

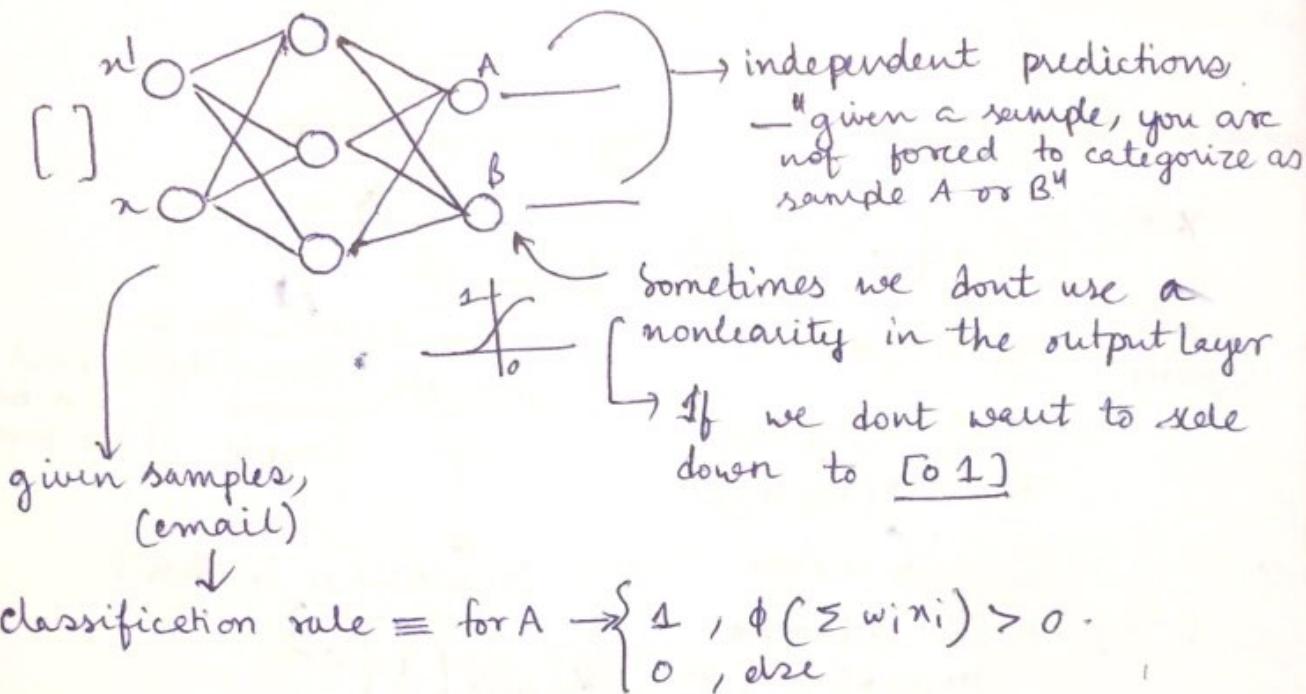
partial derivatives of output with respect to parameters
should be computable.

and
(input)

$$\frac{\partial E}{\partial w_1} = \frac{\partial E}{\partial x_3} \cdot \frac{\partial x_3}{\partial x_2} \cdot \frac{\partial x_2}{\partial x_1}$$

"Chain rule"

then we can apply chain rule to compute E w.r.t any of the parameters.



Suppose, the email gets classified as ~~person~~ by both Spam.

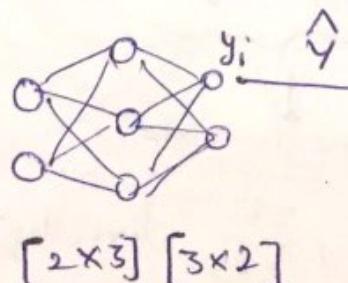
use "softmax"
across neurons

Ex: $\begin{matrix} \text{value} \\ A - 0.49 \\ B - 0.001 \end{matrix}$

* Learn w?

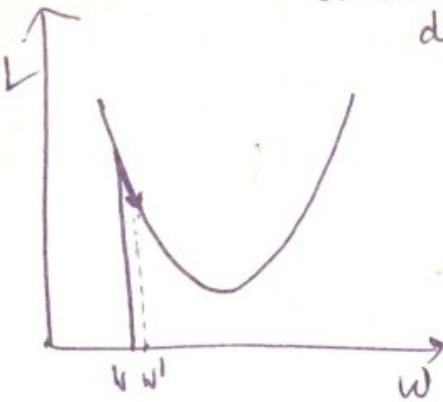
$(x_i, y_i), i=1, N$
Example set

→ Find optimal w?



\equiv Error $= [\hat{y} - y_i]$

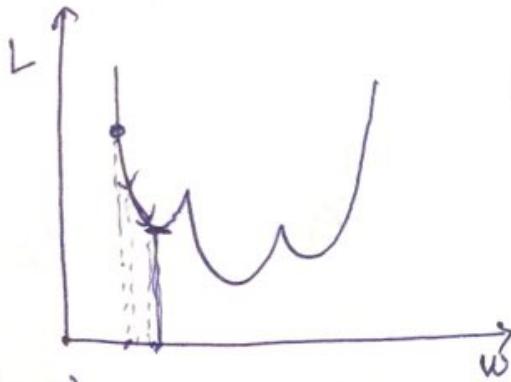
minimise the error
— adjust the weights



Classic gradient descent

$$w^{k+1} \leftarrow w^k - \eta \frac{\partial L}{\partial w}$$

gets modified to

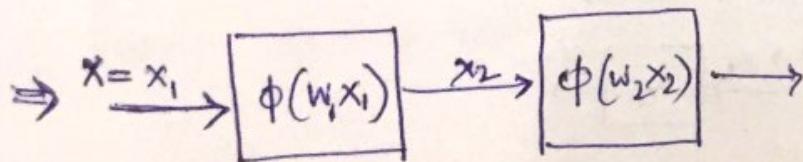


(Error)

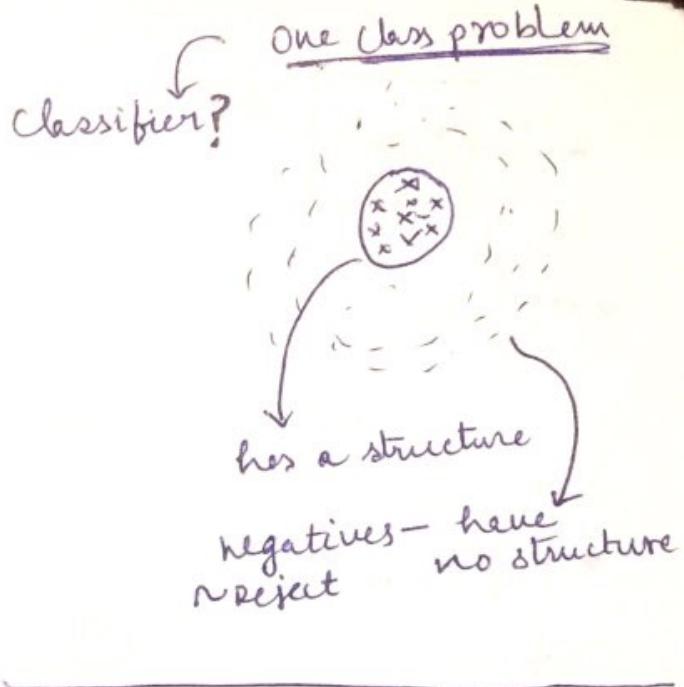
→ E_1 in the output layer indirectly depends on the first layer weights

→ adjust the weights and go to better solution in steps (less error)

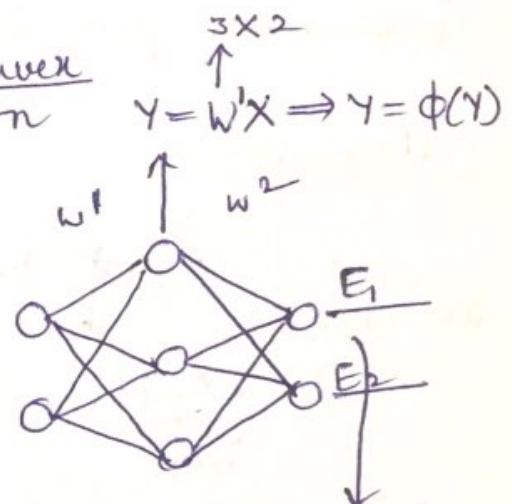
Given input $\Rightarrow (x, y)$



Neural Network

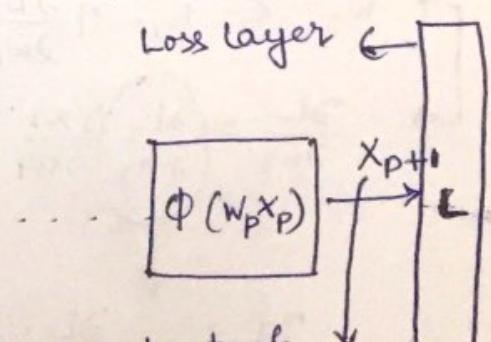


Not a convex optimisation problem

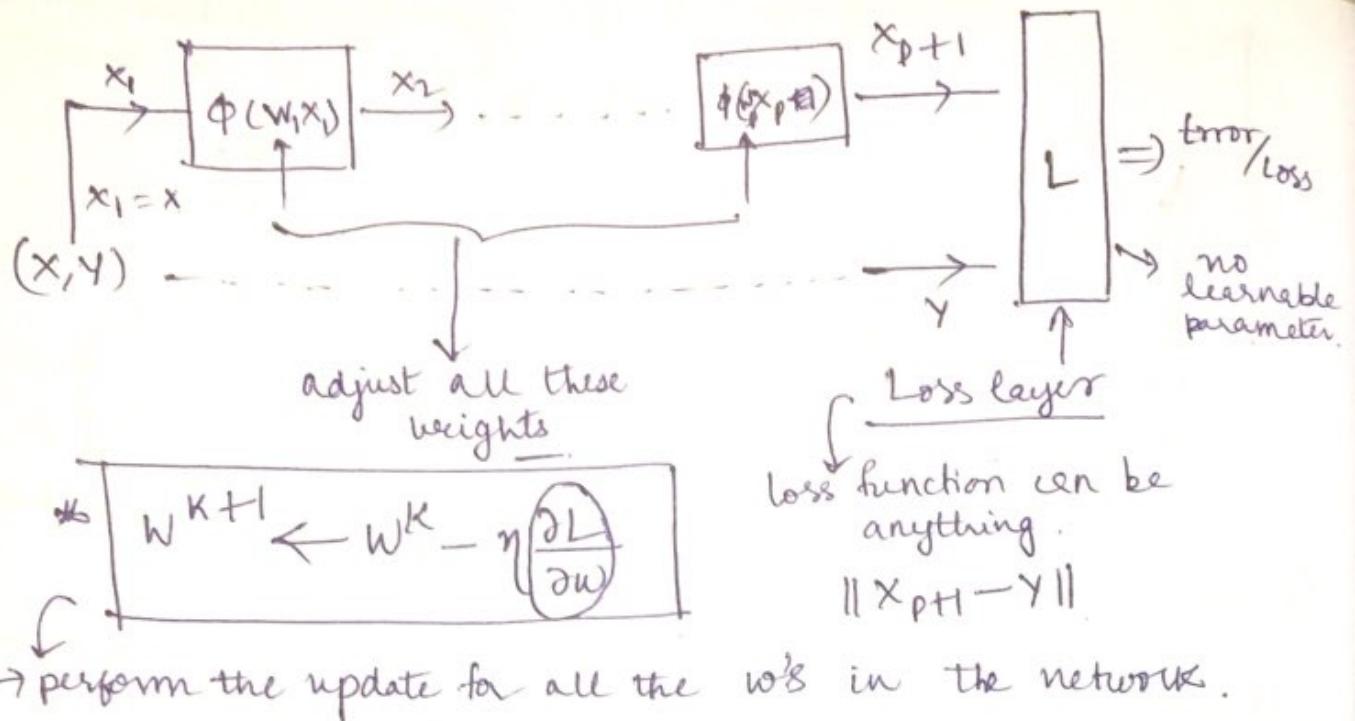


Error here

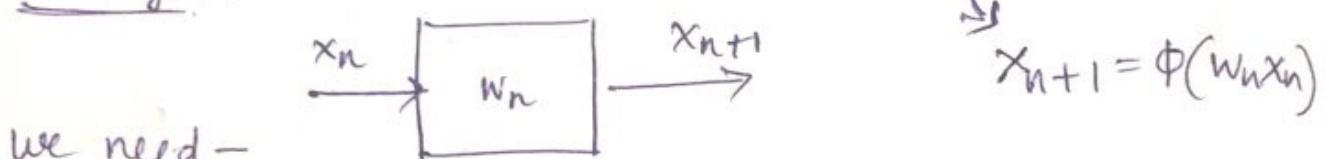
doesn't solely depend on w² weights but also on the outputs of layer 1



