

11ly we can add dummy features to get complex boundary.

But overfitting, F!

Regularisation: "I'm gonna save you"

$$\text{cost}' = \text{cost} + \lambda \sum (m_i)$$

~~Also~~ = high $\lambda \rightarrow$ underfitting

$$L1 = \text{cost} + \lambda \sum |m_i| \quad \text{lasso}$$

$$L2 = \text{cost} + \lambda \sum m_i^2 \quad \text{ridge}$$

* No regularisation on intercept.

Performance measures in classification

(1) Confusion matrix

sklearn metrics \rightarrow confusion matrix

Predicted	Actual	
	T	F
T	TP	FP
F	FN	TN

$$\text{TPR} = \frac{TP}{TP + FN}$$

$$\text{FPR} = \frac{FP}{FP + TN}$$

~~Base~~ For a class (n)

$$\text{Precision} = \frac{TP}{TP + FP}$$

$$\text{Accuracy} = \frac{TP + \overset{T}{\cancel{F}}N}{TP + FN + TN + FP}$$

$$\text{Recall} = \frac{TP}{TP + FN} = \text{TPR} = \text{Sensitivity}$$

$$\text{Sensitivity} = \frac{TP}{TP + FN} = \text{TPR} \quad (\text{ability to detect +ve class})$$

$$\text{Specificity} = \frac{TN}{TN + \cancel{F}P} = 1 - \text{FPR}$$

Bayes Theorem : Naive Bayes :

Bayes Theorem : $P(A|B) = \frac{P(B|A) P(A)}{P(B)}$

$y = a_1, a_2, \dots, a_k$ (output classes)

then $Y = \max_{i=1}^k (P(y = a_i | x = x))$

$= \max_{i=1}^k \left(\frac{P(x = x | y = a_i) P(y = a_i)}{P(x = x)} \right)$

$Y = \max_{i=1}^k (P(x = x | y = a_i) P(y = a_i))$

output class

$P(y = a_i) = \frac{|a_i|}{|y|} = \frac{|a_i|}{|a_1| + |a_2| + \dots + |a_k|}$

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

Independent Events : $P(A \cap B) = P(A) * P(B)$

Given input $x \rightarrow \{x_1, x_2, \dots, x_k\}$ where x_i is value corresponds to feature f_i .

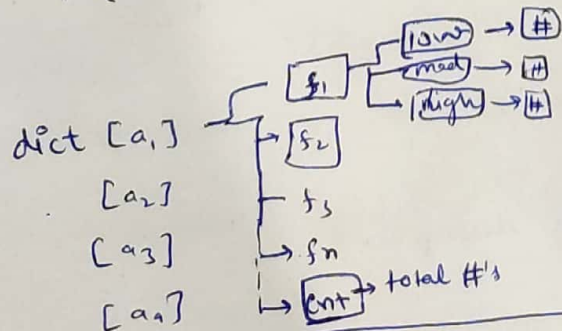
$P(x = x | y = a_i) = P((x_1, x_2, x_3, \dots, x_k) | y = a_i)$

* Assumption \rightarrow All features are independent.

$= P(f_1 = x_1 | y = a_i) * P(f_2 = x_2 | y = a_i) * \dots * P(f_k = x_k | y = a_i)$

$P(x = x | y = a_i) = \prod_{j=1}^k P(f_j = x_j | y = a_i)$

$P(f_j = x_j | y = a_i) = \frac{\#(f_j = x_j \text{ \& } y = a_i)}{\#(y = a_i)}$



Given data x train

Y - train

	f_1	f_2	f_3	\dots	f_k	label
low						a_1
med						a_2
high						a_3
						a_5
						a_n

$P(f_j = x_j | y = a_i) = \frac{\text{dict}[a_i][f_j][x_j]}{\text{dict}[a_i][\text{ent}]}$

Laplace correction: $P()$ can be 0 if there are no value (x_j) present in feature (F_j) given a_i
 so we modify formula

$$= \frac{\text{dict}[a_j][F_j][x_j] + 1}{\text{dict}[a_j][\text{cnt}] + |F_j|}$$

→ values F_j can take

Continuous data?

$$p(x_j = x_j | y = a_i) = \frac{1}{\sqrt{2\pi}\sigma^2} e^{-\frac{(x_j - \mu)^2}{2\sigma^2}}$$

Naive-bayes can be used as ^{text} ~~speech~~ collection

KNN → k-Nearest Neighbour Do feature scaling before KNN.
 You know what it is!

