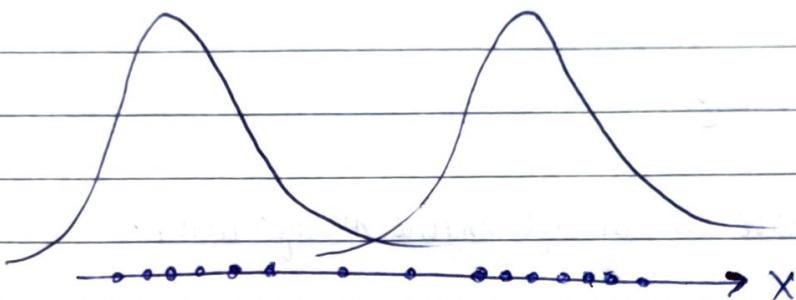
now try to visualise
this in 3D (we will get hills)

There is no standard algorithm model distribution to accomplish all these density distributions.

So, we use mixture of gaussians model to make these distributions.

Mixture of Gaussians model

1-D example



How to know which dataset comes from which Gaussian??

Mixture of Gaussians model :→

suppose there is a latent (hidden/unobserved) random variable z , and $x^{(i)}, z^{(i)}$ are distributed.

$$p(x^{(i)}, z^{(i)}) = P(x^{(i)} | z^{(i)})P(z^{(i)})$$

where $z^{(i)} \sim \text{Multinomial}(1) \quad z \in \{1, 2, \dots, k\}$

$$x^{(i)} | z^{(i)} = j \sim N(\mu_j, \Sigma_j)$$

↓

if we are

representing as
a mixture of k
Gaussians.

$z^{(i)}$ → denotes which gaussian does a point belong to.

how is mixture of gaussian model different from Gaussian discriminant Analysis?

→ In GDA $(x^{(i)}, z^{(i)})$, we know $z^{(i)}$; i.e. to which Gaussian class a data item belongs. On the other hand $z^{(i)}$ is unknown in a mixture of Gaussian model.

If we knew values of $z^{(i)}$ we could use Maximum Likelihood Estimate (MLE) as below.

$$l(\Phi, \mu, \Sigma) = \sum_{i=1}^m \log(p(x^{(i)}, z^{(i)}; \Phi, \mu, \Sigma))$$

upon differentiating to minimise maximise we get $\left(\frac{dl}{d\Phi} = 0 \text{ and so on}\right)$

$$\hat{\Phi}_j = \frac{1}{m} \sum_{i=1}^m \sum \{z^{(i)} = j\}$$

$$\hat{\mu}_j = \frac{1}{m} \sum_{i=1}^m \sum \{z^{(i)} = j\} x^{(i)}$$

$$\sum_{i=1}^m \sum \{z^{(i)} = j\}$$

$$\hat{\Sigma}_j = \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \hat{\mu}_{z^{(i)}})(x^{(i)} - \hat{\mu}_{z^{(i)}})^T$$

↳ verify from lecture notes

EM (expectation maximisation)

has 2 steps :-

1) E-step (Gives value of $z^{(i)}|x$)

$$\text{set } w_j^{(i)} = P(z^{(i)}=j | x^{(i)}; \Phi, \mu, \Sigma) \text{ by bayes rule}$$

$$w_j^{(i)} = \frac{P(x^{(i)} | z^{(i)}=j) \cdot P(z^{(i)}=j)}{\sum_{k=1}^K P(x^{(i)} | z^{(i)}=k) \cdot P(z^{(i)}=k)}$$

prob. that $z^{(i)}=j$ given $x^{(i)}$

$$\mathcal{N}(\mu_j, \Sigma_j) = \frac{1}{(2\pi)^{n/2} |\Sigma_j|^{1/2}} \exp\left(-\frac{1}{2} (x^{(i)} - \mu_j)^T \Sigma_j^{-1} (x^{(i)} - \mu_j)\right)$$

$$\begin{aligned} z &\sim \text{multinomial} \\ P(z^{(i)}=j) &= \Phi_j \end{aligned}$$

2) M-step

$$\Phi_j := \frac{1}{M} \sum_{i=1}^m w_j^{(i)}$$

(when we knew $z^{(i)}$)
earlier formula had

$$\mu_j := \frac{\sum_{i=1}^m w_j^{(i)} \cdot x^{(i)}}{\sum_{i=1}^m w_j^{(i)}}$$

$\sum_j \Phi_j = \dots \dots$ (refer notes)

$$2 \{ z^{(i)} = j \}$$

now we have

$$E[2 \{ z^{(i)} = j \}]$$

expected value of $2 \{ z^{(i)} = j \}$

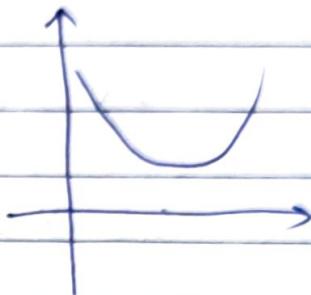
which is the prob. of $z^{(i)} = j$

which is $w_j^{(i)}$ in place of $2 \{ z^{(i)} = j \}$

~~EM step~~ is another EM is a pretty decent estimate of deciding which cluster a point belongs to. Is an alternative for K-means clustering.

Jensen's inequality

Let f be a convex func. $\Rightarrow (f''(x) > 0)$



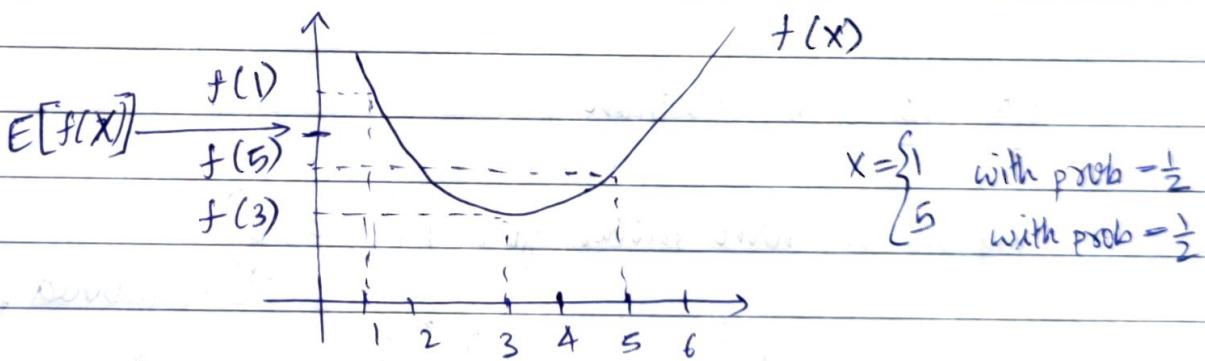
Let X be a random variable
then

$$f(E[X]) \leq E[f(X)]$$

$E[X] \rightarrow$ expected value of X

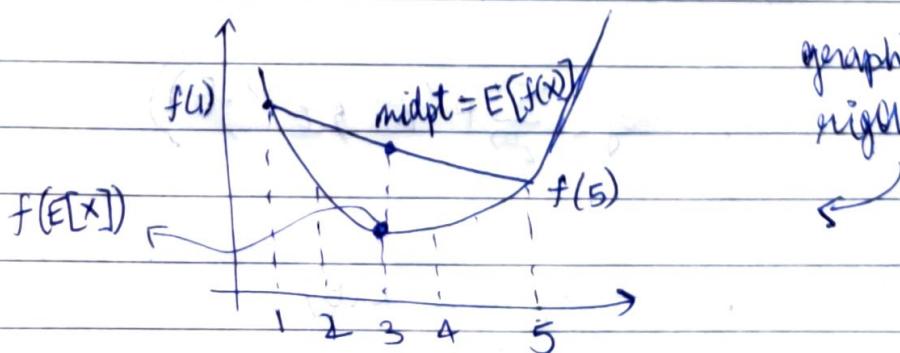
$E[f(X)] \rightarrow$ expected value of $f(X)$

Example: →



$$f(3) = f(E[X])$$

$$E[f(X)] = \frac{1}{2}f(1) + \frac{1}{2}f(5)$$



graphical non-negativity proof

Further if $f''(x) > 0$ (f is strictly convex)

then if

$$E[f(x)] = f(E[x]) \Leftrightarrow x \text{ is a constant}$$

(i.e. $x = E[X]$ with prob. 1)

for concave function

Let f be a concave funct. $\Rightarrow f'' < 0$

Let X be a random variable

$$\text{then } f(E[X]) \geq E[f(X)]$$

Further if $f'' < 0$ then if $f(E[X]) = E[f(X)] \Leftrightarrow X \text{ is constant}$

EM model visualisation

suppose we have model for $p(x, z; \theta)$

on $X \in \{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$

↗ combined parameter

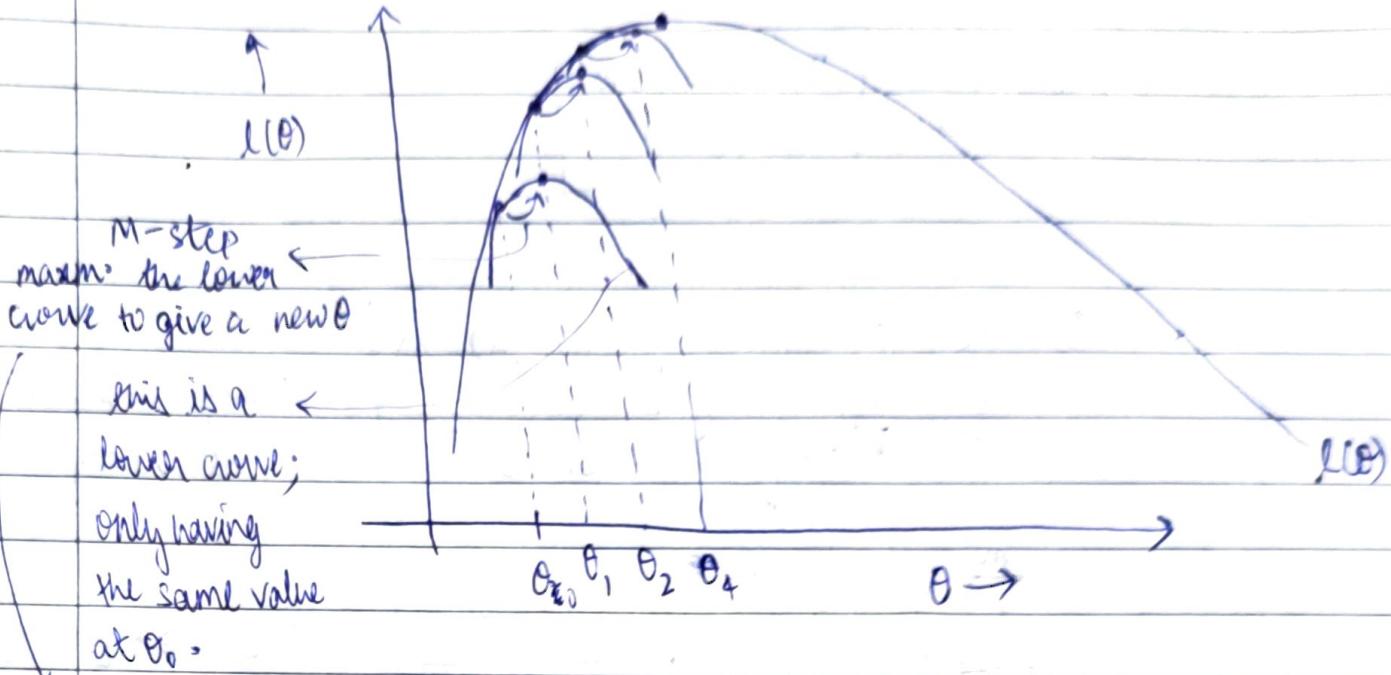
↓
has all values

$$l(\theta) = \sum_{i=1}^m \log(p(x^{(i)}; \theta))$$

μ, Σ, Φ contained

$$= \sum_{i=1}^m \log \left(\sum_{z^{(i)}} p(x^{(i)}, z^{(i)}; \theta) \right)$$

what we want : $\arg \max_{\theta} l(\theta)$



doing this construction and maximisation we slowly converge to the local maxima.

