

So, L_2 reg. doesn't change the value of w_1 from one iteration to another.

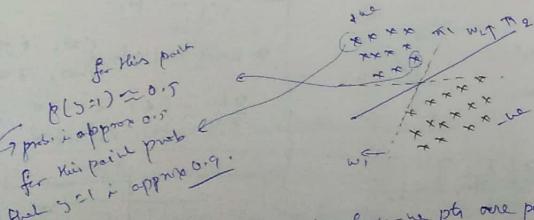
It, however, continues to constantly reduce w_1 towards $w_1 = 0$.

So, choose $\text{let } w_1 \approx 0$ at the end of iterations

Support-Vector-Machine (SVM)

Geometric Intuition:

SVM is a popular ML algo. used in classification / Regression, developed back in 1990's

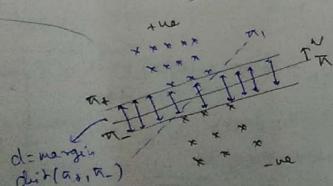


There are many lines that separate the +ve and -ve pts are possible.

key idea of SVM: Choose the line that separates the +ve and -ve pts as widely as possible.

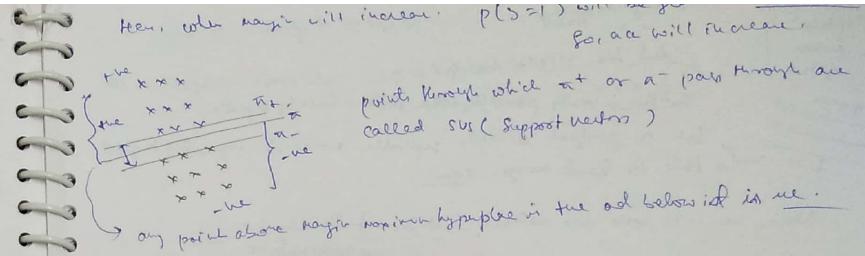
Here a is better than a'
 $\rightarrow a$ is margin maximizing hyperplane

a is parallel to a'
 a is parallel to a''
 a is also parallel.



To find a place a that maximizes the margin = $\text{dist}(a+, a-)$

because, margin ↑, generalization accuracy ↑



⇒ Alternative geometrical intuition of SVM:

one way to do it using convex hull



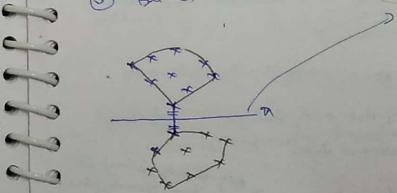
This is a closed path that each point should lie inside it.

B This is not a convex hull, because, like joining two point is order for store file.

① Create Convex Hull for the +ve and -ve pts.

② Find the shortest line connecting these hulls.

③ Direct the line to get margin max hyperplane.



Mathematical Derivation:

α = margin-maximization.

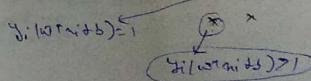
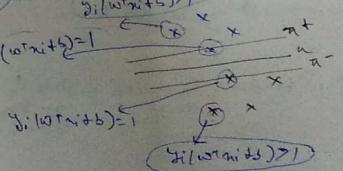
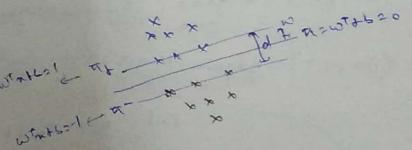
$$w^T x_i + b = 0$$

$$\text{if } \alpha_i \neq 0 \quad w^T x_i + b = 1$$

$$\alpha_i : w^T x_i + b = -1$$

$$\text{overall margin} = d = \frac{2}{\|w\|} \quad \text{This can be derived.}$$

$$(w^*, b^*) = \underset{w, b}{\operatorname{argmax}} \frac{2}{\|w\|} = \text{margin} \quad g_i(w^*x_i + b^*) = 1$$



Constraint optimization problem of SVM

$$\min_{w, b} \frac{1}{2} \|w\|^2$$

subject to $y_i(w^T x_i + b) \geq 1$

True not true \rightarrow This is only possible when data is linearly separable
but a perfect linearly separable model is impossible to find.
 \rightarrow This is hard margin SVM.

Now, we will look at datasets that is not linearly separable.

for non linearly separable data $\exists i$ & repeated

from 1 for the misclassified $y_i(w^T x_i + b) < 1$
pt. ad for point y_i is $= -0.5$
in at π .
But for rightly classified $y_i(w^T x_i + b) = 1 - (1.5)$
pt. it is not separated from π .

$w \rightarrow \pi$
 $g(x) \geq 1$ if $y_i(w^T x_i + b) \geq 1$

Correctly classified pt

if $y_i(w^T x_i + b) \leq 1$ it is equal to some unit of dist away from the correct hyperplane in the incorrect dir.

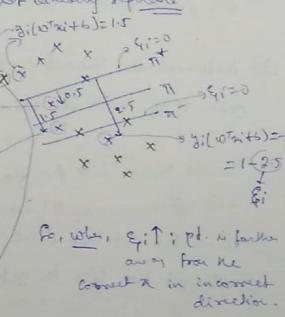
$$(w^*, b^*) = \arg \max_{w, b} \left(\frac{1}{2} \|w\|^2 \right) \text{ subject to } y_i(w^T x_i + b) \geq 1 \quad \text{for correctly classified}$$

$$(w^*, b^*) = \arg \max_{w, b} \left(\frac{1}{2} \|w\|^2 \right) + C \cdot \frac{1}{n} \sum_{i=1}^n \xi_i \quad \text{where } \xi_i = \text{dist of misclassified pt. from correct plane}$$

s.t. $y_i(w^T x_i + b) \geq 1 - \xi_i; \forall i$ } for correctly classified points $\xi_i \geq 0$

minimum errors = minimum misclassification = $\min E_{\xi_i}$
 \rightarrow C is in general regularization term.

on CT, tendency to make mistakes on D_{train} \rightarrow overfit \rightarrow high variance



$C \perp$; underfit \rightarrow high bias

As in logistic-regression \rightarrow $C \perp \rightarrow$ high bias
 $C \parallel \rightarrow$ high var

$$C = \frac{1}{\lambda}$$

Why we take value +1 and -1 for support vector planes

Here just for convenience we take 1 and -1, we can also take any two numbers $+k$ and $-k$
 \rightarrow But these values should be same because, we need a proper bisector of plane.

$$\begin{aligned} w^T x + b \\ w^T x + b = 0 \\ w^T x + b = -1 \end{aligned}$$

$$\|w\| \neq 1$$

\rightarrow so can be any vector it need not to be a unit vector.

Since, our main goal is to maximize the margin

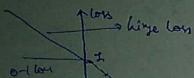
$$\begin{aligned} \therefore w^*, b^* = \arg \max_{w, b} \frac{1}{2} \|w\|^2 &= \arg \max_{w, b} \frac{1}{2} \frac{k^2}{\|w\|^2} & \text{if } k \geq 0 \\ \text{then take } w^* \\ \text{we see for } \underline{\underline{w}} \end{aligned}$$

$$w^T x + b = k$$

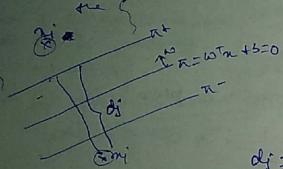
$(w^T x + b)/k = 1 \rightarrow$ since k is the no. so,
 \rightarrow it can't be had side value can be easily change
 $\underline{\underline{k}}$.

Loss Function (Hinge loss) based interpretation.

Because nature of loss in SVM is hinge type.



case 1: $z_i \geq 0$; x_i is correctly classified
 $z_i < 0$; x_i is incorrectly classified.



$$d_i = 1 - z_i / (\omega^T x_i + b) = 1 - z_i$$

$$\text{if } \epsilon_i = \text{dist from } x_i \text{ to the } \pi^\perp = d_i = 1 - z_i$$

when x_i is misclassified.

Since, x_i is misclassified from the ~~positive~~ region so, distance is taken from the negative.

$\epsilon_i = 0 \Rightarrow$ since x_i is correctly classified.

Soft SVM:
 $\min_{\omega, b} \frac{1}{2} \|\omega\|^2 + C \sum_{i=1}^n \epsilon_i$ → loss.
 C ↑ → soft margin
 C ↓ → underfit
 $\omega \rightarrow$ weight
 $b \rightarrow$ bias

reflecting that here because it is convex function of ω .

$$\text{loss-min} \rightarrow \min_{\omega, b} \sum_{i=1}^n \max(0, 1 - y_i(\omega^T x_i + b)) + \frac{1}{2} \|\omega\|^2$$

when $\|\omega\| \geq 0 \Rightarrow$ min value of $\frac{\|\omega\|^2}{2}$ is same as $\min \|\omega\|^2$

So, Soft SVM is similar to loss-min,

SVM = hinge loss + regularization.

Dual Form of SVM Formulation

Suppose if someone gives us dissimilarity matrix, then in soft-margin SVM it is not interpretable, so in \Rightarrow it is converted into its equivalent Dual Form, and worth after dual form is used.

Dual Form:
 $\text{SVM} \left\{ \begin{array}{l} \min_{\omega, b} \frac{1}{2} \|\omega\|^2 + C \sum_{i=1}^n \epsilon_i \\ \text{s.t. } y_i(\omega^T x_i + b) \geq 1 - \epsilon_i, \epsilon_i \geq 0 \end{array} \right.$

This is also known as

Primal of SVM.

On main aim of soft margin SVM is to find $f(x) = \omega^T x + b$

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

- ① for each x_i there is a corresponding value d_i ($x_i \rightarrow d_i$)
- ② x_i 's only occur in the form of $x_i^T \omega$

$$\textcircled{3} f(x) = \sum_{i=1}^n d_i x_i^T \omega + b$$

For Support vector (SVs) $\rightarrow \epsilon_i = 0$

for non-SVs $\rightarrow \epsilon_i > 0$

So, for $f(x)$ \Rightarrow only point that matters is SVs, if we add any non-SVs points, it won't effect $f(x)$

$$\Rightarrow \max_{\omega, b} \sum_{i=1}^n \epsilon_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \epsilon_i \epsilon_j y_i y_j x_i^T x_j$$

$$\text{Set } \epsilon_i \geq 0 \rightarrow \text{for } \epsilon_i = 0 \Rightarrow \text{SVs}$$

$$\sum_{i=1}^n \epsilon_i = 0 \quad \epsilon_i > 0 \Rightarrow \text{non-SVs}$$

$$\text{So, } \omega^T x_i = x_i^T \omega = \text{Correlation}(x_i, \omega) \rightarrow \text{if } \|\omega\|=1 \Rightarrow \|x_i^T \omega\| = 1$$

Correlation of

$$x_i^T \omega$$

So, $x_i^T \omega \Rightarrow$ Can be thought of as kernel function.

$$k(x_i, \omega)$$

G. $k(x, \omega)$ = kernel function.

So, during sum time:

$$f(\mathbf{x}_q) = \sum_{i=1}^n \hat{y}_i k(\mathbf{x}_i, \mathbf{x}_q) + b$$

$\hookrightarrow k(\mathbf{x}_i, \mathbf{x}_q)$

$$\hookrightarrow \sum_{i=1}^n \hat{y}_i k(\mathbf{x}_i, \mathbf{x}_q) + b.$$

Kernel Trick

fill more we see that $(\mathbf{x}_i^T \mathbf{x}_q)$ can be changed to any similarity kernel function.

$$\hookrightarrow k(\mathbf{x}_i, \mathbf{x}_q)$$

↳ the most important idea in SVM is Kernel Trick.

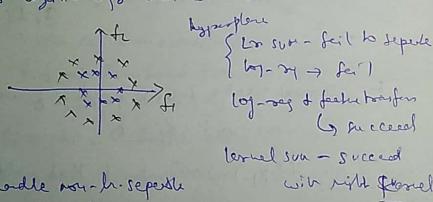
Since, soft margin-hyperplanes can only be applied to linearly separable data.

But Kernel trick or feature transformation is logistic regression, transform non-linear data and separate them using hyperplane.

The only difference L/R logistic separator and SVM is that SVM uses weight vector.

$$\text{linear sum} = \mathbf{x}^T \mathbf{w}$$

$(\text{kernel sum} \rightarrow k(\mathbf{x}_i, \mathbf{x}_j))$



So, kernelization → SVM handle non-linearly-separable datasets.

Polynomial kernels

This dataset is not linearly separable so, kernelization is used.

so we use feature transformation in log-regression $(f_1, f_2) \xrightarrow{\text{FT}} (f_1^2, f_2^2)$

$$\log-\mathbf{x}_q$$

$$k(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{x}_i^T \mathbf{x}_j + c)^d \quad d \text{ can be any number.}$$

$$\hookrightarrow k(\mathbf{x}_i, \mathbf{x}_j) = (1 + \mathbf{x}_i^T \mathbf{x}_j)^2$$

$\hookrightarrow \text{chebyshev kernel}$

$$k(\mathbf{x}_i, \mathbf{x}_j) = (1 + \mathbf{x}_i^T \mathbf{x}_j)^2 \rightarrow \mathbf{w} = \langle \mathbf{x}_1, \mathbf{x}_2 \rangle \xrightarrow{\text{2D data}} \\ \mathbf{x}_2 = \langle \mathbf{x}_{21}, \mathbf{x}_{22} \rangle \\ = (1 + \mathbf{x}_{11}\mathbf{x}_{21} + \mathbf{x}_{12}\mathbf{x}_{22})^2 \\ = 1 + \mathbf{x}_{11}^2\mathbf{x}_{21}^2 + \mathbf{x}_{12}^2\mathbf{x}_{22}^2 + 2\mathbf{x}_{11}\mathbf{x}_{21} + 2\mathbf{x}_{12}\mathbf{x}_{22} + 2\mathbf{x}_{11}\mathbf{x}_{12}\mathbf{x}_{21}\mathbf{x}_{22}$$

$$\begin{aligned} &\text{Let } [1, \mathbf{x}_{11}^2, \mathbf{x}_{12}^2, \sqrt{2}\mathbf{x}_{11}, \sqrt{2}\mathbf{x}_{12}, \sqrt{2}\mathbf{x}_{11}\mathbf{x}_{12}] : \mathbf{y}' \\ &\text{taking } \xrightarrow{\text{FT from }} [1, \mathbf{x}_{11}^2, \mathbf{x}_{12}^2, \sqrt{2}\mathbf{x}_{11}, \sqrt{2}\mathbf{x}_{12}, \sqrt{2}\mathbf{x}_{11}\mathbf{x}_{12}] : \mathbf{x} \\ &\text{taking } \xrightarrow{\text{2D}} \mathbf{x} = (\mathbf{x}_1')^T (\mathbf{x}_2') \end{aligned}$$

↳ 6-Dimension.

So, kernelization: $\xrightarrow{\text{d1 FT implicitly}} \text{d1} \xrightarrow{\text{d2}} \text{d1} \xrightarrow{\text{d2}} \text{d1} \xrightarrow{\text{d3}} \text{d1}$

$$\text{Kernel } 2D \xrightarrow{\text{FT}} 6D.$$

So, kernelization is an implicit process.

Mercer's Theorem:

↳ It shows how any kernel, a non-separable data can be converted into linearly separable.

$$d \rightarrow d'$$

RBF Kernel → Radial Basis Function.

SVM is the most popular/general-purpose: RBF

$$K_{RBF}(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{2\sigma^2}\right)$$

↳ Hyperparameter.

In soft-margined

kernel sum

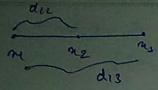
$\sum_i \hat{y}_i k(\mathbf{x}_i, \mathbf{x}_i) = \sum_i \hat{y}_i \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_i\|^2}{2\sigma^2}\right)$

These distances are related to similarity.

σ is also a hyperparameter.

$$\text{where } d_{12} \uparrow ; K(\mathbf{x}_1, \mathbf{x}_2) \downarrow$$

$$\begin{cases} \sigma \uparrow \text{and} \\ d \uparrow, d^2 \uparrow \\ \sigma \uparrow \text{and} \\ \frac{1}{\sigma^2} \uparrow \end{cases}$$

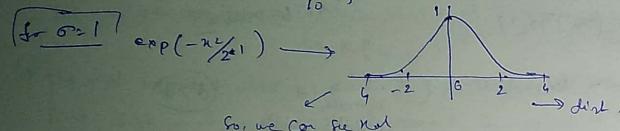


$$\text{So, } k(x_1, x_2) > k(x_1, x_3)$$

For, distances b/w x_1 and x_2 are greater than similarity b/w x_1 and x_3 .

\Rightarrow Hence, similarity depends on σ . and σ in RBF are similar to k in SVM.

Suppose, we have. $\sigma = \begin{cases} 1 \\ 0.01 \end{cases}$



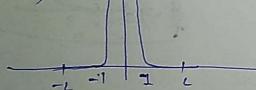
So, we can see that

when $d \rightarrow \infty \rightarrow \exp(-\frac{x^2}{2\sigma^2}) \rightarrow$ reaches to 0

$$\exp(-\frac{x^2}{2\sigma^2}) = 0 \rightarrow \text{when } d = \infty.$$

for, the similarity b/w the different points here when about in b/w -2 to 2.

for $\sigma=0.01$ $\exp(-\frac{x^2}{2\sigma^2})$



when $d \rightarrow 1, k \approx 0$.

So, here, the similarity b/w the dist -1 to 1.

So, we can see that $\sigma=10$

when, $\sigma \uparrow$, Similarity \uparrow $\sigma \downarrow \rightarrow$ 1

(See in SVM when $\sigma \uparrow$

the range b/w points \uparrow

\hookrightarrow

So, $(\sigma \uparrow \text{ in RBF}) \Rightarrow (\sigma \uparrow \text{ in SVM})$

But in SVM we had memory issue $O(n^2)$

But in RBF SVM it only depends on support vectors, which is smaller than all points in SVM.

So, it can easily calculate to get $f(x) \rightarrow$ close corresponding to any query point.

So, when we have no any idea about which kernel to use, then we use RBF SVM \rightarrow it will work good.

$(C, \sigma) \rightarrow$ Can be set by grid search, Random Search.

Domain Specific's kernels.

When we have domain knowledge + Feature Transformation \rightarrow Kernel-tricks

Like in text classification \rightarrow ~~linear~~ string kernel is used.

So, we should have a good domain knowledge.

Training and Runtime Complexity of SVM.

Best library for training SVM is LIBSVM, but we can use sklearn anyway.

Training is done by sequential minimal optimization (SMO).

Runtime complexity of SVM is $\rightarrow O(n^2)$ for (normal-SVM).

\hookrightarrow so for large dataset, use kernel search operator.

SVM can never be used.

$$f(x) = \sum_{i=1}^n \alpha_i y_i k(x_i, x) + b$$

$\alpha_i = 0$ for non SVs.

$$\text{for Support Vectors (SVs), } \alpha_i = 1 \quad \hookrightarrow \quad 0 \leq \alpha_i \leq n$$

nu-SVM Control and Support Vectors.

In C-SVM → that is the original form, we have $(C > 0)$

but we have an alternative formulation of SVM

i.e. nu-SVM, where $(C \rightarrow \infty)$

Hypoparameter ν , fraction of errors

$\nu \leq$ fraction of SVs

$SVs \geq 1/\nu$ of n

if $n = 100,000$ $SVs \geq 1000$

But then much support vectors could lead to large time complexity.

Suppose we are training our data with SVM algo.
and Error can't be more than $1 - \epsilon$.
Then, $\nu n = 0.05$

1 Cases

→ In SVM instead of doing feature engineering & feature transform, we find the right kernel i.e. "RBF"

→ Decision Surfaces ↑ when we have lin. SVM → hyperplane

Kernel SVM is used in case of non-linear Surface

→ when we have similar features/distance func.

↳ Kernel SVM will benefit from it i.e. (n^2)

→ SVM can't do good Interpretability and can't find a ~~good~~ feature importance.

→ SVM is very less affected by outliers since it depends on support vector only.

↳ Bias-Variance: $C \uparrow \rightarrow$ overfit \rightarrow high-variance
underfit \downarrow
 $C \downarrow \rightarrow$ underfit \rightarrow high-bias

\downarrow \hookrightarrow \hookrightarrow \hookrightarrow \hookrightarrow \hookrightarrow \hookrightarrow

["SVM Regression"]

As SVM is used for classification purpose, it is also used for Regression problem.

$$\min_{w, b} \frac{1}{2} \|w\|^2$$

$$\text{subject to } y_i - (w^T x_i + b) \leq \epsilon \quad \text{hyperplane}$$

$$y_i - (w^T x_i + b) \geq -\epsilon$$

$$(w^T x_i + b) - y_i \leq \epsilon$$

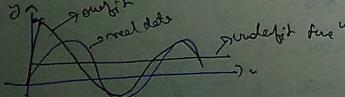
This can also be

linearized to fit in our linear

data also.

$\epsilon \downarrow \Rightarrow$ error in loss in training data but overfitting ↑

$\epsilon \uparrow \Rightarrow$ error in Dtrain ↑ \rightarrow Underfit ↑



error ↓

$\epsilon \downarrow \Rightarrow$ error in loss in training data but overfitting ↑

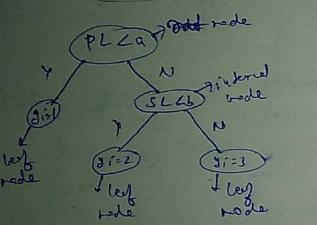
$\epsilon \uparrow \Rightarrow$ error in Dtrain ↑ \rightarrow Underfit ↑

Decision Trees

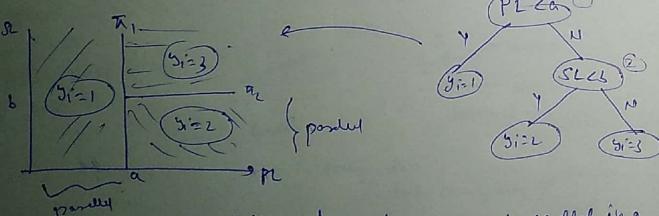
Geometric Intuition of decision tree Axis parallel hyperplanes

Decision trees
 ↗ it is simply a vertical ref of else statement.
 ↗ (base) learn, NB, log reg, dec reg, SVM → kernel.
 ↗ (impure based) probabilistic method
 ↗ (geometric, hyperplane)

Assume I have dataset, $\{x_i\}_{i=1}^n$, $x_i = (x_1, p_1, s_1, p_2)$



nodes at vertex are tree, any node that is not a root node or a leaf node
 is an internal node



all these hyperplanes are axis parallel in a decision tree.