Cross validation:

low value of k leads to overfitting

high value of k leads to underfitting

Make sure we have relevant data and independent.

→ Assign weights to feature (w_i for f_i)

$$\sum_{i=1}^n w_i (x_1^i - x_2^i)^2$$

→ Feature selection (Backward elimination)

Handling labeled data:

→ use 0, 1 in binary data. (2 values)

→ we can't assign 0, 1, 2, 3... to discrete data bcz distance would be effected.

	Red	Blue	Green
Red	0	0	1
Blue	0	0	1
Green	1	0	0
	0	1	0

→ Other KNNs

→ Brute force → compare with all datapoint

→ KD Trees something like BST

→ Ball Trees

Pros of KNN → Easy to understand and code (lol)

(2) Works for multiclassification

also - not - ...

Cons of KNN:

(1) Testing time is huge.

(2) It can be biased, if training data is biased.

(3) Curse of dimensionality.

SVM

The shortest distance b/w the observations and the threshold is called margin.

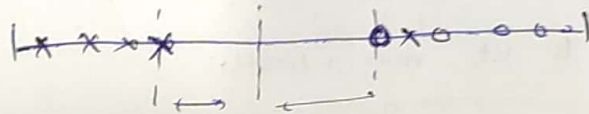
Maximal margin classifier (MMC): threshold that gives us the largest margin.

MMC is not optimal (think of outliers)

We do cross validation to find the best soft margin aka Support Vector Classifier.

The support vector classifier is a hyperplane in n -dimension.

So SVC can handle outliers nicely using cross validation.



But what about

so many overlaps...

Support Vector machine: We add extra dimensions to points and then try to fit hyperplane, ~~sup~~ SVC.

How to transform the data? SVM uses kernel function. We can find good value of (d) dimension via cross validation.

Other kernel is RBF (finds SVM in infinite dimension).

Kernels just calculate the σ b/w points in high dimension without actually transforming them.

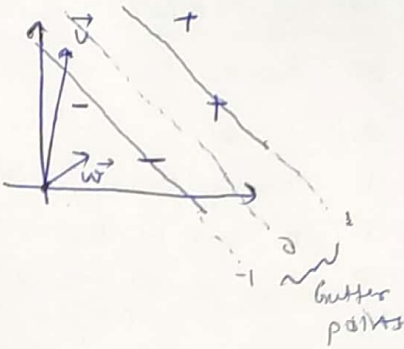
$$\text{let } r = \frac{1}{2} d^2$$

$$\text{Polynomial kernel} = (a \cdot b + r)^d$$

$$= a \cdot b + a^2 b^2 + \frac{1}{4}$$

$$= \left(\underset{x}{a}, \underset{y}{a^2}, \underset{z}{\frac{1}{2}} \right), \left(\underset{x}{b}, \underset{y}{b^2}, \underset{z}{\frac{1}{2}} \right)$$

MIT OCW



is b to margin.
be unknown sample point
 $\vec{w} \cdot \vec{x} \geq c$

i.e. $\vec{w} \cdot \vec{x} + b \geq 0$ for + sample
but we still don't know \vec{w} and b .