our goal is to $\max_{\theta} l(\theta)$

$$= \max_{\theta} \sum_{i=1}^m \log(P(x^{(i)}; \theta))$$

$$= \sum_{i=1}^m \log \left(\sum_{z^{(i)}} P(x^{(i)}, z^{(i)}; \theta) \right)$$

$$= \sum_{i=1}^m \log \left(\sum_{z^{(i)}} Q_i(z^{(i)}) \left[\frac{P(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} \right] \right)$$

where $Q_i(z^{(i)})$ is a probability distribution

$$(i.e. \sum_{z^{(i)}} Q_i(z^{(i)}) = 1)$$

$$= \sum_{i=1}^m \log \mathbb{E}_{z^{(i)} \sim Q_i} \left[\frac{P(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} \right]$$

(Expectation step)

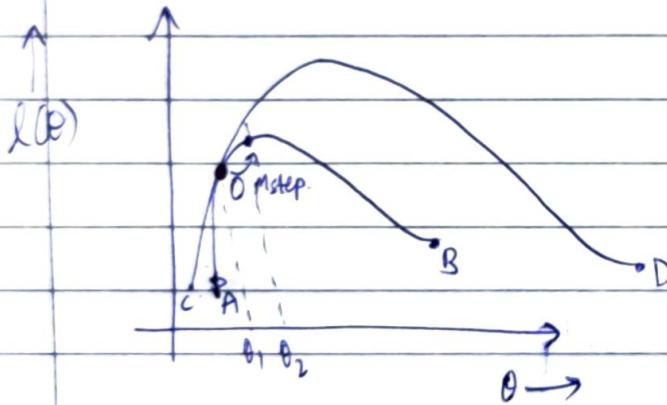
$$\geq \sum_{i=1}^m \mathbb{E}_{z^{(i)} \sim Q_i} \left[\log \frac{P(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} \right]$$

using Jensen's
inequality

$$= \sum_{i=1}^m \sum_{z^{(i)}} Q_i(z^{(i)}) \log \frac{P(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})}$$

$$f(\bar{x}) \geq E[f(x)]$$

$\uparrow f(x) = \log x$



On a given iteration of EM with current param. θ
we want : →

$$\textcircled{i} \rightarrow \log \mathbb{E}_{z^{(i)} \sim Q_i} \left[\frac{P(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} \right] = \mathbb{E}_{z^{(i)} \sim Q_i} \left[\log \left[\frac{P(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} \right] \right]$$

for lower curve to touch upper curve at 0.

(i) is true when random variable is constant. (By Jensen's inequality)

$$\text{i.e. } \frac{p(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} = \text{constant}$$

$$\text{hence, } Q_i(z^{(i)}) \propto p(x^{(i)}, z^{(i)}; \theta)$$

$$\text{but since } \sum_{z^{(i)}} Q_i(z^{(i)}) = 1$$

hence;

$$Q_i(z^{(i)}) = \frac{p(x^{(i)}, z^{(i)}; \theta)}{\sum_{z^{(i)}} p(x^{(i)}, z^{(i)}; \theta)}$$

$$Q_i(z^{(i)}) = p(z^{(i)} | x^{(i)}; \theta) \quad [\text{Refer "Lecture notes"}]$$

E-step therefore involve equality at 0 (point in graph)

Set, \checkmark these were taken as $w_j^{(i)}$ earlier on

$$Q_i(z^{(i)}) \doteq p(z^{(i)} | x^{(i)}; \theta)$$

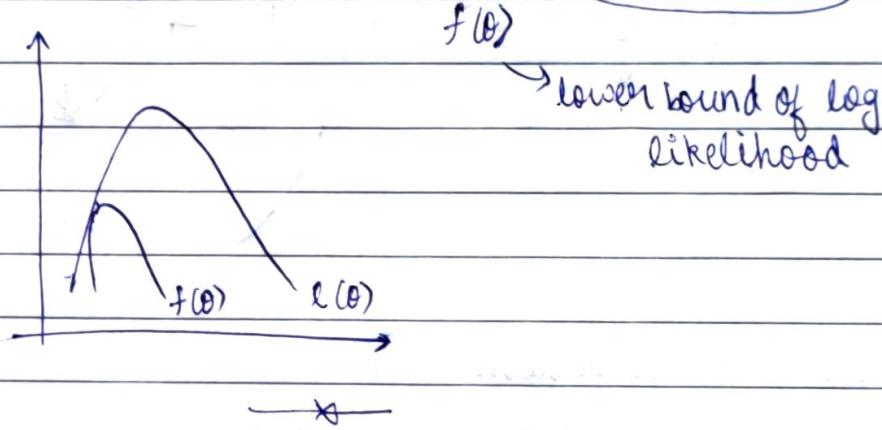
M-step involves finding new θ where $f(\theta)$ is maxm. ($f(\theta)$ is the lower curve here)

$$\theta := \arg \max_{\theta} \sum_i \sum_{z^{(i)}} Q_i(z^{(i)}) \log \left(\frac{p(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} \right)$$

E-M convergence

$$E\text{-step: } Q_i(z^{(i)}) = p(z^{(i)} | x^{(i)}; \theta)$$

$$M\text{-step: } \theta := \operatorname{argmax}_{\theta} \sum_i \sum_{z^{(i)}} Q_i(z^{(i)}) \log \left(\frac{p(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} \right)$$



For mixture of Gaussians model we had: →

$$p(x^{(i)}, z^{(i)}) = p(x^{(i)} | z^{(i)}) p(z^{(i)})$$

$$z^{(i)} \sim \text{Multinomial}(1)$$

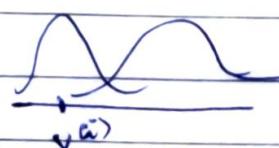
$$\text{i.e. } [p(z^{(i)}=j) = \Phi_j]$$

$$(x^{(i)} | z^{(i)}=j) \sim N(\mu_j, \Sigma_j)$$

E-step for mixture of gaussian

$$\omega_j^{(i)} = Q_i(z^{(i)}=j) = p(z^{(i)}=j | x^{(i)}; \Phi, \mu, \Sigma)$$

$$"p(z^{(i)}=j)"$$



$\omega_j^{(i)}$ determines what fraction of each gaussian combine to give $x^{(i)}$

M-step

$$\max_{\{\mu, \Sigma\}} \sum_i \sum_{z^{(i)}} Q_i(z^{(i)}) \log \left(\frac{P(x^{(i)}, z^{(i)}; \mu, \Sigma)}{Q_i(z^{(i)})} \right)$$

collectively
work taken
as θ

$$= \sum_i \sum_j w_j^{(i)} \cdot \log \left(\frac{1}{(2\pi)^{1/2} |\Sigma|^{1/2}} \exp \left(-\frac{1}{2} (x^{(i)} - \mu_j)^\top \Sigma_j^{-1} (x^{(i)} - \mu_j) \right) \times \Phi_j \leftarrow p(z^{(i)}) \right)$$

$w_j^{(i)}$

We want to maxm above expr.

so, we take derivative w.r.t. param.

strength with which $x^{(i)}$ is assigned
jth gaussian

$$\nabla_{\mu_j} (\dots) \stackrel{\text{set}}{=} 0 \Rightarrow \mu_j = \frac{\sum_i w_j^{(i)} X^{(i)}}{\sum_i w_j^{(i)}} \quad "P(z^{(i)} = j | X^{(i)}, \dots)"$$

$$\nabla_{\Phi_j} (\dots) \stackrel{\text{set}}{=} 0 \Rightarrow \Phi_j = \frac{\sum_i w_j^{(i)}}{\sum_i \sum_j w_j^{(i)}}$$

$$\nabla_{\Sigma_j} (\dots) \stackrel{\text{set}}{=} 0 \Rightarrow \Sigma_j = \dots \quad (\text{refer notes})$$

Factor analysis model

We will be using mixture of Gaussian model for this but, instead

instead of $z^{(i)}$ being discrete it is continuous
 $z^{(i)} \sim N$

$$J(\theta, Q) = \sum_i \sum_{z^{(i)}} Q_i(z^{(i)}) \log \left(\frac{P(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} \right)$$

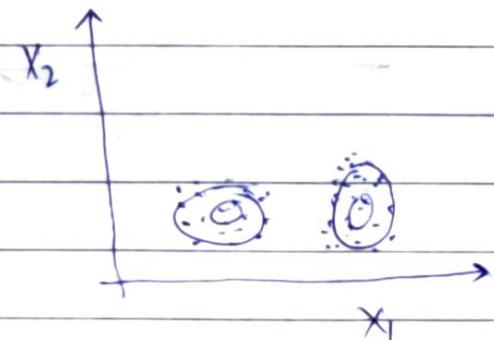
we proved via Jensen's,

$$l(\theta) \geq J(\theta, Q)$$

for any θ, Q

E-step: maximise J wrt - Q } set it such that equality occurs at initial point
M-step: maxim. J wrt. θ } coordinate ascent } we have multiple features we want to maxim. func. on
 and we go back and forth b/w them to maxim. other func.

Mixture of gaussians example → say $n=2$ and $m=100$



$m \gg n$ here

we will not use mixture of gaussians when $m \approx n$ or $m < n$

$$\text{say } m = 30, n = 100$$

try to
suppose we model above case as a single gaussian

$$x \sim \mathcal{N}(\mu, \Sigma)$$

$$\text{MLE: } \mu = \frac{1}{m} \sum_i x^{(i)}$$

$$\Sigma = \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \mu)(x^{(i)} - \mu)^T$$

If $m < n$, then Σ is "singular / non-invertible"

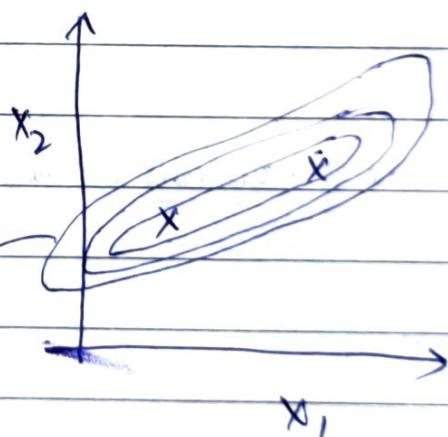
$$\frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp(-\dots) (\Sigma^{-1}) (\dots)$$

problem as → hence we can't model as
 Σ is singular. a single Gaussian

illustration of above example.

$$n=2 \quad m=2$$

contours as essentially
straight lines with zero width
and infinite length

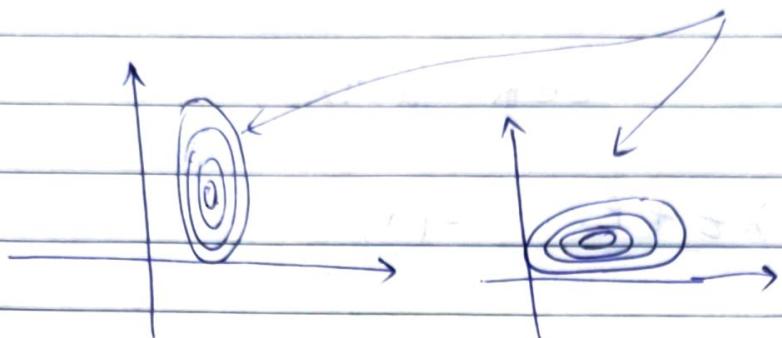


Hence we need a new model

Option 1: constraint Σ to be diagonal

$$\Sigma = \begin{bmatrix} \sigma_1^2 & \dots & 0 & 0 \\ 0 & \sigma_2^2 & \dots & 0 \\ 0 & 0 & \ddots & \vdots \\ 0 & 0 & \dots & \sigma_n^2 \end{bmatrix} \rightarrow \sigma_j^2 = \frac{1}{m} \sum_{i=1}^m (x_j^{(i)} - \mu_j)^2$$

→ this constrains our gaussian to be along axis



MLE: is same ; Above model is limiting as it assumes any two features are completely "uncorrelated".



we won't use this model

Option 2:

Constrain Σ to be $\Sigma = \sigma^2 I$

i.e. $\Sigma = \begin{bmatrix} \sigma^2 & 0 & 0 \\ 0 & \dots & 0 \\ 0 & 0 & \sigma^2 \end{bmatrix}$ → gives circular gaussians.