Sample Decision Tree

Refer Poly provided in last of this pdf

Building a Decision Tree (DT) Entropy

Entropy is often used in information theory & electronics, and physics.

Suppose we have random variable

$$Y \rightarrow y_1, y_2, y_3, \dots, y_k$$

$$\text{Entropy } H(Y) = - \sum_{i=1}^k P(y_i) \log_b (P(y_i))$$

b can be
 $b=2$
 $b=e = 2.718$

$$P(y_i) = P(Y=y_i)$$

$$\log_2 = \ln$$

Suppose we have a dataset where there is probability that the person will play tennis = not

$$\text{if } P(Y_i=y_{\text{yes}}) = \frac{9}{14} \text{ then } P(Y_i=\text{No}) = \frac{5}{14}$$

not play
femin

$$\begin{aligned} \text{After calculation: } \\ H(Y) &= -\frac{9}{14} \log\left(\frac{9}{14}\right) - \frac{5}{14} \log\left(\frac{5}{14}\right) \\ &= 0.94 \end{aligned}$$

Properties:- $Y = y_+, y_-$ (2-class, 2 categories)

$$\text{Case 1: } \begin{cases} y_+ \rightarrow 99\% \\ y_- \rightarrow 1\% \end{cases} \quad \begin{cases} H(Y) = -0.99 \log(0.99) - 0.01 \log(0.01) \\ = -0.085 \end{cases}$$

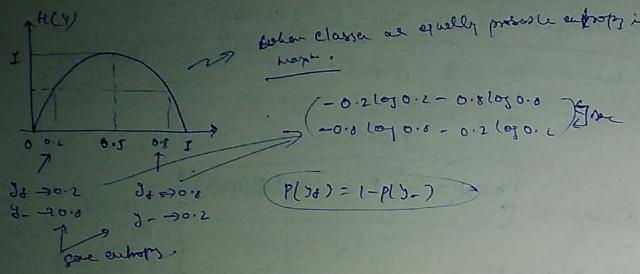
$$\text{Case 2: } \begin{cases} y_+ \rightarrow 50\% \\ y_- \rightarrow 50\% \end{cases} \quad \begin{cases} H(Y) = -0.5 \log(0.5) - 0.5 \log(0.5) \\ = 1 \end{cases}$$

$$\text{Case 3: } \begin{cases} y_+ \rightarrow 0\% \\ y_- \rightarrow 100\% \end{cases} \quad \begin{cases} H(Y) = 0. \end{cases}$$

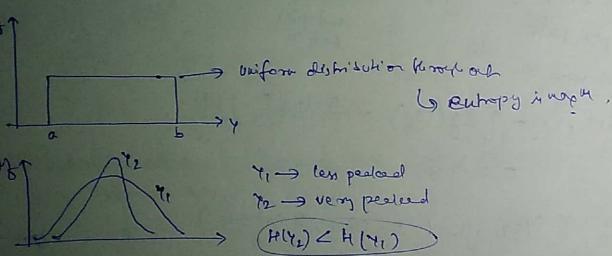
So, when classes are equally probable then Entropy is maximum.

Here, at 50% y.

Entropy decrease when there is dissimilarity in class.



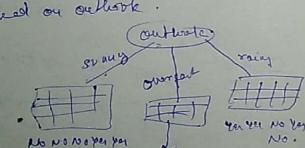
So, $Y \rightarrow S_1, S_2, S_3, \dots, S_K$
 \hookrightarrow Equally probable \rightarrow Entropy is max.
if $S_1 \rightarrow$ most probable
 $S_2, S_3, \dots, S_K \rightarrow$ less probable \rightarrow Entropy is min.



Building a decision tree formation (minimizing entropy)

Suppose we have made our decision tree based on outcome.

We previously have calculated the entropy for outcome $H_D(\text{outcomes}) = 0.94$



$H_{D_1}(Y) = 0.92$
 $H_{D_2}(Y) = 0$
 $H_{D_3}(Y) = 0.97$

Entropy for play (sum of weighted entropies)
 $I_G(Y, \text{outlook}) = \left(\frac{5}{14} \times 0.97 + \frac{6}{14} \times 0 + \frac{5}{14} \times 0.97 \right) = 0.94$

weighted entropy after D_1, D_2, D_3 .

$$= \left(\frac{5}{7} \times 0.97 - 0.94 \right) = 0.6$$

So, we have

parent D has child nodes D_1, D_2, D_3

$I_H = \left(\frac{|D_1|}{|D|} H_{D_1}(Y) + \frac{|D_2|}{|D|} H_{D_2}(Y) + \frac{|D_3|}{|D|} H_{D_3}(Y) \right) - H_D(Y)$

total points in D_1 H_{D_1} pt $\sim D$

Building a decision tree (Gini Impurity)

This is similar to entropy. But calculation of entropy involves log function and calculation of Gini Impurity involves square term. And we know that log takes more computation time than square. So, Gini Impurity is much often used in sklearn.

Gini Impurity $\hookrightarrow I_{Gin}(Y) = 1 - \sum_{i=1}^K (P(S_i))^2$

Case 1: $P(S_+) = 0.7$
 $P(S_-) = 0.3$

$I_{Gin}(Y) = 1 - (0.7^2 + 0.3^2) = 0.5$

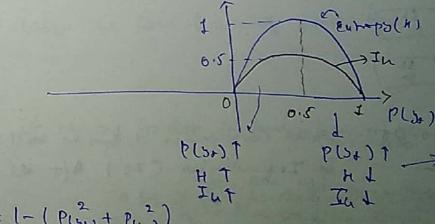
$H(Y) = 1$

$Y \rightarrow Y_1, Y_2, Y_3, \dots, Y_K$

Case 2: $P(Y_+) = 2$
 $P(Y_-) = 0$

$I_{Gin}(Y) = 1 - (1^2) = 0$

$H(Y) = 0$



$I_G = 1 - (P(S_+)^2 + P(S_-)^2)$

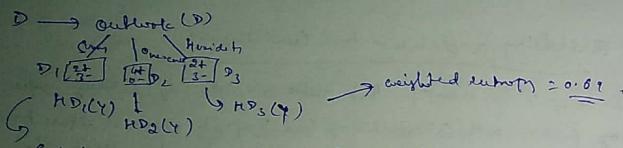
$H(Y) = -P(S_+) \log(P(S_+)) - P(S_-) \log(P(S_-))$

$P(S_+) \uparrow \rightarrow$ So entropy and Gini Impurity between similar.

Building a decision tree/Constructing a DT

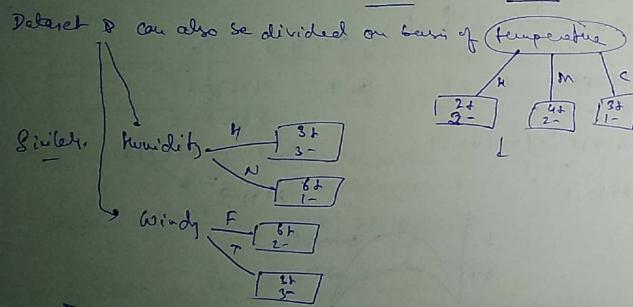
Suppose we have a dataset with 9 (true y_+) and 5 (false y_-). So, Entropy. $H(Y) = 0.94$.

This dataset is further divided on basis of outlook.



Entropy based on dataset D_1, D_2, D_3

$$\text{So, Information Gain } Ig = 0.25 = 0.94 - 0.69$$



$$I_a(Y = \text{outlook}) = 0.25$$

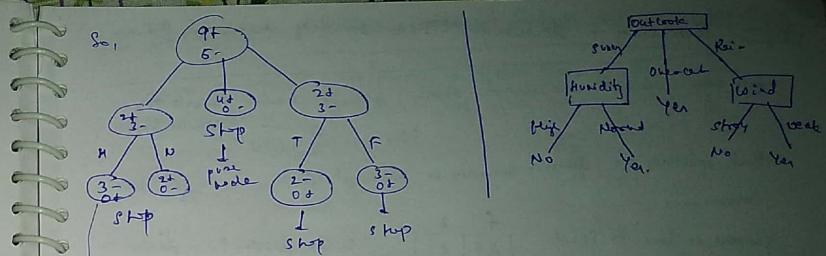
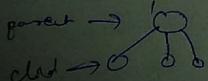
$$I_a(Y = \text{temp}) = -$$

$$I_a(Y, \text{humidity}) = -$$

$$I_a(Y, \text{wind}) = -$$

proves Correspondingly to this we will have to Information Gain (Ig) A/T all dataset,
So, we will choose that ~~one~~ on root node which has highest Ig .

$$I_a(Y_f) = \text{entropy at parent level} - \text{weighted entropy at child level.}$$



So, recursively breaking each node using Ig as the criterion.

So when parent node is decided, to decide the next node again, Ig is used.

Once we reach at a point when we have all the same point in a node left → we stop growing our tree.
After we stop going further.

So, we will stop making decision tree when,

(1) pre-node → stop growing the node

(2) when we have lot of point then stop growing nodes

(3) If we are too deep in tree, then we can stop growing nodes

So, now after going through a lot of if else condition,
Suppose originally we have $n=10$ k

at last point found = 2

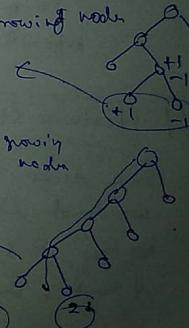
Now it's noise, this can be due because Considering this we are overfit.

→ depth of tree ↑ ; overfit ↑ (few points)

depth is small → underfit

DT-hyperparameter = depth

Can be found by -
cross validation.



Building a decision Tree splitting numerical features.

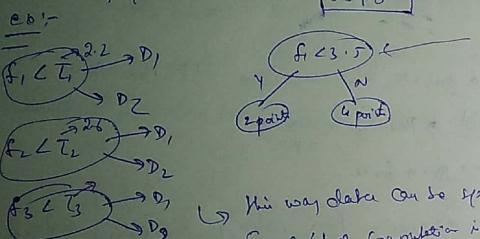
till now we have discussed about discrete random variable or Categorical feature.
This was easy to split.

But splitting the dataset according to numerical feature is a bit like comparing.

Suppose we have to split

f_i	y
2.2	1
2.6	0
3.5	0
3.8	0
4.6	1
5.3	0

Here we have to split data at every point.



This way data can be splitted into ways, so a lot of computation is required.

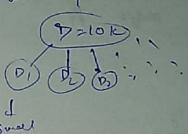
So, one we did it, we take the feature with most information gain in node

Feature Standardization.

Since decision tree is not distance based so unlike kNN, SVM, we don't do feature standardization or feature normalization here. We here in DT each value of feature is taken as threshold and if number is sorted is not, DT only care about k^* .

Building a decision Tree Categorical features with many possible values

Suppose, we have Zipcode/Pincode as feature then at this point making a DT over pincode is not good because there will be a lot of sparsity.



Zipcode	$y \in \{0, 1\}$
B1	0
B2	0
B3	0
B4	1
B5	1

So, if we convert these zipcodes into categorical even then won't solve the problem.

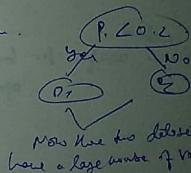
So, one task is change zipcode into numerical feature as it is

$$\text{Suppose } P(\text{3}) = (\text{P}_j) = \frac{19}{20}$$

$$\text{prob of get 3 when zipcode} = P_j$$

this will convert zipcode into discrete numerical value.

So, Now, we can repeat and do further process.

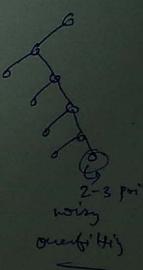
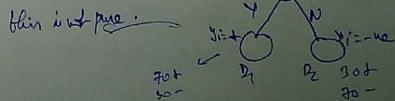


Now we do this for all have a large number of values

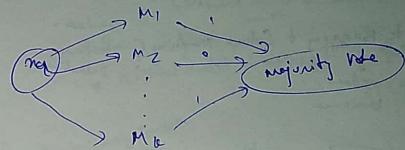
Overfitting and Underfitting

when depth of tree increases (\uparrow) \rightarrow possibility of very few points at leaf nodes increases. \rightarrow Interpretability of model also decreases because there will be a lot of if else condition depth $\uparrow \rightarrow$ overfitting increases.

when, depth = 1 = decision stump (like shooting cricket)



Here, each model M_i has seen a different subset of data.



Suppose when we have binomial classification, when different model predicts different value, we take the majority vote value. When, it is a regression problem, we take the mean or median of all the predicted output by different models.

Bootstrapping (Bagging) is used to reduce Variance.

Here, if we reduce the number of samples from total data then it won't effect much.

For, Bagging can reduce Variance in a model without impacting the bias.

$$\text{model error} = \text{Bias}^2 + \text{Var}$$

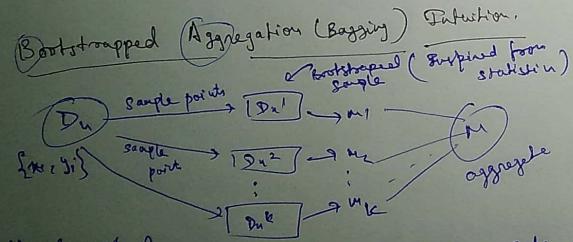
If base model (M_i) = low bias, high var model reduced by bagging.

Ex., in DT of depth (large)

high variance
low bias

Can be solved by bagging

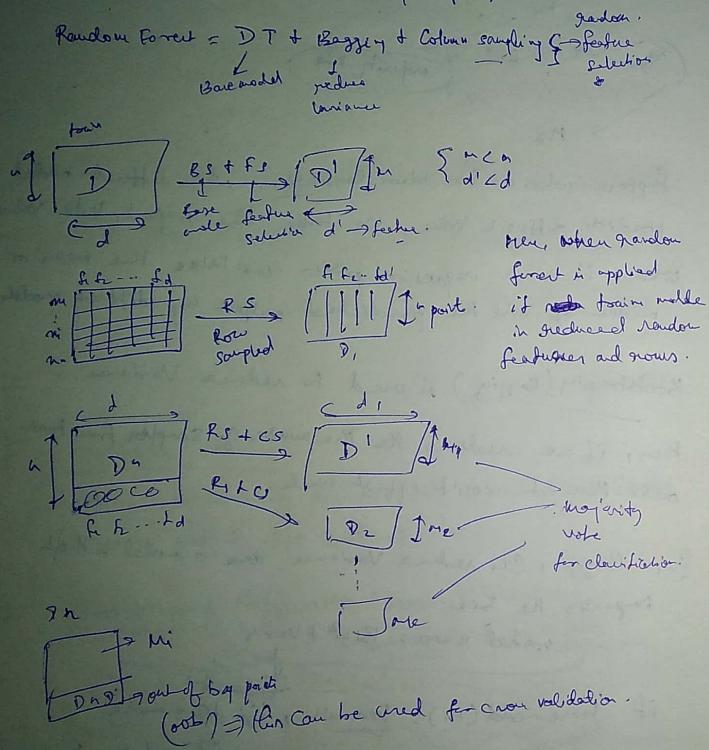
RF



Here, the concept of bagging is → Sample D_n will ~~be~~ different sets with replacement, then train all ~~the~~ the sets with different models, and then finally aggregate the model.

Random Forest and their Construction:

Random forest is also a bootstrap sampling.



Bias-Variance tradeoff

As discussed RF reduces variance.

→ less bias because base learners (M_i 's) are low bias.

actual model is aggregate of all the models.

$$M = \text{agg}(M_1, M_2, M_3, \dots, M_K)$$

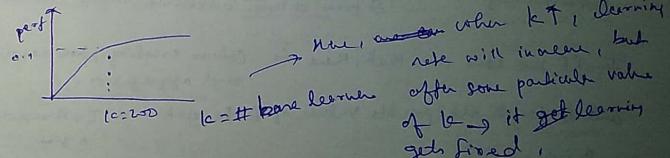
Since, bias of ~~base~~ discrete models (M_i) is small \therefore bias of stacked model is also small $\text{bias}(M) = \text{bias}(M_i)$

$D_n \xrightarrow{\substack{RS (20\%) \\ CS (20\%)}} (D')$

Row sampling ratio = $\frac{R_S}{CS} = \frac{M}{n} = \text{Rows SR}$.

Column sampling ratio = $\frac{CS}{d} = \frac{C}{d} = \text{Col. CR}$.

When we sample rows and columns later say for 20% data and 40% data, then there is very less chance that model will change when training data changes. \therefore Col. CR $\downarrow \rightarrow$ low variance
rows CR $\downarrow \rightarrow$ var \downarrow

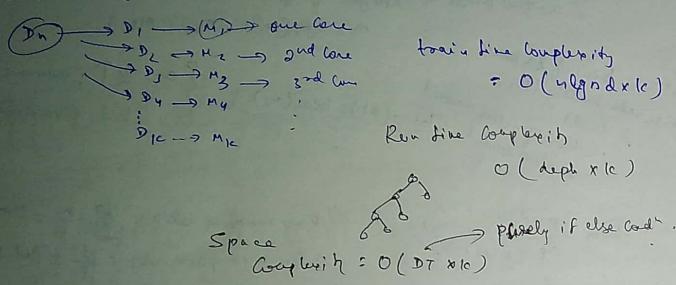


For, $(Rows SR, Col. CR)$ can be calculated by Grid Search.

Train and Runtime complexity.

RF with k -base learners (DT)

Since most computers have different cores, so, RF is designed to perform ~~parallelly~~ parallelly on individual core.



Extremely randomized trees

It is somewhat similar to Random forest.

→ ExtraTree Classifier \Rightarrow from sklearn.model_selection
ExtraTree Regressor \Rightarrow from sklearn.ensemble import ExtraTreeClassifier

In RF we have seen that, there is column sampling, row sampling, and aggregation.

and it finds all possible values of f_i to determine $T \rightarrow$ threshold

Part f_i

1	2	3	4
2	3	4	5
3	4	5	6
4	5	6	7

check for every threshold

But in ExtraTree:
There is column sampling + S.S + agg + randomization when selecting "T"
it selects threshold randomly.

reduces variance
better than
RF.

Random Tree cases.

\hookrightarrow DT \rightarrow don't handle large feature dimension.
RF

In DT \rightarrow bias variance tradeoff depends on depth.
But here in RF \rightarrow it can be controlled by $k = \#$ base learners.

For feature importance

DT \rightarrow It gives feature importance by overall reduction in Entropy & min entropy (J_m) because of ~~it's~~ only one tree in DT.
(model)

But in RF \rightarrow overall reduction in Jm because of feature variance level of each of M_i 's.

Boosting Intuition.

In Bagging we have high variance, low bias base model
+ randomization + aggregation
 \hookrightarrow CS (RS) \hookrightarrow reduce variance.

But in Boosting, we have low variance, high bias + additively combine.
 \hookrightarrow reduces bias while keeping the var low.

Core idea \Rightarrow reduce bias.

$$D_{\text{train}} = \{x_i, y_i\}_{i=1}^n$$

high bias \rightarrow high train error.

\hookrightarrow e.g. DT which has we fit one tree a model, then we get tree errors 2 or 3 depth.

② $D_{\text{train}} \rightarrow$ Model 1 \rightarrow bye train error.
 $(x_i, y_i)_{i=1}^n \hookrightarrow y_i = h_0(x)$

③ $y_i - h_0(x_i) = \text{error}_1 \rightarrow$ simple difference error

Step ④ \hookrightarrow train again a model with feature x_i and error_1 position.

$$(h_1(x_i))_{i=1}^n \hookrightarrow y_i - h_0(x_i) = \text{error}_2$$

$$\text{error}_2 = y_i - h_0(x_i)$$

model of end stage 1 = $F_1(x) = \alpha_0 h_0(x) + \alpha_1 h_1(x)$
 ↳ weighted sum of two base models.

$$(3) \{x_i, \text{error } e_i\} \rightarrow (\underline{\alpha_2}, \underline{\alpha_2})$$

$$\hookrightarrow y_i - F_1(x_i)$$

$$F_2(x) = \alpha_0 h_0(x) + \alpha_2 h_2(x) + \alpha_1 h_1(x)$$

and model

So, at the end of stage k .

$$F_k(x) = \sum_{i=0}^k \alpha_i h_i(x)$$

→ trained to fit the residual error at the end of previous stage
 ↳ additive weight model

$$h_i(x) \in \{x_i, \text{error}\}$$

↳ residual error at end of stage $(i-1)$

So, with reducing training error in each stage, we end up reducing the residual train error. So, \downarrow error

Gradient Boosted DT (GBDT) → also used in internet companies.

Residuals, Loss function and Gradients

$$F_k(x) = \sum_{i=0}^k \alpha_i h_i(x) \quad \text{where } h_i(x) \text{ depends on residual error.}$$

$$\text{residual at end of stage } k \rightarrow e_i = y_i - F_k(x_i)$$

residual

Loss minimization can be done by different ways for different model

logistic-reg → classif

Linear Reg → Squ. loss → visualization of Hinge loss

SVM → hinge loss. Because will take this.

$$L(y_i, F_k(x_i)) = (y_i - F_k(x_i))^2 \quad \{ \text{let } f_k(x_i) = z_i \}$$

$$\begin{aligned} \text{Slope, } \frac{\partial L}{\partial F_k(x_i)} &= \frac{\partial}{\partial z_i} = \frac{\partial}{\partial z_i} (z_i - z_i)^2 \\ &= (-1) + 2z_i(z_i - z_i) \\ &= -2(y_i - z_i) \end{aligned}$$

$$\frac{\partial L}{\partial F_k(x_i)} = \frac{\partial}{\partial F_k(x_i)} (y_i - z_i)$$

negative gradient → residual

So, negative gradient ≈ residual.

↳ pseudo residual

↳ so, residual can be replaced by negative gradient

So, pseudo residual lets us have any loss function which is differentiable.

But Gradient Boosting can minimize any losses which is differentiable.