Output of the network.



one layer:-



we need -

$$1.) \frac{\partial x_{n+1}}{\partial x_n}$$

$$2.) \frac{\partial x_{n+1}}{\partial w_n}$$

* if we want to update w_1 , we will need

$$\left[\frac{\partial L}{\partial w_1} \right]$$

$$w_1^{K+1} \leftarrow w_1^K - \eta \frac{\partial L}{\partial w_1}$$

$$\frac{\partial L}{\partial w_1} = \left(\frac{\partial L}{\partial x_2} \right) \frac{\partial x_2}{\partial w_1} \quad (\text{from the block})$$

$$\left[\frac{\partial L}{\partial x_{p+1}} \right] \leftarrow \text{can be computed directly}$$

$$\frac{\partial L}{\partial x_2} = \frac{\partial L}{\partial x_{p+1}} \cdot \underbrace{\frac{\partial x_{p+1}}{\partial x_p} \cdot \frac{\partial x_p}{\partial x_{p-1}} \cdots \frac{\partial x_3}{\partial x_2}}_{\text{diff output wrt input } \rightarrow \text{from the block.}}$$

Computing loss =

$$L = \sum_{i=1}^N [x_{p+1}^i - y_i]^T [x_{p+1}^i - y_i]$$

↓
sum of squared differences of the predicted and actual output.

To update a weight, we need to know the derivatives of the blocks before.

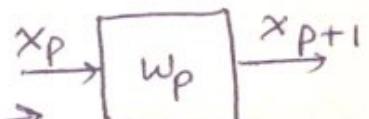
Ex:

$$w_p^{k+1} \leftarrow w_p^k - n \frac{\partial L}{\partial w_p}$$

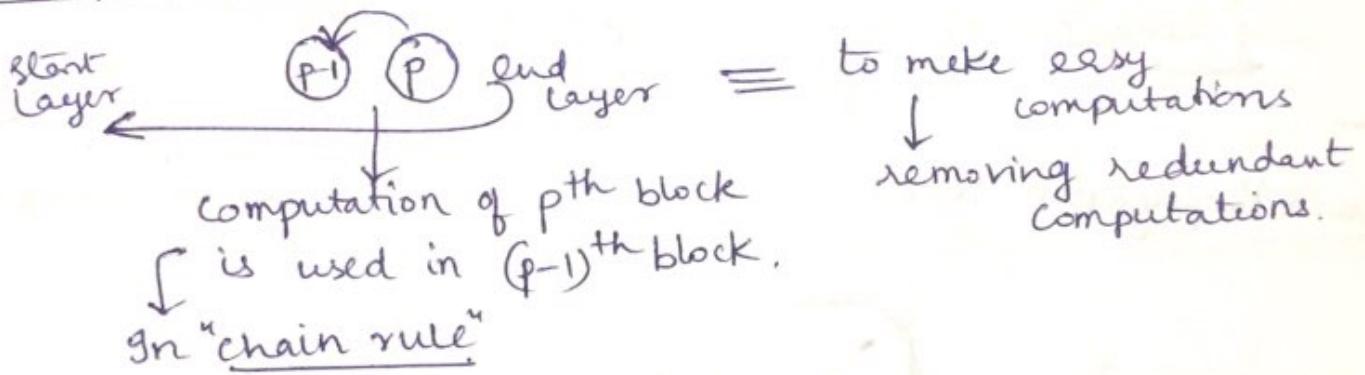
$$\frac{\partial L}{\partial w_p} = \frac{\partial L}{\partial x_{p+1}} \cdot \frac{\partial x_{p+1}}{\partial w_p}$$

can be
computed
directly.

last block



Back propagation — order $\equiv \{ \text{back to front} \}$



Pseudo code \equiv

- >1. Initialise $w_1 \dots w_p$ randomly
- >2. Compute loss for all N samples. (Forward propagation)
- >3. Update w_i ;

$$w_i^{k+1} \leftarrow w_i^k - n \frac{\partial L}{\partial w_i} \quad \begin{cases} i = p, \dots, 1 \\ \text{(or)} \\ i = 1, \dots, p \end{cases} \quad \begin{array}{l} \text{(backward pass),} \\ [\text{wastage of computation}] \times \end{array}$$

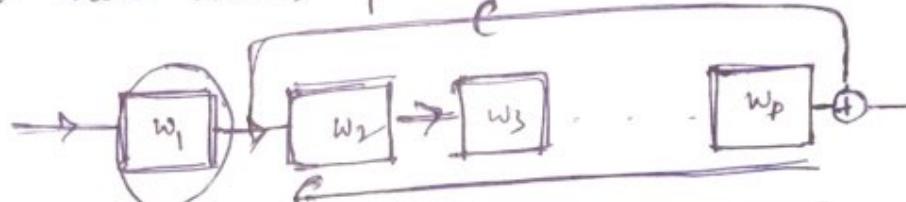
product of partial derivatives.

- >4. Repeat step 2-4 until $\underbrace{\text{convergence}}_{\text{No significant change in loss/weight}}$

- * Vanishing gradient problem } → as you propagate, the gradient might diminish
- * Exploding gradient problem

for these issues; provide skip connections

skip connections
↓
Not dynamic architecture.



both
changing this $w_1 \rightarrow$ depends on the channels

If one of the gradients in chain rule becomes zero
very dangerous → because output doesn't depend on the input and also the previous blocks
might be due to convergence.

Back Propagation Tricks:

$$\frac{1}{1+e^{-ax}} = \text{curve changes w.r.t. } a.$$

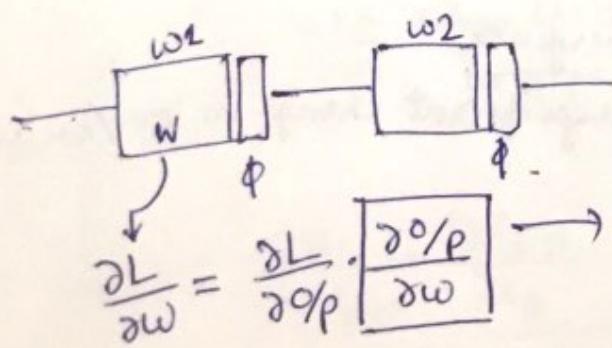
we can have -

learnable parameters in the non-linearity

J/L → Loss function.

$$\text{if } \frac{\partial L}{\partial \phi} \text{ is available } \rightarrow \text{then, } \frac{\partial L}{\partial w} = \frac{\partial L}{\partial \phi} \cdot \boxed{\frac{\partial \phi}{\partial w}}$$

↓
should be known.



after computing this,
update rule =

$$w^{k+1} \leftarrow w^k - \eta \frac{\partial L}{\partial w}$$

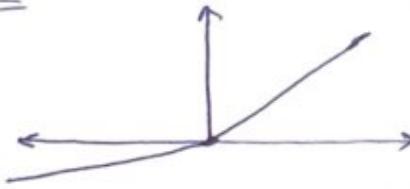
sigmoid and tanh activation functions

↓
the derivatives of such functions can be expressed in terms of the original function.

New activation function = ReLU \Rightarrow continuity?
derivative is easy to compute.
SVM Loss function \rightarrow [enforcing margin.] using "subgradients"
hinge loss. derivate evaluated at x
 $\phi'(x)$.

Max (0, 1 - $w^T x$). \sim similar to ReLU.
loss can be computed at any point

Leaky ReLU =



Stochastic

$\frac{\partial L}{\partial w}$ → loss on the samples that we have access to.

$\frac{\partial L}{\partial w}$ → so the estimate can be calculated over a subset

Assume: out of billion samples → take random 100

so we look at a subset $\left[\begin{array}{l} \text{then there is no point in considering loss over all billion samples.} \\ \text{if the loss is same for any random 100} \end{array} \right]$

$$\boxed{\frac{\partial L}{\partial w}}$$

"over a subset"

"batch" - extremely reliable but slow

"single sample" - not reliable

"minibatch"

$\frac{\partial L}{\partial w} \rightarrow$ computing derivative for every sample and adding

$\hookrightarrow \frac{\partial}{\partial w} \sum_i (\cdot) \rightarrow \sum_i \frac{\partial (\cdot)}{\partial w}$ \approx summation of each derivative.

that's why we look
at mini batch

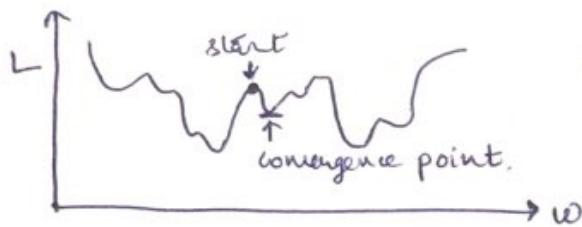
'batch'
↑
computationally
not attractive.

Algorithm =

- ① Initialise
- ② Forward pass and compute "Loss"
- ③ Back propagate and update "w"
- ④ Repeat 2-4 until [convergence.]

Problem = Non convex in
loss function
↓ nature

MSE, $L \geq 0$
 $L^0, L^1, L^2, L^3, \dots, L^i, L^{i+1}$
moving across negative
gradient
 $L^i > L^{i+1}$



1.

Initialisation - "don't know where to start?"

\hookrightarrow pre training - similar but different problem

↓ can give a good hint of the starting
 \hookrightarrow picking "w" of appropriate sizes. point.

\hookrightarrow Random Initialisation -



— network can have many inputs / outputs
"depends on number of inputs / outputs"

\hookrightarrow Xavier's Initialisation (heuristics)
He's Initialisation.

Update Rule :-

$$\theta^{k+1} \leftarrow \theta^k - \eta \frac{\partial L}{\partial \theta} \quad (\text{gradient based update})$$

η (eta) can be controlled.

Momentum

"idea": ball on the surface, moving in the same direction = pickup velocity.



↓ Previous change is added to the update rule.