~~MLE~~ $\sigma^2 = \frac{1}{m} \frac{1}{N} \sum_i \sum_j (x_j^{(i)} - \mu_j)^2$

This is even more limiting than option 1 and so isn't used

Consider a model where;

$$p(x, z) = p(x|z) \cdot p(z)$$

z is hidden

(example)
 $d = 3$
 $n = 100$
 $m = 30$

$$z \sim N(0, I) \quad z \in \mathbb{R}^d \quad (d < n)$$

$$x = \mu + \lambda z + \epsilon. \quad \text{--- (i)}$$

$$\text{where } \epsilon \sim N(0, \Psi)$$

parameters: $\rightarrow \mu \in \mathbb{R}^n$; $\lambda \in \mathbb{R}^{n \times d}$; $\Psi \in \mathbb{R}^{n \times n}$ and Ψ is diagonal

$\rightarrow x|z \sim N(\mu + \lambda z, \Psi)$ is equivalent to (i)

Example: \rightarrow

$$z \in \mathbb{R}^1 \quad x \in \mathbb{R}^2 \quad d=1$$

$\# n=2$
 $m=7$

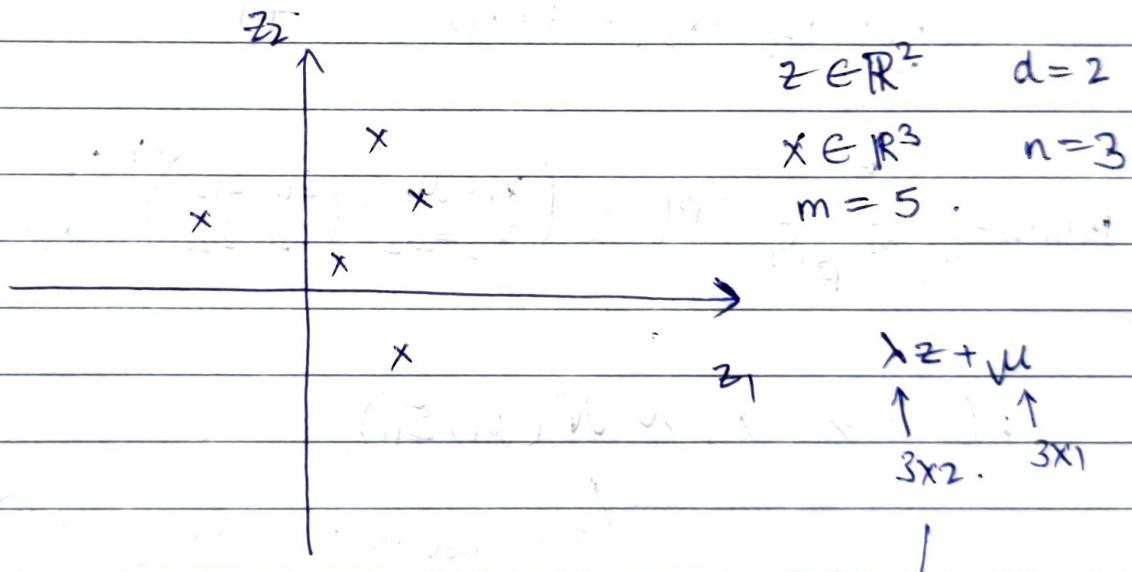
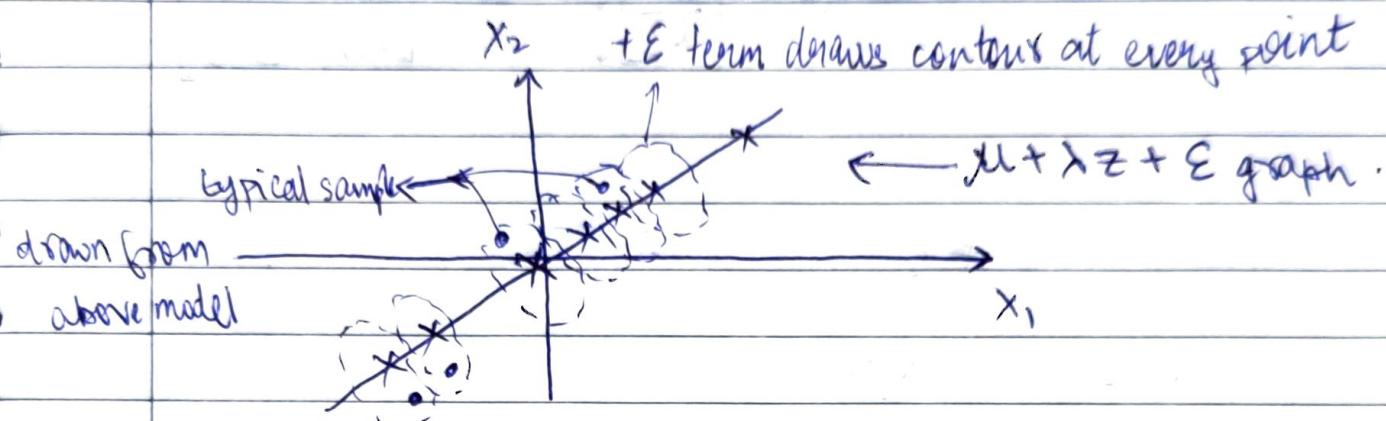
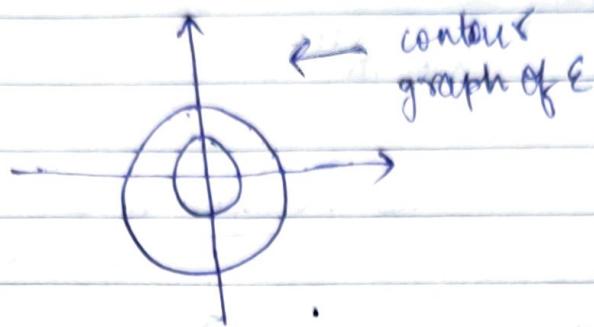
$$\text{--- } \times \times \times \times \times \times \times \rightarrow z^{(i)}$$

say $\lambda = \begin{bmatrix} 2 \\ 1 \end{bmatrix}$ $\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$

$$\lambda z + \mu \in \mathbb{R}^2$$



say $\Psi = \begin{bmatrix} 1 & 0 \\ 0 & 2 \end{bmatrix}$



$\Psi \in \mathbb{R}^{3 \times 3}$

so we can imagine like spherical contours.

this maps our points onto a 3-D plane (still ^{all} planar points).

Multivariate Gaussian

(Refer notes)

$$x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \quad \begin{array}{l} \text{↑ or components } x_1 \in \mathbb{R}^n, x_2 \in \mathbb{R}^g \\ \text{↓ } \dots \quad \quad \quad x \in \mathbb{R}^{n+g} \end{array}$$

$x \rightarrow$ is "partition" vector.

$$x \sim N(\mu, \Sigma) \Rightarrow \mu = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix} \quad \begin{array}{l} \text{↑ } n \\ \text{↑ } g \end{array} \quad \Sigma = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}$$

Marginal : $P(x_1) = ?$

$$P(x) = P(x_1, x_2)$$

$$\int_{x_2} P(x_1, x_2) dx_2 = P(x_1). \quad \textcircled{1}$$

$$P(x_1, x_2) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2} \begin{pmatrix} x_1 - \mu_1 \\ x_2 - \mu_2 \end{pmatrix}^\top \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix}^{-1} \begin{pmatrix} x_1 - \mu_1 \\ x_2 - \mu_2 \end{pmatrix}\right).$$

$$P(x_1) \Rightarrow x_1 \sim N(\mu_1, \Sigma_{11}).$$

↳ to show this we integrate ①

Conditional : $P(x_1 | x_2) = ?$

$$\frac{P(x_1, x_2)}{P(x_2)}$$

$$x_1 | x_2 \sim N(\mu_{1|2}, \Sigma_{1|2})$$

$$\mu_{1|2} = \mu_1 + \Sigma_{12} \Sigma_{22}^{-1} (x_2 - \mu_2)$$

$$\Sigma_{1|2} = \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21}$$

deriving E-M algo. using multivariate gaussians

1. Derive $p(x, z)$

[Refer notes]

$$\begin{pmatrix} z \\ x \end{pmatrix} \sim N\left(\mu_{x,z}, \Sigma\right)$$

$$z \sim N(0, I) \quad E_z = 0 \rightarrow \text{as mean is } 0 \text{ for } z.$$

$$x = \mu + \lambda z + \varepsilon \quad E_x = E[\mu + \lambda z + \varepsilon] = E[\mu] = \mu.$$

$$\mu_{x,z} = \begin{bmatrix} 0 \\ \vdots \\ \mu \end{bmatrix} \begin{array}{l} \uparrow d \\ \downarrow n \end{array}$$

$$\Sigma = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix} = \begin{bmatrix} E[(z - E_z)(z - E_z)^T] & E[(z - E_z)(x - E_x)] \\ E[(x - E_x)(z - E_z)^T] & E[(x - E_x)(x - E_x)] \end{bmatrix}$$

$$\Sigma_{22} = E[(x - E_x)(x - E_x)^T] = E[(\lambda z + \mu + \varepsilon - \mu)(\lambda z + \mu + \varepsilon - \mu)^T]$$

$$= E[(\lambda z + \varepsilon)(\lambda z + \varepsilon)^T]$$

$$E[\varepsilon \varepsilon^T] = \Psi \rightarrow \text{defn.}$$

$$\Sigma_{22} = E[\lambda z z^T \lambda^T + \lambda z \varepsilon^T + \varepsilon z^T \lambda + \varepsilon \varepsilon^T]$$

$$E[z z^T] = I \text{ as variance is } I.$$

$$\Sigma_{22} = E[\lambda z z^T \lambda^T] + E[\varepsilon \varepsilon^T]$$

value as
z and ε^T have 0
expected value and
are not correlated

$$\Sigma_{22} = \lambda E[z z^T] \lambda^T + \Psi$$

$$\Sigma_{22} = \lambda I \lambda^T + \Psi = \lambda \lambda^T + \Psi.$$

deriving similarly for all we get

$$\Sigma = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix} = \begin{bmatrix} I & \lambda^T \\ \lambda & \lambda\lambda^T + \Psi \end{bmatrix}$$

so we derived,

$$\begin{bmatrix} z \\ x \end{bmatrix} \sim N \left(\begin{bmatrix} 0 \\ \mu \end{bmatrix}, \begin{bmatrix} I & \lambda^T \\ \lambda & \lambda\lambda^T + \Psi \end{bmatrix} \right) \quad \#$$

~~we can~~ $p(x^{(i)})$ will be drawn from above Gaussian density

we can find the log likelihood for $p(x^{(i)})$ and differentiate it w.r.t parameters but we find out there is no closed form solution.

so, we use ~~the~~ E-M algorithm.