$$\vec{w} \cdot \vec{x}_+ + b \geq 1 \quad \dots (1)$$

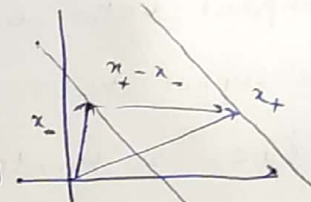
$$\vec{w} \cdot \vec{x}_- + b \leq -1 \quad \dots (2)$$

introduce y_i such that $y_i = 1$ +ve sample
-1 -ve sample
eq 1 and 2 become

$$y_i (\vec{x}_i \cdot \vec{w} + b) \geq 1$$

$$\begin{cases} y_i (\vec{x}_i \cdot \vec{w} + b) - 1 \geq 0 \\ y_i (\vec{x}_i \cdot \vec{w} + b) - 1 = 0 \end{cases} \rightarrow \text{For buffer points}$$

$$\begin{aligned} \text{width} &= (\vec{x}_+ - \vec{x}_-) \cdot \frac{\vec{w}}{\|\vec{w}\|} \\ &= \frac{\vec{x}_+ \cdot \vec{w}}{\|\vec{w}\|} - \frac{\vec{x}_- \cdot \vec{w}}{\|\vec{w}\|} = \frac{1-b}{\|\vec{w}\|} + \frac{1+b}{\|\vec{w}\|} = \frac{2}{\|\vec{w}\|} \end{aligned}$$



we are maximising $\frac{1}{\|\vec{w}\|}$ to get max width.

$$\boxed{\text{minimise } \|\vec{w}\| \sim \text{minimise } \frac{1}{2} \|\vec{w}\|^2} \quad \text{Primal Form}$$

we have to use Lagrange multipliers.

$$L = \frac{1}{2} \|\vec{w}\|^2 - \sum \alpha_i [y_i (\vec{w} \cdot \vec{x}_i + b) - 1] \quad (3)$$

$$\frac{\partial L}{\partial \vec{w}} = \vec{w} - \sum \alpha_i y_i \vec{x}_i = 0 \Rightarrow \vec{w} = \sum \alpha_i y_i \vec{x}_i$$

$$\frac{\partial L}{\partial b} = \sum \alpha_i y_i = 0$$

substitute in eq 3.

$$L \Rightarrow \frac{1}{2} (\sum \alpha_i y_i \vec{x}_i) \cdot (\sum \alpha_j y_j \vec{x}_j) - \sum \alpha_i [y_i (\sum \alpha_j y_j \vec{x}_j \cdot \vec{x}_i + b) - 1]$$

$$L \Rightarrow \frac{1}{2} (\sum \alpha_i y_i \vec{x}_i) \cdot (\sum \alpha_j y_j \vec{x}_j) - (\sum \alpha_i y_i \vec{x}_i \cdot (\sum \alpha_j y_j \vec{x}_j) + \sum \alpha_i b) = \sum \alpha_i y_i b + \sum \alpha_i$$

$$\boxed{L = \sum \alpha_i - \frac{1}{2} \sum_i \sum_j (\alpha_i \alpha_j y_i y_j (\vec{x}_i \cdot \vec{x}_j))} \rightarrow \text{Max of this.} \quad (4)$$

depend on dot product.

$$\sum a_i y_i \underbrace{x_i \cdot \vec{u}}_{\text{dot product}} + b \geq 0 \quad +ve \text{ scalar.}$$

Won't work with non linear decision boundary.

transformator: $\phi(\vec{x}) = \phi(\vec{x}_i) \cdot \phi(\vec{x}_j)$ to max
 $\phi(x_i) \cdot \phi(x_j)$

$$k(x_i, x_j) = \phi(\vec{x}_i) \cdot \phi(\vec{x}_j)$$

↳ kernel function. we actually don't have to do ~~dot~~ transform
 so ~~$\phi(\vec{x})$~~ , all depends on dot product.

some kernels: (1) $(\vec{u} \cdot \vec{v} + c)^d$ Polynomial kernel.
 (2) $e^{-\frac{\|\vec{x}_i - \vec{x}_j\|^2}{2}}$ Radial basis kernel

Neural Network

Activation Function → sigmoid f^n , ReLU, Soft Plus.

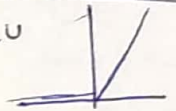
We get x-axis coordinate

from the weight & and use activation fun to get Y.

Weights come from Activation Function.

Act. (weight + biases) → new perception.

ReLU



$$f(x) = \begin{cases} x & x \geq 0 \\ 0 & x < 0 \end{cases}$$

$$f(x) = \max(0, x)$$

Maam PPT's after mid sem.