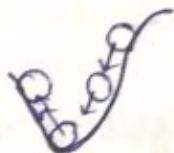
$$\theta^{k+1} \leftarrow \theta^k - \eta \frac{\partial L}{\partial \theta} + \gamma \underline{\theta^k} \quad (\text{Previous})$$

→ if previous direction of motion is same, we add them else, they cancel each other and become slow.

⇒ Reinforces the continuity

Nesterov Momentum — look ahead prediction of gradient.

↓ issue: skips the minima and climbs on.



Adaptive Gradient Change

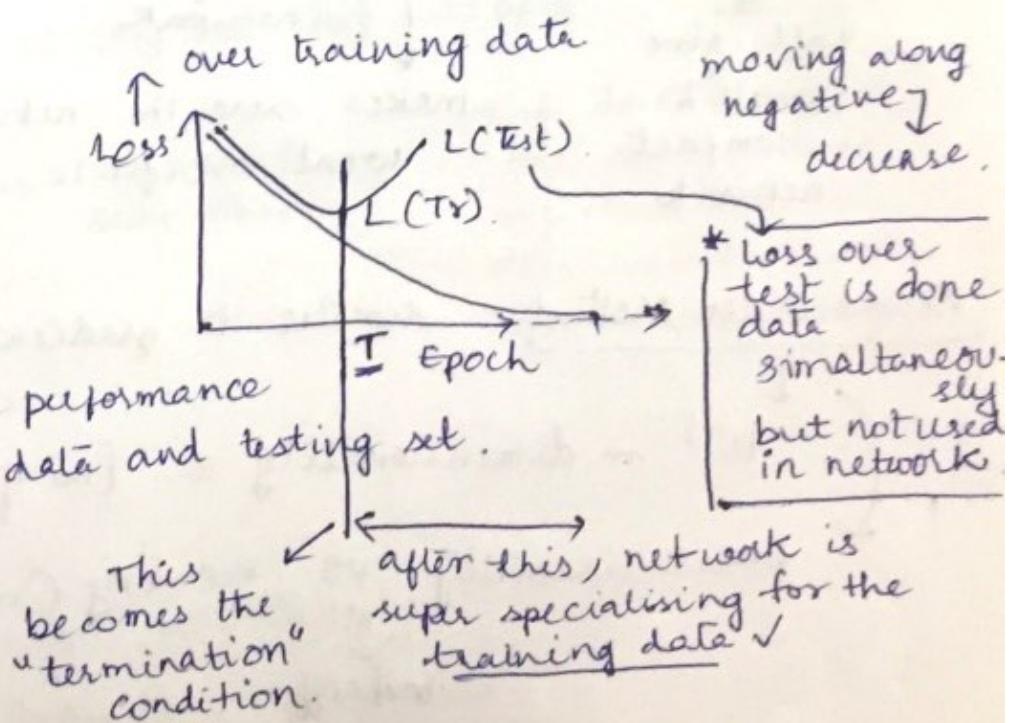
- Ada delta
- Ada grad
- RMS delta
- Adam.

Termination :-

— overfitting

* should give same performance on both training data and testing set.

$L(\text{test})$
need not be smooth.



Overfitting \rightarrow neural network - is "overparametrised".

$L(\text{Test}) \rightarrow$ not a minimum, because it is not a smooth curve.
↓
has no distinct property (might have)
zig-zags.

has extra weights.

to counter "overparametrised" \Rightarrow

add a norm of weight $J' = J + \|\theta\|$ learnable parameters.

Small network K - cannot by heart, so we reduce the parameters

$$\frac{\partial J'}{\partial \theta_i} = \frac{\partial J}{\partial \theta_i} + \frac{\partial \|\theta\|}{\partial \theta_i}$$

$$J' = J + \|\theta\|$$

L0 (nonzero) L1 L2 \Rightarrow norms.

not easily differentiable.

\rightarrow weight decay - or Regularisation.

will give small & compact network.

remove small θ s in the network

pruning the network

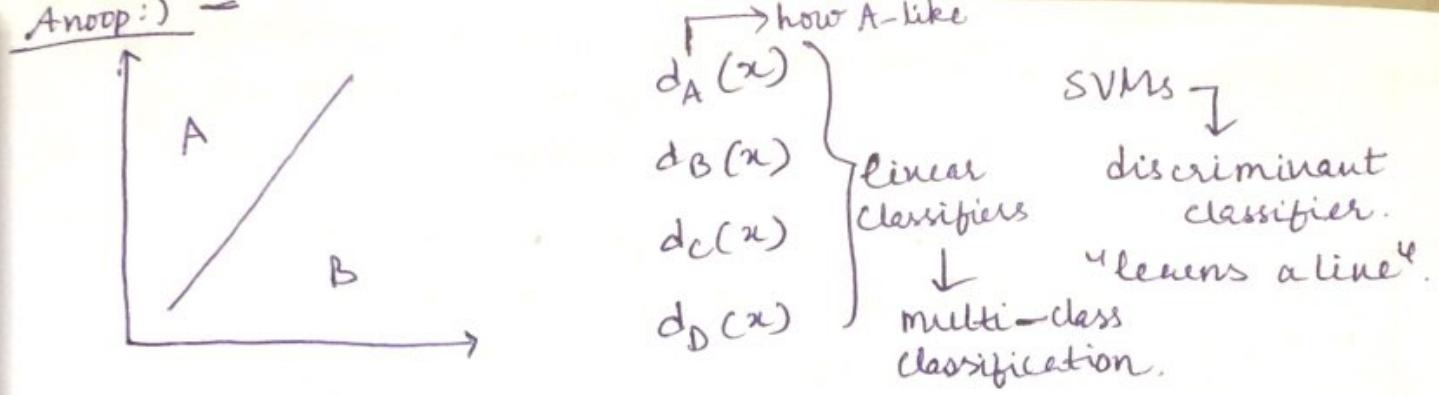
makes sure the network wont overfit.

* if w is very small, it can be removed from the network.
Not much change.

Second Order Method - similar to gradient descent.

H^{-1} - dimensionality $\equiv [\text{no. of parameters}]^2$

parallelisability vs memory (very high)
↓
tradeoff

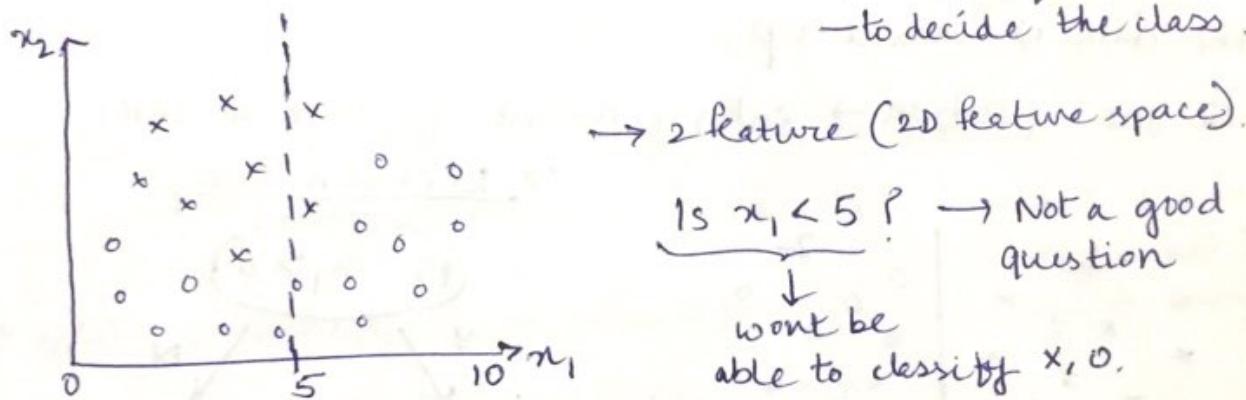


multi layer perception - generative model.

decision trees (~ 20 questions game).

for cat/non-cat classification [length of the tail $\leq 5\text{cm}?$]

answers of these questions
→ to decide the class.



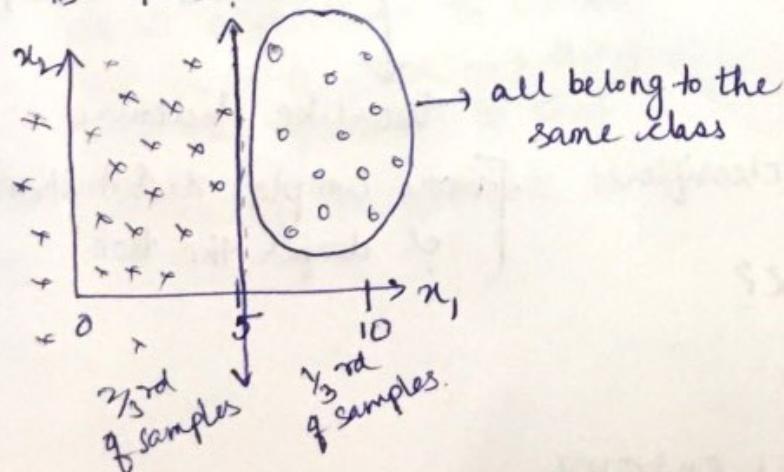
[Question = answer] → reduction in space of set of samples.

(best question).

Is $x_1 \geq 5$? → not exactly reducing into half

"optimal"

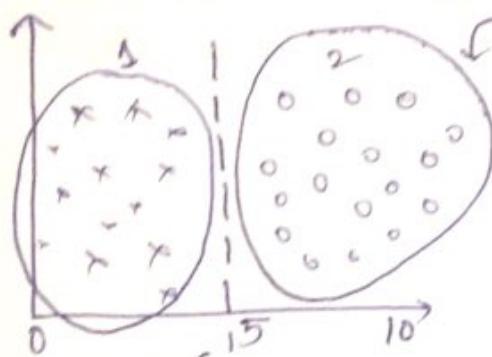
→ if due to a split,
one side contains
one class and the
other contains the other
class - then done!



Entropy - randomness / uncertainty.

$$H(x) = - \sum_i P_i \log_2 P_i \quad \rightarrow \text{class.}$$

high entropy - high uncertainty - when a sample is picked from the whole set of samples.



when a sample is taken from a side, entropy is zero, as all the samples on one side belong to the same class

on side-2 =

$$H = -P_1 \log_2 P_1 + -P_2 \log_2 P_2 = 0 \text{ (Entropy)}$$

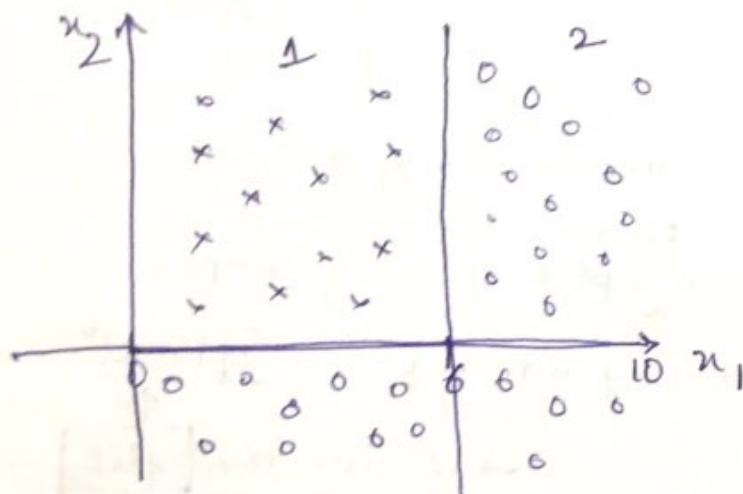
0×0 1×0

these type of classifiers are called "decision trees".