E-step:

$$Q_i(z^{(i)}) = p(z^{(i)} | x^{(i)}; \theta). \quad z^{(i)} \rightarrow \text{continuous here for factor analysis}$$

we can't use $w_j^{(i)}$ here as $z^{(i)}$ is continuous.

↳ we will need infinitely many $w_j^{(i)}$'s

$$z^{(i)} | x^{(i)} \sim N(\mu_{z^{(i)}|x^{(i)}}, \Sigma_{z^{(i)}|x^{(i)}}) \quad \rightarrow \text{store all}$$

$$\text{where, } \mu_{z^{(i)}|x^{(i)}} = \vec{0} + \lambda^T (\lambda\lambda^T + \Psi)^{-1} (x^{(i)} - \mu)$$

$$\Sigma_{z^{(i)}|x^{(i)}} = I - \lambda^T (\lambda\lambda^T + \Psi)^{-1} \lambda.$$

we got above equations using our formula for multivariate gaussians.

M-step: (Refer lecture notes)

mathematical trick for derivation: \rightarrow

$$\theta_i(z^{(i)}) = \frac{1}{(2\pi)^{\frac{d_h}{2}} | \dots |} (\exp(-\frac{1}{2}(\dots)))$$

if we get something expr. as below then,

$$\int_{z^{(i)}} \theta_i(z^{(i)}) z^{(i)} dz^{(i)} \rightarrow \text{don't subn. } \cancel{\theta_i(z^{(i)})} \text{ and } \theta_i(z^{(i)})$$

$$E_{z^{(i)} \sim \theta_i} [z^{(i)}] = \mu_{z^{(i)} | x^{(i)}}$$

as θ_i is gaussian
so mean is

$$\theta := \underset{\theta}{\operatorname{argmax}} \sum_i \int_{z^{(i)}} \theta_i(z^{(i)}) \log \frac{p(x^{(i)}, z^{(i)})}{\theta_i(z^{(i)})} dz^{(i)}$$

$$= \sum_i E_{z^{(i)} \sim \theta_i} \left[\log \frac{p(x^{(i)}, z^{(i)})}{\theta_i(z^{(i)})} \right]$$

plug in gaussian density.

$$= \sum_i E_{z^{(i)} \sim \theta_i} \left[\text{quad. func.} \right]$$

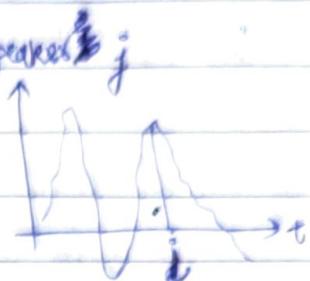
$$\nabla_{\mu} (\dots) = 0. \text{ solve for } \mu.$$

please refer lecture notes.

Independent Component analysis (ICA) part - 2

suppose we have n speakers, $s \in \mathbb{R}^n$. speaker j

$s_j^{(i)} \rightarrow$ speaker j at time i



$$(n, n) \xleftarrow{\text{microphones}} X^{(i)} = A s^{(i)} \xrightarrow{\text{record}} (n, n), X \in \mathbb{R}^n$$

here no. of microphones record. $\Rightarrow (n, n)$
 $= n$ - no. of speakers

Goal: Find $\omega = A^{-1}$ so,

$$s^{(i)} = \omega X^{(i)} \rightarrow \text{so we can tell each speaker individually}$$

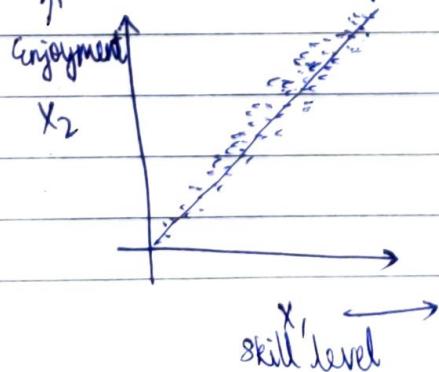
about

ICA, continued later

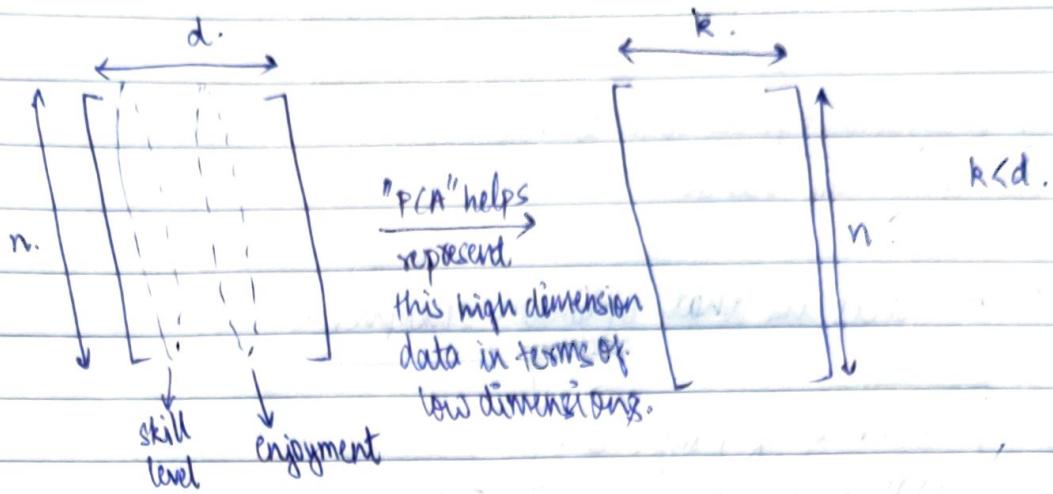
Principal Component Analysis

$$\{X^{(i)}\}_{i=1}^n, X^{(i)} \in \mathbb{R}^d$$

consider two features x_1 and x_2 which indicate skill level and enjoyment a pilot has during flight. These are other features too



$\therefore x_1$ and x_2 are correlated.



Steps in PCA.

- i) Standardize your dataset.

$$x_j^{(i)} = \frac{x_j^{(i)} - \mu_j}{\sigma_j}$$

basically each column is updated in such a manner that it is centered around its mean \Rightarrow new mean = 0 and is scaled such that $\Sigma_{new} = 1$

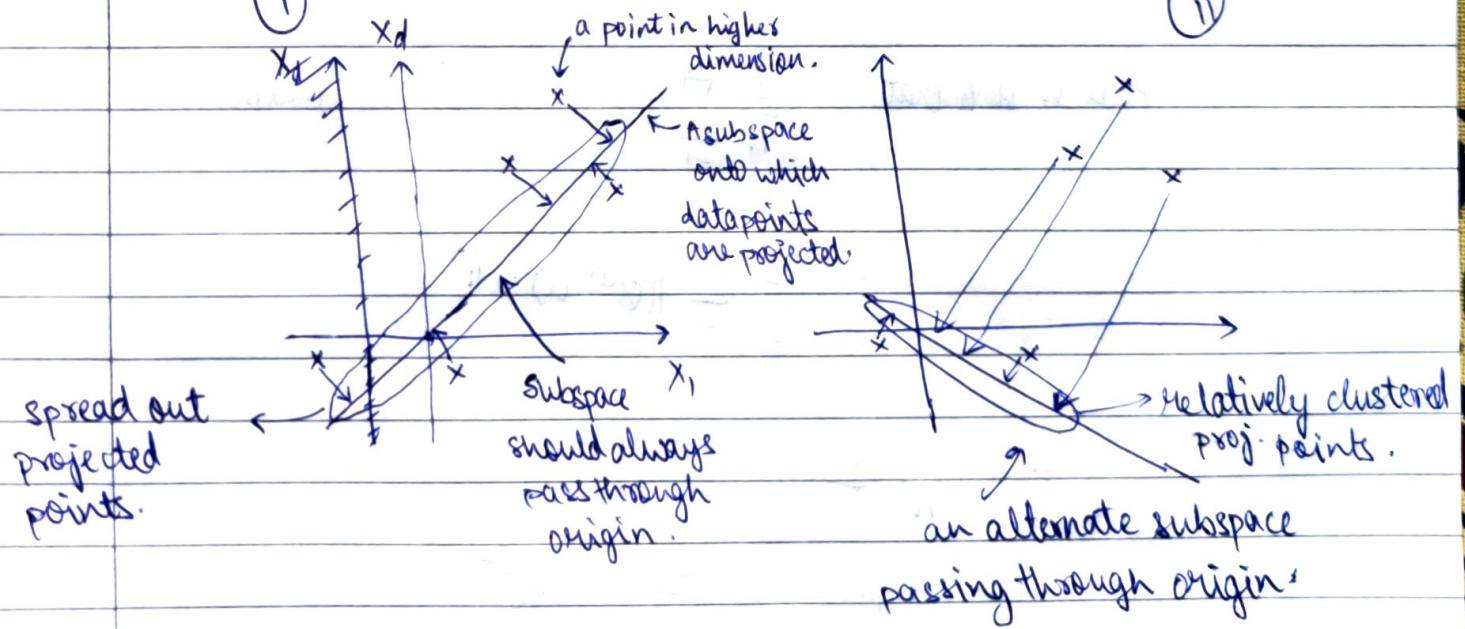
We do so because units may be in different units.

↳ example \rightarrow height and weight

↓ ↓
160-190 cm 60-90 kg.

⑨

⑩



in (i) projected points are sparse

in (ii) they are dense.

PCA tells us that optimum subspace.

variance of projected points is maximised → "sparser"

Idea is that we want to lose only the dimensions while retaining the variance.

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

unit length of ~~subspace onto~~ which we want to project the data.

