

FDP on Data Analytics in Healthcare

Traditional Learning Algorithms for Healthcare

Presented by

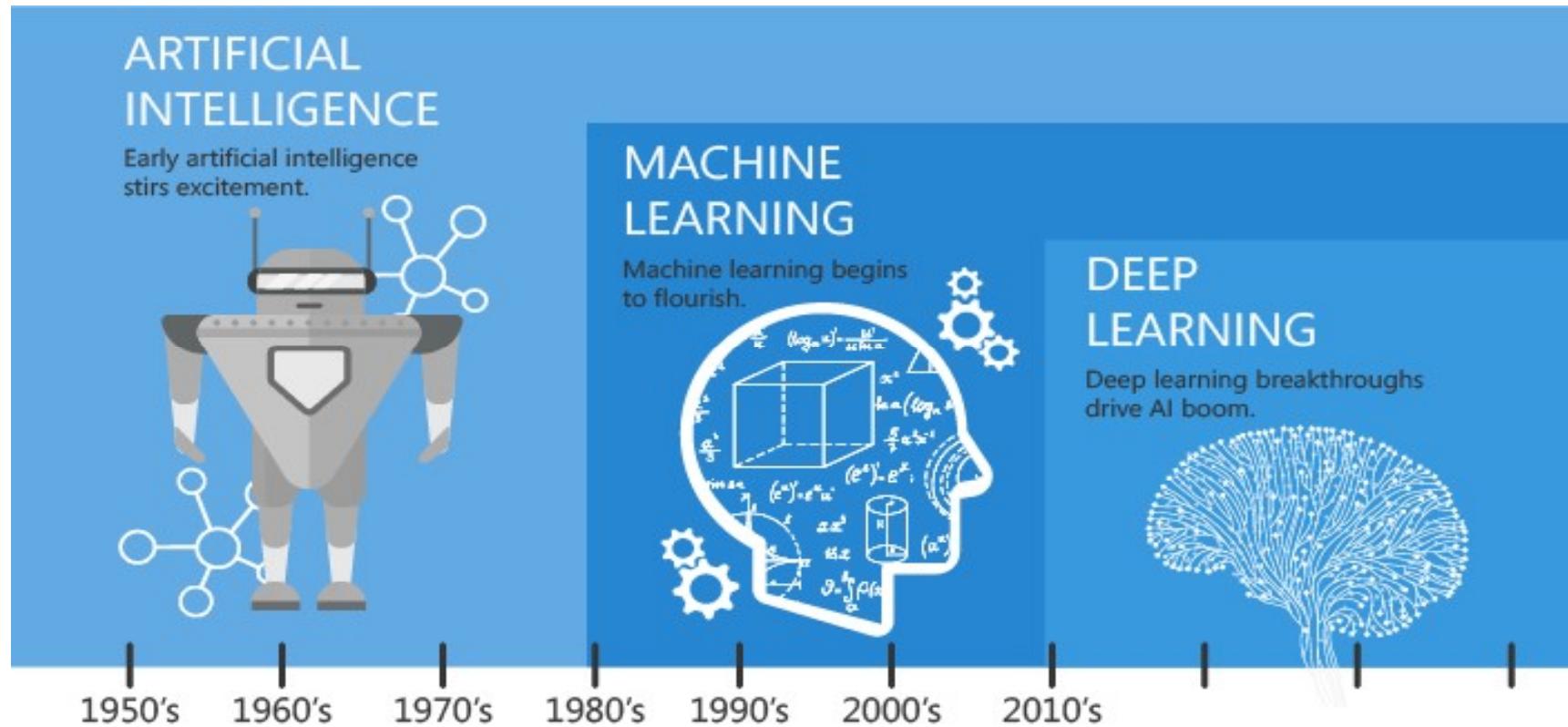
Dr. S. Kavitha

Associate Professor

Department of CSE, SSN CE

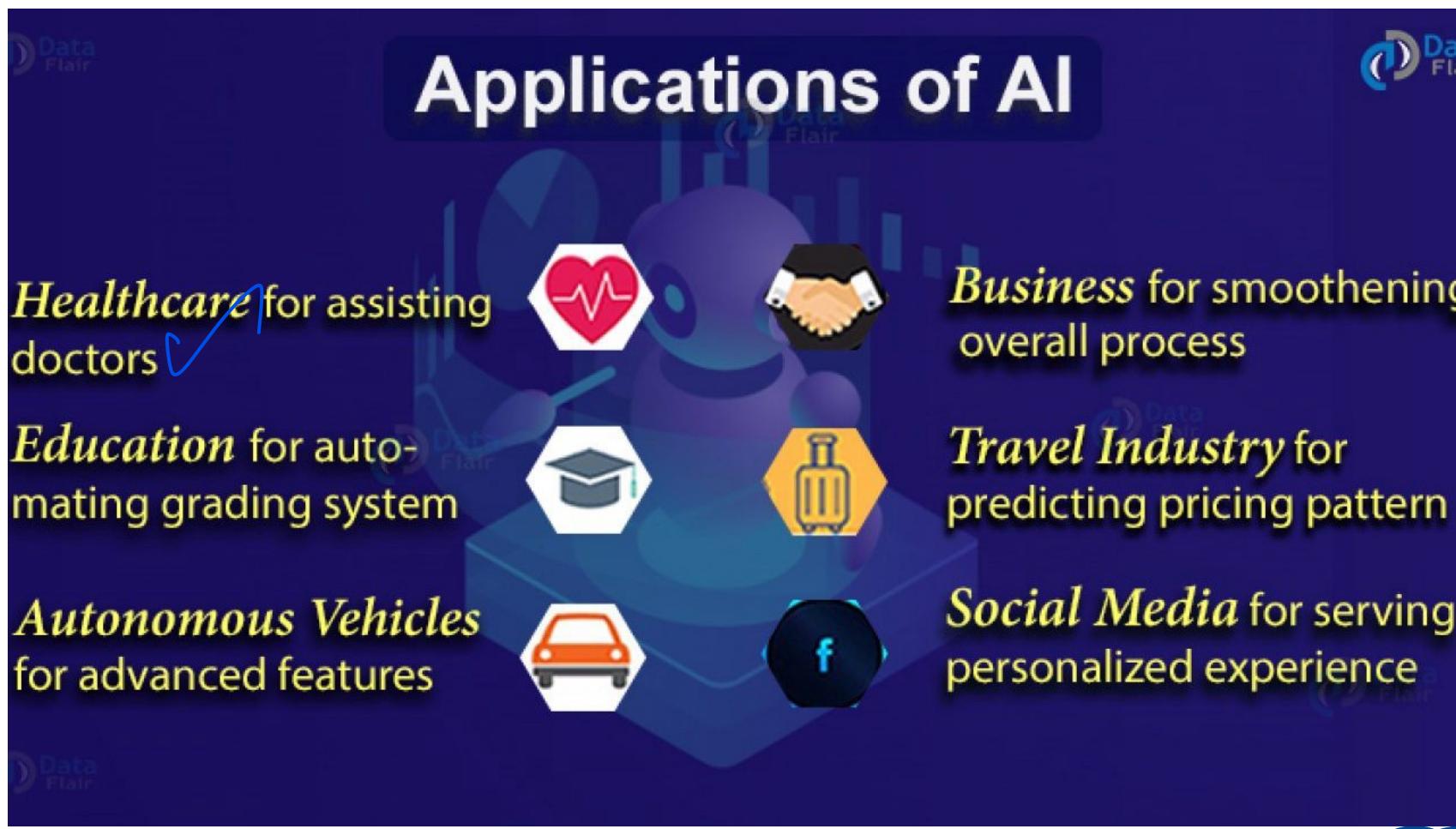