KNN → pattern recognition

What k to choose? $k < \sqrt{n}$

→ While scaling we can add weight to the important features. w_k
we can learn w_k using cross-validation.

Ez:)

Kmeans Algo : (Unsupervised learning)

Set quest PPT

- (1) Select a k (clusters to identify) ($k=3$)
- (2) Randomly select k distinct dp. (~~centroids~~)
- (3) Measure distance of remaining points to clusters and classify them to nearest cluster.
- (4) Calculate mean of each cluster (Centroid)
- (5) Repeat. (Measure distance along the mean this time)
- (6) ~~Repeat~~ When it does change its result and calculate the variance.

Centroid can be real or imaginary dp

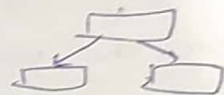
Repeat whole with different k and diff. starting points.
The best is when we have a gradual change in ~~variance~~ variance.

resource: Set Quest

PPT

Ada Boost → contains stumps

The tree with just one node and two leaves is stump.



~~Multiple~~ Forest of stumps.

A tree uses all features to make decision, but a stump would take one feature to make decision.

So 'Stumps' are weak learners.

In RF, each tree has equal vote to classification.

In Adaboost, each stump has different vote to classification.

In RF, each tree is independent.

In Adaboost, order is important, The error of first stump influences how second stump is made and soon.

Steps

(1) ~~make~~ stumps of

(1) Assign sample weight (this will tell which dp is important)
(initially for all = $1/\text{no. of samples}$)

(2) Make stumps using each feature and calculate gini index.
the one with lowest would be the first stump.

(4) Now we find the ~~weight~~ ^{vote} of this stump (by calculating error).

$$\frac{\text{Amount of vote}^{\text{log}}}{\frac{1}{2} \log \left(\frac{1 - \text{Total Error}}{\text{Total Error}} \right)} \left\{ \begin{array}{l} \# \text{ of incorrectly} \\ \text{classified with their} \\ \text{sample weight.} \end{array} \right.$$

(5) Now we modify sample weight.

We increase the sample weight of incorrectly classified weight of that feature, so that next time we try to classify it correctly and decrease all other sample weights.

$$\uparrow \text{Net sample weight} = \text{sample weight} \times e^{\text{error of vote}}$$

↓ and calculate new sample weight. (also we normalize later) and find rest of stumps.

Ensemble Technique → combining multiple models.

Need to add from Main

→ Bagging: Random Forest
(Bootstrap aggregation)

→ Boosting: ADABOOST, GRADIENT BOOSTING, XG BOOST.

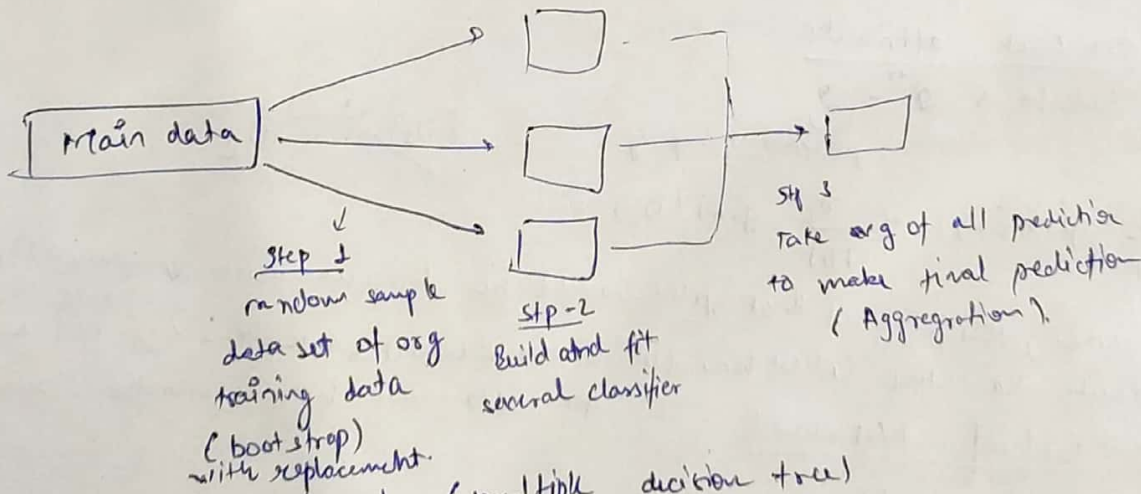
Simple learn sklearn ✓

Error → Bias + variance.