So, projection of x on u will be = $\text{proj}(u) \cdot \vec{x}$

$$= \frac{u u^T}{u^T u} \vec{x} \quad u^T u = 1$$
$$= ((x^{(i)})^T u) \vec{u}.$$

Find u such that

$$\frac{1}{n} \sum_{i=1}^n \|\text{proj}(u) \cdot x\|^2 \text{ is max.}$$

$$= \frac{1}{n} \sum_{i=1}^n \|((x^{(i)})^T u) u\|^2$$

$$u = \frac{1}{n} \sum_{i=1}^n x^{(i)}$$

$$u = \underset{u}{\operatorname{argmax}} \frac{1}{n} \sum_{i=1}^n (x^{(i)^T} u)^2$$

$$u = \underset{u}{\operatorname{argmax}} \quad u^T \left(\frac{1}{n} \sum_{i=1}^n x^{(i)} (x^{(i)})^T \right) u$$

↓

sample covariance
matrix as mean is
zero

solution to an eqn. of form; $u = \underset{u}{\operatorname{argmax}} \quad u^T A u$

is $u = \text{eigenvector of largest eigenvalue of } A$

proof: $\rightarrow u^T A u$ has to be maxm. s.t. $u^T u = 1$

let $g(u, \lambda) = u^T A u - \lambda(u^T u - 1) \rightarrow$ using Lagrange multiplication technique

critical points: $\rightarrow \frac{\partial g}{\partial u} = 2A u - 2\lambda u = 0$

$A u = \lambda u$

we see u is an eigenvector corresponding to an eigen value λ .

$$f(\lambda) = \underset{u}{\operatorname{argmax}} \quad u^T A u = u^T \lambda u = \lambda u^T u = \lambda$$

so λ will be the largest eigen value of A
and u will be λ 's corresponding eigen vector.

so, in PCA we first have data as X so first we calculate

$$X^T X \quad \begin{array}{l} \rightarrow (\lambda_1, u_1) \\ \rightarrow (\lambda_2, u_2) \\ \rightarrow (\lambda_3, u_3) \\ \vdots \\ \rightarrow (\lambda_k, u_k) \end{array}$$

find k s.t. $\sum_{i=1}^k \lambda_i / \sum_{i=1}^n \lambda_i = 95\% \text{ or some percentage}$

$\rightarrow k$ largest eigenvalues.

$$\frac{\sum_{i=1}^k \lambda_i}{\sum_{i=1}^n \lambda_i} = 95\% \text{ or some percentage}$$

eigen value of $X^T X$ \rightarrow corresponding eigen vector

we retained 95% of variance in our data even though we went from d -dimensions to k dimensions.

Suppose we have a matrix M

if M is square & symmetric $\Rightarrow M$ has orthogonal eigen vectors & real eigen values.

if M is PSD $\Rightarrow M$ is square and symmetric $\Rightarrow \parallel \parallel$

in our case, we know $X^T X$ is PSD

Eigen decomposition on $X^T X \stackrel{\sim}{\equiv}$ singular value decomposition on X
 \downarrow
is equivalent
to
 $\longrightarrow *$

ICA part-1

Suppose we have a cocktail party with

d speakers

d microphones

what speakers
are saying at
an instance

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

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

$s_j^{(i)}$ \rightarrow j th speaker say at i th time

$x_j^{(i)}$ \rightarrow recording in j th mic at i th time.

what mics are
picking up at an
instance

$$X = A \cdot S$$

up

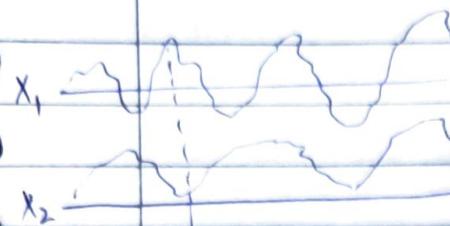
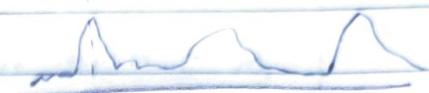
Mixing matrix

Example: →

(S₁)

(S₂)

S₁



Q

X₁

Q

X₂

S₂

time

$$x^{(i)} = (x_1^{(i)}, x_2^{(i)}) \in \mathbb{R}^d; d=2$$

$$s^{(i)} = (s_1^{(i)}, s_2^{(i)}) \in \mathbb{R}^d$$

where d=2

$$x^{(i)} = A s^{(i)}$$

\downarrow
(d,d)

$$w = A^{-1} \quad (\text{unmixing matrix}) \quad [s^{(i)} = w x^{(i)}]$$

Assumptions : → i) # S = # x

ii) S = w x → (linearly associated)

iii) S_j ⊥ S_k j ≠ k

R.V are independent

not necessarily true in a convo

iv) S_j is not Gaussian

p(x)

suppose we have x ~ Unif [0, 1]



p(y)

↑

y

↑

y

↑

y

$$p_y(y) = p_x(y/2) \cdot \frac{1}{2}$$

$$y = 2x$$

$$\begin{matrix} y \\ \mathbb{R}^d \\ \downarrow \\ \mathbb{R}^{d \times d} \\ \downarrow \\ \mathbb{R}^d \end{matrix}$$

$$y = w x \quad \text{say } p_y(y) = p_x(w^{-1}y) \cdot \frac{1}{|w|}$$

↓ known as Jacobian

in ICA we have

$$p(x) = \prod_{j=1}^d p_s(w_j^T x) \cdot |w|$$

$$W = \begin{bmatrix} -w_1 & - \\ \vdots & \end{bmatrix} \text{ unmixing matrix}$$

$$\begin{bmatrix} -w_j & - \\ \vdots & \end{bmatrix} \begin{bmatrix} x \\ \vdots \end{bmatrix} = s_j$$

$p_s \sim$ logistic distribution \rightarrow pdf of logistic func $\rightarrow \frac{1}{1+e^{-x}}$

$$\ell(w) = \sum_{i=1}^n \left[\left(\sum_{j=1}^d \log [\sigma(x^{(i)}) \cdot (1 - \sigma(x^{(i)}))] \right) + \log |w| \right]$$

do gradient ascent till convergence

$$w := w + \alpha \begin{bmatrix} 1 - 2\sigma(w_1^T x^{(i)}) \\ 1 - 2\sigma(w_2^T x^{(i)}) \\ 1 - 2\sigma(w_3^T x^{(i)}) \\ \vdots \end{bmatrix} x^{(i)\top} + (w^\top)^{-1}$$

Reinforcement learning

in games like chess or go, common reward function R is

$$R(s) = \begin{array}{ll} +1 & \text{for win} \\ -1 & \text{for lose} \\ 0 & \text{for draw} \end{array}$$

reward ↓ ↓
func. state

Credit assignment problem :-

i) Consider a chess game which was lost at move 50. Suppose, the blunder step was at move 20. How does the computer know this?

ii) Consider a car crash. Before crashing (-ve outcome) most often brakes are applied. But blunder/fault occurs earlier on.

~~How does~~

helps figure out what to do more of in +ve example
and what exactly to do less of in -ve example.

Markov Decision Process (MDP)

- helps model RL problems
- notation for modelling how world works

MDP involves 5-tuples →

$$(S, A, \{P_{sa}\}, \gamma, R)$$

S → set of states

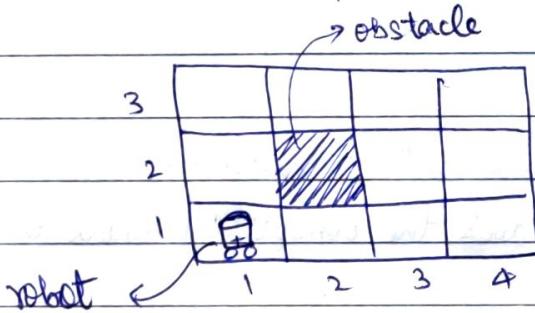
A → set of actions

P_{sa} → state transition probabilities ($\sum_{s'} P_{sa}(s') = 1$)

γ → discount factor, $\gamma \in [0, 1]$

R → reward function

Consider an example to show MDP



① here the robot can take 11 states

② Actions robot takes are {N, S, E, W}

③ suppose we press N command for our robot
0.8 → 0.1 0.1 → 0.1
so if giving N has 0.8 prob., 0.1 for east and west

So,
 $P_{(3,1)N}((3,2)) = 0.8$; $P_{(3,1)N}((3,1)) = 0.1$, $P_{(3,1)N}((2,1)) = 0.1$
states $s \downarrow$ action $a \downarrow$ states s'

$P_{(3,1)N}(\text{anything else}) = 0 \rightarrow$ so we have P_{sa} for all states.

③ Reward for this example :→

3			+1
2			-1
1	2	3	4

$$R(4,3) = +1 \rightarrow \text{incentivise going there}$$

$$R(4,2) = -1 \rightarrow \text{punish for going there.}$$

if we want robot to reach $R(4,3)$ ASAP → we put a small negative reward for all other states (ex → -0.01) $R(s) = -0.01$

punishes robot
for taking long.

Our robot starts at s_0 and picks an action a_0 .

s_0

choose action a_0

Get to $s_1 \sim P_{s0a0}$

choose action a_1

Get to $s_2 \sim P_{s1a1}$

choose action a_2

and so on

in finance terms

helps relate to time value

of money.

some amount.

get debt later

γ here is discount factor

→ helps take positive actions sooner
→ and delays negative rewards to later

→ generally chosen to be like 0.99 and such.

γ is useful as it helps total payoff converge.

get profit/loan

sooner (inflation/invest)

Goal: choose actions over time to maximise EXPECTED total payoff

$$E[R(s_0) + \gamma R(s_1) + \gamma^2 R(s_2) + \gamma^3 R(s_3) \dots]$$

policy
(controller) $\pi : S \mapsto A$ (π maps states to actions)
when in state s take action $\pi(s)$

Optimal policy for MDP in example (solved on PC by prof not by hand)

3	→	→	→	+1
2	↑		↑	-1
1	↑	←	←	←

$\pi((3,1)) = \omega$

How to compute optimal policy?

Define: V^π, V^*, π^*

for a policy π , $V^\pi : S \mapsto \mathbb{R}$ is s.t. $V^\pi(s)$ is the expected total payoff for starting in state s and executing π .

$$V^\pi(s) = E[R(s_0) + \gamma R(s_1) + \dots | \pi, s_0 = s]$$

"value func for policy π "

Bellman's equations

$$V^\pi(s) = R(s_0) + \gamma \sum_{s'} P_{\pi}(s')(s') V^\pi(s')$$

"immediate reward"

$$V^\pi(s) = E \left[R(s_0) + \gamma \underbrace{(R(s_1) + \gamma R(s_2) + \gamma^2 R(s_3) + \dots)}_{V^\pi(s_1)} \right] \mid \pi, s_0=s$$

$$V^\pi(s) = R(s_0) + \gamma \sum_{s'} P_{s,\pi(s)}(s') \cdot V^\pi(s')$$

Given π , get a linear system of eqns. in terms of $V^\pi(s)$

3				+1
2				-1
1	1	2	3	4

$\pi(3,1) = \text{North}$

$$V^\pi(3,1) = R((3,1)) + \gamma (0.8 * V^\pi(3,2) + 0.1 * V^\pi(2,1) + 0.1 * V^\pi(4,1))$$

→ unknowns.

we will get 11 such linear eqns. here which we have to solve

V^* defn: → V^* is the optimal value function

$$V^*(s) = \max_\pi V^\pi(s)$$

Bellman's eqn. for V^* :

$$V^*(s) = R(s) + \max_{\pi(s)} E[V^*(s')]$$

$$\Rightarrow V^*(s) = R(s) + \max_a \gamma \sum_{s'} P_{s,a}(s') \cdot V^*(s')$$

if we take action a .

$$\Rightarrow \pi^*(s) = \arg \max_a \gamma \sum_{s'} P_{s,a}(s') \cdot V^*(s')$$

optimal policy

$$\text{property: } \rightarrow V^*(\delta) = V^{\pi^*}(\delta) \geq V^\pi(\delta)$$

Strategy for finding π^* (optimal policy): →

- 1) find V^*
- 2) use argmax eqn. to find π^* .