* If the data is mixed up,

→ with one question → entropy doesn't go down to zero.

⇒ Decision Tree.

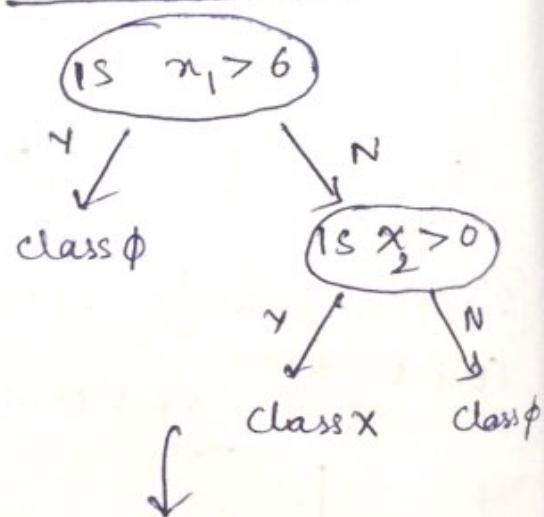


series of questions n classifier

- * what questions to ask?
- * when to stop?

Compute total entropy

n choose questions such that it reduces the entropy



tree-like structure

[more complex distribution]
* deeper the tree

Mathematically, entropy cannot go higher.

half-half samples (x, o) \Rightarrow

$$\rightarrow \text{Entropy, } H = -[0.5 \times -1 + 0.5 \times -1] \quad (\text{Total Entropy})$$
$$= 1 \quad (\text{before splitting})$$

\rightarrow after splitting, ($2/3^{\text{rd}} x, 1/3^{\text{rd}} o$)

$$\begin{cases} H_0 = -1 \log 1 + 0 \log 0 \\ = 0. \end{cases}$$

on one side

$$\begin{cases} H_x = -2/3 \log 2/3 - 1/3 \log 1/3 \\ = 0.72 \end{cases}$$

on the other side

$$H_{\text{total}} = 0.72$$

Entropy reduced from $1 \rightarrow 0.72$. (with just one question).

\rightarrow How to decide which question to ask?

\downarrow * how many questions? (infinite). \rightarrow but answers might not be infinite. ~~same~~.

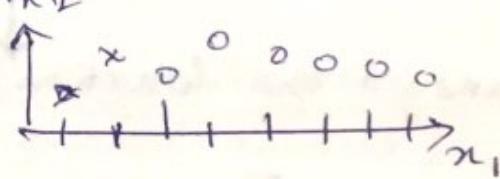
\rightarrow splits possible? (Not infinite)

— all possible values of x ,

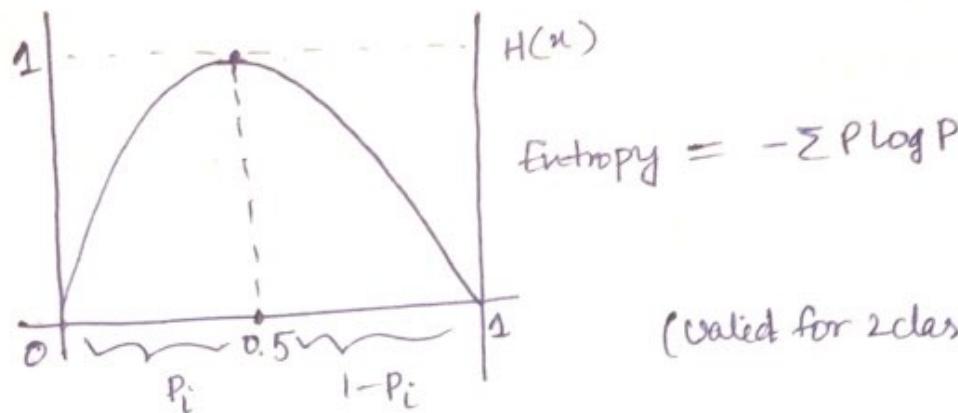
if n samples in the database, $n-1$ splits

* decision trees are the worst algorithms to train but once trained \rightarrow they are the fastest to test

* [extremely inefficient to train, extremely efficient to test]



H vs P_i : (change in entropy function w.r.t P_i).



Impurity measure [if pure, low value; if highly mixed up, high value]

Entropy function $= -\sum P \log P$ ← calculating entropy for every question
computationally expensive ↓ hard to compute → became difficult.

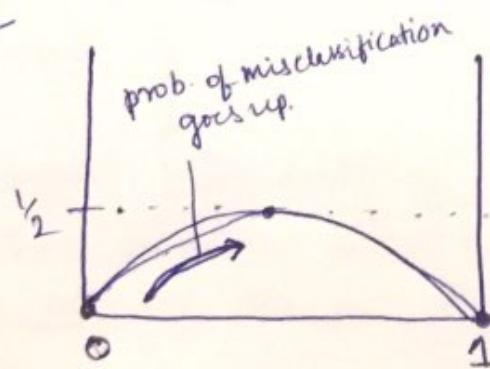
$\downarrow P \log P$?
probability of sample being picked.

changed the function \Rightarrow GINI IMPURITY.

$$= \sum_i P_i \sum_{k \neq i} P_k$$

not \times Entropy.
prob. of misclassifying = Prob. doesn't belong to i
= $1 - P_i$

highly pure, 0 mixed up, $\frac{1}{2}$



(Valid for 2 classes)

$$\sum_i p_i \sum_{k \neq i} p_k = \sum_i p_i (1-p_i).$$

$$\text{gini impurity} = \sum_i p_i - \sum_i p_i^2 = 1 - \sum_i p_i^2$$

\downarrow
computed millions of times (for each question)

\downarrow
We can learn the set of questions to be asked.

\downarrow
 \sim much better than the computation of $\log p_i$

\sim other classifiers assumed euclidean space, decision trees don't care [categorical data].

\downarrow
impurity can be computed.

* Advantages :-

- Efficient in practice
- Handles categorical data very well.
- Tells you the reasons (\sim answers), about the classification.

[this is the problem of deep neural network.]

generally)
categorical data
converted to
numerical data
 \sim most of the classifiers.

human interpretable
(decisions-answers).

in: length of the tail, whiskers or not?

- Tells you what features to look for.

| feedback about what features to collect]

↳ importance of features

* Limitation :-

- Training time (expensive to train)
- Not good with "Numerical" data

(split cannot give)
(a line → ↘
only axis parallel.)

→ Extreme case of overfitting can happen → decision trees can fit any complex distribution.

→ Stop overfitting? — stop the growth of the tree.

→ Early stopping (threshold on the entropy).

→ Pruning (chop off the bottom till you reach simplicity)

* Missing data can create issues in decision trees (you cannot skip questions) as you won't be able to go ahead.

missing data = average of features
(or)
value of previous data.

↓
how? training-validation

when the accuracy of validation starts to drop [due to overfitting of training data] = you stop!

semi-testing
[testing data while training]

* Random Forest Classifiers → built on decision trees.

multiple classifiers (neural network, SVM, decision tree, bayesian ...)

$$\hookrightarrow \begin{array}{l} h_1(x) \\ h_2(x) \\ h_3(x) \end{array} \quad \left\{ \begin{array}{l} 3 \text{ different classifiers for an input, } x. \end{array} \right.$$

how to learn from all of the classifiers (even though few classifier's performances are bad).

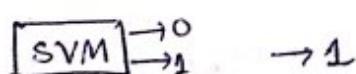
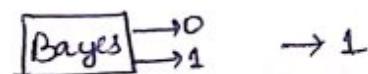
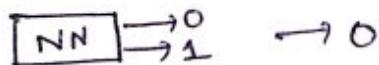
→ let suppose the accuracies are given;

$$\begin{array}{l} - h_{NN}(x) = 92\% \\ - h_{Bayes}(x) = 90\% \\ - h_{SVM}(x) = 88\% \end{array} \quad \left\{ \begin{array}{l} \text{one choice is to take the best out of all the classifiers.} \\ (\text{high accuracy}) \equiv \underline{\text{Majority Vote}}. \end{array} \right.$$

→ Majority Vote — due to its simplicity, it works with any classifier.

→ Weighted Voting — mixing the accuracy with the results.

Ex:



$$y_K = f(y_1, y_2, y_3)$$

$$= \left\{ \frac{92}{270} \times 0 + \frac{90}{270} \times 1 + \frac{88}{270} \times 1 \right\}$$

$$= \left\{ \frac{178}{270} > 0.5 \right\} \equiv 1.$$

In general; $y_K^1 = f(y_1, y_2, \dots, y_n, \epsilon_1, \dots, \epsilon_n)$

where y_1, y_2, \dots, y_n — outputs of n classifiers

$\epsilon_1, \epsilon_2, \dots, \epsilon_n$ — errors of classifiers.

$$\equiv \frac{\sum_K y_K (1 - \epsilon_K)}{\sum K \epsilon_K} < 0.5 \quad \begin{array}{l} \uparrow \text{if} \\ \downarrow \text{else.} \end{array} \quad \begin{array}{l} \text{for classification.} \\ \text{use.} \end{array}$$

→ Classifier not only gives the output, but also the confidence measure = (a single number).
if the classifier can output it.

$$y \Rightarrow f(y_k, c_k) \Rightarrow \frac{\sum_k y_k c_k}{\sum_k c_k} \sim \text{determines the result.}$$

↑
confidence
measure

→ Feature level Fusion = (information fusion for classification)

each classification is trained on a different representation of features,
can learn the function.

$$\begin{aligned} * & y = f(y_k, c_k) \\ & \left[\begin{array}{l} \text{train a second level classifier} \\ \text{Stack all the outputs and train a classifier based on it.} \\ [c_1, c_2 \dots c_k] \\ \text{into a feature vector.} \end{array} \right] \end{aligned}$$

* Will combining many classifiers improve the result?

→ Accuracy doesn't go much beyond the best classifier accuracy. (best of the experts).

→ It works good with a group of laymen.

→ Strong classifier
→ Weak classifier.

* Combining a large no. of weak classifiers can give a better result than a group of strong classifiers (ML justification).

- Methods to combine a group of classifiers to give a single classifier \equiv Ensemble methods.

\rightarrow Bagging : Bootstrap Aggregation.