But models like RF can only minimize hinge loss.

$$\text{So, } (x_i, \text{error}) \rightarrow (M_i, \text{error})$$

↳ pseudo residual $\frac{\partial L}{\partial F_k(x)}$

Gradient Boosting

Input training set $\{(x_i, y_i)\}_{i=1}^n$, a differentiable loss function $L(y_i, f(x))$

No. of iterations M .

Algorithm (Coordinate Descent)

① Initialize model with a constant value.
 zeroth model $\leftarrow F_0(x) = \arg\min_{\gamma} \sum_{i=1}^n L(y_i, \gamma)$
 just γ is mean of y_i

about γ is mean of y_i

$m=k$

② From $m=1$ to M

③ Compute so-called pseudo-residuals.

$$g_{im} = -\left[\frac{\partial L(y_i, F_{m-1}(x_i))}{\partial F_m(x_i)} \right] \quad \text{for } i=1, \dots, n$$

$$F_m(x) = F_{m-1}(x)$$

④ Fit a base learner (e.g. tree) $h_m(x)$ to pseudo-residuals.
 i.e. train it using the training set $\{(x_i, g_{im})\}_{i=1}^n$.

⑤ Compute multiplier γ_m by solving the following one-dimensional optimization problem.

$$\gamma_m = \arg\min_{\gamma} \sum_{i=1}^n L(y_i, F_{m-1}(x_i) + \gamma h_m(x_i))$$

⑥ Update the model:

$$F_m(x) \leftarrow F_{m-1}(x) + \gamma_m h_m(x)$$

⑦ Compute $F_M(x)$

Since V is very small, so $\underbrace{V \gamma_m h_m(x)}$

\hookrightarrow is also small so effect of this on model won't increase the variance.

Then γ and V can be searched by coordinate-mining CV.

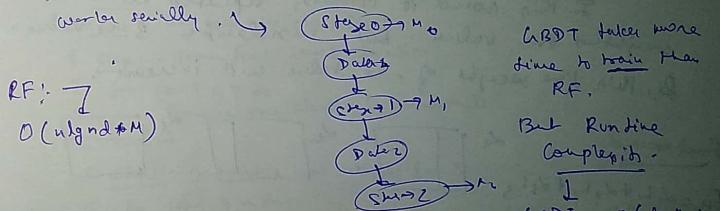
when M is large \hookrightarrow model is more likely to overfit

$\text{say } V=200$

Train and Runtime Complexity

We have seen that we can grow different sets of models parallelly in RF.

But in gradient boosting we can't parallelize it, because it works serially.



GBDT takes more time to train than RF.

But Runtime Complexity -

$O(\text{depth} * M)$

\downarrow
 GBDT $\rightarrow O(\text{depth} * M)$
 \downarrow
 Space time $\leftarrow O(\text{depth} * M + \text{depth} * M)$
 \downarrow
 $O(\text{depth} * M + \text{depth} * M)$

So, slow learning algorithm.

most of the internet protocols

Companies like Google, use GBDT for low space complexity.

Regularization by Shrinking.

In gradient boosting there is a very high chance of overfitting a model.

So, by shrinking we saves the model from overfit.

$$F_m(x) = h_m(x) + \sum_{m=1}^M \gamma_m h_m(x)$$

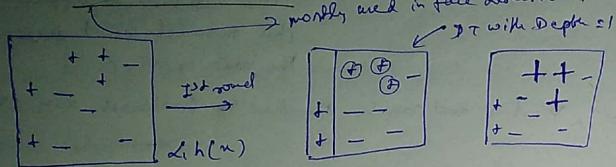
bare-model $\uparrow \rightarrow$ overfit $\uparrow \rightarrow$ variance \uparrow

∴ Shrinkage is added

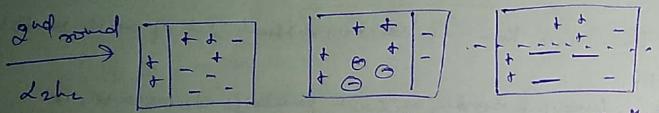
$$F_m(x) = F_{m-1}(x) + V \cdot \gamma_m h_m(x) \quad \text{where } 0 < V \leq 1$$

(Learning rate).

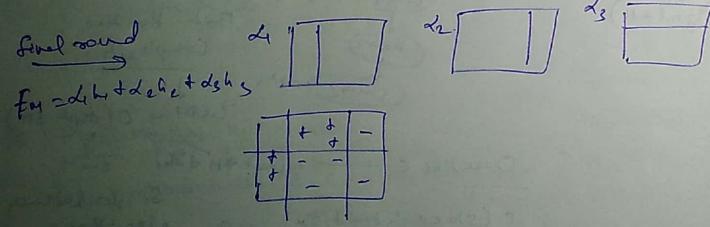
AdaBoost Geometric Intuition.



Working of AdaBoost - At first it will separate 2 positive pts.
So, if misclassifies 3 +ve pts. So, next time it increases
the weights of 3 +ve pts.

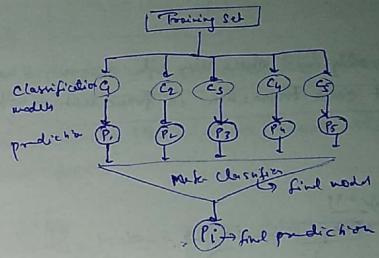


In this round it correctly separates all three weighted +ve values. but misclassifies 3 -ve value.
So, their true weight of -ve value will increase.



Stacking Models

In boosting we build models based on errors, but here in stacking different models are trained independently (parallelly).
more different the model, more better is stacking.



Suppose we have to train a model, and classify \mathbf{x}_i .
Then, \mathbf{x}_i will be passed to all the different parallel models, then,
the predicted value for (P_1, \dots, P_6) is passed to meta classifier
then final prediction is made.

Algorithm:

Input → training data $D = \{x_i, y_i\}_{i=1}^n$ ($x_i \in \mathbb{R}^m$, $y_i \in \mathcal{Y}$)
Output → An ensemble classifier H .

1. Step 1 → learn first level classifier
2. for $i=1$ to T do
 3. learn a base classifier h_i based on D
 4. end for
 5. Step 2 → Construct a new datasets from D
 6. for $i=1$ to n do
 7. Construct a new data set $K_{i,i}$ contains $\{x_i, y_i\}$, where $x_i^j = \{h_1(x_i), h_2(x_i), \dots, h_T(x_i)\}$
 8. end for
 9. Step 3 → learn a second-level classifier
 10. learn a new classifier h' based on the newly constructed dataset
 11. $H_{i,i}(x) = h'(h_1(x), h_2(x), \dots, h_T(x))$

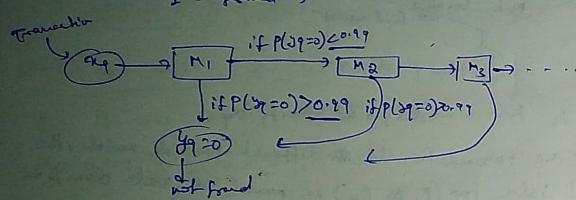
Pipe instead of extend

for stacking.

Cascading Models

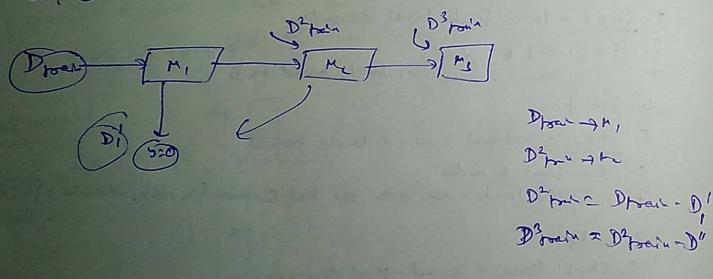
↳ Cascading model is generally used where very high accuracy is needed.
like in ~~the~~ credit card transaction is fraudulent or not.

Suppose $\{0 \rightarrow \text{not found}, 1 \rightarrow \text{found}\}$



When the model fails to classify then a final call is made to customer about the transaction.

↳ Cascading is used when the cost of making a mistake is high.
Ex: Cancer or not Cancer.



Deep learning

History of Networks and Deep learning.

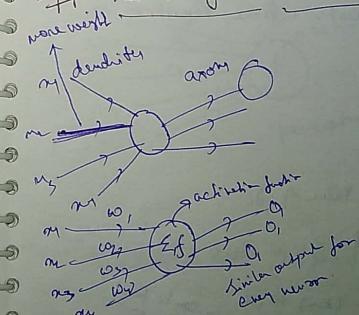
Perceptron \rightarrow 1957 $\xrightarrow{\text{by Rosenblatt}}$

In 1986 Hinton & others added Backpropagation algo to NN
This is a simple chain rule.

After 1986 there was a lot of hype about AI. But due to less computation power and data, NN was not so much popular.
But SVM, logistic regt was.

NN gets its popularity in 2012 when it perform extremely well in ImageNet Competition.

How biological Neuron Work



{ activated }
fire

When in ~~the~~ neuron, there is enough data as input then it fires.

Data with more weight are more important.

$$o_i = f(w_1x_1 + w_2x_2 + \dots + w_nx_n)$$

$$o_i = f(\sum_{i=1}^n w_i x_i)$$

output activation weight input

Growth of Biological neural networks.

When a child gets birth it has less interconnection in brain, but as he grows up from 0 to 6 year old. It learns most of the things, like, walking, seeing different things etc. As he grows older, again, interconnection of NN in brain decreases.

The neuron which is thicker has higher weight.



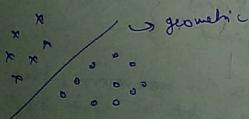
Learnings happen when, neuron connects with edges.

Given, a Japanese student ^{learns} Japanese and an Indian learns Hindi, so, this also creates different weights in neurons.

As we are studying at night right now, it also creates some weights -

Diagrammatic Representation of Logistic Regression and Perceptron.

We have seen LR earlier, it creates a linear separation.



LR: $x_i \rightarrow \hat{y}_i \rightarrow \text{predicted value of } y_i$

$$\hat{y}_i = \text{Sigmoid}(w^T x_i + b) \quad x_i \in \mathbb{R}^d$$

$$D = \{(x_i, y_i)\}$$

Train LR $\rightarrow w, b$ can be found by SGD

WEIRD
HELP

In logistic regression sigmoid acts as an Activation function.

But in perceptron if $w^T x_i + b > 0 \rightarrow 1$
 $w^T x_i + b < 0 \rightarrow 0$ Activation function

In perceptron, if $w^T x_i + b > 0 \rightarrow 1$ reward network fires otherwise not.

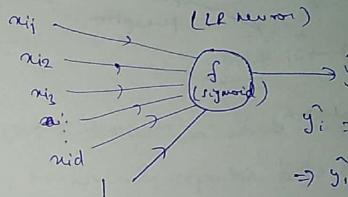
$$\text{In LR, } \hat{y}_i = \text{Sigmoid}\left(\sum_{j=1}^d w_j x_{ij} + b\right)$$

$$x_i = [x_{i1}, x_{i2}, x_{i3}, \dots, x_{id}]$$

$$w = [w_1, w_2, w_3, \dots, w_d]$$

In perceptron.

$$o = f\left(\sum_{j=1}^d o_j x_{ij}\right)$$

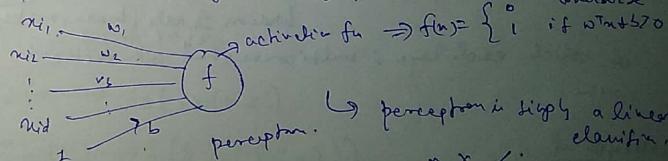


$$\hat{y}_i = f(w_1 x_{i1} + w_2 x_{i2} + \dots + w_d x_{id} + b)$$

$$\Rightarrow \hat{y}_i = f(w^T x_i + b)$$

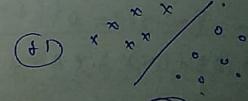
So, LR can be represented as neural network.

For, train a NN, weights on edges/vertices.



perception.

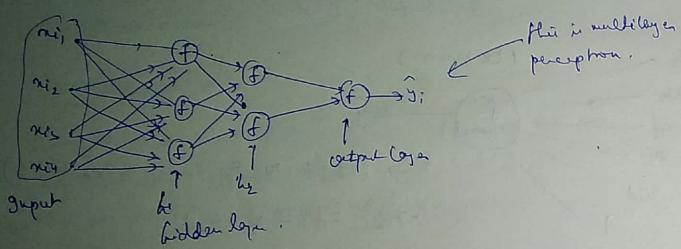
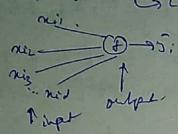
perception is highly a linear classifier.



But in LR
Squashing is done by sigmoid

Multi-Layered Perceptron (MLP)

Till now we saw Single layer perceptron \rightarrow LR



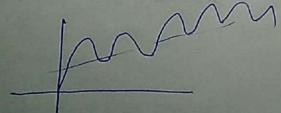
Q Why should we care about MLP?

C Because it is ~~inspired~~ inspired from neuroscience, brains of rats, human, monkeys, in which each layer is multi-connected with different layers.

Mathematical representation $\{x_i, y_i\} = D$
 $y = f(x) \text{ SFR}$

Suppose we have to make prediction $y = f(2 \sin(x^2) + \sqrt{5}x^5)$

So, if we do LR or single perceptron
 Then, this model will make a lot of errors.



But, this model can be solved easily by multilayer perceptron.

$$\text{Given } F(x) = 2 \sin(x^2) + \sqrt{5}x^5$$

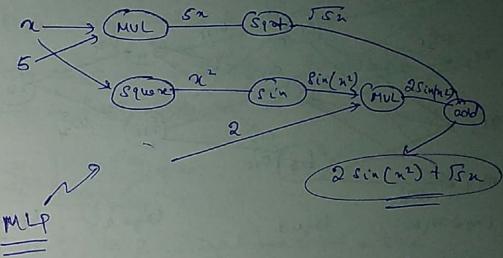
$$f_1 \rightarrow \text{add}()$$

$$f_2 \rightarrow \text{sqrt}()$$

$$f_3 \rightarrow \text{sqrt}()$$

$$f_4 \rightarrow \sin()$$

$$f_5 \rightarrow \text{multi}()$$



MLP is trained when we need to make prediction on non-linearly separable data then MLP can be easily used, but this can easily overfit.

MLP is basically a composite function, $f_1 \circ f_2 \circ f_3 \circ \dots$

$$F(x) = 2 \sin(x^2) + \sqrt{5}x^5$$

$$\sqrt{5}x^5 \rightarrow f_3(f_5(x, 5))$$

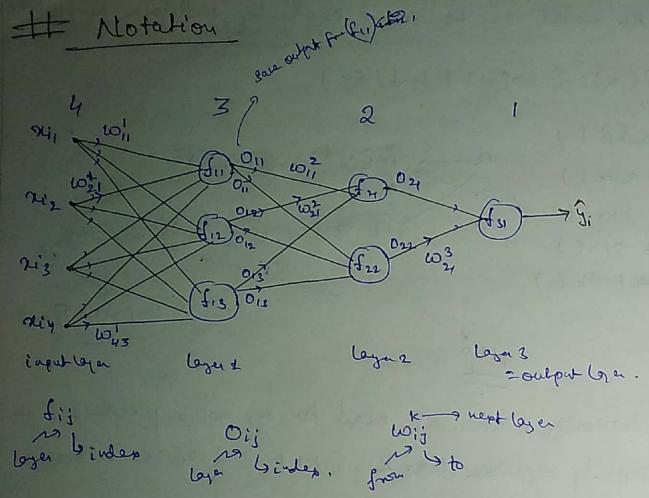
$$2 \sin(x^2) \rightarrow f_1(2, f_4(f_2(x)))$$

$$\therefore F(x) = f_1(2, f_4(f_2(x)), f_3(f_5(x, 5)))$$

⇒ By using multi-layered structure we can come up with complex math formulae to solve your problem.

MLP → graphical way to represent $f \circ g \circ h$
 ⚡ powerful model

Notation



This is a fully connected layer.

Since, from 1st to 2nd layer, there are total $4 \times 3 = 12$ connections.

Similarly, from 2nd to 3rd $\rightarrow 3 \times 2 = 6$ connections.

So, these weights can be represented as matrices.

$$W^1 = \begin{bmatrix} w_{11}^1 & w_{12}^1 & w_{13}^1 \\ w_{21}^1 & w_{22}^1 & w_{23}^1 \\ w_{31}^1 & w_{32}^1 & w_{33}^1 \\ w_{41}^1 & w_{42}^1 & w_{43}^1 \end{bmatrix}$$

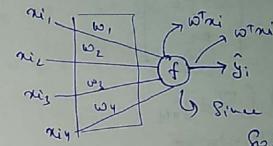
Similarly,
 w^2 can be made.

Training a single-neuron model

Here with training we find the best edge weight using D_{train} .

Perceptron & Logistic Reg → Single Neuron models for classification.
Linear Reg → Single Neuron model for Reg.

As we know in Linear Reg → function f is identity function because input given is same as output.



Since input and output is same
so, f is identity function.

$$\text{Linear Reg: } \hat{y}_i = \sum_{j=1}^d w_{ij} x_{ij} \quad \left\{ \begin{array}{l} \text{linear optimization, } m \in \mathbb{R}^d \\ \hat{y}_i = w^T x_i \end{array} \right.$$

$$\text{Optimization: } \min_{w_i} \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \text{regularization}$$

$$\min_{w_i} \sum_{i=1}^n (y_i - w^T x_i)^2 + \|w\|_2^2$$

Steps to train

(1) Define a loss-function.

$$L = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \text{reg} \quad \rightarrow \text{for all training data}$$

$L_i = (y_i - \hat{y}_i)^2 \rightarrow \text{sum of single train data}$

(2) Optimization.

$$\min_{w_i} \sum_{i=1}^n (y_i - w^T x_i)^2 + \text{reg} \quad \rightarrow \text{neglecting regularization for simplicity.}$$

$\hat{y}_i = f(w^T x_i)$

↳ Identity for linear reg.
Logistic for log. reg.

$$\omega^* = \underset{\omega}{\operatorname{argmin}} \sum_{i=1}^n (y_i - f(\omega^T x_i))^2 + \lambda \|\omega\|_2^2$$

(3) Solve the optimization problem.

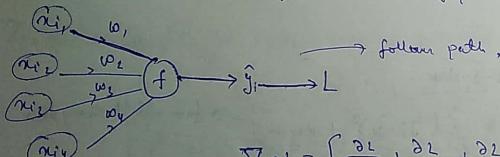
(a) initialize weight $\omega_0 \rightarrow$ random

(b) Compute $\nabla_{\omega} L = \begin{bmatrix} \frac{\partial L}{\partial \omega_0} \\ \frac{\partial L}{\partial \omega_1} \\ \vdots \\ \frac{\partial L}{\partial \omega_n} \end{bmatrix}$ $\omega \in \mathbb{R}^n$

(c) $\omega_{\text{new}} = \omega_{\text{old}} - \eta (\nabla_{\omega} L)_{\text{old}}$
 update weight ω_i \rightarrow learning rate
 $(\omega_i)_{\text{new}} = (\omega_i)_{\text{old}} - \eta \left[\frac{\partial L}{\partial \omega_i} \right]_{\text{old}}$
 for $i = 1 \text{ to } n$

Gradient Descent $\nabla_{\omega} L \rightarrow x_i's \& y_i's \rightarrow$ all

SGD $\rightarrow \nabla_{\omega} L = \frac{\text{for small batch of points SGD.}}{\text{for small batch of points SGD.}}$



$$\nabla_{\omega} L = \left[\frac{\partial L}{\partial \omega_1}, \frac{\partial L}{\partial \omega_2}, \frac{\partial L}{\partial \omega_3}, \frac{\partial L}{\partial \omega_4} \right]^T$$

$$\frac{\partial L}{\partial \omega_i} = \frac{\partial L}{\partial f} \frac{\partial f}{\partial \omega_i} \rightarrow \text{chain rule of differentiation.}$$

We have, $L = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda \|\omega\|_2^2$.

$$L = \sum_{i=1}^n (y_i - f)^2$$

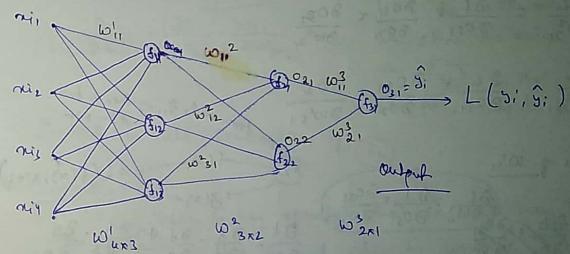
$$\therefore \frac{\partial L}{\partial f} = - \sum_{i=1}^n 2(y_i - f) \omega_i \quad \left| \frac{\partial f}{\partial \omega_i} = \omega_i \right. \quad (f = \omega^T x)$$

$$\therefore \frac{\partial L}{\partial \omega_i} = \sum_{i=1}^n 2x_i(y_i - \hat{y}_i)$$

Training a MLP Chain rule.

In previous video we learned about training a single neuron model.

Given, $D = \{(x_i, y_i)\}$



Now, we need to determine

$12 + 6 + 2 = 20$ weights and update them.

$$(1) L = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda \|\omega\|_2^2 \quad \left(\begin{array}{l} \text{actual} \\ \text{predicted} \\ \text{sum} \end{array} \right) \sum_{i,j=1}^n (\omega_{ij})^2$$

(2) SGD or GD

(a) initialize ω_{kij} randomly

(b) $(\omega_{kij})_{\text{new}} = (\omega_{kij})_{\text{old}} - \eta \left(\frac{\partial L}{\partial \omega_{kij}} \right)$
 learning rate.