PPT ✓

using Bagging we can decrease variance
using Boosting bias

Bagging → bootstrap aggregation reduces variance of an estimate by taking mean of multiple estimates.



eg Random forest (multiple decision tree)

Boosting reduces bias by training weak learners sequentially, each trying to correct predecessor.

eg: → ADA BOOST (multiple stumps)
→ Gradient Boosting (GBM)
→ XG BOOST

{ 2, 1, 5 : 3 }

confidence if these rules
can be formed on not.

{ 5 } → { 1 }

SVM

$$\text{Hyperplane} = w^T \phi(n) + b = 0$$

↳ error
(bias)

$$\text{Distance of a point } (x_0) \text{ from hyperplane} = \frac{|w^T \phi(x_0) + b|}{\|w\|_2} = d_H \phi(x_0)$$

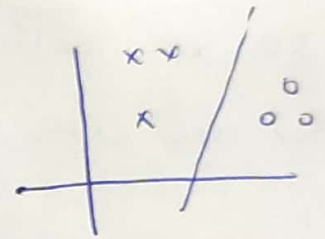
we aim to find hyperplane that has maximum distance from the closest point - (maximise the minimum distance)

$$\text{i.e. } w^* = \arg \max_w [\min_n d_H \phi(x_n)]$$

$$\text{For } x \rightarrow w^T \phi(n) + b \geq 0$$

$$0 \rightarrow w^T \phi(n) + b < 0$$

$$\rightarrow y_n [w^T \phi(n) + b] \geq 0 \quad \{\text{correct}\}$$
$$< 0 \quad \{\text{incorrect}\}$$



$$w^* = \arg \max_w \left[\min_n \frac{w^T \phi(x_0) + b}{\|w\|_2} \right] \quad \because \text{perfect separation}$$

$$w^* = \arg \max \frac{1}{\|w\|_2^2}$$

$$\min \frac{1}{2} \|w\|_2^2$$

~~Dual~~
Primal Form

Can 2 Non-perfect separation.

New Primal Form of SVM

$$\xi_n \geq 0$$

$$\min \frac{1}{2} \|w\|_2^2 + c \sum_n \xi_n$$

↑ penalty for incorrect classification

$c=0$: Less complex boundary

$c=\infty$: More complex boundary.

We want to get rid of $\phi(n)$.

We use kernels to get dual form.

Method of Lagrangian multipliers

Steps

(1) Obtain new primal form = $\min \frac{1}{2} \|w\|_2^2 + c \sum \xi_n$

(2) Determine Lagrange.

(3) $\frac{dL}{dw} = 0$.

(4) Express in form of duals

(5) Back substitution.

(6) Depends to dot product.

Dual form

↳ in terms of multiplication of 2 $\phi(n)$.

which is simplified by kernels

Hyperparameters in kernelisation:

(1) Margin

(2) Kernel

(3) Regularisation (c) [how much you want to avoid misclassification]

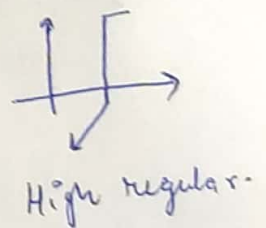
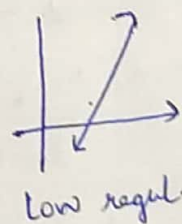
(4) Gamma.

low gamma

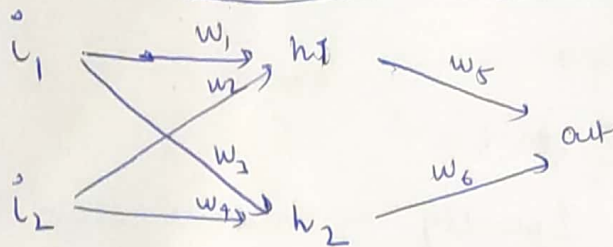
high gamma

↳ Hyperplane is influenced by far distant points

↳ Hyperplane is influenced by nearby points only.