\downarrow
Sampling on the training data

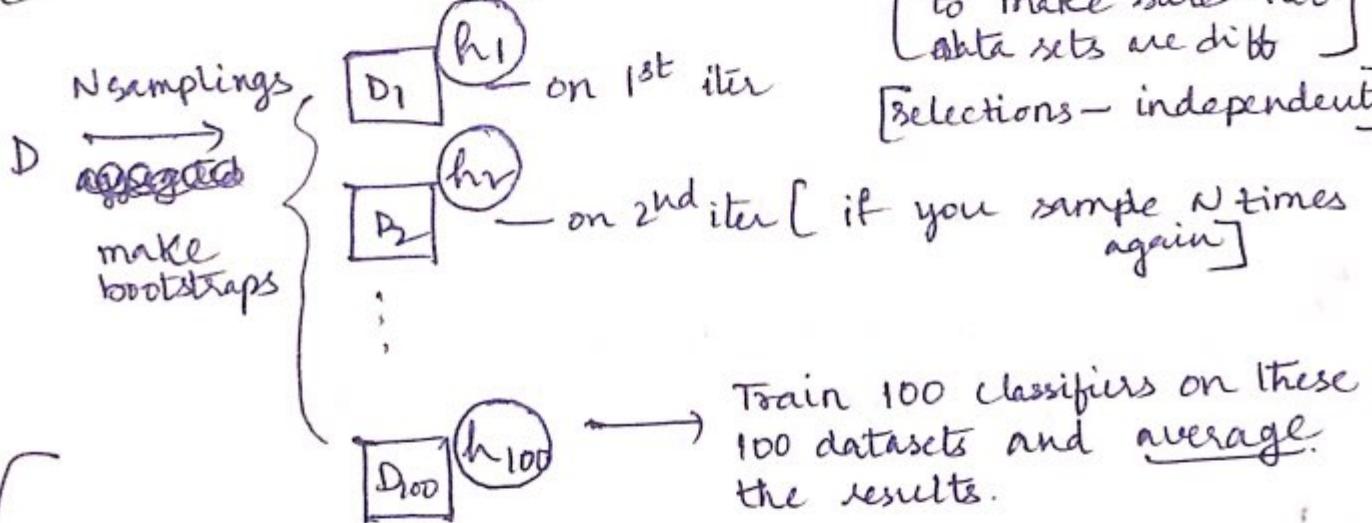
\Rightarrow Since we can't train all the classifiers on the same data (give similar results)

\rightarrow created a subset.

[called as bootstrapping]

fair
(random sampling) \rightarrow sample N times,
with replacement

[to make sure two
data sets are diff
selections - independent]



\Rightarrow if we sample less than N times, variability will be high.

\Rightarrow if we sample 100 times, the expected no. of times each sample will be selected is 100.

N done to create variations.

\rightarrow Boosting [works better in diversity] (not uniform)

\hookrightarrow weak classifiers but with complementary strengths are combined.

\rightarrow After training $\hookrightarrow h_1$ on $\boxed{D_1}$, test it on the training data \boxed{D}

→ Now, testing on D_1 might give classification errors.
— modify the sampling for D_2 ✓

↳ we assign some weights to the samples

$$\delta_i(x) = \gamma_N + i$$

should \downarrow be Normalised

$$\delta_i(x) = \begin{cases} e^{-\alpha_t} & \text{if misclassified, weight goes up} \\ e^{\alpha_t} & \text{if classified correctly, weight goes down} \end{cases}$$

→ Now for the next step, sample the data such that it contains misclassified samples [for the next classifier to test on] (weighted) decides

↓ This is a sequential process [not a parallel one]

* So, h_2 is a complement of h_1

h_{100} is a complement of h_1, h_2, \dots, h_{99} .

→ D_{100} has only those samples on which the other classifiers have failed on. [h_{100} doesn't work that well on the whole dataset as it has the complex samples]

⇒ So instead of directly averaging, they have to be weighted.

• Combine the classifiers based on weights \Rightarrow

$$\alpha_t = \frac{1}{2} \log \left[\frac{1 - \epsilon_t}{\epsilon_t} \right], \quad \epsilon_t = \text{error on the whole set given by the classifier.}$$

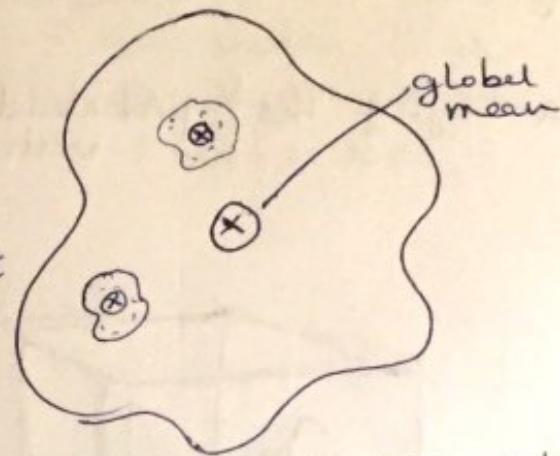
$$\downarrow \text{Output} \Rightarrow H = \sum \alpha_t h_t(x)$$

↳ called as ADA Boost \equiv Adaptive boosting.

* Variance of estimate -

→ as we take small subsets, mean tends to jump around.

variance of the ^{subset} estimate → variance of the whole set
underestimate.



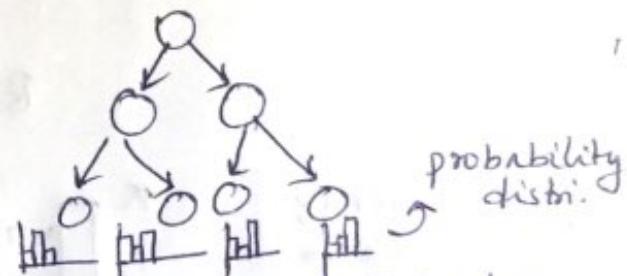
→ Ensemble methods work well even in the case of ~~bad~~ week classifiers.

↑
decision trees are not good for ensemble methods.

→ weaken the decision tree [as decision tree is a strong classifier]
[to a weak classifier.]

↓
→ ignoring few features [access restricted] → * Select a subset of features.
[for ex [pick up 1000 out of 5000]]
→ causes significant reduction of feature space

⇒ Stop the tree growth → stumps
— short decision trees.
"limiting the height".



⇒ Bagging (selecting a subset of samples → build a decision tree on it).

↓
doing such simplifications - underfit

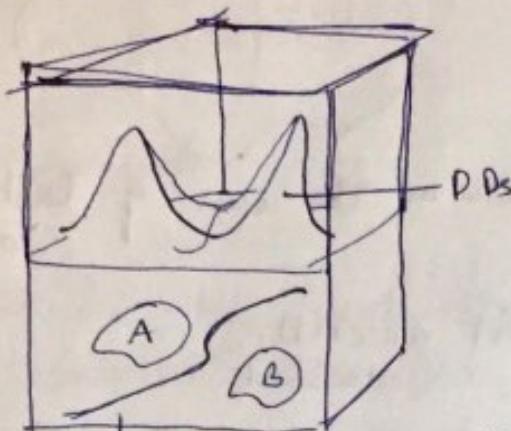
→ generate such ^{decision} _{weak} trees = T_1, T_2, \dots, T_{200} .

↓
and average the output [mixture of classes]

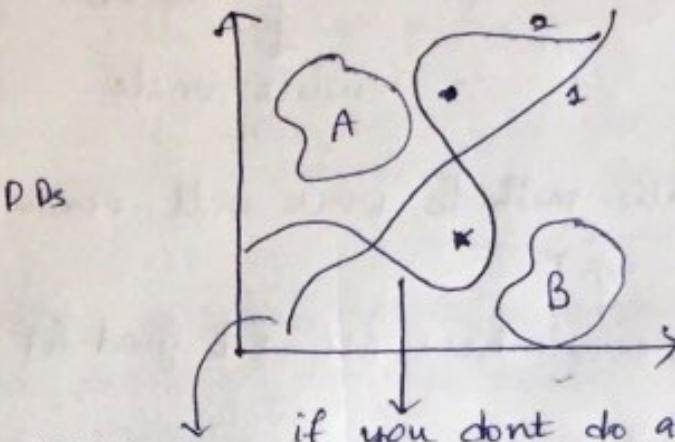
↓
⇒ the output generated by the trees is not a perfect class [as we cut it off], will have a PD of classes.

Height is limited \rightarrow so ensemble of trees is efficient.

→ averaging the probabilities for all the trees vector [ensemble]



decision boundaries
are smooth with
more samples



many sample sets miss out those outliers

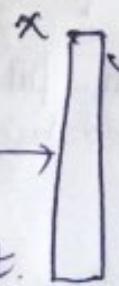
if you don't do any subset sampling, and because decision trees are extremely strong, your decision boundary turns out this way.

random features are selected along with random sampling

↓
Random Trees \Rightarrow Random Forests

Deep learning:
is growing \Rightarrow

raw data \rightarrow



learned feature vector

instead of handcoded features like [Random trees, SVMs]

doesn't overfit.

→ Random Forests \rightarrow easy to train as each tree is very small
might be faster than decision trees.

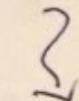
- extremely parallelizable
- efficient to train/test
- doesn't overfit

- can model any complex distribution

\times net \rightarrow runs a random forest classifier. Microsoft
 [using very simple features — like whether it is a shoulder or not?
 (based on depth)]
 computes depth & $g_{rgb} \gamma = \underline{rgbd}$

classification = extremely fast done on a weak processor.

- Random Forests \rightarrow can be used for regression also.

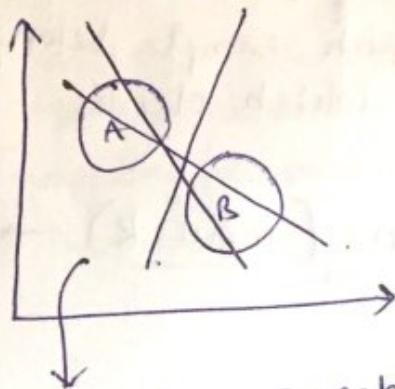


most important features
 "feature selection".

* gene data + cancer patient data
 ↓ Random Forest

give you the genes that are important in classification

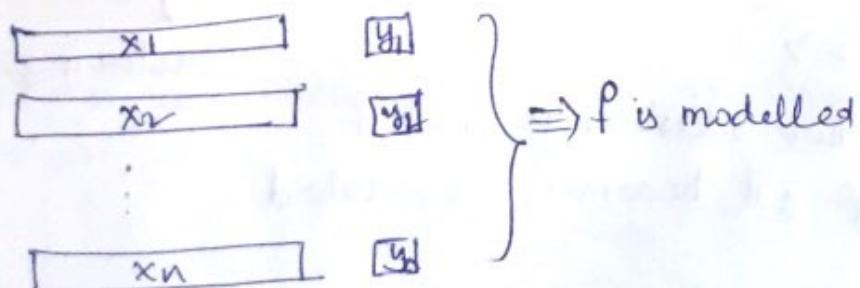
= { can also give non-linear boundaries } \leftarrow random lines [weak classifiers]



↓ decision trees.

Unsupervised Learning \rightarrow

+ supervised learning $\equiv y \leftarrow f(x)$
 $\{x, y\}$ pairs of data is given to you. mapping is learnt



Unsupervised learning

↓
data is given → don't know what to learn?

find patterns in the data

into multiple clusters.

* clustering - group the data ^{with similar patterns}.

↳ "ill-defined problem". \equiv Find clusters in the data

+ multiple unknowns -

- no. of clusters

- which sample belongs to
which cluster.

K-Means (input: K) \rightarrow

↓
grouping into no. of groups \equiv unknown
can be 2 or 3

if $K=2$, clusters $\equiv 2$

↳ on defining $K=2$,

cluster 1 cluster 2

find clusters in the data

↓ converted to

There are 2 clusters,
find out those

* Which sample belong to
which cluster?

↳ each sample could be put into $y =$ No. of ways to
cluster-1 or cluster-2 split into 2 sets
 \Downarrow

$$2^{n-1}$$

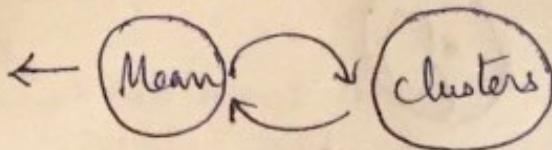
clusterings

* Naive approach \Rightarrow

- Compare any two clusterings, but if n becomes large, it becomes complicated.

* Assume, you know the means of the clusters
(close samples belong to the cluster)
to the mean - A

these are
hard
problems.