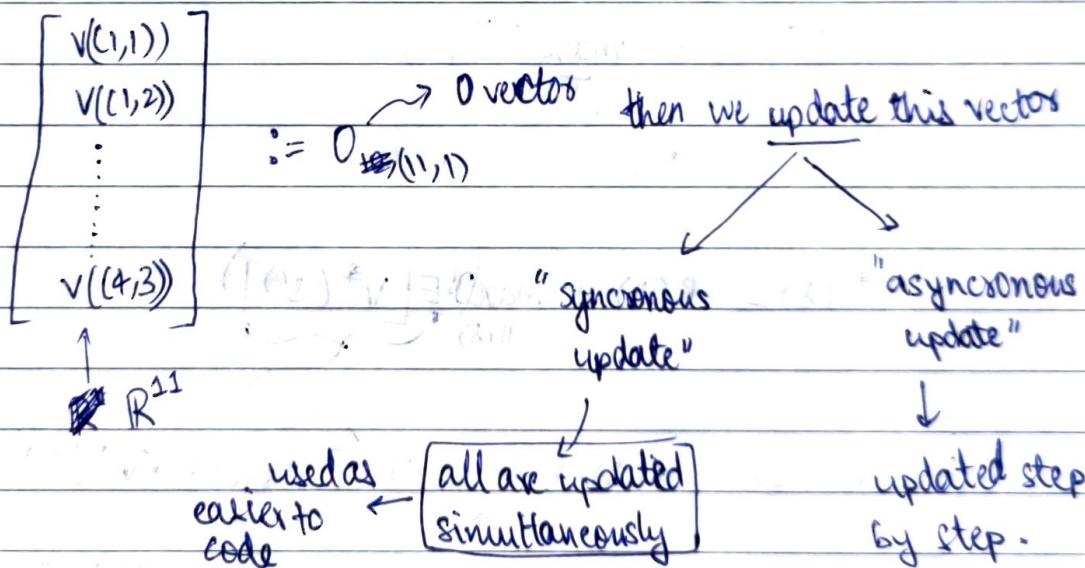
Finding V^* : →

Value iteration: → algo to find V^*

↳ Initialise $V(s) := 0$ for every s .

for every s , update :

$$V(s) := R(s) + \max_a \gamma \sum_{s'} P_{s,a}(s') \cdot V(s')$$

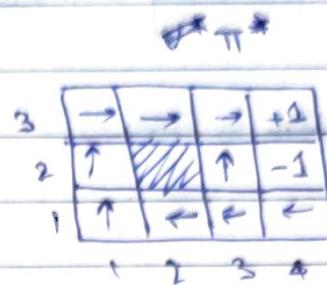
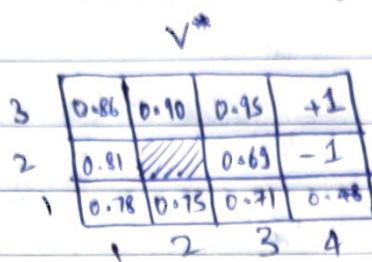


Bellman backup operator:

$$V \leftarrow B(V) \quad \text{as each new vector updated from old one.}$$

how to get π^* from V^* ?

consider the solution from example



$$\text{use eqn: } \pi^*(s) = \underset{a}{\operatorname{argmax}} \sum_{s'} p_{sa}(s') \cdot V^*(s')$$

Suppose we want to find $\pi^*((z_1))$

If we suppose we picked the action w here

$$W \Rightarrow \sum_{s'} P_{s,a}(s') V^*(s') = 0.8 \times 0.75 + 0.1 \times 0.69 + 0.1 \times 0.71$$

west ↓ go south
 north and bounce
 off wall

$= 0.74$

for North it comes $N \Rightarrow 0.676$

so, going west is a better course of action at $((3,1))$.

Policy iteration : \rightarrow

→ Initialise π randomly

computationally expensive if large set of states.

Repeat : \rightarrow Set $V := V^\pi$ (solve bellman's eqns to get V^π) \rightarrow using linear eqn. solving as we saw before-
 till convergence
 set $\pi(s) := \underset{a}{\operatorname{argmax}} \sum_{s'} P_{sa}(s') V(s')$ \rightarrow we update $\pi(s)$ for each state using the values we obtained.

Continuous state MDPs

consider the following example,

suppose we have a car,



6 dimensions.

state $\rightarrow (x, y, \theta, \dot{x}, \dot{y}, \dot{\theta})$

suppose we have an helicopter,



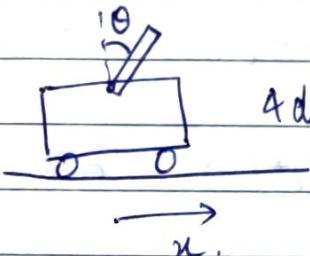
12 dimensional.

roll pitch yaw
state $\rightarrow (x, y, z, \dot{x}, \dot{y}, \dot{z}, \dot{\theta}, \dot{\phi}, \dot{\psi})$

consider case of an inverted pendulum,

we have a car with a rod with full hinge.

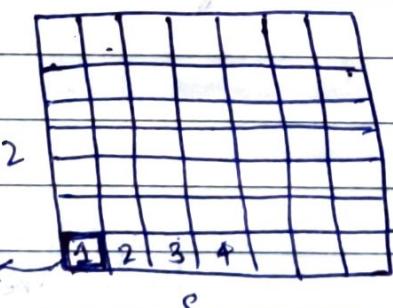
} we want to move car in such a way that rod ~~is balanced~~ is balanced.



4 dimensions $\rightarrow (x, \theta, \dot{x}, \dot{\theta})$

suppose our problem has state space as n dimensional $\rightarrow s \in \mathbb{R}^n$

Discretisation:



① \rightarrow all values in square 1 are taken as one state.

suppose we have a two dimensional state space s_1 and s_2

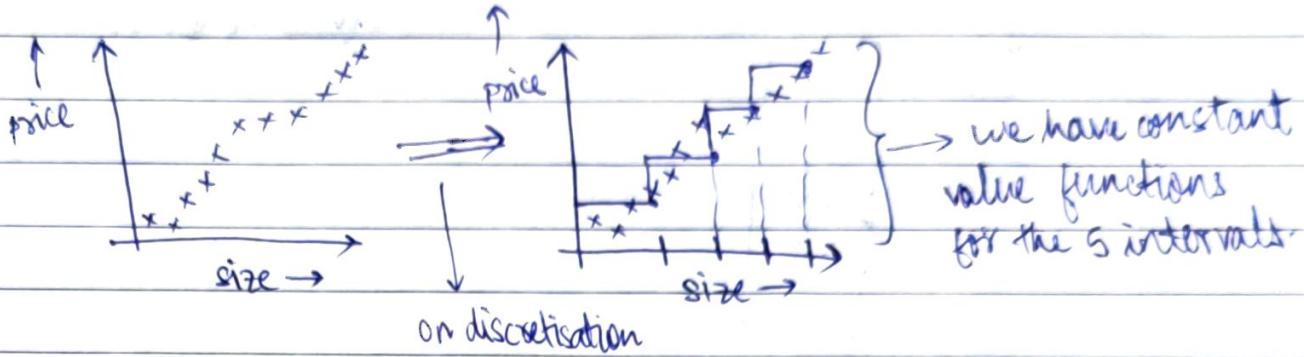
we discretise each state and convert it from continuous to discrete MDP problem.

This is a good method for low dimensional state space.

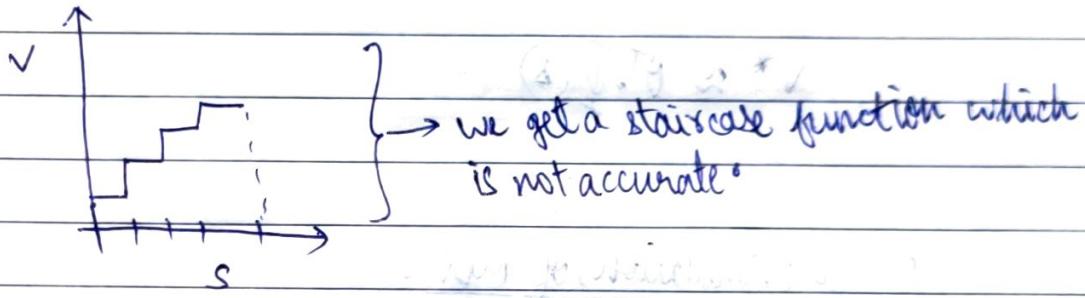
Disadvantages of discretisation :-

→ Naive representation for V^* and π^* .

consider an analogy to predicting housing price based on size.



for RL :-



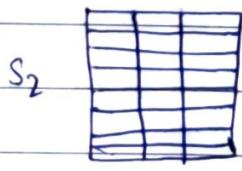
→ Curse of dimensionality

$S \in \mathbb{R}^n$, and we discretize each dimension into k values; then we get k^n discrete states.
↳ problematic for high dimensional state space.

2 to 3 dimensions → discretisation okay

4 to 6 dimensions → discretise carefully.

7 to 8 dimensions → discretisation won't work



→ suppose s_2 has more influence on output so it is discretised finely, but s_1 has larger intervals as we want to minimise # of states.

Approximate V^* directly without discretisation

consider another analogy to linear regression;

$$y \approx \theta^T \cdot \Phi(x)$$

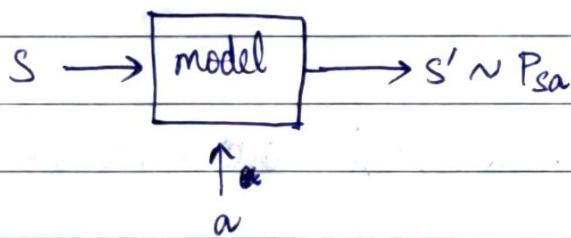
y can be approximated as linear func. of $\Phi(x)$.

$$\Phi(x) = \begin{bmatrix} x_1 \\ x_2 \\ x_1x_2 \\ x_1^2 \\ x_2^2 \\ \vdots \end{bmatrix}$$

so, here for RL we try to approx V^* as linear func. of $\Phi(s)$

$$V^* \approx \theta^T \cdot \Phi(s)$$

Model (simulation) of MDP:



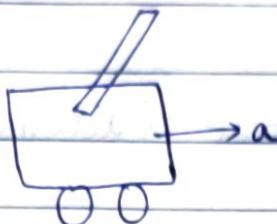
here we are assuming, that we can discretise action space.

In most cases action space has very low dimensions as compared to state space.

Hence, here action space is discretised.

How to get a model?

→ Physical simulation → consider the inverted pendulum problem.



$$S = \begin{pmatrix} x \\ \dot{x} \\ \theta \\ \dot{\theta} \end{pmatrix}$$

$$\dot{S} = \begin{pmatrix} \dot{x} \\ a - L\ddot{\theta} \cos(\theta)/m \\ \dot{\theta} \\ \ddot{\theta} \end{pmatrix}$$

simulator said so.

$$S' = S + (\Delta t) \dot{S}$$

$\Delta t = 0.1$ seconds or so on

→ Learn model from data → consider a helicopter flying problem.

$$S_0^{(1)} \xrightarrow{a_0^{(1)}} S_1^{(1)} \xrightarrow{a_1^{(1)}} S_2^{(1)} \xrightarrow{a_2^{(1)}} S_3^{(1)} \xrightarrow{\dots} S_T^{(1)}$$

basically above example → we took our helicopter out to the field and recorded the states every 10 seconds or so and the action taken to get to that state.

we do above step M times.

$$S_0^{(2)} \xrightarrow{a_0^{(2)}} \dots$$

$$\vdots$$

$$S_0^{(M)} \xrightarrow{a_0^{(M)}} \dots$$

} basically generating data.