(c) perform update till convergence

Now, to update weight.

$$\frac{\partial L}{\partial w_{11}^2} = \frac{\partial L}{\partial o_{21}} \times \frac{\partial o_{31}}{\partial w_{11}^2} \quad \leftarrow \text{chain rule}$$

$$\frac{\partial L}{\partial w_{21}^3} = \frac{\partial L}{\partial o_{31}} \times \frac{\partial o_{31}}{\partial w_{21}^3}$$

$$\frac{\partial L}{\partial w_{11}^2} = \frac{\partial L}{\partial o_{21}} \times \frac{\partial o_{21}}{\partial o_{22}} \times \frac{\partial o_{21}}{\partial w_{11}^2}$$

$$\frac{\partial L}{\partial w_{21}^3} = \frac{\partial L}{\partial o_{21}} \times \frac{\partial o_{21}}{\partial o_{22}} \times \frac{\partial o_{21}}{\partial w_{21}^3}$$

$$\frac{\partial L}{\partial w_{31}^4} = \frac{\partial L}{\partial o_{21}} \times \frac{\partial o_{21}}{\partial o_{22}} \times \frac{\partial o_{21}}{\partial w_{31}^4}$$

for update of w_{11}^1 ,

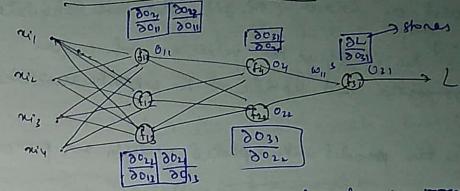
$$\frac{\partial L}{\partial w_{11}^1} = \frac{\partial L}{\partial o_{21}} \times \frac{\partial o_{31}}{\partial w_{11}^1}$$

$$\begin{aligned} \frac{\partial L}{\partial w_{11}^1} &= \frac{\partial L}{\partial o_{21}} \times \frac{\partial o_{31}}{\partial w_{11}^1} \\ &\quad \text{from node } f(u) \rightarrow h(f(u), g(u)) \\ &\quad \text{then in the division} \end{aligned}$$

$$\frac{\partial L}{\partial w_{11}^1} = \frac{\partial o_{31}}{\partial o_{21}} \cdot \frac{\partial o_{21}}{\partial o_{22}} \cdot \frac{\partial o_{21}}{\partial w_{11}^1} + \frac{\partial o_{31}}{\partial o_{22}} \cdot \frac{\partial o_{22}}{\partial o_{21}} \cdot \frac{\partial o_{21}}{\partial w_{11}^1}$$

$$\therefore \frac{\partial L}{\partial w_{11}^1} = \frac{\partial L}{\partial o_{31}} \cdot \left\{ \frac{\partial o_{31}}{\partial o_{21}} \cdot \frac{\partial o_{21}}{\partial o_{22}} \cdot \frac{\partial o_{21}}{\partial w_{11}^1} + \frac{\partial o_{31}}{\partial o_{22}} \cdot \frac{\partial o_{22}}{\partial o_{21}} \cdot \frac{\partial o_{21}}{\partial w_{11}^1} \right\}$$

Training an MLP memorization



Memorization is a technique from dynamic programming.

As we have seen earlier, while applying chain rule,

$$\frac{\partial L}{\partial w_{11}^1} = \frac{\partial L}{\partial o_{21}} \cdot \frac{\partial o_{31}}{\partial w_{11}^1}$$

This diagram has used a lot of time. So, with memorization, we store the value of $\frac{\partial L}{\partial o_{31}}$ in f_3 . So that we need not to recompute it each time.

Similarly, we will save $\frac{\partial o_{31}}{\partial o_{21}}$ to f_4 .

Memorization

↳ if there is any operation that is used many times repeatedly, it's a good idea to compute it once, store it and reuse it.

↳ it boosts up the speed. (uses slightly more memory)

For here, we need to store only f values, and can use a lot of time.

Chain rule + memorization
Back propagation.

Back Propagation

Finally we train our model on dataset $\{x_i, y_i\}$. Then we calculate the loss $L(y - \hat{y})$, if these two values are not similar then we backpropagate and update all the weights associated with neurons.

We again train the model with updated weights and process

Continues, till we get ($W_{\text{new}} = W_{\text{old}}$)

Step 1 Data $\Rightarrow D\{x_i, y_i\}$

- ① initialize W_{ij}^k, b^k
- ② for each x_i in D
 - a) pass x_i forward through the network \rightarrow forward propagation
 - b) Compute $L(y_i, \hat{y}_i)$
 - c) Compute all the derivatives chain rule & memoization
 - d) Update the weights from end of the network to start.

Cat say, from previous network.

$$(w_{ij}^k)_{\text{new}} = (w_{ij}^k)_{\text{old}} - \eta \left(\frac{\partial L}{\partial w_{ij}^k} \right)$$

Q: what even activation function we take it should be easily differentiable, because if not then, we can't use chain rule.

- ③ repeat step 2 till convergence

Convergence occurs when

$$(w_{ij}^k)_{\text{new}} \approx (w_{ij}^k)_{\text{old}}$$

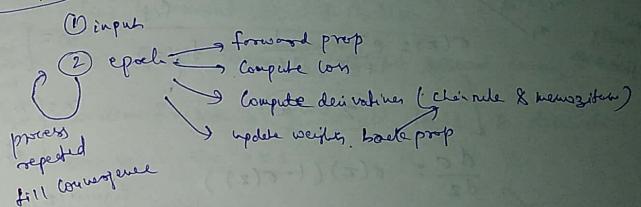
Since, suppose, we send the input data and get the output and updated the weight

\hookrightarrow This is one epoch.

Now, we again send the updated weight forward and this step

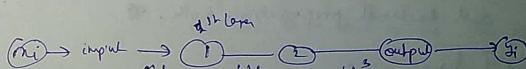
Continues \hookrightarrow This is 2nd epoch, then 3rd epoch...

Back prop \rightarrow chainrule & memoization



P.V.T \Rightarrow In backprop activation func must be differentiable.

It speeds up the training of the NN.



For, on a input we can use full data $\overset{\text{at a time}}{\text{also, but can take a}}$ batch of neurons. (mini-batch SGD)

Or, we can take each point at a time, this is (SGD)

Or, we can take batch of points at a time (mini-batch SGD)

Suppose, 1000 points in D \downarrow most popular approach.

mini-batch = 100

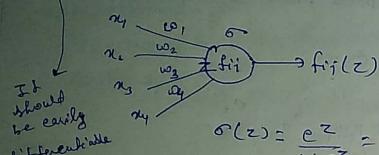
$$\text{epoch} = \frac{100000}{100} = 1000$$

- ② for each epoch mini batch of size 100,
do, forward prop + loss + $\frac{\partial L}{\partial w_{ij}}$ + update & backprop.

Activation Function.

The concept of activation function was first introduced in 1980 & 90's.

Generally, Sigmoid & tanh are mostly used.



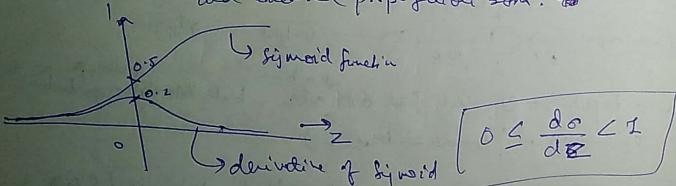
If it should be easily differentiable

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

$$z = w^T x$$

$$\frac{d\sigma}{dz} = \sigma(z)(1-\sigma(z))$$

Here we can see that derivative of sigmoid can be used as a function of sigmoid. So it can be used for forward and backward propagation both.

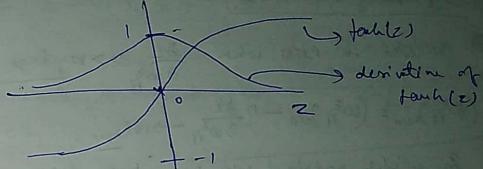


Similarly, for

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

$$\frac{da}{dz} = 1-a^2$$

Derivative of $\tanh(z)$ can also be represented as a function of $\tanh(z)$. So, easy in computation.



Vanishing Gradient Descent Problem.

During 80's and 90's problem of vanishing gradient descent arises in NNs.

As the activation function is sigmoid and tanh, so derivative becomes same value b/w 0 to 1.

So, during weight update.

$$\frac{\partial L}{\partial w_{ij}} = \frac{\partial L}{\partial o_{ij}} \left[\frac{\partial o_{ij}}{\partial o_{ii}} \times \frac{\partial o_{ii}}{\partial o_{ii}} + \frac{\partial o_{ij}}{\partial o_{ii}} \times \frac{\partial o_{ii}}{\partial o_{ii}} \times \frac{\partial o_{ii}}{\partial w_{ij}} \right]$$

This will come out to be very less, nearly zero.

So, $w_{ij, \text{new}} \approx w_{ij, \text{old}}$

So, weight are not updated when size of hidden layer are large.

$$(w_{ij, \text{new}} = (w_{ij, \text{old}})_{\text{old}} - \eta \frac{\partial L}{\partial w_{ij}})$$

$$\text{calc} [0.2 \times 0.1 \times 0.05] = 0.0010 \ll 1$$

So, $\frac{\partial L}{\partial w_{ij}} \rightarrow \text{v.v. small}$
because of chain rule.

Exploding gradient

It occurs in case when $\frac{\partial L}{\partial w_i} \gg \text{very large}$.

$$(w_i)_{\text{new}} = (w_i)_{\text{old}} - \eta \frac{\partial L}{\partial w_i} \rightarrow \text{large}.$$

For which update will have large difference from previous weight.

↳ It also occurs due to chain rule.

$$\frac{\partial L}{\partial w_i} = \frac{\partial L}{\partial z_1} \times \frac{\partial z_1}{\partial w_i} \times \frac{\partial z_2}{\partial w_i} \times \frac{\partial z_3}{\partial w_i} \gg \text{large}.$$

Bias - Variance Tradeoff

In MLP, there is high chance of overfitting.
So, regularization is added.

Hyperparameter tuning is done in
↳ 1 and layers.

① when layer $L \rightarrow$ more weights/parameters
↳ high chance of overfitting.
↳ high variance.

② when layer L (less) \rightarrow higher chance of underfitting,
high bias.

So, in MLP \rightarrow add multiple layers
↳ for optimization \rightarrow regularization is added

$$w \rightarrow L = \sum_{i=1}^n \text{loss}(y_i, \hat{y}_i) + \lambda \sum_{j=1}^L \|w_j\|^2$$

(λ is added to loss)
with λ , regularization creates a loss
of sparsity.

$$\text{So, } L = \sum_{i=1}^n \text{loss}_i + \lambda \sum_{j=1}^L \|w_j\|^2 \text{ as weight}$$

Layer $i \rightarrow$ lesser weight
 $\lambda \uparrow \rightarrow$ variance ↓
layer $\uparrow \rightarrow$ variance ↑.

Deep Multilayer perceptrons 1980s to 2010s

from 1980s to now \rightarrow people name many DL models.
↳ problem was training time
↳ too little data \rightarrow high chance of overfitting
↳ to little computing \rightarrow taking too much time

So, due to these problems, training
in deep net was tough. Since, it
was impossible to handle lots of weights and then to update it.

So, after 2010
↳ we get a lot of data \rightarrow by Google, Facebook etc
(stored as unlabeled)
↳ Computing power increased
↳ a lot of new ideas and algorithms
↳ GPU (GPU for displaying nearly 100 times faster than CPU).

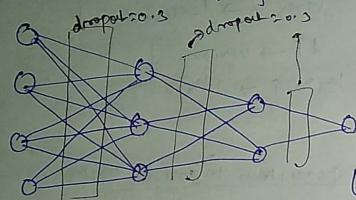
Important
Classical ML was based on theory and then experiments.
But Modern DL is based on my first experiment then theory.

Dropout & Regularization.

The Deep NN \rightarrow easily overfit
 ↳ many layers \rightarrow many weights

So, regularization (L_1 & L_2) can be added to reduce overfitting.

But in 2012 Nitish Srivastava gave the idea of Dropout, which also reduces overfitting in a much more elegant way.



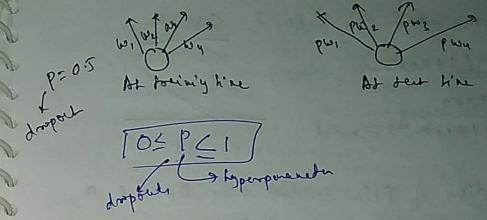
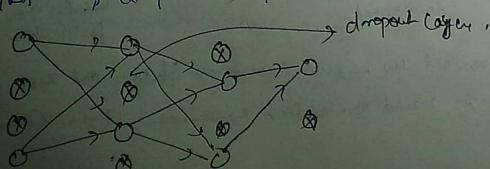
(Let's say dropout = 0.3 means, reduce 30% of weight)

- So, the idea of dropout is as follows:
- randomly reduce the no. of ~~weights~~^{weights (layer)} in each iteration.
 - Then train the model and backpropagate.
 - Then again randomly reduce the no. of connected weights, and the process is repeated.

→ This way we are not training on whole data, but a different sets of reduced ~~weights~~^{weights (layer)}. So, reduces overfitting, like, we train a model in RF on different set.

(ii) → Then during testing we multiply each weight by 0.3 because it was reduced to 30% during training.

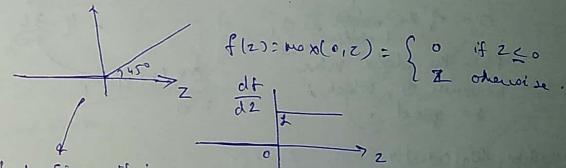
This way dropout works.



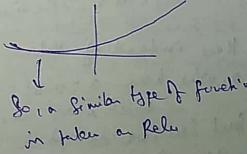
Rectified Linear Units (ReLU)

Since, with sigmoid and tanh function, we had a problem of Gradient Descent.

So, in 2012 ReLU was developed, best activation function till now.



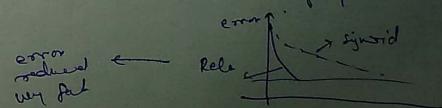
but since it is not differentiable at $z=0$.



So, a similar type of function is taken as ReLU

↳ This derivative is either 0 or 1.
 So, there is no vanishing gradient, or gradient boosting condition with ReLU.
 ↳ when $z > 0$, derivative is 0
 This is called the dead ~~ReLU~~ activation.

In experiment people have seen that with ReLU as activation function, error reduces with less no. of epochs.

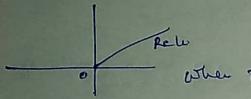


A smooth approximation of ReLU is Softplus.

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

$$f'(x) = \frac{e^{px}}{1 + e^{px}} \rightarrow \text{sigmoid}$$

ReLU Since, derivative of ReLU involves only 0 and 1 so, it is faster, unlike Sigmoid and Tanh.



When $z < 0$, $f(z) = 0$

$$\frac{df}{dz} = 0$$

So, in weight update: $\Delta w_{ij}^k = (w_{ij}^k)_{\text{new}} - (w_{ij}^k)_{\text{old}}$
one of the ~~weights~~ is 0, derivative:
then the whole value becomes 0, which is no good.

$$= 0 \times |x_1| x_1 = 0$$

So, to our concern

we can use Leaky ReLU



$$f(x) = \begin{cases} x & \text{if } x > 0 \\ 0.01x & \text{otherwise} \end{cases}$$

Here instead of Value 0, we have a value which is not 0, but here, we can face problem of vanishing gradient.

- potential problem with ReLU
- Non-differentiable at 0
- Non-zero centered
- Unbounded

but still better than ~~Sigmoid~~ Sigmoid and Tanh.

Weight Initialization.

As we have seen in logistic regression, initialization of weight w_{ij} is random.

$$w_{ij} \sim N(0, \sigma)$$

uniform random

During initialization of weight of MLP

we should follow the following steps.

① initialize $w_{ij}^k = 0$ or any constant value $\forall i, j, k \rightarrow$ v.v. bad.

{ because all the neurons compute the same thing
→ and same gradient update.

② initialize $w_{ij}^k = \text{large -ve number}$

because, for -ve z in ReLU, $f(z) > 0$

* So data normalization is mandatory.

$w_{ij}^k = z - \text{ave value}$
so if $w_{ij}^k = 0$ then $z = \text{ave}$
 w_{ij}^k is generally normalized
in $b/10$ and 1
near centering
variance scaling

So the idea is

① weights should be small (not too small)

↳ all zeros or constant value

↳ good variance

$$\text{So, } \{w_{ij}^k \sim N(0, \sigma)\}$$

↳ so weight should be normalized before
with 0 mean and
constant variance σ^2 .

↳ small

gaussian or
normal distributed

→ So, with lots of research experiments, scientists have come up with lots of different initialization techniques, using fan-in and fan-out (but there is no concrete agreement amongst all researchers).

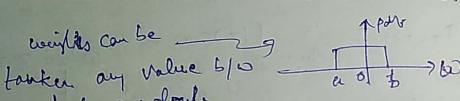


$$\text{fan-in} = 4, \text{ fan-out} = 2.$$

② Uniform Initialization

$$w_{ij}^k \sim \text{Uniform} \left[-\frac{1}{\sqrt{\text{fan-in}}}, \frac{1}{\sqrt{\text{fan-out}}} \right]$$

weights can be taken any value b/w
a and b randomly.



③ Xavier/Glorot Initialization (2010) → better for sigmoid

$$\text{④ Normal } w_{ij}^k \sim N(0, \sigma_{ij}^2) \quad \sigma_{ij}^2 = \frac{2}{\text{fanin} + \text{fanout}}$$

$$\text{⑤ Uniform } w_{ij}^k \sim U \left[-\frac{\sqrt{6}}{\sqrt{\text{fanin} + \text{fanout}}}, \frac{\sqrt{6}}{\sqrt{\text{fanin} + \text{fanout}}} \right]$$

⑥ He Initialization (2015) & ReLU

$$\text{⑦ Normal } w_{ij}^k \sim N(0, \sigma) \quad \sigma = \sqrt{\frac{2}{\text{fanin}}}$$

$$\text{⑧ Uniform } w_{ij}^k \sim U \left[-\frac{\sqrt{6}}{\text{fanin}}, \frac{\sqrt{6}}{\text{fanin}} \right]$$

Batch Normalization.

Suppose we have a fully connected layer.

$$n \rightarrow L_1 \rightarrow L_2 \rightarrow L_3 \rightarrow L_4 \rightarrow L_5 \rightarrow \text{output} \rightarrow f$$

fully connected layer

Since n is generally preprocessed

\rightarrow mean centering

$$D = \{x_i, y_i\}$$

$$\left\{ \hat{x}_i = \frac{x_i - \mu}{\sigma} \right\}$$

for a small change in batch input can lead to large change in deep NN layers.

Since, at each stage, input changes to small small fraction.

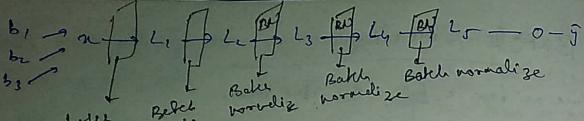
$$\begin{aligned} \therefore \text{if } & 1.01 \leftarrow \text{changed value} \\ \text{first value } & 1.01 \times 1.0201 \times 1.01 \\ & \vdots \\ & = 1.0406 \end{aligned}$$

So, this can cost us a lot (known as internal covariance shift)

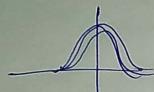
So, we do ~~batch~~ normalization, before every layer to get the same normalized value, to train much better.

$$\begin{aligned} n & \rightarrow L_1 \rightarrow L_2 \rightarrow L_3 \rightarrow L_4 \rightarrow L_5 \rightarrow \text{output} \\ \{x_i\}_{\text{batch}} & \xrightarrow{\text{Batch}} \xrightarrow{\text{Batch}} \xrightarrow{\text{Batch}} \xrightarrow{\text{Batch}} \xrightarrow{\text{Batch}} \xrightarrow{\text{Batch}} \end{aligned}$$

changes in distribution at stage 5.



Now for every batch each layer will encounter similar distribution.



Alg's: Values of α over a graph mini-batch: $B = x_1, \dots, x_m$.

parameters to be learned: γ, β

Output $\rightarrow \frac{1}{m} \sum_{i=1}^m x_i \rightarrow \mu$ mini batch mean

$\sigma^2_B \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu)^2$ minibatch Variance

$\hat{x}_i \leftarrow \frac{x_i - \mu}{\sqrt{\sigma^2_B + \epsilon}}$ normalize

$\gamma_i \leftarrow \frac{1}{\hat{x}_i^2 + \epsilon}$ scale
added.
to escape ∞ .

$y_i \leftarrow \gamma_i + \beta = BN(\gamma_i, \sigma(x_i))$

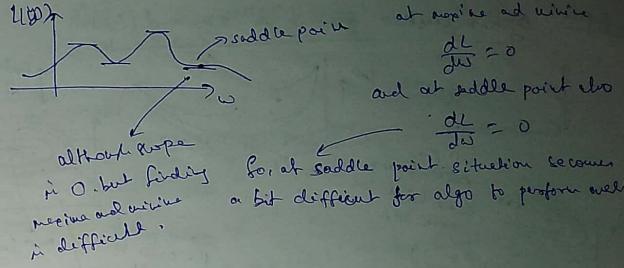
scale
and shift

for batch normalization,

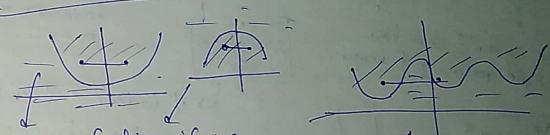
(1) Convergence is fast

(2) Reduces inter-class Covariance shift.

(3) BN + drop.



Convex function & Non Convex function



At Convex function, if we take any two point and join them, it will be inside the function.

At non-convex function, points may lie outside of function.

At non-convex function, there is lot of noise and noise.