— just 2 comparisons
for every sample
[one with mean - A
and other with
mean - B]

given clusters, means can be computed
a given means, clusters can be computed.

↓ with both as unknowns (chicken and egg problem).

* Trick: Randomly pick means. and assign each sample
to the cluster
(computing distance)

— once you have clusters, compute the means of the clusters.

↓
Recompute the clustering

↓
Recompute the means.

Algorithm: KMeans (very powerful).

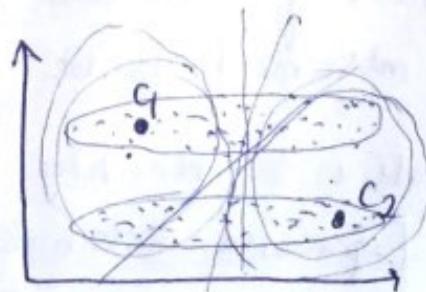
→ Termination:-
No change in
assignment, stop.

deep neural
learning ↗ improved
unsupervised
learning
(unlabelled data).

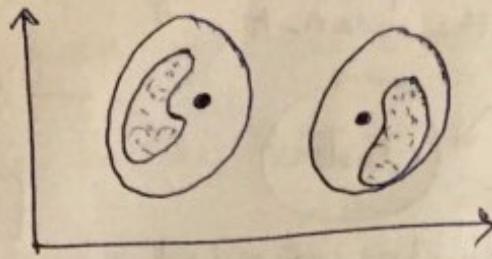
* Data clustering
by Jain &
Dubes
(classic algorithm)

KMeans - computes the distance to the centers
(euclidean).

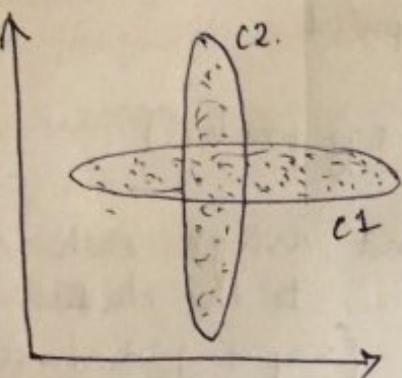
shapes are circular, spherical in shape.
clusters.



boundary between the clusters \rightarrow far apart



(a)



To Solve -

\rightarrow group them into more clusters and can later combine the clusters

- Mahalanobis distance

\rightarrow [distances in different directions are different] instead of Euclidean distance.

- Euclidean distance -

$$\sqrt{\sum_i (x_i - y_i)^2}$$

- Mahalanobis distance - depends on the distribution
[covariance matrix]
distances are "weighted".
based on the distribution of points.

$$D = (n-\mu) \sum^{-1} (x-\mu)$$

covariance matrix [of the cluster]

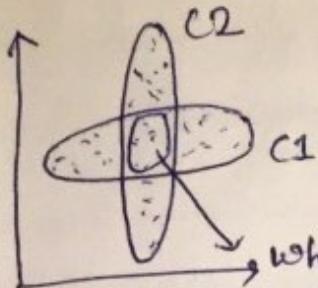
P.C. Mahalanobis - Indian Mathematician

\Downarrow Mahalanobis dist. instead of Euclidean dist.

[• clusters are no more assumed to be spherical]

if circular - covariance matrix = identity matrix. ($\Sigma = I$)

• clusters \rightarrow Elliptical [gaussian distribution]



what about these points? where do they belong?

- ① → cluster shapes are made not spherical [elliptical]
- modelled as gaussians.

$$C [M_c, \Sigma_c]$$

- ② → soft Assignments = [every point is assigned to each cluster with some probability (different likelihoods).]

$$\text{cluster} = N(x; M_c, \Sigma_c)$$

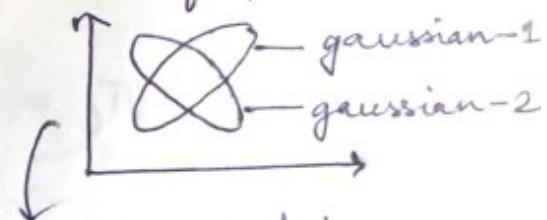
Combination of

→ Multiple gaussians model the distribution of points.

$$\rightarrow \text{size of gaussian} \equiv \pi_c$$

* samples assigned
to cluster c

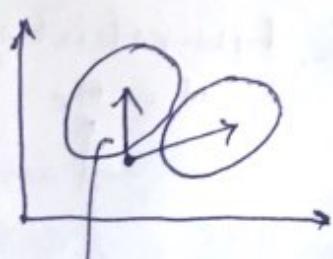
becomes = density estimate.



model = sum of weighted gaussians.

[based on size]

$$\boxed{\text{E-step}} \equiv \text{Expectation} \quad (\text{of assignments}) \\ \equiv \pi_c \cdot N(x_i; M_c, \Sigma_c)$$



Not Euclidean distance, instead likelihood estimate

TO sum it up to one \Rightarrow
sum it over all the clusters

* responsibility of i^{th} sample for cluster c

$$\gamma_{i,c} = \frac{\pi_c \cdot N(x_i; M_c, \Sigma_c)}{\sum_j \pi_{cj} \cdot N(x_i; M_{cj}, \Sigma_{cj})}$$

$r_{ic} \rightarrow$ computed for $N \times C$
 ↓
 ↓ no. of samples no. of clusters.

represented by a matrix

	C	R	rows
1st	1	... 1	sum up to one
N			row contains for each sample-i soft assignment to all the clusters.

For cth cluster, all the responsibilities together \Rightarrow

$$m_c = \sum_i r_{ic}$$

$$\downarrow \pi_c = \frac{m_c}{m} \quad (\text{updating } \pi_c)$$

you are the cluster, m_c = all the responsibilities of data points towards you by sum of responsibilities to all the clusters.

M-step

updating, $\mu_c, \Sigma_c \Rightarrow$

$$\mu_c = \frac{\sum_i r_{ic} \cdot x^i}{m_c}$$

maximum likelihood estimate of μ

\therefore equivalent of $N = m_c$
 \uparrow
 soft assignment

$$\Sigma_c = \frac{1}{m_c} \sum_i (\hat{x}_i - \hat{\mu})^T (\hat{x}_i - \hat{\mu}) r_{ic}$$

MLE of Σ

$$\frac{\sum x_i}{N} = \frac{\sum r_{ic} x_i}{m_c}$$

K-means algorithms

↳ 2 iterations =

- non-spherical clusters
- soft assignments (instead of hard assignments)

→ computing π_{ic} [cluster assignments]

→ μ_c

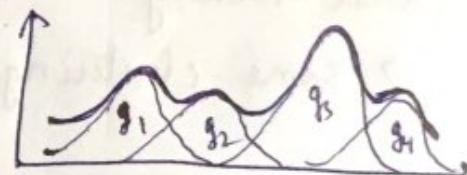
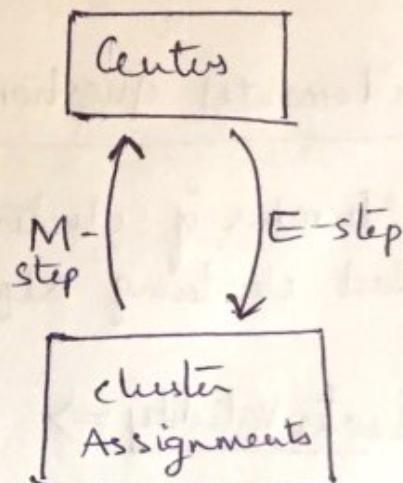
→ updating π_c

→ updating μ_c, Σ_c [cluster centers]

Complex distributions = can be

effectively modelled by a combination of multiple gaussian distributions.

↓
Gaussian Mixture Model (GMM)



E-M algorithm:

↓
for generic
chicken & egg
problems

(Z) (unobserved variables)

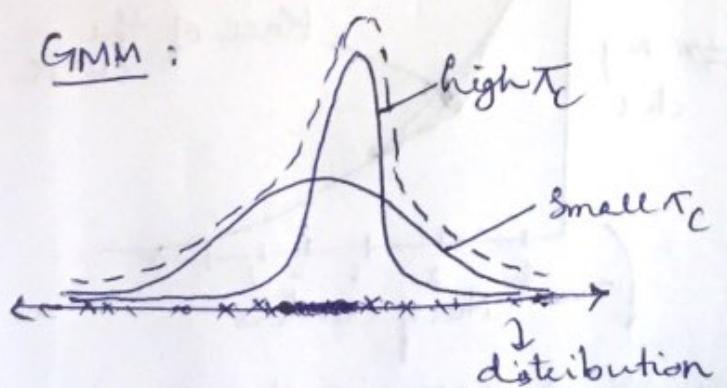
they are called "Latent Variables".

guaranteed
to converge.

- gaussians can grow larger / smaller by the varying π_c

↓
not the spread of points but the no. of points.

GMM:



Gaussian Mixture Model = Density Estimate

[Generative Model]

↳ can also generate samples to decide which cluster it belongs to.

- K-means - doesn't get stuck in local minima, but that's not the case with GMM Estimation

GMM Estimation = Latent Variable in here is "Pic"
 K-Means = Latent Variables = centers, assignments.

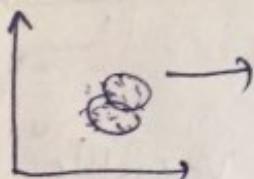
Fundamental Questions in Clustering :-

- * Number of clusters?
- * Best clustering algorithm.

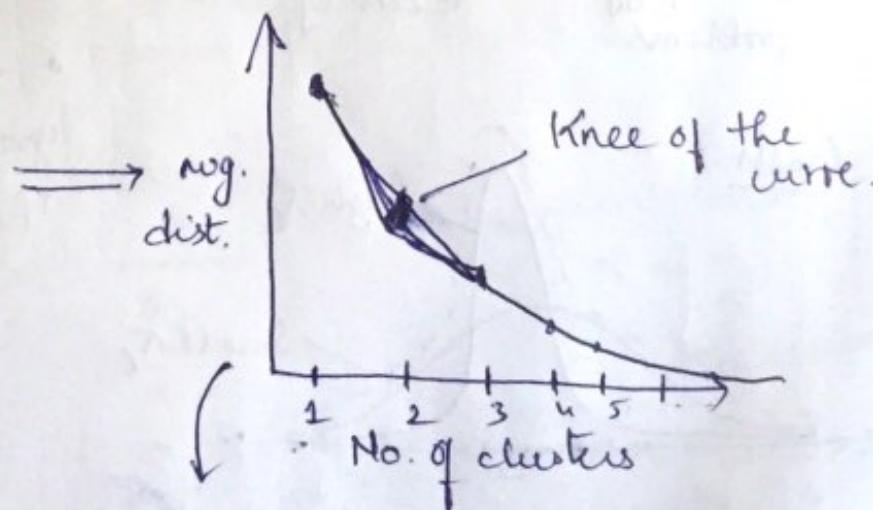
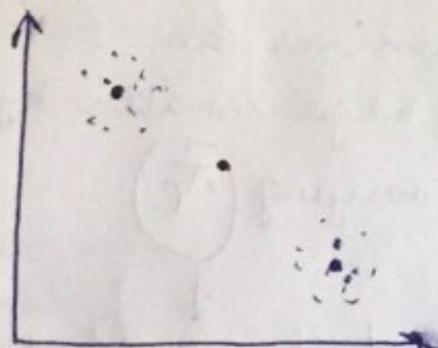
Cluster Validity \Rightarrow Is the clustering valid or not?