Feed Forward Neural Network



$$h_1 = i_1 w_1 + i_2 w_2$$

$$h_2 = i_1 w_3 + i_2 w_4$$

$$\text{predicted out} = w_5 h_1 + w_6 h_2 = (i_1 w_1 + i_2 w_2) w_5 + (i_1 w_3 + i_2 w_4) w_6$$

$$\text{Error} = \frac{1}{2} (\text{out} - \text{predicted})^2$$

$$\text{new } w_x = w_x - \alpha \frac{d \text{Error}}{d w_x}$$

$$\begin{aligned} w'_6 &= w_6 - \alpha \frac{d \frac{1}{2} (\text{out} - \text{pre})^2}{d \text{pre}} \times \frac{d \text{pre}}{d w_6} \\ &= w_6 - \alpha (\text{out} - \text{pre}) \times (i_1 w_3 + i_2 w_4) \\ &= w_6 - \alpha \Delta \times h_2 \end{aligned}$$

$$w'_6 = w_6 - \alpha \Delta h_2$$

$$w'_5 = w_5 - \alpha \Delta h_1$$

$$\begin{aligned} w'_1 &= w_1 - \alpha \frac{d \text{Error}}{d w_1} \\ &= \alpha \frac{d \frac{1}{2} (\text{out} - \text{pre})^2}{d \text{pre}} \times \frac{d (\text{pre})}{d h_1} \times \frac{d (i_1 w_1 + i_2 w_2)}{d w_1} \\ &= \alpha \Delta w_6 \times i_1 \end{aligned}$$

$$w'_1 = w_1 - \alpha i_1 \Delta w_6$$

$$w'_2 = w_2 - \alpha i_2 \Delta w_6$$

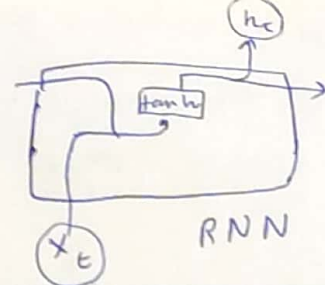
$$w'_3 = w_3 - \alpha i_1 \Delta w_6$$

$$w'_4 = w_4 - \alpha i_2 \Delta w_6$$

No memory in
Feed Forward Network.

- RNN models sequence data
- short term value.

long short term memories.

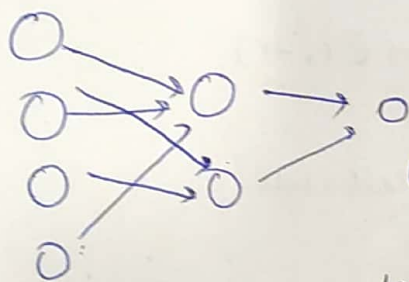


- CNN is a feed forward neural network that is generally used to analyse visual image by preprocessing data with grid like topology.
- Every image is represented in form of arrays of pixel values.

→ layers in CNN

- ① Convolution layer → has no of filters that perform convolution operation ^{on} input → filters → feature map
- ② ReLU layer → Move convolute value to ReLU (Rectified feature map)
- ③ Pooling layer → Pooling is a down sampling operation that reduces dimensionality of feature map. (~~2 dimension arrays~~) (pooled feature map)
- ④ Fully connected layer → then we do flattening them into 1D.

↓
The flattened matrix from the pooling is fed as input to fully connected layer. (Forward fed Neural network)