* In logistic reg & log-reg sum function is only one local minima \Rightarrow Global minima

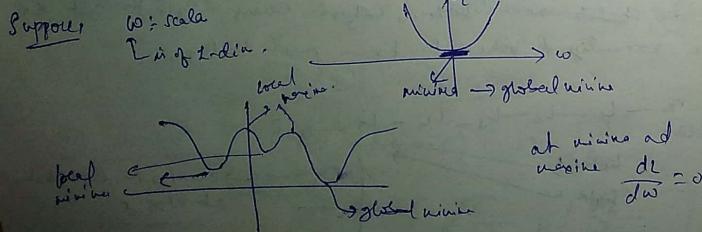
\hookrightarrow low function \rightarrow convex in nature

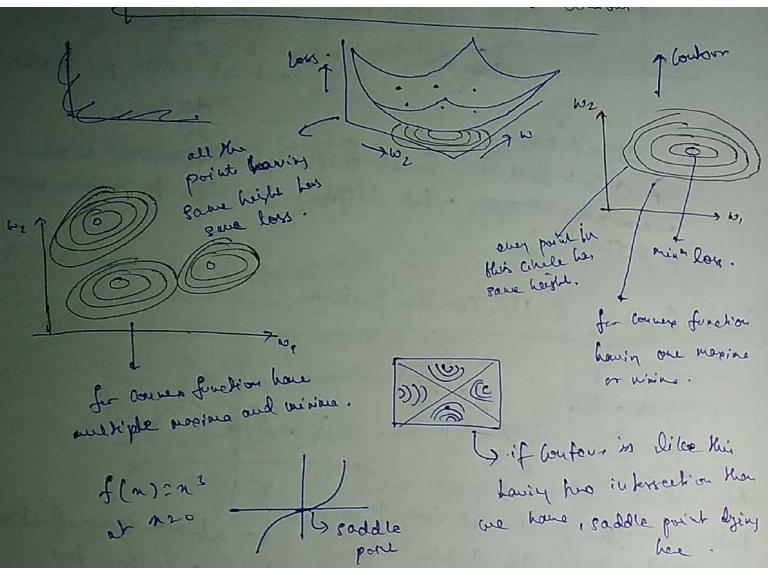
But in $\frac{dL}{dw}$ local function is non-convex because there is multiple minima and maxima.

\rightarrow It all depends on where we have initialized the weight. If it is around local minima, after reaching local minima, it becomes hard to climb up.

Optimizer Hill-descent analogy in 2D

With optimization we mean to find the optimal value of w^* to which gain minimum loss. L (in LR, linear Regress)





$\left(\frac{\partial L}{\partial w}\right)_{MB-SGD} \approx \left(\frac{\partial L}{\partial w}\right)_{GD}$

MB SGD \hookrightarrow greater noise.

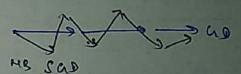
update of weight w_i (avg)

loss L

MB SGD \hookrightarrow a lot of noise.

although it is moving in same direction but with noise

\hookrightarrow MB SGD is estimation of gradient descent with SGD
 \hookrightarrow each of the updates are more noisy in SGD than G.D.



So, for the convergence of MB SGD, decaying of learning rate is done

SGD Recap

where $(w_{ij}^*)_{new} = (w_{ij}^*)_{old} - \eta \left[\frac{\partial L}{\partial w_{ij}} \right]_{(w_{ij}^*)_{old}}$

weight update function.

 $w_{ij}^* = w_{ij} - \eta \left[\frac{\partial L}{\partial w_{ij}} \right]_{w_{ij}^*}$

for all $D = \{x_i, y_i\}_{i=1}^n$

update $\frac{\partial L}{\partial w}$ \rightarrow using all the n points in $D \rightarrow$ GD

\hookrightarrow using only one pt $x_i @$ random \hookrightarrow SGD

\hookrightarrow using a random subset of D pts in \rightarrow mini-batch SGD

mini-batch GD doesn't give ~~exact~~ each value on w , it gives approximate values.

Batch SGD with momentum.

Suppose take any random numbers.

$$a_1 \quad a_2 \quad a_3 \quad a_4 \quad a_5 \quad a_6 \dots$$

$$t=1 \quad t=2 \quad t=3 \dots$$

so, at time t , \rightarrow say $v_1 = a_1$,

$$t=2 \Rightarrow v_2 = \gamma v_1 + a_2 \quad \text{where, } 0 < \gamma < 1$$

$$\text{for, take } \gamma=1 \rightarrow v_2 = v_1 + a_2 = a_1 + a_2$$

$$\gamma=0.5 \rightarrow v_2 = 0.5v_1 + a_2 \\ = 0.5a_1 + a_2$$

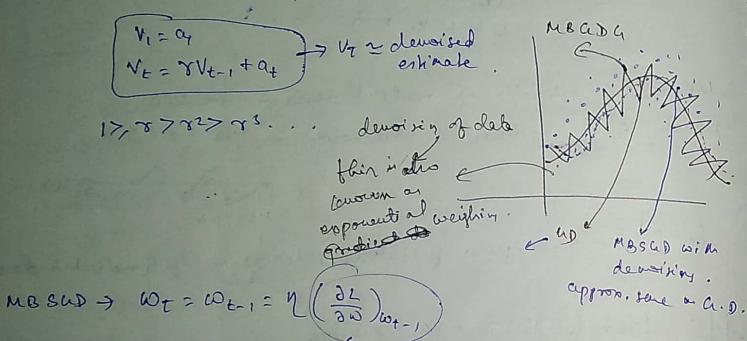
$$\gamma=0 \rightarrow v_2 = a_2$$

here, we can see that here output v_2 is always dependent on previous output a_1 , but a_2 has more weight than a_1 .

$$\begin{aligned} \Rightarrow v_1 &= a_1 \\ v_2 &= \gamma v_1 + a_2 \\ v_3 &= \gamma v_2 + a_3 = \gamma(\gamma v_1 + a_2) + a_3 \\ &= \gamma^2 a_1 + \gamma a_2 + a_3 \\ \text{let } \gamma = 0.5 &= 0.25 a_1 + 0.5 a_2 + 1.0 a_3. \\ f=1 & \quad f=2 \quad f=3. \end{aligned}$$

Weight of $a_3 > a_2 > a_1, \dots$

So, decay in weight update will reduce in MASCD.



$$MASCD \rightarrow w_t = w_{t-1} = \eta \left(\frac{\partial L}{\partial \theta} \right)_{w_{t-1}}$$

$$w_t = w_{t-1} - \eta g_t$$

exponentially weighted

$$\begin{aligned} V_t &= \gamma V_{t-1} + \eta g_t \\ w_t &= w_{t-1} - V_t \end{aligned}$$

$$\text{Initially at } \gamma=0, V_t = h_g, \quad w_t = w_{t-1} - (\gamma V_{t-1} + \eta g_t)$$

Usually γ is taken

as 0.9.

$$\gamma = 0.9 \rightarrow w_t = w_{t-1} - (0.9 V_{t-1} + \eta g_t)$$

$$\text{so, } w_t = w_{t-1} - [\gamma V_{t-1} + \eta g_t]$$

friction is called
momentum because
it depends on previous value
 $g_{t-1}, g_{t-2}, g_{t-3}, \dots$

g_t

$$V_t = 1(\eta g_t) + \gamma(\eta g_{t-1}) + \gamma^2(\eta g_{t-2}) + \dots$$

γV_{t-1} momentum
to find this will move
in this direction
if gradient is going in
this direction
if we move
in going in that
direction

so initially, it will see in
which direction update is moving,
after knowing that it will move
converge more rapidly.

Nesterov Accelerated Gradient (NAG)

Earlier we have seen, SGD + momentum,
momentum $\rightarrow w_t \rightarrow$ where find direction travel in the
 w_{t-1} gradient vector sum of gradient + momentum.

$$w_t = w_{t-1} - (\gamma V_{t-1} + \eta g_t)$$

we first travel in the direction of momentum (from w_{t-1})
then we travel in the direction of gradient (point $(w_{t-1} - \gamma V_{t-1})$)

so, here point w_t has changed from $(w_{t-1} - \gamma V_{t-1})$.

NAG gives better result.

$$\begin{aligned} \eta g_t &\rightarrow \text{gradient term} \\ \gamma V_{t-1} &\rightarrow \text{momentum} \end{aligned}$$



when $g_t \neq g$
since starting point
has changed.

Optimizers Adaline

In SGD and SADT, we see that learning rate was constant ($\eta = \text{const}$) → same for each weight.

But here, idea is:

→ each weight/pair has a different η

Because, when we apply BOW, features become sparse, and for sparse features, constant learning rate is no good.

So, AdaGrad, η should change such that $\eta \rightarrow$ needs to reduce for better model.

SGD $\rightarrow w_t = w_{t-1} - \eta g_t$ same for all weights

Adagrad $\rightarrow w_t = w_{t-1} - \eta' g_t$ adaptive.

$$\text{So, } \eta' = \frac{\eta}{\sum_{i=1}^{t-1} g_i^2}$$

different η' for each weight

$$\text{So, } \eta' = \frac{\eta}{\sum_{i=1}^{t-1} g_i^2} \rightarrow \text{small the norm (to avoid divisibility by zero)}$$

$$g_t = \sum_{i=1}^{B-1} g_i^2 \rightarrow \left(\frac{dL}{dw} \right)_{w_{t-1}}$$

already +ve → $g_1, g_2, g_3, \dots, g_{t-1}$

So, after time(t), x_{t-1} , $\theta_0, (\eta')$ ↓

$$\eta' \geq \eta_{t-1}$$

So, in Adagrad)

(+) → no need of manually tuning η → weight, time/iter

(+) sparse features and dense one can be easily handled

(-) d_{t-1} can become very large, as $t \uparrow$

So, when, $\alpha_{t-1} \uparrow$ with $t \uparrow$.

↳ convergence rate will be slower.

Optimizers Adadelta and RMSprop

↳ idea of adadelta is to use exponential decreasing avg. So to reduce, $\alpha_{t-1} = \sum_{i=1}^{t-1} g_i^2$ is used in Adadelta.

$$\text{Since, } w_t = w_{t-1} - \eta' g_t$$

$$\eta' = \frac{\eta}{\sum_{i=1}^{t-1} g_i^2 + \epsilon}$$

exponential decaying avg.

$$\text{where, } eda_{t-1} = \gamma eda_{t-2} + (1-\gamma) g_t^2$$

$$\text{Param, } \gamma = 0.95$$

$$eda_{t-1} = 0.95 eda_{t-2} + 0.05 g_{t-1}^2$$

$$eda_{t-1} = (0.95)^2 g_{t-1}^2 + 0.95 \cdot 0.05 g_{t-2}^2$$

for, this value is never zero, + 0.95 eda_{t-2}

So, $\eta' \rightarrow$ will never be zero.

It has all the features of Adagrad with eda_{t-1} cannot be very large with $t \uparrow$.

Adam (Adaptive moment estimation)

In Adadelta we saw it stores exponential decreasing avg (eda) $\rightarrow g_t^2 \rightarrow$ learning rate ($\eta' t$)

Idea of adam is it stores eda of g_t also.

It has moment in it because

In statistics, mean collected 1st order moment
Variance " 2nd ... "

$$\begin{aligned} \text{edt} &\leftarrow \hat{g}_t \leftarrow m_t = \beta_1 m_{t-1} + (1-\beta_1) g_t \quad (0 \leq \beta_1 \leq 1) \\ \text{edt} &\leftarrow \hat{g}_t^2 \leftarrow v_t = \beta_2 v_{t-1} + (1-\beta_2) g_t^2 \quad (0 \leq \beta_2 \leq 1) \end{aligned}$$

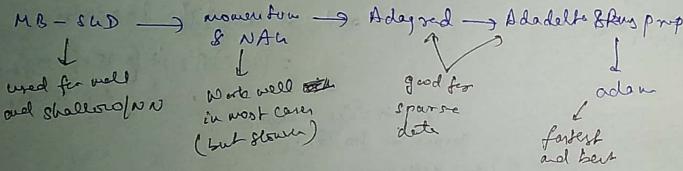
Note: $\left\{ \hat{m}_t = \frac{\hat{m}_t}{1-(\beta_1)^t}, \hat{v}_t = \frac{v_t}{1-(\beta_2)^t} \right\}$

$$\left\{ w_t = w_{t-1} - \frac{\hat{m}_t}{\sqrt{\hat{v}_t + \epsilon}} \right\} \quad \text{if } \beta_1 = 0 \rightarrow \text{Adadelta.}$$

jellyfish

This is the best optimization technique among all, it converges much faster.

Which Algorithm to choose when?



Gradient Checking and Clipping

In a normal way we should always keep an eye on gradient on how it is changing. Because, updating gradient is the most important thing in D2.

So, monitor gradient for each epoch, weights and layers.

because, not doing so, we can face vanishing gradient or exploding gradient.

One solution for exploding gradient is Clipping.

weights in network $w = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 & 1 \end{bmatrix}$

$G = \text{gradient} = \begin{bmatrix} \frac{\partial L}{\partial w_1} & \frac{\partial L}{\partial w_2} & \dots & \frac{\partial L}{\partial w_6} \end{bmatrix}$ (let 200)

So, gradient can be so high, so, we clip it by l2 normalization clipping.

When, $l_2 \text{ norm} = \frac{G}{\|G\|_2} * T$ $T=2 \rightarrow \text{clip}$

$G = \frac{G}{\sqrt{\sum G_i^2}} * T$ $\|G\|_2 = \sqrt{G_1^2 + G_2^2 + \dots + G_n^2}$

gradient $\approx T$

If we are vanishing gradient descent, change Sigmoid to ReLU.

Softmax and Cross-entropy for multi-class classification.

As we have seen logistic regression is used for binary classification. If one want to do multi-class classification, then One Vs Rest method is used.

But a new method is developed using same concept of Logistic Regression.

\therefore Log-reg. + Multiclass = Softmax.

In Logistic Regm

$$\begin{aligned} \text{Input: } & \begin{bmatrix} w_1 \\ w_2 \\ w_3 \\ w_4 \end{bmatrix} \rightarrow z = \sum w_i x_i \\ \text{Output: } & \hat{y} = p(y_i=1|x_i) = \sigma(z) = \frac{1}{1+e^{-z}} \\ & z = \sum w_i x_i \quad = \frac{e^z}{e^z + 1} \end{aligned}$$

↳ here, we have only two possibilities.

Softmax classifn.

$$D = \{(x_i, y_i) \mid y_i \in \{1, 2, 3, \dots, k\}\}$$

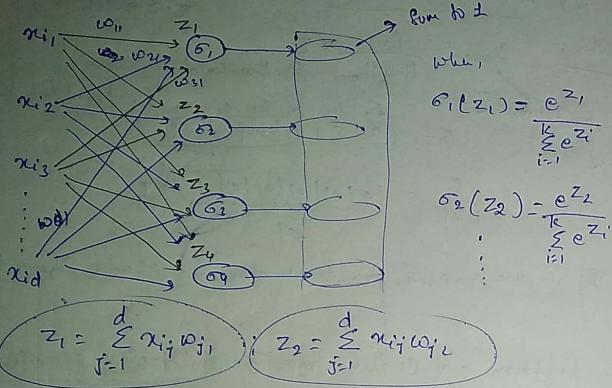
↳ multiclass.

$$\begin{aligned} x_i \rightarrow \text{model} \rightarrow & P(y_i=1|x_i) \\ & P(y_i=2|x_i) \\ & P(y_i=3|x_i) \\ & \vdots \\ & P(y_i=k|x_i) \end{aligned}$$

all sum up to 1.

Suppose

$D = 4$ softmax



Now

$$\sigma(z_1) + \sigma(z_2) + \dots + \sigma(z_k) = \frac{e^{z_1}}{\sum e^{z_i}} + \frac{e^{z_2}}{\sum e^{z_i}} + \dots + \frac{e^{z_k}}{\sum e^{z_i}} = \frac{\sum e^{z_i}}{\sum e^{z_i}} = 1$$

Put up dot.

So, softmax is the generalization of LR to multiclass level.

Softmax minimizes multiclass log-loss

$$\sigma_i(z_i) = \frac{e^{z_i}}{\sum_{j=1}^k e^{z_j}}$$

$p(y_i=1|n_i)$

$$\log\text{-loss} = -\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^k y_{ij} \log p_{ij}$$

$\left\{ \begin{array}{ll} 1 & \text{if } y_{ij} \neq 0 \\ 0 & \text{otherwise} \end{array} \right\}$

Support given $y_{ij}=1$

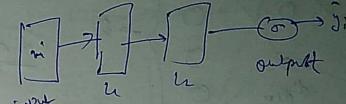
1	2	3	4	..	k	
0	0	1	1	0	..	0

$$x_i \rightarrow \text{model} \rightarrow \begin{cases} 0.2 & \rightarrow p(y_i=1|n_i) \\ 0.1 & \rightarrow p(y_i=2|n_i) \\ 0.1 & \rightarrow p(y_i=3|n_i) \\ 0 & \\ 0 & \end{cases}$$

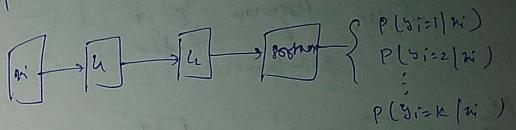
$\therefore \text{log-loss} = -\sum y_{ij} \log p_{ij} \rightarrow \text{one putting on formula}$

MLP

2-class



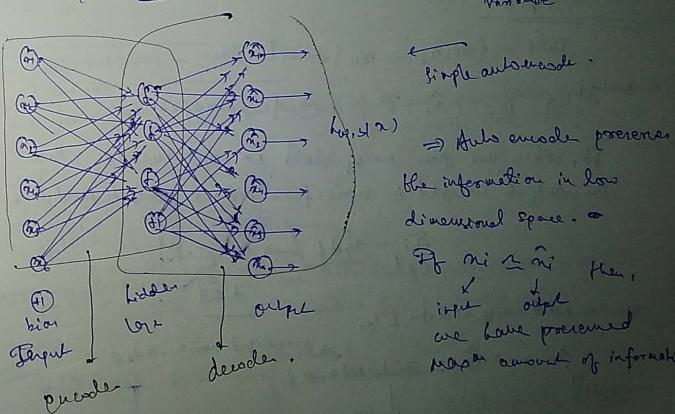
10-class

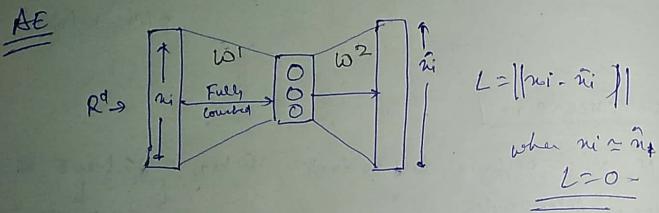
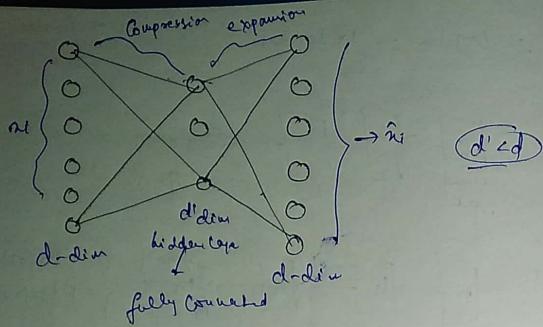


Auto Encoder.

→ MN which perform dimension reduction like PCA, t-SNE is based on neighborhood

Simple Autoencoder





Deriving AE

Data = $\{x_1, x_2, \dots, x_n\} \leftarrow$ actuatable

weights $\mathcal{D} = \{\hat{x}_1, \hat{x}_2, \dots, \hat{x}_n\} \rightarrow$ noisy.

So, when we pass noisy data to autoencoder, then the output we get is noisefree, more robust model.

$$x_i \xrightarrow[\text{Eqd}]{AE} x'_i \in R^d$$

$d'cd$

We can add noise to data like

$$x_i \rightarrow \tilde{x}_i = x_i + N(0, \sigma)$$

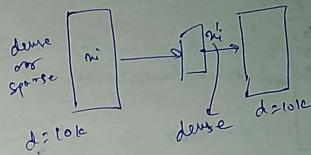
normalized noise added

$$\xrightarrow{\text{AE}} x'_i \in R^{d'}$$

$d'cd$

Sparse Autoencoder

To add sparsity, we add L_1 regularization to loss.

$$|L_1| + \text{loss} \rightarrow \text{sparse}$$


* If linear activations are used or only a single sigmoid hidden layer, the optimal solution to an autoencoder is strongly to PCA.