① Cluster Stability.

\rightarrow same clustering result for multiple initialisations.

 won't be stable (Not red clusters)
 * cluster centers will be moving for every iteration.

② Average distance (MSD) to ClusterCenter.



You see a knee in the curve if the data is well separated.

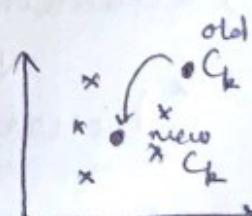
Hierarchical Clustering

[Slides]

- K-Means = simple clustering algorithm.
 - Randomly select k samples as cluster centers, c_k
 - Repeat until convergence =
 - Assign each sample to nearest cluster center (c_k)
 $\& \text{of } x_i^k = \arg\min \|x_i - c_k\|$
 - Update cluster center as mean of samples (x_i^k) assigned to that cluster.
- $$c_k = \frac{1}{n_k} \sum_i x_i^k$$

Convergence?

- Distance of sample to the cluster, over all the samples
 $\downarrow G_k = \sum_i (x_i - c_k)^2$ if goodness \downarrow , it is more compact.
- Goodness, $G = \sum_k G_k$ (overall clusters)
- Reassignment $\rightarrow G$ decreases monotonically.
 - x_i , if it has to change clusters, then distance from x_i to c_p has to be more than the new reassigned cluster distance, x_i to c_k distance
 - \Rightarrow Hence, it is decreasing.
- Updating the cluster centers -
 - on computing average, distances go down from that point.
- At every step of K-Means, G goes down, but it has to either
 - Bound on the step size.
 - No. of G values is finite [permutations = all possible grouping subsets]
 $\text{possible } G \text{ values} \leftarrow \{2^k\}$



>Data $\xrightarrow{\text{PCA}}$ Transformed Data $\xrightarrow{\text{KMeans}}$?

KMeans ↴

Highly sensitive to
outliers because
of squared distance.

define a distance threshold.

→ you ~~will~~ get proper clustering
if there is ~~no~~ loss
[~~sufficient~~ principal components]

Value of K (No. of clusters)?

G_K goes down,
when no. of clusters
is increased.

- ↳ generally prefer lower K
- ↳ balance between, K vs G_K

* Instead of square distances in G_K , we can use cosine distances [easier to compute].

> In GMM, you end up with elliptical clusters

K-Means " " spherical clusters.

comes to a hard stop, after some time [convergence]

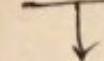
* Clusters \rightarrow sub-clusters.

so, clustering algorithm \Rightarrow output a hierarchy of entities

* Hierarchical Clustering.

① Top Down (Divisive) Clustering = Apply K-means repeatedly on each cluster.

② Bottom Up (Agglomerative) clustering [more efficient]



bottom up grouping \Rightarrow start with each sample in a diff cluster and then join the clusters till you get a single cluster

* Dendrogram — Hierarchy of clusters represented as a tree.

as we increase the distance threshold, we keep grouping into clusters. [inter-cluster distance]

Similar to Kruskal's for minimum spanning tree

$$\sim \text{Time Complexity} \equiv O[n \log n] \equiv O[E \log E] \quad E = m^2 \text{ edges}$$

$$O[m^2 \log m] \quad m = \text{points}$$

Inter-cluster distance =

$\min[\text{distance between the sample points in the clusters}]$

$$E = m^2$$

→ Hierarchical Algorithms \rightarrow don't know where to cut?

for appropriate clusters

{ just like the m in K-Means }

lengths of lines in the tree \Rightarrow represents "lifeline" of cluster

[birth to death]

* Valid cluster — long lifeline

→ Inter-Cluster distance =

* single link clustering algorithm / MST clustering =

$$D(C_p, C_q) = \min_{x_i \in C_p, x_j \in C_q} (d(x_i, x_j))$$

$C_p, C_q \rightarrow$ are the clusters.

as you cross the smallest distance, a link is created.

- Complete Link — max of distances
- Average Link — average of distances.
- Centroid based.

⇒ Single Link Clustering Algorithm:

→ prefers long chain clusters.

⇒ Complete Link = $O[N^3]$ ~ Time complexity.

→ spherical clusters.

* Centroid based — compute the centroid of cluster

$$C_p \quad C_q \quad \text{y} \equiv \text{centroid} = \frac{n_p C_p + n_q C_q}{n_p + n_q}$$

implemented easily
and efficiently.

⇒ How good is a clustering algorithm?

Comparing clustering algorithms =

→ feed "standard data / labelled data / ground truth".

Let's
• Standard data = C gold standard classes

• Clustering algorithms = K clusters, $\{w_1, w_2, \dots, w_K\}$ with n_i members

→ measure the Accuracy?

⇒ Purity metric [cluster-level] → "Biased"

$$\text{Purity}[w_i] = \frac{\max_j (n_{ij})}{n_i} \quad j \in C$$

if majority in the cluster are apples, then how many
impurity ← oranges are in it?

→ Rand Index [overall] — (2×2 matrix)

↳ If 2 samples are together in ground truth, — are they together in the clusters or similarly, verify for the inverse also. on clustering^{not?}

↳ depends defined on pairs, not single data points

Ground truth clustering	Same cluster	Different cluster
Same class [ground truth]	A	B
Different class [ground truth]	C	D

$$\text{Rand Index} = \frac{A + D}{A + B + C + D}$$

related to precision metric of retrieval problems

• SemiSupervised Learning — partially labelled data.

↳ partial labels (not for every pair).

(such hard constraints — are usually softened).

[* adding a penalty term to the optimisation problem]

↳ Many variations of partial labelling can be present.

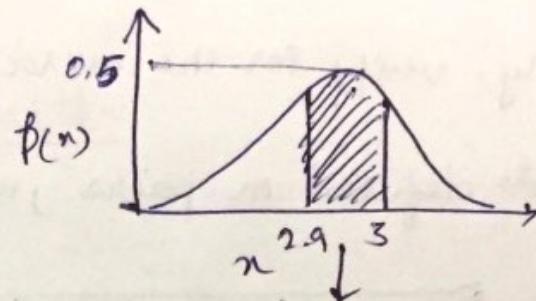
Probabilistic Graphical Models — Daphne Koller. (Stanford) (CPGM).

Probability density function

→ area is the probability, not the height

Bayes Rule

$$\Rightarrow P(A/x) = \frac{f(x/A) \cdot P(A)}{P(x)}$$



can be ignored → area here = Probability that x lies between 2.9 & 3.

- if this is known, (density is unknown)
→ given a class, likelihood of \bar{x} vector ..
- then bayes becomes accurate (the best).

d-dimensional vector $\rightarrow O[Kd^2]$ complexity

$d \mu_s, d^2 \Sigma$ parameters \downarrow for K GMM model. \equiv parameters have to be estimated.

Can be solved using —

Naive Bayes [where independent features]

\sim correlation = zero

correlation matrix \equiv diagonal matrix

\Rightarrow Parameters change to $= O[Kd]$ from $O[Kd^2]$

$$P(\bar{x}/A) = \underbrace{P(x_1/A) \times P(x_2/A) \dots P(x_d/A)}$$

{ multiple features \sim need not be same distribution

[not all gaussian]

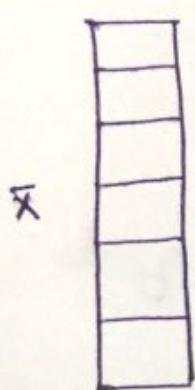
one can be normal density, other can be exponential density.

Probabilistic Graphical Models =

* completely independent features - not always true.

Model what is dependent and what is not dependent..

→ Let x be vector — then we have a 6-dimensional feature vector
features



$\rightarrow 0, 1$, No. of possible values $\equiv 2^6$

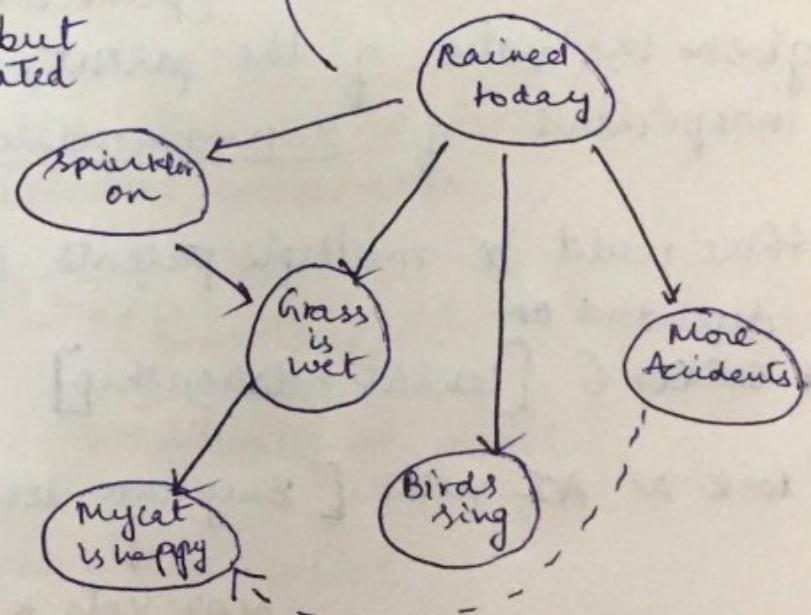
[without any probabilistic relationship information]

\rightarrow Compute $p(x)$, density function

Lines indicate causal relationships.

not causal but correlated

→ There's no direct relation / line between sprinkler on and more accidents, but they might be correlated.



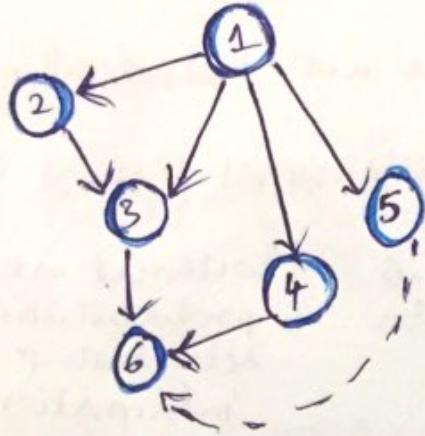
defined over huge no.
of observations [statistically]

Causality vs Correlation.

Only correlated, not causal.

When independent, can use naive bayes but
when dependent, use bayesian belief networks.

Conditional Independence =



$6, 5$ - not completely independent
 $6, 5$ - do not have a causal relationship.
 Knowing the value of 6 can affect the value of 5.

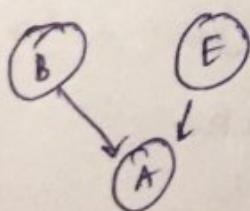
$\Rightarrow 2, 5$ - conditionally independent, given the value of its parent (1)
 (predecessors).