Apply supervised learning to estimate S_{t+1} as function of S_t, a_t .
(s') (s, a)

Eg: → Linear regression version:

$$S_{t+1} = A \cdot S_t + B * a_t \quad \left. \right\} \text{okay model for slow speed} \\ \mathbb{R}^{n \times n} \quad \text{helicopter}$$

for above lin. reg model ,

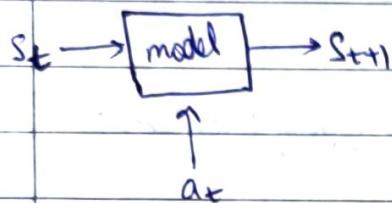
$$\min_{A, B} \sum_{i=1}^M \sum_{t=0}^T \| s_t^{(i)} - (As_t^{(i)} + Ba_t^{(i)}) \|^2$$

suppose if our model is $s_{t+1} = As_t + Ba_t$ (deterministic)

if it was $s_{t+1} = As_t + Ba_t + \epsilon_t$ (stochastic)

$\epsilon_t \sim N(0, \sigma^2 I)$ helps simulate random

noise like wind
and such .



All Above was Model based Reinforcement Learning.

Model free RL \rightarrow we let the robot free and let it crash
and such.

\rightarrow good for videogame RL but not for
real life RL as we don't want crashed.

Fitted value iteration:

choose feature $\Phi(s)$ of state s

$$V(s) = \theta^T \Phi(s)$$

$$\Phi(s) = \begin{bmatrix} x \\ \dot{x} \\ x^2 \\ \vdots \\ x_0 \\ \vdots \end{bmatrix} \left. \right\} \text{example for inverted pendulum.}$$

we can generate lots of $\Phi(s)$ from
our model we obtained
earlier.

Value iteration for discrete states:

$$V(s) := R(s) + \gamma \max_a \sum_{s'} P_{sa}(s') V(s')$$

$$V(s) := R(s) + \gamma \max_a E_{s' \sim P_{sa}} [V(s')]$$

Fitted value iteration: \rightarrow in lin-reg we find $x \rightarrow y$.
 here we find $s \rightarrow V(s)$
 we take some s and try to find $V(s)$ lin-regression

step 1) Sample $\{s^{(1)}, s^{(2)}, \dots, s^{(m)}\} \subseteq S$ randomly.

\rightarrow not same m as modelling one.

step 2) initialise $\theta := 0$

step 3) Repeat {

for deterministic case we can do this as s_{t+1} will take a fixed value. $k=1$ for deterministic

For $i = 1, 2, \dots, m$ {

For each action $a \in A$ {

we take k . s' because

sample $s'_1, s'_2, \dots, s'_k \sim P_{sa}$ of stochastic model

$$\text{set } q(a) = \frac{1}{k} \sum_{j=1}^k [R(s^{(i)}) + \gamma V(s'_j)]$$

$$S_{t+1} = AS_t + BA_q + \epsilon_t$$

we obtain this using

$$\theta^T \Phi(s_j)$$

\downarrow
 $\theta^T \neq \theta$ from prev. iteration

$$\text{set } y^{(i)} = \max_a q(a)$$

because of this we take

k samples as each s'

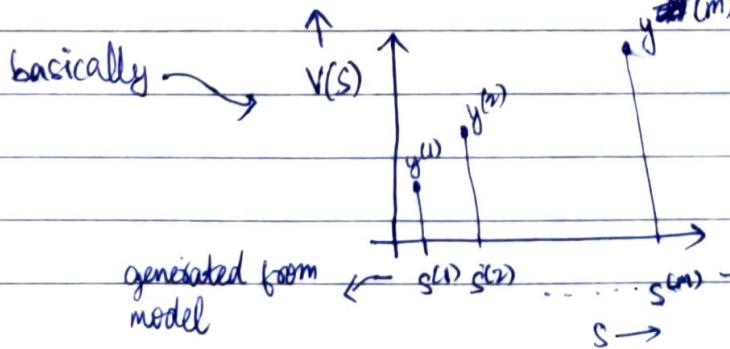
is different because of Gaussian term.

\downarrow
 then we average them.

fit lin-regression

$$y^{(i)} = \max_a E[R(s) + \gamma V(s')]$$

$\theta^T \Phi(s^{(i)}) \approx y^{(i)}$
 we get θ^T from here and we can get
 $V(s)$ for any s using $V(s) = \theta^T \Phi(s)$



(value iteration)

Fitted VI give approximation to V^*

We can get π^* from fitted VI

as,

$$\pi^*(a) = \arg \max_a E_{S' \sim P_{Sa}} [V^*(S')]$$

We can draw k states from for this to approximate state

$s'_1, s'_2, \dots, s'_k \sim P_{Sa}$ to approximate expectation
but this is like a random number generator.

Say our model is of the form

$$s_{t+1} = f(s_t, a_t) + \epsilon_t$$

$$f(s_t, a_t) = As_t + Ba_t$$

so, during runtime we set,

$\epsilon_t = 0$ and $k=1$ \rightarrow we make it deterministic

so, when in state s , pick a s.t.

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

simulation without noise

State action rewards : →

$$R: S \times A \mapsto \mathbb{R}$$

→ reward is a function of S and A

payoffs

$$\text{Total payoff} = R(s_0, a_0) + \gamma R(s_1, a_1) + \gamma^2 R(s_2, a_2) \dots$$

Bellman's eqns : →

$$V^*(s) = \max_a \left(R(s, a) + \gamma \sum_{s'} P_{sa}(s') V^*(s') \right)$$

↓
immediate
reward ↓
future reward

we can do value iteration as we did prev. : → $V(s) = R(s)$ (from above eqn)

$$\pi^*(s) = \arg\max_a \left[R(s, a) + \gamma \sum_{s'} P_{sa}(s') V^*(s') \right]$$

Finite Horizon MDP : →

$$(S, A, \{P_{sa}\}, T, R)$$

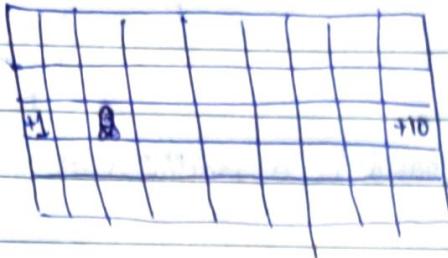
(Horizon time T)
(no γ here)

example can be helicopter
with limited fuel

$s_0, a_0 \rightarrow s_1, a_1 \rightarrow \dots \rightarrow s_T, a_T$. only T transitions
happen

$$\text{Total payoff} = R(s_0, a_0) + R(s_1, a_1) + \dots + R(s_T, a_T)$$

goal is to maximize $E[R(s_0, a_0) + R(s_1, a_1) + \dots + R(s_T, a_T)]$



action of robot here depends on horizon time.

If there are 10-12 steps it'll go to +10 ← if there are 2-3 steps left it'll go to +1.



So, the policy here depends on Horizon time T

$\pi^*(s) \rightarrow$ is a function of t here



$\pi_t^*(s) \rightarrow$ non-stationary policy → as it changes over time

Non stationary state transitions:

$$S_{t+1} \sim P_{state}^{(t)}$$

the state transitions here are a function of time here

example: → driving car from A to B → depends on time i.e. rush hour and so on.

Non-stationary rewards:

$R_{s,a}^{(t)} \rightarrow$ reward is a function of time

MDPs

when to use non-stationary models?

- in airplanes, 1/3rd weight is fuel, so it is changing dynamic → lighter when landing compared to takeoff
- weather forecast → depending on whether it'll rain and all.
- Industrial automation → It is easier to get labour at 3pm than 3am → so this too is non stationary.

$$V_t^*(s) = E \left[R(s_t, a_t) + R(s_{t+1}, a_{t+1}) + \dots + R(s_T, a_T) \mid \pi^*, s_0 = s \right]$$

Expected total payoff starting in state s_t at time t executing π^*

Value iteration: (Dynamic programming)

$$V_t^*(s) = \max_a \left[R(s, a) + \sum_{s'} P_{sa}(s') V_{t+1}^*(s') \right]$$

$$\pi^*(s) = \operatorname{argmax}_a (\quad \text{"} \quad)$$

defines V_t as function of V_{t+1}

at time T final V^* is V_T^*

$$V_T^*(s) = \max_a (R(s, a)) \rightarrow \text{as we can't take further steps after time } T$$

so, we calculate V_T^* first, then we find

$$V_{T-1}^*, V_{T-2}^*, \dots, V_0^*$$

then we find π_t^* using V_t^* where $0 \leq t \leq T$

Linear Quadratic Regulation (LQR) :-

condn. for LQR :- finite horizon

$$(S, A, \{P_{sa}\}, T, R)$$

$$w_t \sim N(0, \Sigma_w)$$

$$S = \mathbb{R}^n \quad A = \mathbb{R}^{d \times d} \quad P_{sa} : S_{t+1} = AS_t + Ba_t + w_t$$

$$R(s, a) = -(s^T U s + a^T V a) \quad \text{where } U \in \mathbb{R}^{n \times n}, V \in \mathbb{R}^{d \times d}$$

$U, V \geq 0$ (positive semi-definite)

$$s^T U s \geq 0, a^T V a \geq 0$$

If we want $s \approx 0 \rightarrow$ i.e. want our helicopter to hover

and if we choose $U = I, V = I$.

$R(s, a) = -(||s||^2 + ||a||^2)$ \rightarrow basically penalise state s and action a from deviating
penalises large \leftarrow from zero.
jucky motions

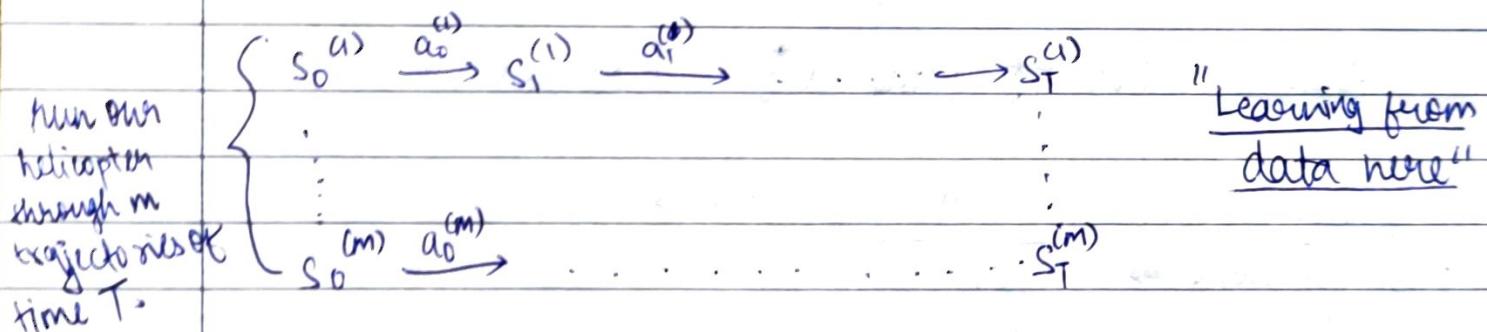
if we want non-stationary transitions then;

$$A, B \xrightarrow[\text{with}]{\text{replace}} A^{(t)}, B^{(t)}$$

if we want non stationary Rewards;

$$U, V \xrightarrow[\text{with}]{\text{replace}} U^{(t)}, V^{(t)}$$

Where to get A, B ?



$$s_{t+1} \approx As_t + Bs_a$$

$$\min_{A, B} \sum_{i=1}^m \sum_{t=0}^{T-1} \|s_{t+1} - (As_t + Bs_a)\|^2$$

↓

Like linear regression

Second method to get A, B \rightarrow linearize a non linear model :-

consider our case of inverted pendulum :-

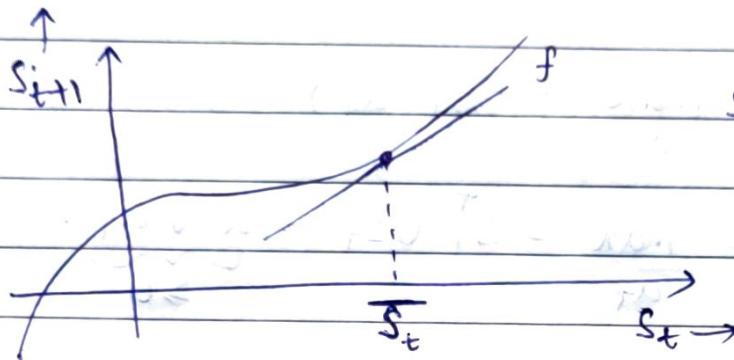


$$s_{t+1} = f(s_t, a_t)$$

↓

^{suppose}
our physics friend helped
with this

$$\begin{pmatrix} x_{t+1} \\ \dot{x}_{t+1} \\ \theta_{t+1} \\ \dot{\theta}_{t+1} \end{pmatrix} = f\left(\begin{pmatrix} x_t \\ \dot{x}_t \\ \theta_t \\ \dot{\theta}_t \end{pmatrix}, a_t\right)$$



$$s_{t+1} \approx f'(\bar{s}_t) \cdot (s_t - \bar{s}_t) + f(\bar{s}_t)$$

\bar{s}_t here is constant.