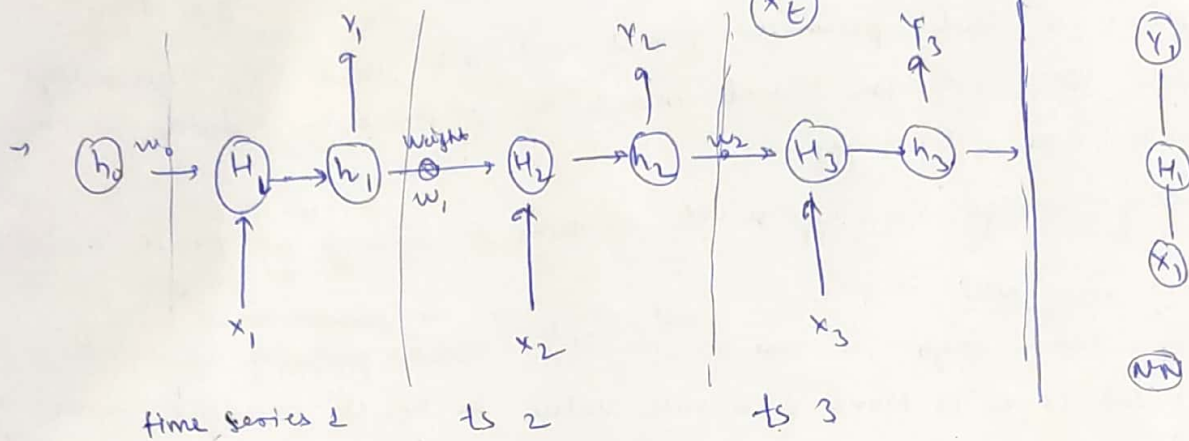
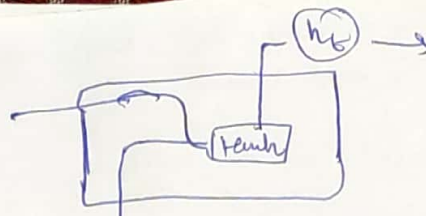
We get multiple feature map as we slide filters over the image.
There are multiple features.
[We can see this is not fully connected graph] ~~[# not]~~

- (1) So in convolution layer we have multiple filters slides over image so we get multiple feature matrix
- (2) Pass them through ReLU to eliminate -ve values.
- (3) Pooling ~~steps~~ layers, take a window of some size n and reduce dimensionality of a ~~the~~ rectified feature map (we take max among all values in window)

Now we stack up all the feature matrix, take each matrix and shrink to 1D matrix and append and feed it to Fully connected layer (Forward feed network).

Recurrent Neural Network

- model for sequential data
- has memory. (state)



$$h_t = f(h_{t-1} \times x_t)$$

the function is tanh

$$y_t = (\text{some output weight}) \times h_t$$

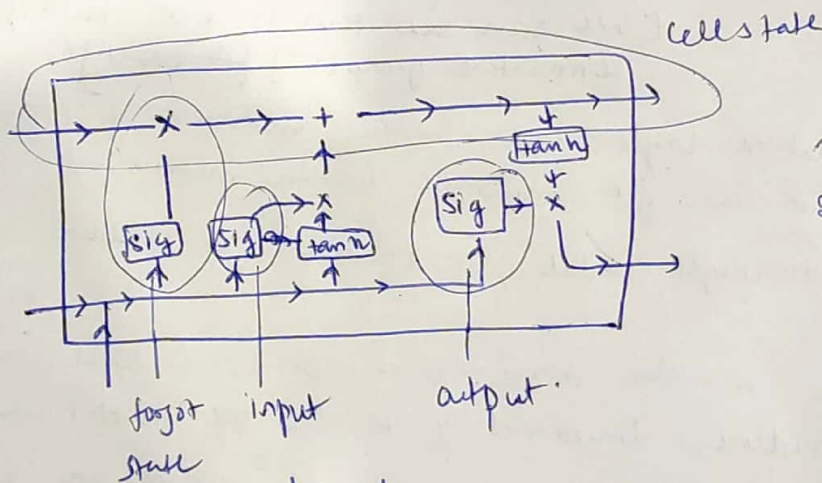
taylor series

Previous states starts vanishing due to backpropagation so we use LSTM (vanishing gradient)

LSTM (Long Short Term Memory)

→ use gates to store info

tanh makes sure value are b/w [-1, 1] (h_t)



3 sigmoid & tanh

sigmoid → $\begin{cases} 1 \leftarrow \text{retained} \\ 0 \leftarrow \text{rejected} \end{cases}$

used in speech recognition.