- given the value of the parent variable is conditionally independent of "non-descendants".
- there could be multiple parents (like (1,2) for 3)
- 4 ^{dependent on} 6 [causal relationship]

look at AI notes [Bayesian belief networks]

↓
represents a joint probability distribution of variables.

$$P(x_1, x_2, \dots, x_n) = \prod_i P(x_i / \text{parent}(x_i))$$

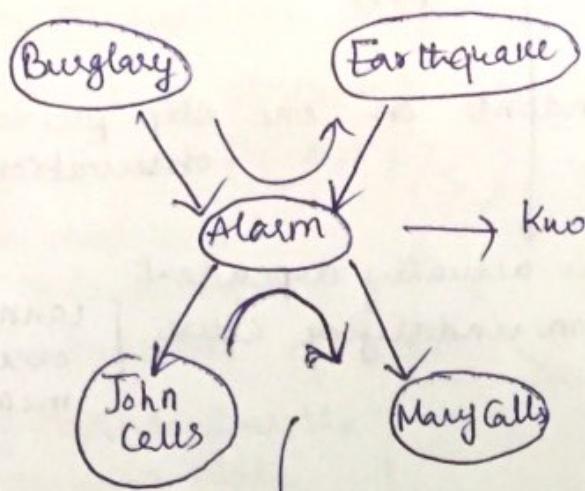


$$P(A/B) = P(A/B, E)P(E) + P(A/B, \neg E)P(\neg E)$$

$P(B/A) \Rightarrow$ use Bayes' rule

$$\frac{P(A/B) P(B)}{P(A/B) P(B) + P(A/\neg B) P(\neg B)}$$

Inter-Causal =



→ Knowing the values of common descendants, there exists an inter-causal influence.

When alarm goes off, (not knowing this = makes them dependent)
John calling → probability of Mary calling ↑

Hidden Markov Model →

$P(x_1, x_2, \dots, x_d)$ — Joint density can be factorised into independent components

multiple features = observed at time t $\pi P(x_i / p_{\text{parent}}(x_i))$

↓
conditioned on parent.

$x_t, t=0, 1, 2, \dots, T$

A single entity measured over a time period.

→ $P(x_0, x_1, \dots, x_T)$

↳ there exists some dependency

Correlation between what happens today / tomorrow

we make a markovian assumption.

→ dependency is only 1 step [simplest]

$p(x_i) \rightarrow$ depends on $p(x_{i-1})$ only.

on the previous one,
not on anything
in the future.

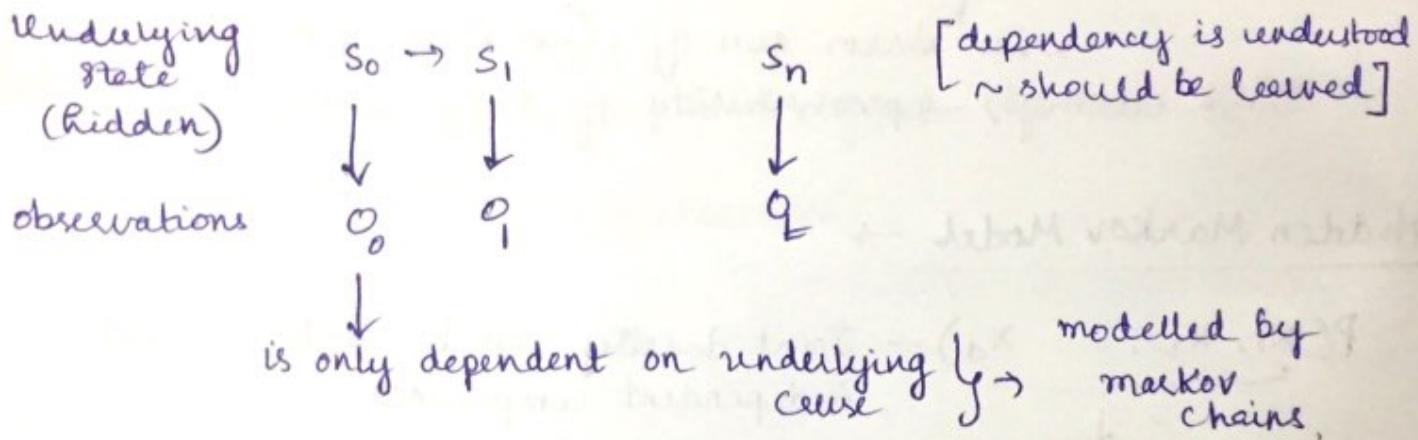
$$P(x_0, x_1, \dots, x_T) = p(x_0) \cdot p(x_1/x_0) \cdot p(x_2/x_1) \cdot \dots \cdot p(x_T/x_{T-1})$$

$$= P(x_0) \cdot \prod_{t=1}^T P(x_t/x_{t-1})$$

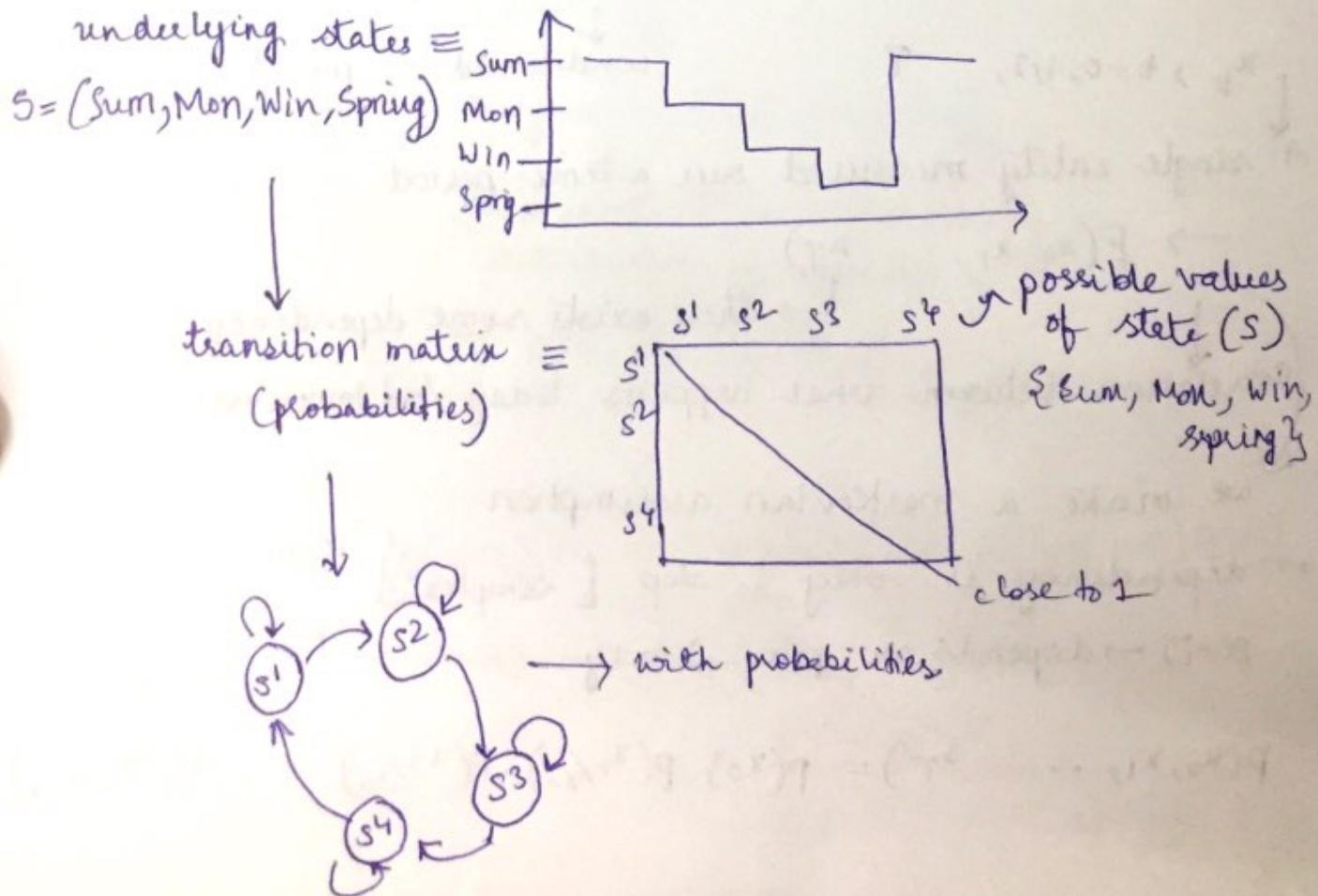
\Rightarrow every observation is dependent on one step previous observation.
 [first-order markovian process]

it is actually dependent on an underlying cause. [cannot be observed/measured]

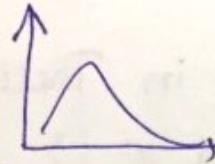
Underlying Cause



For Ex: for predicting the weather,



- we should know -
 - * state transitions $s_0 \rightarrow s_1 \dots s_T$ [unobserved]
 - * likelihood of observing an observation given the underlying state.
- $P(O = y | s)$ = observation probability.
 can be modelled as a distribution function
 [purely dependent on the underlying state]



$$P(O_0, O_1, \dots, O_T | s_0, s_1, \dots, s_T) = P(s_0) \cdot P(O_0 | s_0) \cdot \dots \cdot \prod_{t=1}^T P(s_t | s_{t-1}) P(O_t | s_t)$$

↓ available

① Given O , find most likely sequence of s_i .

$$P(O) = \sum_s P(O_1, O_2, \dots, O_T, s_0, s_1, \dots, s_T)$$

↓ summing over all possible sequences.

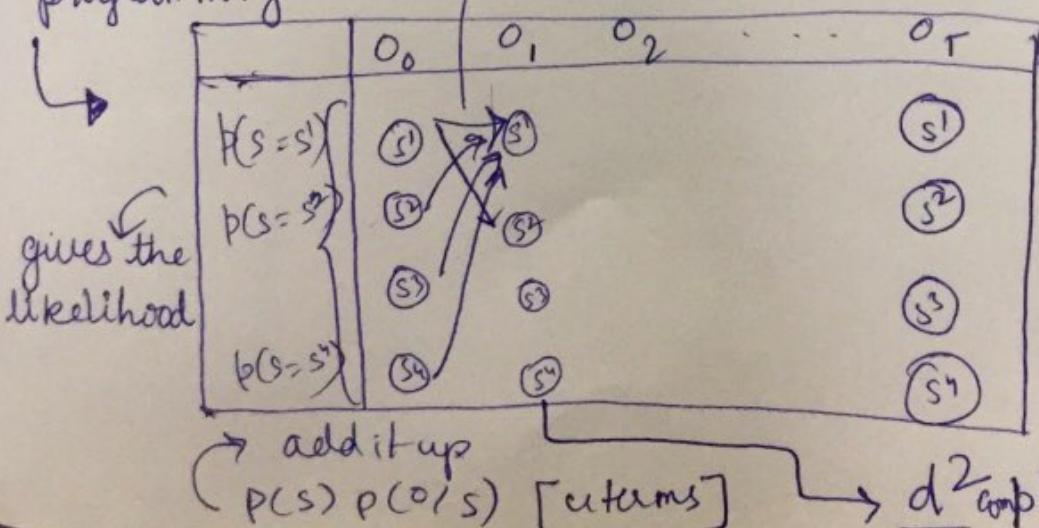
for ex: observation = O

it can be \hookrightarrow hidden or \rightarrow kitten.

out of all sequences, find the best.

(can be solved using dynamic programming.)

$P(s_t | s)$ transition probability (have to be known) multiply and add it.



Trellis diagram / decoding
 [which sequence do we pass thru?]

Computations to be done = $T \cdot d^2$

$A^n \rightarrow T \cdot d^2$ [Dynamic Programming]
(Reduction in complexity)

Most likely path sequence in Trellis diagram

→ use "argmax" instead of Σ

↓ Viterbi decoding / algorithm

→ Learning problem - given O , learn transition probs.

* left to right - forward computation / algorithm

right to left - backward computation / algorithm

↓ Forward-backward algorithm