AI-ML-DL



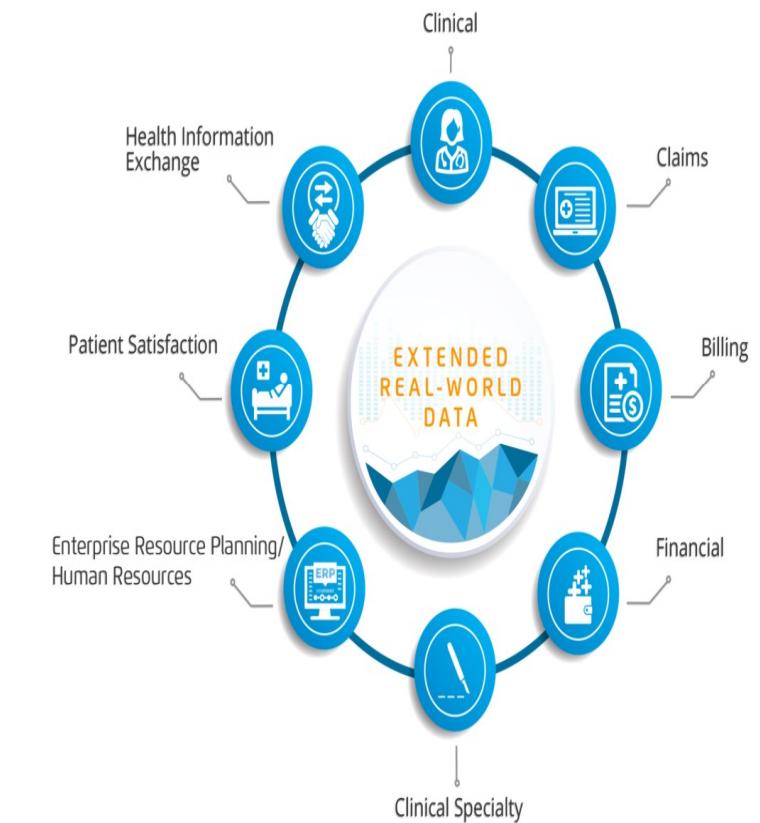
Since an early flush of optimism in the 1950's, smaller subsets of artificial intelligence - first machine learning, then deep learning, a subset of machine learning - have created ever larger disruptions.