MNIST \rightarrow 784dim $\xrightarrow{\text{Sdlin}}$ 50dim \rightarrow feature representation

| Word2Vec - CBOW

We have seen earlier in Amazon food reviews \rightarrow how to convert words to vector using CBOW, TFIDF, w2v using softmax

↳ c_1 found by c_2 ↳
The Cat sat on the Wall
~~~~~ Context

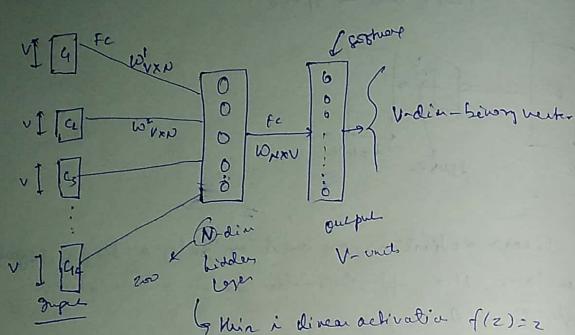
↳ in w2v  $\rightarrow$  we take the help of context words to predict focused words.

CBOW ↳ (1) Create dictionary/vocabulary  
 $V = \text{length of vocab.}$

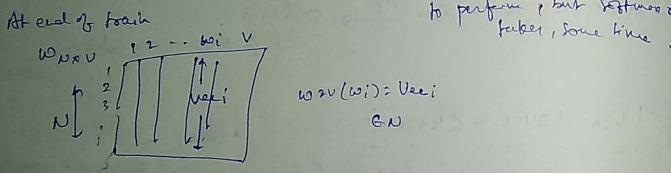
(2) Do one hot encoding of each word.

WEV: binary vec of V-dim.

Now, idea is using all the vectors of context words, providing the focused vector using softmax.



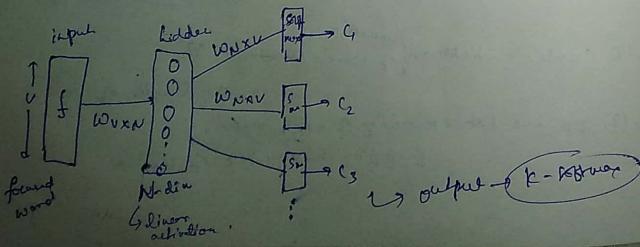
$\hookrightarrow$  Here  $i$  is linear activation  $f(z) = z$   
This does not take much time to perform, but softmax takes longer, some time.



## # Skip-gram (Word2Vec)

Earlier we have seen that given context words CBOW predicts focused word.

But in Skip-gram situation is reversed, here given focused words we have to predict context words, which is quite difficult.



So, No. of weights in CBOW & skipgram  $\rightarrow (C+1)(N \times V)$

CBOW  $\rightarrow k$  softmax

Skipgram  $\rightarrow k$  softmax

{ Skipgram is computationally more expensive.

CBOW

- (+) faster than skipgram
- (+) better for frequently occurring words.

Skipgram

- (-) can work with smaller amount of data
- (+) well for infrequent words.

So, in CBOW and skipgram -

we have ' $k$ ' no. of context words.

so when  $k \uparrow \rightarrow N$ -dim for each point is better

$$\text{Weights} = (C+1)(N \times V)$$

$$\downarrow$$

$$6 \times (20 \times 10k)$$

$$= \cancel{+2 \times 10k} + 1.2 \text{ more}$$

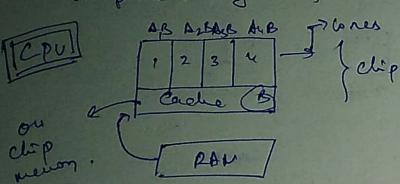
So this can take forever to train.

So, optimizer is used to train it more efficiently.

## # Word2Vec Algorithmic Optimizations.

## # CPU vs GPU for Deep learning.

Amplifier was generally developed by NVIDIA for learning purpose.



In CPU it is general of 2<sup>48</sup>,  
i.e., in 4 cores or 8 cores,  
each operates ~~independently~~  
parallelly.

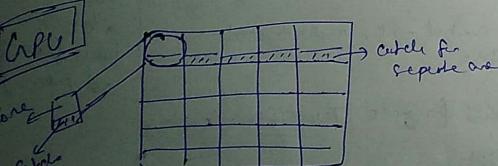
like, one for music,  
one for video,

$$\text{So, Suppose we have two computations} \quad \begin{bmatrix} \leftarrow A_1 \\ \leftarrow A_2 \\ \leftarrow A_3 \\ \vdots \end{bmatrix} \begin{bmatrix} B \\ \downarrow \end{bmatrix} = \begin{bmatrix} A_1 B \\ A_2 B \\ A_3 B \\ \vdots \end{bmatrix} \quad \text{both are parallel.}$$

Now, each matrix multiplication is independent of each other.

So, can operate in parallel.

Cache is used to keep some part  
of data so that interaction of  
core with data becomes easier and fast.



Each GPU Core Core is much  
more slower than CPU

→ 400 MHz.

$$\begin{array}{c|c|c|c} \text{Computation} & A_1 B & A_{10} B & \dots \\ & A_2 B & A_{11} B & \dots \\ & \vdots & \vdots & \ddots \\ & A_n B & A_{n+1} B & \dots \end{array}$$

In GPU no. of cores  
can be  
100  
1024  
2048

So, each core operates  
parallel to, computation  
becomes much easier.

Since here each core has separate cache so computation is fast.

$$1\text{GB/RAM} \rightarrow 32 \times 10^9 \text{ Ti} \quad \begin{array}{l} \text{VRAM} \\ \text{video} \end{array}$$

## # [49.6] Software classifier on MNIST dataset

### [ Biological Inspiration Visual Cortex ]

⇒ In this chapter we will learn about (CNN, ConvNets)

so, CNN is basically used for visual tasks.

As MNIST data is also a small  $\rightarrow$  visual task.

In the last ~~topic~~

Research on ML/AL + ~~Image~~ / video  $\rightarrow$  has evolved a lot  
and for ~~object~~ object recognition.

The credit of all this research people goes to a research paper  
in 1981 by Hubel & Wiesel (Nobel prize Medicine),  
where he performed an exp. on cat on how he actually detects  
image.

While performing the exp. he found that 8 neuron gets fired  
when it sees some ~~of~~ similar pattern.

like →

- neuron fired
- neuron little fired
- neuron little fired
- neuron doesn't fire
- doesn't fire

and our brain has three different layers ( $V_1, V_2, V_3, \dots$ )  
each layer has different tasks.

$V_1 \rightarrow$  detects motion  
 $V_2 \rightarrow$  stereo  
 $V_3 \rightarrow$  color

$V_4 \rightarrow$  feature representation  
 $V_5 \rightarrow$  classification  
etc.

key findings ↗

⑤ Some neurons in the visual cortex that fire when it sees depth perceived from a specific angle.

⑥  $\Rightarrow$  primary visual cortex  $\rightarrow$  detects edge.

⑦ brain has different layers to detect different things, like edge detector, motion detector, depth, color detector etc.

each layer operates to detect a particular thing



Convolution Edge Detection on image

→ here in the next four layers to detect edge.

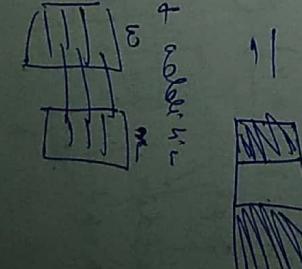
$6 \times 6$

$\star$  Convolution

$3 \times 3$

$\star$  max pool

$\star$  max pool



|   |   |   |   |   |   |
|---|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 |

|    |    |    |  |  |  |
|----|----|----|--|--|--|
| 41 | +2 | +1 |  |  |  |
| 0  | 0  | 0  |  |  |  |
| -1 | -2 | -1 |  |  |  |
| 0  | 0  | 0  |  |  |  |
| 0  | 0  | 0  |  |  |  |

|       |       |       |       |       |       |
|-------|-------|-------|-------|-------|-------|
| 0     | 0     | 0     | 0     | 0     | 0     |
| -1020 | -1020 | -1020 | -1020 | -1020 | -1020 |
| -1020 | -1020 | -1020 | -1020 | -1020 | -1020 |
| -1020 | -1020 | -1020 | -1020 | -1020 | -1020 |
| 0     | 0     | 0     | 0     | 0     | 0     |

|   |   |   |   |   |   |
|---|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 |

|       |       |       |       |       |       |
|-------|-------|-------|-------|-------|-------|
| 0     | 0     | 0     | 0     | 0     | 0     |
| -1020 | -1020 | -1020 | -1020 | -1020 | -1020 |
| -1020 | -1020 | -1020 | -1020 | -1020 | -1020 |
| -1020 | -1020 | -1020 | -1020 | -1020 | -1020 |
| 0     | 0     | 0     | 0     | 0     | 0     |

|       |       |       |       |       |       |
|-------|-------|-------|-------|-------|-------|
| 0     | 0     | 0     | 0     | 0     | 0     |
| -1020 | -1020 | -1020 | -1020 | -1020 | -1020 |
| -1020 | -1020 | -1020 | -1020 | -1020 | -1020 |
| -1020 | -1020 | -1020 | -1020 | -1020 | -1020 |
| 0     | 0     | 0     | 0     | 0     | 0     |

|       |       |       |       |       |       |
|-------|-------|-------|-------|-------|-------|
| 0     | 0     | 0     | 0     | 0     | 0     |
| -1020 | -1020 | -1020 | -1020 | -1020 | -1020 |
| -1020 | -1020 | -1020 | -1020 | -1020 | -1020 |
| -1020 | -1020 | -1020 | -1020 | -1020 | -1020 |
| 0     | 0     | 0     | 0     | 0     | 0     |

|       |       |       |       |       |       |
|-------|-------|-------|-------|-------|-------|
| 0     | 0     | 0     | 0     | 0     | 0     |
| -1020 | -1020 | -1020 | -1020 | -1020 | -1020 |
| -1020 | -1020 | -1020 | -1020 | -1020 | -1020 |
| -1020 | -1020 | -1020 | -1020 | -1020 | -1020 |
| 0     | 0     | 0     | 0     | 0     | 0     |

$$+25 \times 1 + 25 \times (-2) + 25 \times (-1) = -125$$

↳ we do computation with matrix multiplication & addition and get a row vector  $\hat{w}_4$  having max value. Since 0 is min and 25 is max.

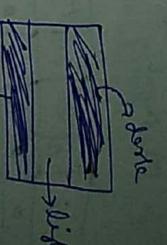
so  $\hat{w}_4$   $\rightarrow$  0 means Line - 1020  
so  $\hat{w}_4$   $\rightarrow$  25 means Line - 1020

so  $\hat{w}_4$   $\rightarrow$  0 is min and 25 is max.

so  $\hat{w}_4$   $\rightarrow$  0 is min and 25 is max.

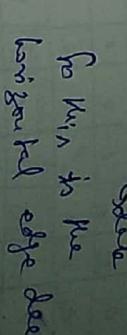
|    |    |    |    |
|----|----|----|----|
| 25 | 25 | 25 | 25 |
| 0  | 0  | 0  | 0  |

|    |   |    |   |
|----|---|----|---|
| 0  | 0 | 0  | 0 |
| -1 | 0 | -2 | 0 |



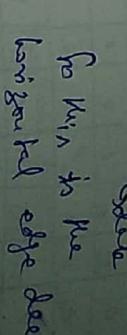
|    |    |    |    |
|----|----|----|----|
| 25 | 25 | 25 | 25 |
| 0  | 0  | 0  | 0  |

|    |   |    |   |
|----|---|----|---|
| 0  | 0 | 0  | 0 |
| -1 | 0 | -2 | 0 |



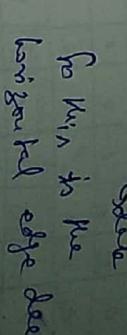
|    |    |    |    |
|----|----|----|----|
| 25 | 25 | 25 | 25 |
| 0  | 0  | 0  | 0  |

|    |   |    |   |
|----|---|----|---|
| 0  | 0 | 0  | 0 |
| -1 | 0 | -2 | 0 |



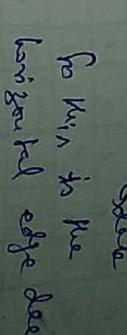
|    |    |    |    |
|----|----|----|----|
| 25 | 25 | 25 | 25 |
| 0  | 0  | 0  | 0  |

|    |   |    |   |
|----|---|----|---|
| 0  | 0 | 0  | 0 |
| -1 | 0 | -2 | 0 |



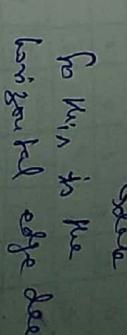
|    |    |    |    |
|----|----|----|----|
| 25 | 25 | 25 | 25 |
| 0  | 0  | 0  | 0  |

|    |   |    |   |
|----|---|----|---|
| 0  | 0 | 0  | 0 |
| -1 | 0 | -2 | 0 |



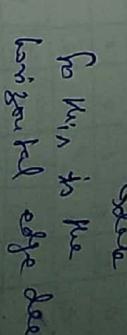
|    |    |    |    |
|----|----|----|----|
| 25 | 25 | 25 | 25 |
| 0  | 0  | 0  | 0  |

|    |   |    |   |
|----|---|----|---|
| 0  | 0 | 0  | 0 |
| -1 | 0 | -2 | 0 |



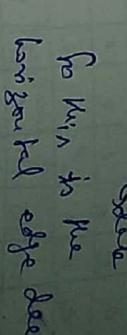
|    |    |    |    |
|----|----|----|----|
| 25 | 25 | 25 | 25 |
| 0  | 0  | 0  | 0  |

|    |   |    |   |
|----|---|----|---|
| 0  | 0 | 0  | 0 |
| -1 | 0 | -2 | 0 |



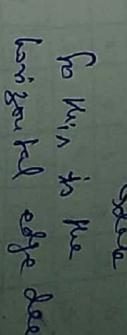
|    |    |    |    |
|----|----|----|----|
| 25 | 25 | 25 | 25 |
| 0  | 0  | 0  | 0  |

|    |   |    |   |
|----|---|----|---|
| 0  | 0 | 0  | 0 |
| -1 | 0 | -2 | 0 |



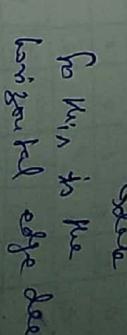
|    |    |    |    |
|----|----|----|----|
| 25 | 25 | 25 | 25 |
| 0  | 0  | 0  | 0  |

|    |   |    |   |
|----|---|----|---|
| 0  | 0 | 0  | 0 |
| -1 | 0 | -2 | 0 |



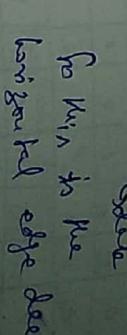
|    |    |    |    |
|----|----|----|----|
| 25 | 25 | 25 | 25 |
| 0  | 0  | 0  | 0  |

|    |   |    |   |
|----|---|----|---|
| 0  | 0 | 0  | 0 |
| -1 | 0 | -2 | 0 |



|    |    |    |    |
|----|----|----|----|
| 25 | 25 | 25 | 25 |
| 0  | 0  | 0  | 0  |

|    |   |    |   |
|----|---|----|---|
| 0  | 0 | 0  | 0 |
| -1 | 0 | -2 | 0 |



|    |    |    |    |
|----|----|----|----|
| 25 | 25 | 25 | 25 |
| 0  | 0  | 0  | 0  |

<table border="1

## # Convolution padding and stride.

As seen earlier, the output we get won't have dimension  
than output.

$$\text{If input} \rightarrow 6 \times 6 \text{ matrix} \rightarrow \text{output} = n - k + 1 \\ = 6 - 3 + 1 \\ = 4$$

$$\text{So when stride 1} \rightarrow 6 \times 3$$

$$= 4$$

But when with padding we can achieve output of same dimension  
as input.

|   |    |    |    |    |    |   |   |
|---|----|----|----|----|----|---|---|
| 0 | 0  | 0  | 0  | 0  | 0  | 0 | 0 |
| 0 | 0  | 0  | 0  | 0  | 0  | 0 | 0 |
| 0 | 0  | 0  | 0  | 0  | 0  | 0 | 0 |
| 0 | 25 | 25 | 25 | 25 | 25 | 0 | 0 |
| 0 | 25 | 25 | 25 | 25 | 25 | 0 | 0 |
| 0 | 25 | 25 | 25 | 25 | 25 | 0 | 0 |
| 0 | 25 | 25 | 25 | 25 | 25 | 0 | 0 |
| 0 | 25 | 25 | 25 | 25 | 25 | 0 | 0 |

So, here we want output  
as  $6 \times 6$  dimension.  
So, in  $n - k + 1 = 6$   
 $\therefore n = 8$

$$k = 3$$

So we add one extra layer  
to make extra  $2 \times 2$  cells  
to make input  $8 \times 8$ .

And fill that extra layer  
with 0 or value zero to it.

$$\text{So, } (n \times n) \xrightarrow{(k \times k)} (h - k + 1) \times (w - k + 1)$$

and

padding:

$$(n \times n) \xrightarrow{\text{left padding}} (n - k + p - 1) \times (n - k + p - 1)$$

Strides: cells in input we traversed each step.  
So, with stride 2 we traversed 2 step.

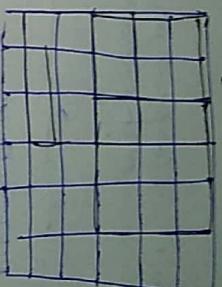
$$\text{Stride: } (n \times n) \xrightarrow{\text{S=2}} \left(\frac{n-k}{2}\right) \times \left(\frac{n-k}{2}\right)$$

$$6 \times 6$$

$$\xrightarrow{k=3}$$

$$2 \times 2$$

$$3 \times 3$$



### # Convolution over RGB images.

Till now we have been convolution in grayscale image. (2D)  
Now we will see convolution in 3-D image.



Volume

2D

matrix

1-D

vector

3-D

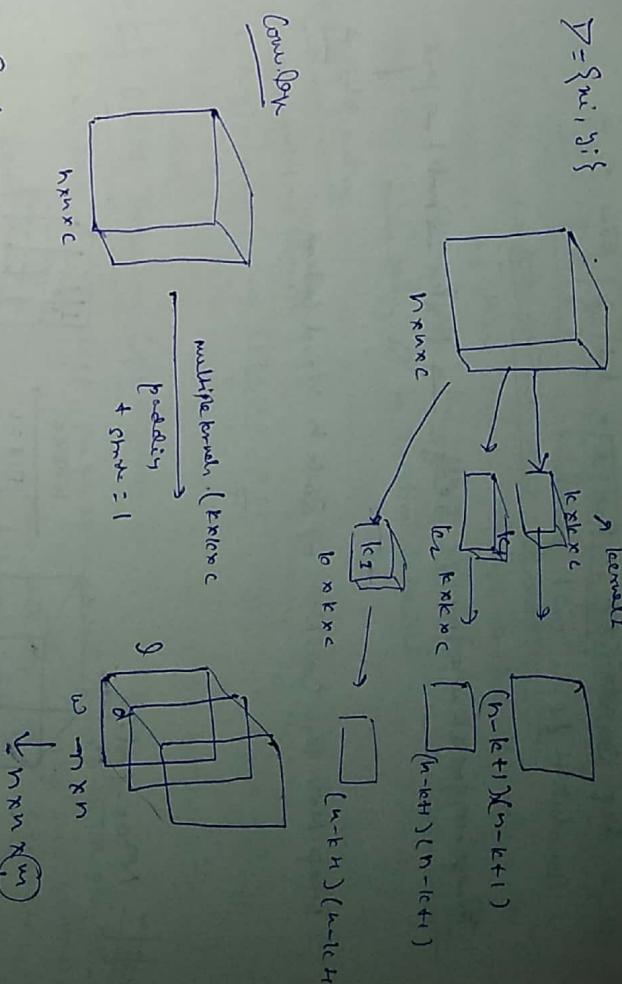
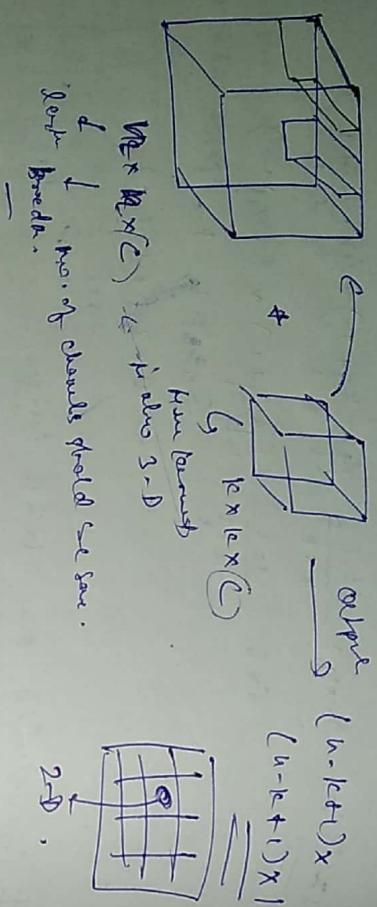
matrix

3-D

As we have seen that

convolution in 2D matrix happens with  $\rightarrow$  component wise multiplication.

So, similarly,  $3D \rightarrow \dots$



## Convolution layer

$\rightarrow$  Earlier we saw how a MLP worked

$$\text{Step } L \rightarrow \text{at } n = 2$$

$$\text{Step } L \rightarrow \text{ReLU}(z) = y$$

So, convolution layer also works kind of similarly.

① Conv. layers are inspired from ~~biologically~~ from  $V_1, V_2, V_3, \dots$

② As in the first stage of seeing there is multiple edge detectors,

so in conv. layer,  $1 - \text{edge detector} \rightarrow 1 \text{ kernel}$

multiple edge detector  $\rightarrow$  multiple kernel.

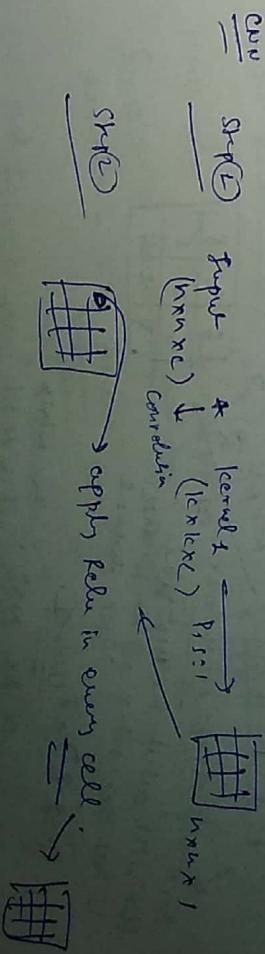
The CNN also kernels learn from backprop.

The MLP and weights were updating but in CNN kernels will update.

weight update

MLP  $\rightarrow$  learn  $w_i \rightarrow w_i \leftarrow w_i + \text{dot product}$

CNN  $\rightarrow$  learn learned matrices  $\rightarrow$   $I \otimes \text{Matrix} \rightarrow$  learned input by convolution



In a Deep ~~Conv~~ CNN, there is a layer of pair of convolution & pool unit paired.