Suppose for our inv. pendulum, we assume bar almost upright and low speeds, \bar{s}_t will be zero. \leftarrow suppose we want our helicopter to hover, so \bar{s}_t will be 0 here \rightarrow "typical value" as s_{t+1} will be close to \bar{s}_t .

$$\text{in general, } s_{t+1} \approx f(\bar{s}_t, \bar{a}_t) + (\nabla_s f(\bar{s}_t, \bar{a}_t))^T (s_t - \bar{s}_t) + (\nabla_a f(\bar{s}_t, \bar{a}_t))^T (a_t - \bar{a}_t)$$

$$s_{t+1} = A s_t + B a_t$$

we may need to add an intercept term to $s_t \rightarrow \begin{pmatrix} 1 \\ x \\ \dot{x} \\ \theta \\ \dot{\theta} \end{pmatrix}$ because of constants.

problem we want to solve:

$$S_{t+1} = A S_t + B a_t + \omega_t.$$

$$R(s_t) = -(s^T U s + a^T V a)$$

$$N(0, \Sigma_w)$$

noise

} if these condns. are met
we can compute exact
 v^*

$$\text{total payoff} = R(s_0, a_0) + R(s_1, a_1) + \dots + R(s_T, a_T)$$

Dynamic programming solution to find v^*

$$v_T^*(s_T) = \max_{a_T} R(s_T, a_T)$$

$$= \max_{a_T} -s_T^T U s_T - \underbrace{a_T^T V a_T}_{\geq 0}$$

so max. \rightarrow we choose

$$= -s_T^T U s_T$$

$$a_T = 0$$

$$\pi_T^*(s_T) = \vec{0}$$

Inductive step: \rightarrow

$$\text{suppose, } v_{t+1}^*(s_{t+1}) = s_{t+1}^T \Phi_{t+1} s_{t+1} + \psi_{t+1}$$

where, $\Phi_{t+1} \in \mathbb{R}^{n \times n}$, $\psi_{t+1} \in \mathbb{R}$

$$v_t^*(s_t) = \max_{a_t} \left(R(s_t, a_t) + E_{s_{t+1} \sim P_{s_t, a_t}} [v_{t+1}^*(s_{t+1})] \right)$$

this format as state transition
is continuous.

$$V_t^*(s_t) = \max_{a_t} \left(-s_t^T U s_t - a_t^T V a_t + \underbrace{\left[\underbrace{s_{t+1}^T \Phi_{t+1} s_{t+1}}_{s_{t+1} \sim N(A s_t + B a_t, \Sigma_w)} + \Psi_{t+1} \right]}_{\text{comes out to be a big quadratic}} \right)$$

take derivatives ← func. of a_t
and solve for a_t

$$a_t = \underbrace{(B^T \Phi_{t+1} B - V)^{-1} B^T \Phi_{t+1} A \circ s_t}_{L_t}$$

$$\pi_t^*(s_t) = L_t \cdot s_t$$

∴ optimal action is a linear function of s_t . and

$$V_t^*(s_t) = s_t^T \Phi_t s_t + \Psi_t$$

where,

$$\Phi_t = A \left(\Phi_{t+1} - \Phi_{t+1} B (B^T \Phi_{t+1} B - V)^{-1} B^T \Phi_{t+1} \right) A^T - U$$

$$\Psi_t = -\lambda \left(\Sigma_w \Phi_{t+1} \right) + \Psi_{t+1}$$

To summarize,

$$\text{initialize } \Phi_T = -U \text{ and } \Psi_T = 0$$

recursively calculate, Φ_t, Ψ_t using Φ_{t+1}, Ψ_{t+1} (for $t = T-1, T-2, \dots, 0$)

calculate L_t

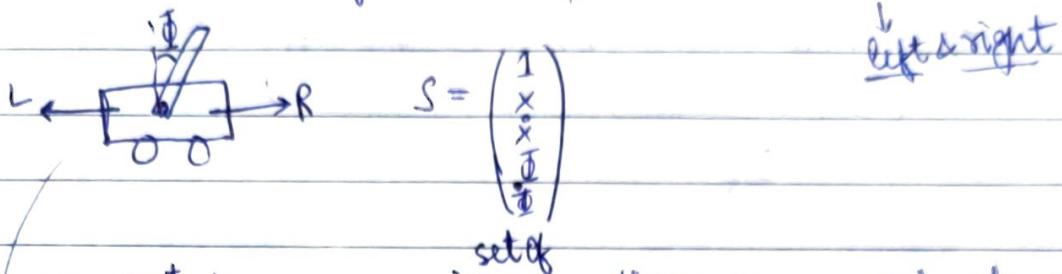
$$\pi_t^*(s_t) = L_t \cdot s_t$$

Policy search & POMDPs

("direct policy search") \rightarrow try to find good π directly.

normally in RL we first find V^* then we find π^*

consider the case of inverted pendulum \rightarrow with two actions



we want to come up with functions to approximate the policy-set of

$$\pi_\theta(s) = \frac{1}{1 + e^{-\theta^T s}}$$

say/assume like how we did for logistic func.

Stochastic policy \rightarrow function $\pi: S \times A \mapsto \mathbb{R}$ where $\pi(s, a)$ is the prob. of taking action a in state s .

$$\left(\sum_a \pi(s, a) = 1 \right)$$

basically on each timestamp give prob. that you want to take action a .

for this example consider something like

$$\pi_\theta(s, "R") = \frac{1}{1 + e^{-\theta^T s}}$$

$$\pi_\theta(s, "L") = 1 - \frac{1}{1 + e^{-\theta^T s}}$$

this is a reasonable policy as consider $\theta = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$, so $\theta^T s = 1$

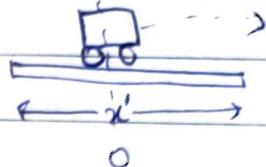
if $1 >> 0 \rightarrow$ it moves more aggressively to right

if $1 << 0 \rightarrow$ move " " " left.

however, this isn't a good policy as it ignores features other than θ ,

$$\text{if, } \theta = \begin{bmatrix} 0 \\ -0.5 \\ 0 \\ 1 \\ 0 \end{bmatrix} \quad \theta^T s = [-0.5x] + 1 \rightarrow \text{helps stay cast in an interval}$$

if position move to right
then more prob. of going left



goal: → Find θ so that we execute $\pi_\theta(s, a)$, we maxim.

$$\max_{\theta} E[R(s_0, a_0) + \dots + R(s_T, a_T) | \pi_\theta] \quad \text{or we assume initial state has a fixed distribution}$$

here s_0 is fixed initial state. → assumption for direct policy search

↳ in general RL ~~assume~~ s_0 is not a fixed state → we see for all states

$$\max_{\theta} E[R(s_0, a_0) + R(s_1, a_1) + \dots + R(s_T, a_T) | \pi_\theta]$$

↳ using $T=1$ to simplify the math and not write such a long summation
→ all derived results work for longer summation too

$$= \sum_{s_0, a_0, s_1, a_1} P(s_0, a_0, s_1, a_1) [R(s_0, a_0) + R(s_1, a_1)]$$

payoff

all possible state and action sequence.

initial state dist. state transition prob.

$$= \sum_{s_0, a_0, s_1, a_1} P(s_0) \pi(s_0, a_0) P(s_1) \pi(s_1, a_1) \cdot [\text{payoff}]$$

↑ to s_1

↓ ↓

prob. of choosing action a_0 prob. of picking action a_1

We will use gradient ascent as func. of θ to maxim. above funct.

Reinforce algorithm to maxim. total payoff: →

Loop: {

→ random state seq; generated

Sample: $s_0, a_0, s_1, a_1, \dots, s_T, a_T \rightarrow$ run it for T timestamps

$$\text{compute payoff} = R(s_0) + R(s_1) + \dots + R(s_T)$$

$$\text{update: } \theta := \theta + \alpha \left[\frac{\nabla_{\theta} \pi_{\theta}(s_0, a_0)}{\pi_{\theta}(s_0, a_0)} + \frac{\nabla_{\theta} \pi_{\theta}(s_1, a_1)}{\pi_{\theta}(s_1, a_1)} + \dots \right] \times \text{payoff}$$

payoff of the random seq used

multiplied by payoff

this is stochastic gradient ascent

i.e. bounces all over the place while on an avg. converging

← because of the random sequence used.

on an avg. in same dirn. as actual gradient.

$$\nabla_{\theta} \sum_{s_0, a_0, s_1, a_1} p(s_0) \underline{\pi_{\theta}(s_0, a_0)} \underline{p_{s_0, a_0}(s_1)} \underline{\pi_{\theta}(s_1, a_1)} \cdot [\text{payoff}]$$

$$= \sum_{s_0, a_0, s_1, a_1} p(s_0) \cdot \underline{\pi_{\theta}(s_0, a_0)} \underline{p_{s_0, a_0}(s_1)} \underline{\pi(s_1, a_1)} \left[\frac{\nabla_{\theta}(\pi_{\theta}(s_0, a_0))}{\pi_{\theta}(s_0, a_0)} + \frac{\nabla_{\theta}(\pi_{\theta}(s_1, a_1))}{\pi_{\theta}(s_1, a_1)} \right] \times [\text{payoff}]$$

$$= \sum_{s_0, a_0, s_1, a_1} p(s_0, a_0, s_1, a_1) \left[\frac{\nabla_{\theta}(\pi_{\theta}(s_0, a_0))}{\pi_{\theta}(s_0, a_0)} + \frac{\nabla_{\theta}(\pi_{\theta}(s_1, a_1))}{\pi_{\theta}(s_1, a_1)} \right] \times [\text{payoff}]$$

$$= \sum_{s_0, a_0, s_1, a_1} p(s_0, a_0, s_1, a_1) [\text{gradient update}]$$

direct policy search is better than finding V^* then π^* in 2 settings: →

i)

POMDP

↪ partially observable MDP



at each step, we get a partial (and

potentially noisy) measurements of

the state. Have to choose action according to POMDP

so it's to these incomplete measurements.

ex: →



$$s = \begin{pmatrix} x \\ y \\ \theta \end{pmatrix}$$

→



but suppose we can only measure x , θ

ii)

is π^* simpler or is V^* simpler.

low level control task → can a skilled human do the task.

→ flying a helicopter

→ balancing inverted pendulum

playing chess → not low level → have simpler π^* so better to use this.