Interpretation: Pixel → edge → texture → motif → part

Deep Multi-layer ConvNet → Similar to visual context

↓  
Output  
of two type of pooling  
↳ max-pooling → concrete  
↳ Avg/mean-pooling →  $\frac{1+1+5+6}{4} = \underline{\underline{3.25}}$

for input  $\begin{bmatrix} 1 & 5 \\ 1 & 6 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 5 \\ 1 & 6 \end{bmatrix} \rightarrow \begin{bmatrix} 6 & 5 \end{bmatrix}$   
↳ ~~Concrete value~~ ~~with~~ ~~with~~ if more values change within this matrix, the even then the max value from max pooling = 6.

## # Max-Pooling.



|   |   |   |   |
|---|---|---|---|
| 1 | 1 | 1 | 1 |
| 1 | 5 | 1 | 6 |
| 1 | 1 | 1 | 1 |
| 1 | 6 | 5 | 6 |

when

# CNN Training Optimization  
As we have seen in MLP, back-propagation is applied in all the layers backprop via sub, Adagrad, Adam etc.

a  $\frac{\partial L}{\partial w} \rightarrow L(y, \hat{y})$

$\frac{\partial L}{\partial w} \rightarrow$  derivative. So here we have to check if derivative and max pooling is differentiable or not.

Surely in

differentiable so, low

Layer is also differentiable  
because it has + add.

Elements wise mult + add.

Value in MLP.

Value is differentiable

MLP  $\xrightarrow{\text{differentiable}}$  max pooling  $\xrightarrow{\text{not differentiable}}$  mult + add

|   |   |   |
|---|---|---|
| 1 | 2 | 3 |
| 6 | 7 | 8 |
| 3 | 4 | 5 |

|   |   |   |
|---|---|---|
| 1 | 2 | 3 |
| 6 | 7 | 8 |
| 3 | 4 | 5 |

Now when talking for max pool,  $\xrightarrow{\text{backprop}}$

So when it back propagate ~~at~~ gradient derivative is slice  $\frac{\partial L}{\partial w} = \frac{\partial L}{\partial u} = 1$

directly

and for all other value derivative = 0.

So max pooling is also differentiable.

if we change the value within this matrix, the even then the max value from max pooling = 6.

It will detect the max in 2x2 window.

If (because) and similarly.

It will detect the max in 2x2 window.

and because with stride = 2.

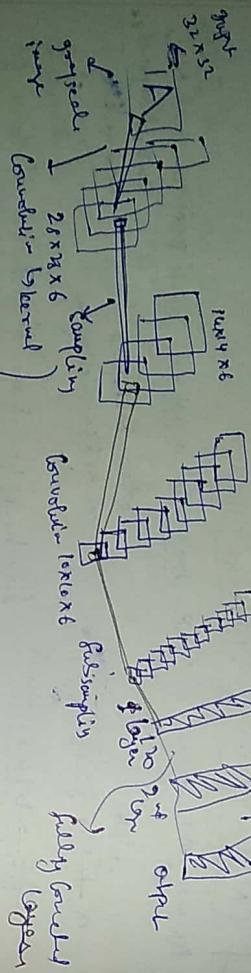
$$\text{map}(x) = \begin{cases} 1 & \text{if } x > 0 \\ 0 & \text{otherwise} \end{cases}$$

so, if for each  $x_i \rightarrow$  there is more than one but same  $y_i$ ,  
then way our dataset is increasing and accuracy  
can also increase.

Can apply backprop to it.

## # Example CNN layer [1998]

Then over paper by LeCun in 1988, since we didn't have  
much of computation power, and data was less, so it can not work.

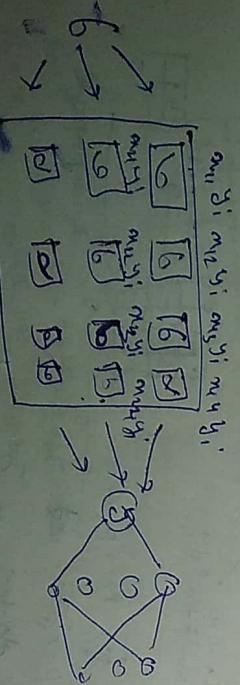


Sampling in convolution.

Since, first time Adam, Polya, were not developed, so used + sigmoid function

## # Data Augmentation

→ it mean to add extra features to data.



So, the concept of residual network was introduced.

Residual Network → In neural Network, it is seen that while  
using skip connection → the error also increases.  
But this should not happen.  
So, the concept of residual network was introduced.

## # Convolution Layer in ResNet

Using in data augmentation method.

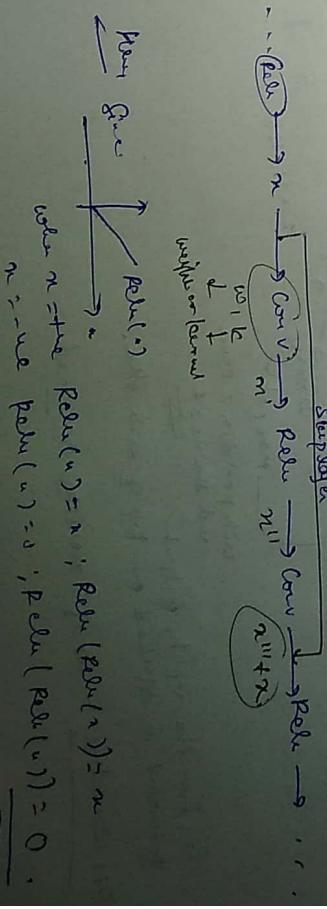
We can take the variation in data, like with rotation,  
zoom, then noise.

→ we can generate large dataset from small one.

→ very popular for image dataset.

## Discriminative

Each input is augmented, rotated, reflected, mirrored,  
sheared, blur etc.



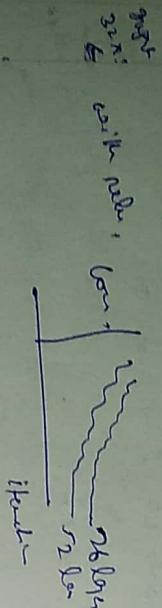
So, the concept of residual network was introduced,  
where  $\text{ReLU}(u) = u$ ;  $\text{ReLU}(\text{ReLU}(u)) = u$   
 $n - \text{one} \quad \text{ReLU}(u) = 0$ ;  $\text{ReLU}(\text{ReLU}(u)) = 0$ .

For when adding some extra layers, if weight or kernel width of some layer reaches closer to zero, then this layer becomes useless so, with respect we can skip those layers which output of next cells is always zero.

before or after, the output before or after some layers remains same.

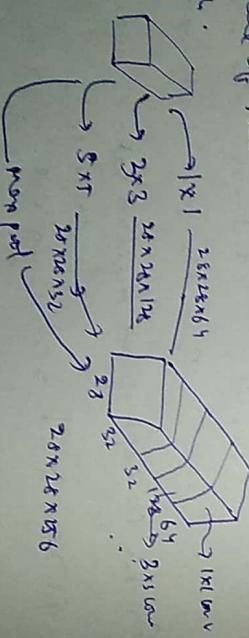
So, this way, our performance increases if each layer is giving some modification to weights, unless zero. even then this layer is good.

So, in both situations model is good.



### # Inception Network

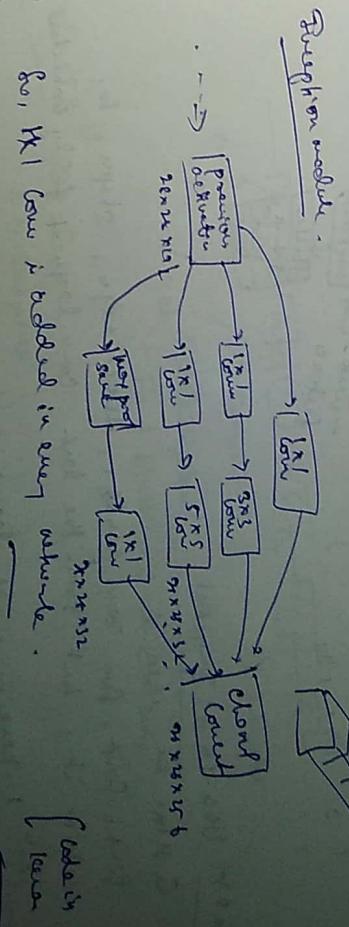
The idea of inception network is instead of deciding whether to use 1x1 convolution or 3x3 - conv or 5x5 conv or pooling layer, why not use all of them. But if we use all of them together at once, then computation becomes very high.



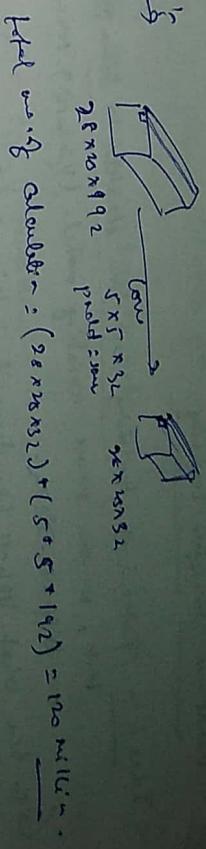
[P.S. Computationally expensive network]

So, before applying 5x5 or more conv, we apply 1x1 conv, and then we add in almost twice.

### Inception module



So, 1 conv is added in every network.



6.

With padding = none  
and stride = 2,  
so, from the input, output generated is stored on top of each other.

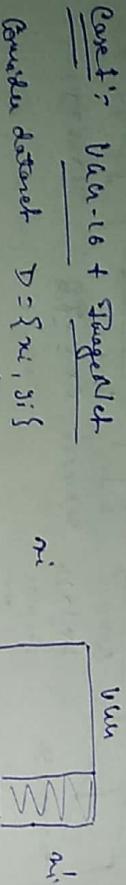
Disadvantage for lot of memory required.

## What is Transfer Learning.

↳ Idea is: Instead of building NN from scratch to solve a task, we can reuse existing models (VGG16) trained on different dataset.

↳ So here, what VGG16 has learned is being transferred to different nodes.

### Case 1: VGG16 + TransferNet

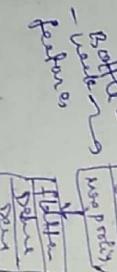


Rather than training in all dataset, we pre-trained weights till low. layer, and then we can make for last fully connected layer.

So far each

$$x_i' \rightarrow x_i' \rightarrow z_i$$

$$D' = \left\{ \begin{array}{l} x_1', x_2', x_3', \dots \\ x_1, x_2, x_3 \end{array} \right\}$$



### Case 2: VGG16 + TransferNet

↳ Here the idea is, first, few conv. layers, slopes, so retain it, but train the last conv. layers + fully connected layer, keeping small learning rate.

Case 2: VGG16 + TransferNet ← initial model.

↳ Learn this on initial model and tune the complete model  $\rightarrow D = \{x_i, y_i\} \rightarrow$  keeping learning rate low.

Case 3: Train Network from scratch.

↳ On Basis of size of Data.

$$D = \{x_i, y_i\}$$

$|D| = \text{size of dataset}$

Case 1:  $D \approx \text{TransferNet} \rightarrow$  data is easier to ingested and  $|D| = \text{small}$

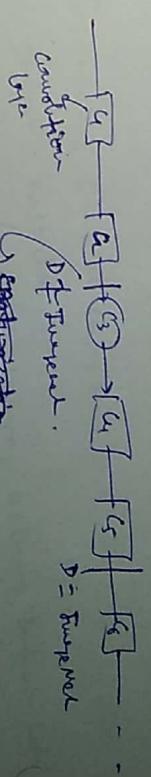
Case 2:  $|D| = \text{large}$   
 $D \approx \text{TransferNet}$

Picture the Complete Network (will bring you)

$|D| = \text{medium size}$  } + Picture the last few layer.

Case 3:  $|D| = \text{small}$

$D \neq \text{TransferNet}$



Case 4:  $\text{FC} + \text{FeatureNet} \rightarrow \text{Complete Model}$

Case 4:  $|D| = \text{large}$

$D \neq \text{TransferNet}$

$\rightarrow$  Train like the model VGG16 + TransferNet

$\rightarrow$  fine the complete model.

Reference  
Handbook  
Case 3:

## Unsupervised Learning

### Clustering

In earlier lectures we understand supervised learning whereas we had feature and labels.

$$D = \{m_1, m_2\}$$

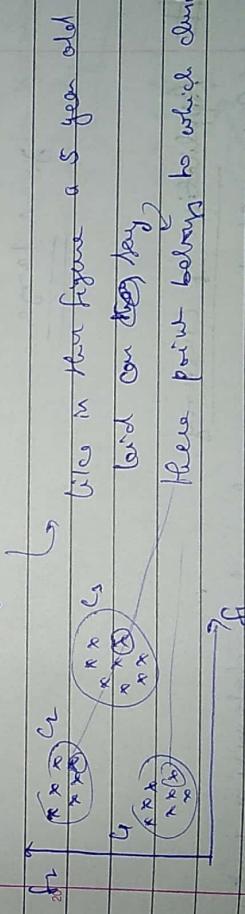
Let  $\{0, 1\} \rightarrow 2$  class classification  
 $0 \in P \rightarrow \rightarrow$  negative

But in clustering,

we have features but no labels.

So with help feature we have to group them

"Similar" cluster points.



- So, aiming to group similar elements  
points are close together,  
④ points in clusters is close together.  
⑤ points in different clusters are far away.
- Next we will see Clustering:-  
↳ Cluster hierarchical clustering,  
↳ BSCAN in detail.

## Unsupervised learning

Clustering is also known as unsupervised learning.

Unsupervised because:

1) Here we don't have (y<sub>i</sub>) to supervise the output.

Supervised learning

2) Here we have to supervise the output.

### Semi-supervised

Here in small sets we have (y<sub>i</sub>, x<sub>i</sub>)

and for large sets we have only (x<sub>i</sub>)

→ this method is used when cost of obtaining

y<sub>i</sub> is large

### Application of Clustering

But clustering was mainly used in data mining work,

but in these days a wide range of applications of clustering is used ~~there~~ in the these days.

use some e-commerce company like, Amazon, Google,

Ebay, flipkart etc. who group similar customers based on their purchase behaviour.

Doing this it becomes easier for them to target a limited no. of people with different needs.  
May right like:

In Amazon this is widely used because suppose if they want to focus the polarity of their products.  
So the comments there can be millions, then could require a lot of manual process to manually read the comments and focus on polarity.

For this reason it uses unsupervised learning & we make clusters of polarity.

Now, to decide the polarity, we take total dataset of same polarity and will check in which cluster will it fit? Then cluster it will fit we give

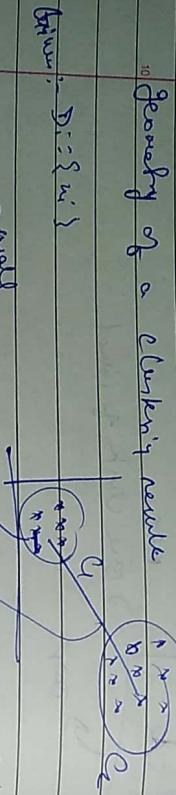
some polarity to that whole cluster.

No, analysing total dataset becomes more than the

## Metrics for Clustering

As we have seen in classification and regression we have some metrics to measure this, one more, say how far the two matrix are from each other (J<sub>i</sub>) or lesser, so can't see much in clustering.

Geometry of a clustering result



small



large

15

inter-cluster = sum (x)  $\leftarrow$   
intra-cluster = sum (x)

loop

20

Here, as a metric we will use Dunn-index -

for our ideal metric 'i' inter-cluster dist  $\rightarrow$  very high

intra-cluster dist  $\rightarrow$  very low

25

perfect

Dunn-index =  $D = \min_{i,j} d(i,j) / \max_{i,j} d(i,j)$

15

near d(i,j)

10

intracluster distance

if  $D \uparrow$  is high good clustering

For clustering we have many metrics to measure but core concept is all in

$\rightarrow$  low inter-cluster dist  
 $\rightarrow$  high intra-cluster dist

20

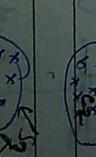
## k-means geometric intuition

Now, this one going to be very simple k-means although there is variants of k-means.

25

How C<sub>1</sub> C<sub>2</sub> C<sub>3</sub> are centroids.  
avg. point is closest to its centroid.

20



15

C<sub>1</sub> C<sub>2</sub> C<sub>3</sub> are sets of points

$$\text{dist}_{ij} = \sqrt{\sum_{k=1}^n (x_{ik} - y_{jk})^2}$$

$$S_1 \cap S_2 = \emptyset, S_2 \cap S_3 = \emptyset, S_1 \cap S_3 = \emptyset$$

if we have

1 cluster: ① ② ③ ④ ...  $\rightarrow$  k

1 centroid:  $C_1 \in S_1 \cup S_2 \cup \dots \cup S_k$

1 set:  $S_1 \cup S_2 \cup S_3 \cup \dots \cup S_k$

$$C_i = \frac{1}{n_i} \sum_{x \in S_i} x \quad \text{mean point to } C_i$$

Now, each point in set is closest to its centroid from centroid in cluster  $i$

k-means centroid based clustering scheme

Big challenge in how to find best  $k$ .

it can be done by Merge method.

Now, each point in set is closest to its centroid from other centroid of other set.

## # K-Means Algorithm (Lloyd's Algorithm)

- ① First step in initialization of k pts from D and call them centroid  $C_1, C_2, \dots, C_k$
- ② Assignment of points in  $D$  to nearest  $C_i$

Assignment

### ② Second step in Assignment

- for each point  $x_i \in D$
- with condition → future new centroid  $C_i$
- in loop for  $i=1, 2, \dots, k$   $\Rightarrow$  next
- end for  $\rightarrow$  add  $x_i$  to set  $S_i$



min  $\sum_{i=1}^k \sum_{x \in S_i} \|x - C_i\|^2 \rightarrow$  sum of dist of point to centroid placed by one. In minimized.



Sum of sq. dist from centroid in cluster.

(2) Recompute centroids / update.

$\Rightarrow$  calculate / update  $C'_j$ 's as follows.

$$C'_j = \frac{1}{n_j} \sum n_i$$

loop.  $\downarrow$

mean point

④ Repeat step 2 & 3 until convergence.

Convergence will take place, when value of old centroids and new centroids are quite similar.

$$\hookrightarrow (c_k - c'_k), \dots, (c_l - c'_l)$$

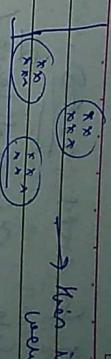
→ approach zero.

↳ final convergence.

⑤ How to initialize k-means +

Chaudhury gave algorithm where we had to choose k-centroids randomly but with constraint, that with centroid initialization there were ~~few~~ many points, hence we chose different sets of centroids, the points are split in different.

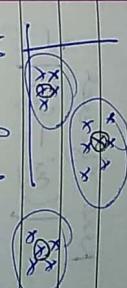
like this.



then in the clustering

cent

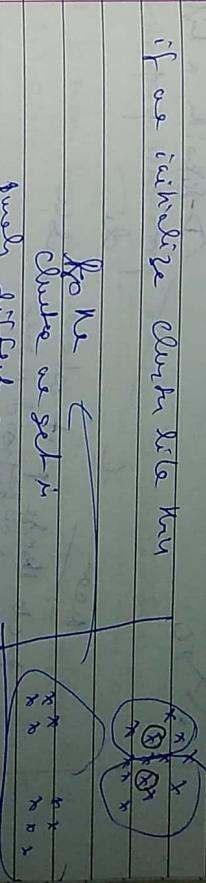
if we initialize cluster like



then cluster we get is

same

if we initialize cluster like this



so we get

cluster

much diff.

⑥ One way to come up with k-means is

① repeat between multiple runs with different initializations and pick the best clustering

based on

→ cluster into set and

→ longer inter distances.

~~→~~  $c$ -means +

② ~~→~~  $c$ -means +

→  $c$ , here we initialize randomly, but smoothly.

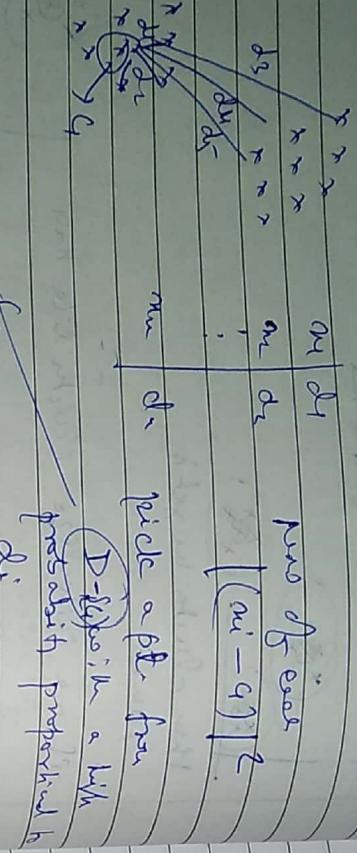
→  $c$  pick one final centroid randomly  $\rightarrow$   $c$ ,  $\forall$  next.

$\rightarrow$  K-Mean is a probabilistic approach.

By doing this approach, we can't choose the centroid point or centroid become bias could be on surface.



so between it does affected by value.



# Limitations of K-mean ~~is different type~~

→ K-mean has problem when clusters are of different sizes

- (1) Size
- (2) Densities

(3) Non-Gaussian shape  $\rightarrow$  non convex shape.

$\Rightarrow$  K-mean has problem when data contain outliers.

① Different size - vary

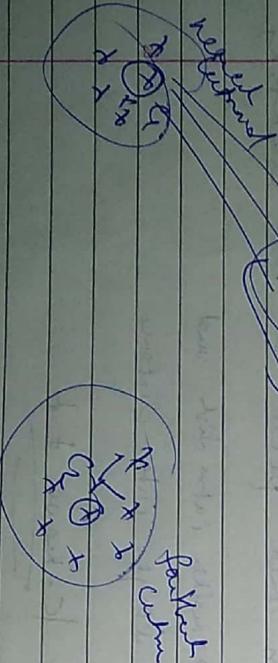


K-means (3 cluster)

mean size change. cluster set also changed, but

shouldn't happen.

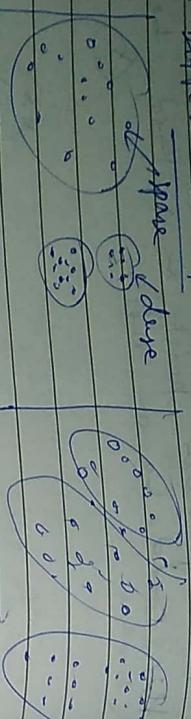
trying to pick pt. which has a far as possible from other centroids but are already picked up.



30

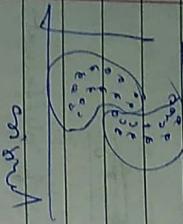
(2) Different Density

Camlin Page  
Date / /

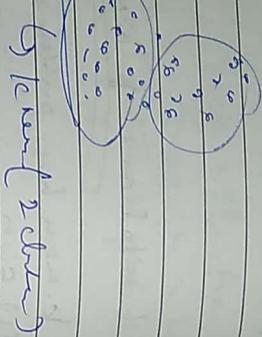


→  
→  
→  
→  
→

→  
→  
→  
→  
→

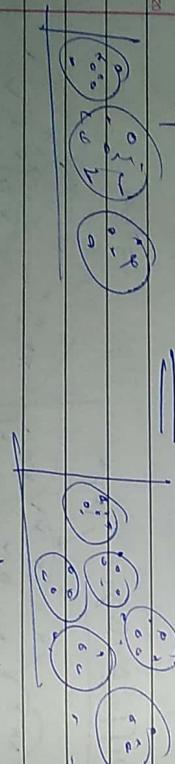


→  
→  
→  
→  
→

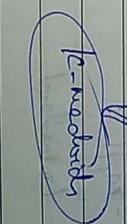


→  
→  
→  
→  
→

How all can be reduced by dividing it  
into very clusters, and put it together, but  
the process is very hard.



→  
→  
→  
→  
→



→  
→  
→  
→  
→

Partitioning around Median (PAM) : k-medoids

Camlin Page  
Date / /

(1) Trivialization → Put all as k-means → probabilistic method  
pick k pts from D.

Camlin Page  
Date / /

But here the cons idea is to take centroid from  
particular obj.

Like if we are working on test data, [what's left...]  
So, if we take only one, then  
But then say, if we take only one, then  
if we take w, will not fall in cluster.

or centroid then then would be  
much worse.

Because we can need previous from (w)  
interpretability  
between centers,

But idea is that if each centroid is a datapoint in  
clustered

k-medoids

→  
→  
→  
→  
→

b = 10

→  
→  
→  
→  
→

Trivialization → Put all as k-means → probabilistic method  
pick k pts from D.

→  
→  
→  
→  
→

(2) An object is closest medoid  
{ mostly if medoid is k closest medoid  
to me }

→  
→  
→  
→  
→

### (3) Update reservoir

Contracted update function is:

$$\hookrightarrow \hat{S}_j = \frac{1}{1 + \epsilon_j} \sum_{i=1}^n n_i$$

But in k-medoids,

(Plan)

(a) we swap each medoid with a non-medoid

points.

(b) if we decrease less we swap some  
medoid with previous medoid.

$$\stackrel{\text{if}}{=} \min_{\pi \in \Pi} \sum_{i=1}^n d(x_i, \pi)$$

non medoids are not

less than two medoids or not

Pair with k-medoids.

( $\hookrightarrow$ ) interoperability ↑

( $\hookrightarrow$ ) generalization, simplification, dist.

# Determining the right k.

$$\stackrel{\text{min}}{=} \sum_{i=1}^n \sum_{j=1}^n ||x_i - m_j||^2$$

if we take one medoid

and all other medoids are remaining

(a)  $\text{Sum} - \text{Sum } m_1 = n_1 \cdot n_2 = n_2 \rightarrow M$

$$\hookrightarrow \frac{\text{Sum} - M}{n_2} = M$$

$$\hookrightarrow \text{Sum} - M = 0$$

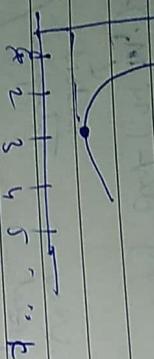
(2) "Elbow Method" or Cone method

$$K = \sum_{i=1}^n \|x_i - C\|^2 \rightarrow \text{minimize}$$

Ans =  
i.e.

Now we plot  $\mu$ ... less and variance  $\Sigma$   
and we choose  $K$  with min loss.

elbow



K-means Selection Documentation.

For 1-dimensional  
we difficult  
points on line  
cluster by

Step 1 → 1 cluster

Step 2 → 2 clusters

Step 3 → 3 clusters

Step 4 → 4 clusters

Step 5 → 5 clusters

Step 6 → 6 clusters

Step 7 → 7 clusters

Step 8 → 8 clusters

Step 9 → 9 clusters

Step 10 → 10 clusters

Step 11 → 11 clusters

Step 12 → 12 clusters

Step 13 → 13 clusters

Step 14 → 14 clusters

Step 15 → 15 clusters

Step 16 → 16 clusters

Step 17 → 17 clusters

Step 18 → 18 clusters

Step 19 → 19 clusters

Step 20 → 20 clusters

Step 21 → 21 clusters

Step 22 → 22 clusters

Step 23 → 23 clusters

Step 24 → 24 clusters

Step 25 → 25 clusters

Step 26 → 26 clusters

Step 27 → 27 clusters

Step 28 → 28 clusters

Step 29 → 29 clusters

Step 30 → 30 clusters

Step 31 → 31 clusters

Step 32 → 32 clusters

Step 33 → 33 clusters

Step 34 → 34 clusters

Step 35 → 35 clusters

Step 36 → 36 clusters

Step 37 → 37 clusters

Step 38 → 38 clusters

Step 39 → 39 clusters

Step 40 → 40 clusters

Step 41 → 41 clusters

Step 42 → 42 clusters

Step 43 → 43 clusters

Step 44 → 44 clusters

Step 45 → 45 clusters

Step 46 → 46 clusters

Step 47 → 47 clusters

Step 48 → 48 clusters

Step 49 → 49 clusters

Step 50 → 50 clusters

Step 51 → 51 clusters

Step 52 → 52 clusters

Step 53 → 53 clusters

Step 54 → 54 clusters

Step 55 → 55 clusters

Step 56 → 56 clusters

Step 57 → 57 clusters

Step 58 → 58 clusters

Step 59 → 59 clusters

Step 60 → 60 clusters

Step 61 → 61 clusters

Step 62 → 62 clusters

Step 63 → 63 clusters

Step 64 → 64 clusters

Step 65 → 65 clusters

Step 66 → 66 clusters

Step 67 → 67 clusters

Step 68 → 68 clusters

Step 69 → 69 clusters

Step 70 → 70 clusters

Step 71 → 71 clusters

Step 72 → 72 clusters

Step 73 → 73 clusters

Step 74 → 74 clusters

Step 75 → 75 clusters

Step 76 → 76 clusters

Step 77 → 77 clusters

Step 78 → 78 clusters

Step 79 → 79 clusters

Step 80 → 80 clusters

Step 81 → 81 clusters

Step 82 → 82 clusters

Step 83 → 83 clusters

Step 84 → 84 clusters

Step 85 → 85 clusters

Step 86 → 86 clusters

Step 87 → 87 clusters

Step 88 → 88 clusters

Step 89 → 89 clusters

Step 90 → 90 clusters

Step 91 → 91 clusters

Step 92 → 92 clusters

Step 93 → 93 clusters

Step 94 → 94 clusters

Step 95 → 95 clusters

Step 96 → 96 clusters

Step 97 → 97 clusters

Step 98 → 98 clusters

Step 99 → 99 clusters

Step 100 → 100 clusters

## Hierarchical Clustering Technique

# Agglomeration & Division, dendrogram.

Hierarchical clustering is divided into two part  
↳ Agglomerative &  
↳ division.

Agglomeration → popular.

initially it consider all the points as cluster.

then, combine the two closest points as cluster, then further repeat the process.

Step 1 → 10 cluster.

Step 2 → 9 cluster.

Step 3 → 8 cluster.

Step 4 → 7 cluster.

Step 5 → 6 cluster.

Step 6 → 5 cluster.

Step 7 → 4 cluster.

Step 8 → 3 cluster.

Step 9 → 2 cluster.

Step 10 → 1 cluster.

Step 11 → 0 cluster.

Step 12 → -1 cluster.

Step 13 → -2 cluster.

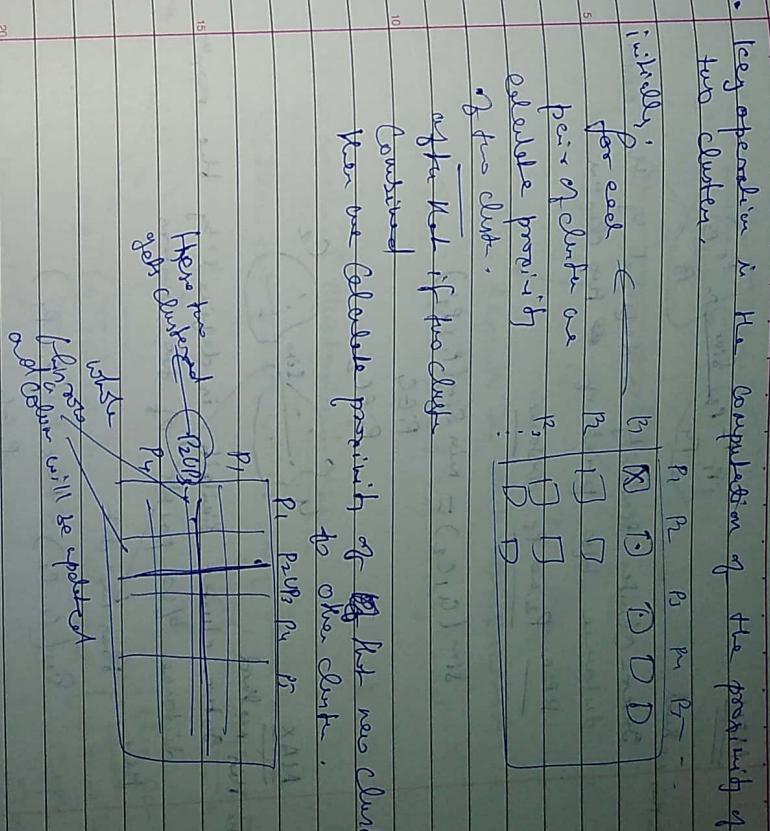
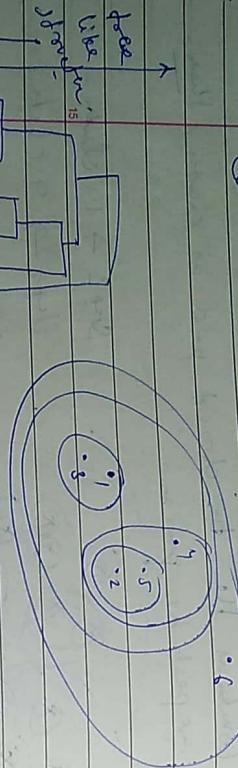
Step 14 → -3 cluster.

Q In how we need to generate the process  
to get different no. of cluster.

But here in use two algo. we don't need  
to assign the whole algo. we can stop at  
any cluster we want.

Agglomerative clustering

Can be visualized as dendrogram (like tree)



Proximity is calculated by sum of dist b/w 2 clusters

### ~~Agglomerative Clustering~~

1. Compute the proximity matrix
2. Let each data point be a cluster
3. Repeat
4. Merge the two closest cluster.
5. Update the proximity matrix
6. Until only a single cluster remains

### ~~Proximity Methods: Advantages and Disadvantages~~

How to define inter cluster proximity?

- ⇒ P1 → New dist
- ⇒ MAX dist
- ⇒ Group avg
- ⇒ Distance b/w Centroids
- ⇒ Other methods defined by own objective

MIN

In this method "similarity" is defined by the min

between cluster similarity is defined by the min

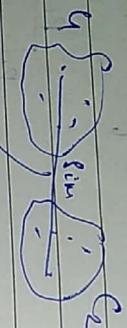
distance b/w points of two clusters.

here P<sub>i</sub> and P<sub>j</sub> distance

$$\text{sim}(C_1, C_2) = \min \text{sim}(P_i, P_j)$$

P<sub>i</sub>, C<sub>2</sub>  
P<sub>j</sub>, C<sub>1</sub>

MAX



In this method "similarity" is defined by the max

distance b/w points of two clusters

here, this method

$$\text{sim}(C_1, C_2) = \max \text{sim}(P_i, P_j)$$

P<sub>i</sub>, C<sub>2</sub>  
P<sub>j</sub>, C<sub>1</sub>

Group Avg

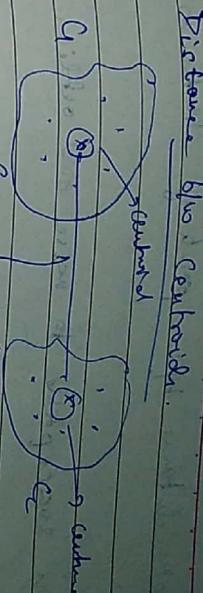
here, similarity b/w two group is defined by avg -  
dist b/w points of two cluster.

$$\text{sim}(C_1, C_2) = \sum \text{sim}(P_i, P_j)$$

P<sub>i</sub>, C<sub>2</sub>  
P<sub>j</sub>, C<sub>1</sub>

|C<sub>1</sub>| \* |C<sub>2</sub>|

here, distance b/w cluster.



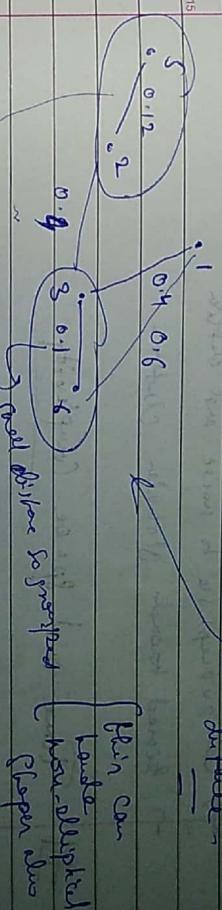
sim depends on dist b/w centroids.

This method is less extensive and

Hierarchical clustering : MIN

↳ Here two points will be group based on min

distance



Hierarchical clustering :

Limitation of MIN

→ sensitivity to noise and outliers.

Hierarchical clustering : MAX.

Since dist b/w pt. small  
its larger than

dist b/w pt. small and  
large 4.  
(S, 2, 3, 4) cannot be in  
same cluster.

Min. mean distance b/w points of cluster is considered.

## # Limitations of Hierarchical clustering

- It is less susceptible to noise and outliers.
- 5 Limitations of HAC
  - ⇒ tends to break large clusters.
  - ⇒ Biased towards globular cluster.

Linear Avg

- It is a compromise b/w Single and Complete link.
- less susceptible to noise and outliers.

- 15 → Biased towards globular clusters.

## # Time and space Complexity

Space =  $O(n^2)$  → since we have to store fin. matrix

Time =  $O(n^2)$  → since we need to iterate n iterations

Time =  $O(n^2)$  → to group 2 clusters.

and update similarity matrix,  $\sim O(n^2)$

$$\therefore O(n^2) = O(n^2)$$

Time complexity

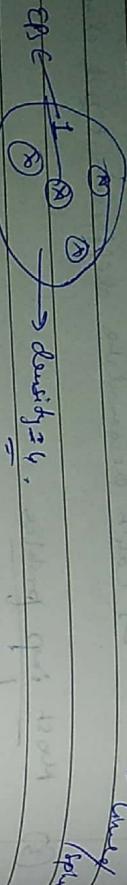
Space complexity

# Minpts and Eps: Density.

↳ Note we have hyperparameter "Eps".

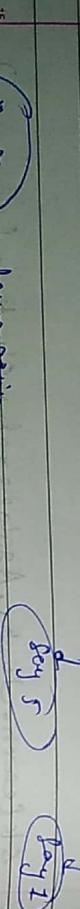
- ① density at a  $P$ : no. of points within a hypersphere of radius  $\epsilon_{\text{pt}}$  around  $P$

→  $\text{core pt}$  if  $\geq \text{minpts}$

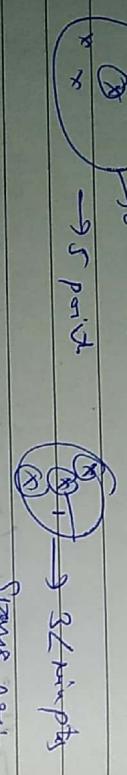


10

- ② 'dense region': a hypersphere of radius  $\epsilon_{\text{pt}}$  that contains atleast minpts points.



15



20

# Core point, Border point & noise point

25

Core D = 2 minpts and minpts, Minpts, Eps and

- ① 'Core point': if  $\text{pt} \geq \text{minpts}$  points in an  $\epsilon_{\text{pt}}$  medium around it

↳ Core point always below  
no dense region.

# Border pt ( $P$ )

①  $P$  is not a core point  $\Rightarrow P$  has  $<$  minpt points in  $\epsilon_{\text{pt}}$  medium

- ②  $P \in \text{Neighbourhood}(\mathcal{Q})$

→ core point

$$\text{dist}(P, Q) \leq \epsilon_{\text{pt}}$$

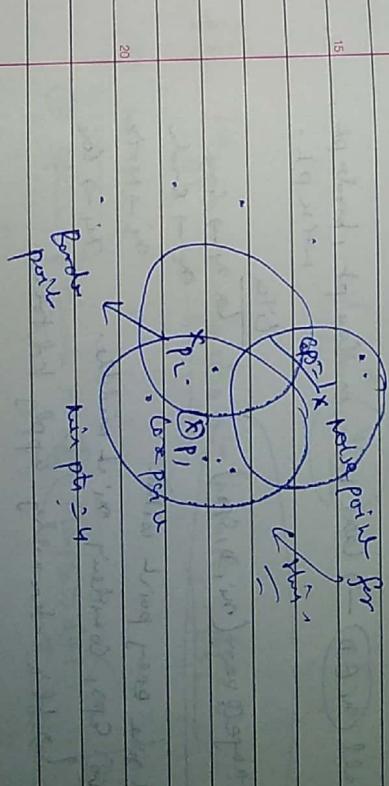
10

- ③ 'Noise point':

↳ neither core pt nor a border pt.



15



20

# Density edge and Density Connected points:

- ① 'Density edge':

↳ if  $P_1 \rightarrow \text{core pt}$       ↓  
and  $\text{dist}(P_1, P_2) \leq \epsilon_{\text{pt}}$       ↓  
 $P_2$  density edge

25

30

## ② Density Connected pts :-

if  $P_{14}$  and  $\rightarrow$  line pt.

$\leq \epsilon_D$   $\leq \epsilon_D$

$P_1$   $P_2$   $P_3$   $P_4$   $P_5$   $P_6$   $P_7$   $P_8$   $P_9$   $P_{10}$   $P_{11}$

line pt

10

10

11

11

11

11

11

11

11

# DSCAN Algorithm.

Step

for all  $n_i \in D \rightarrow$  label mean on line pt, border pt,  $n_i \neq p_2$ .

like

$S_i = \text{range}(n_i, D, \epsilon_D)$

$n_i \rightarrow \text{line}$

$n_i \rightarrow \text{border}$

$n_i \rightarrow \text{inter}$

$n_i \rightarrow \text{line}$

Step

25

remove all noise pts from your data

$\rightarrow$  sparse region  $\rightarrow$  don't assign to any cluster

$n_i \rightarrow \text{line}$

Step

30

for each line  $p_i$ ,  $p_i$  not assigned to a cluster

(a) create a new cluster with  $p_i$

(b) add all pts that are density connected

to  $p_i$  into this new cluster

this can be achieved by,

$\left| p_i - p_j \right| \leq \epsilon_D$

True

Typically,  $\text{Minpts} \approx 2^n$

25

value of

$\text{Minpts} \rightarrow$  dataset is more noisy  $\rightarrow$  reduce  $\text{Minpts}$ .

reduce  $\text{Minpts}$ .

16

$\text{Minpts} \rightarrow$  dataset is more noisy  $\rightarrow$  reduce  $\text{Minpts}$ .

17

$\text{Minpts} \rightarrow$  dataset is more noisy  $\rightarrow$  reduce  $\text{Minpts}$ .

18

$\text{Minpts} \rightarrow$  dataset is more noisy  $\rightarrow$  reduce  $\text{Minpts}$ .

19

$\text{Minpts} \rightarrow$  dataset is more noisy  $\rightarrow$  reduce  $\text{Minpts}$ .

20

$\text{Minpts} \rightarrow$  dataset is more noisy  $\rightarrow$  reduce  $\text{Minpts}$ .

21

$\text{Minpts} \rightarrow$  dataset is more noisy  $\rightarrow$  reduce  $\text{Minpts}$ .

22

$\text{Minpts} \rightarrow$  dataset is more noisy  $\rightarrow$  reduce  $\text{Minpts}$ .

23

$\text{Minpts} \rightarrow$  dataset is more noisy  $\rightarrow$  reduce  $\text{Minpts}$ .

24

$\text{Minpts} \rightarrow$  dataset is more noisy  $\rightarrow$  reduce  $\text{Minpts}$ .

25

$\text{Minpts} \rightarrow$  dataset is more noisy  $\rightarrow$  reduce  $\text{Minpts}$ .

26

$\text{Minpts} \rightarrow$  dataset is more noisy  $\rightarrow$  reduce  $\text{Minpts}$ .

27

$\text{Minpts} \rightarrow$  dataset is more noisy  $\rightarrow$  reduce  $\text{Minpts}$ .

28

$\text{Minpts} \rightarrow$  dataset is more noisy  $\rightarrow$  reduce  $\text{Minpts}$ .

29

$\text{Minpts} \rightarrow$  dataset is more noisy  $\rightarrow$  reduce  $\text{Minpts}$ .

30

$\text{Minpts} \rightarrow$  dataset is more noisy  $\rightarrow$  reduce  $\text{Minpts}$ .

(2) DB can choose min-pt by domain experts.

(2)  $\epsilon_{DB}$  let  $\text{min-pt} = 4$ .

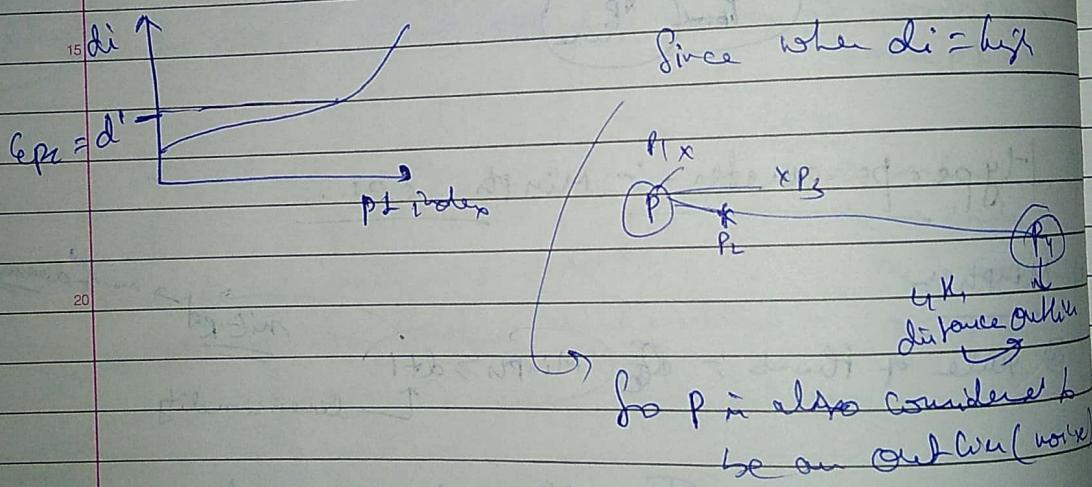
(a)  $x_i \rightarrow d_i$ : dist from  $x_i$  to the 4<sup>th</sup> neighbor of  $x_i$

(b) sort  $d_i$  in increasing order:

$n \propto L$

$n \propto L$

Now, Select  $\epsilon_{DB}$  by Elbow method.



Advantage and Limitations of DBSCAN.