Applications of AI



Healthcare Dimensions

- Data promises different dimensions in health care - Prediction of Chronic Diseases, Personalized Medicine, Precision Medicine, Drug Discovery, Remote Health Monitoring, Genomic Analysis, Patient's Data Management – Fraudulent Medical Claims.
- Medical imaging - improves medical decision making and reduces unnecessary medical procedures

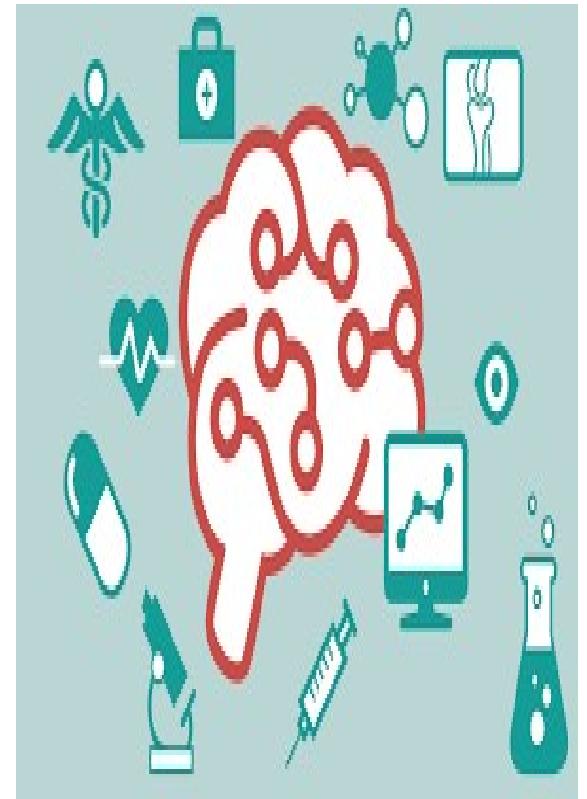


Data Sources in Healthcare



Handshake of AI with Healthcare

- Healthcare faces a overflow of rich data that is enables to operate with greater insight and effectiveness, and deliver better service.
- Advanced AI techniques coupled with the explosion of data in healthcare to uncover
 - leading clinical practices,
 - minimize research discovery time,
 - offer new personalized engagement paradigms at an industrial scale that align people's decisions and actions in ways that improve outcomes.



Objectives

reliability
model

- Why?
 - AI and ML needs to provide explanation in healthcare
- What?
 - Type of explanation do we need in healthcare
- How?
 - To select the right ML algorithm for explanation
- Where?
 - Settings and different types of interpretable ML models
- When?
 - The past, present and future of explainable AI in healthcare

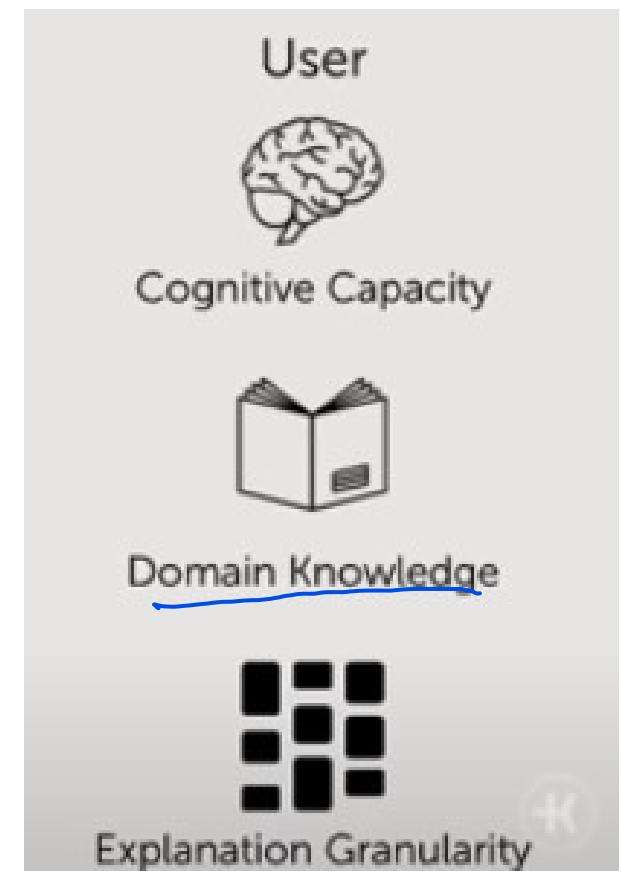
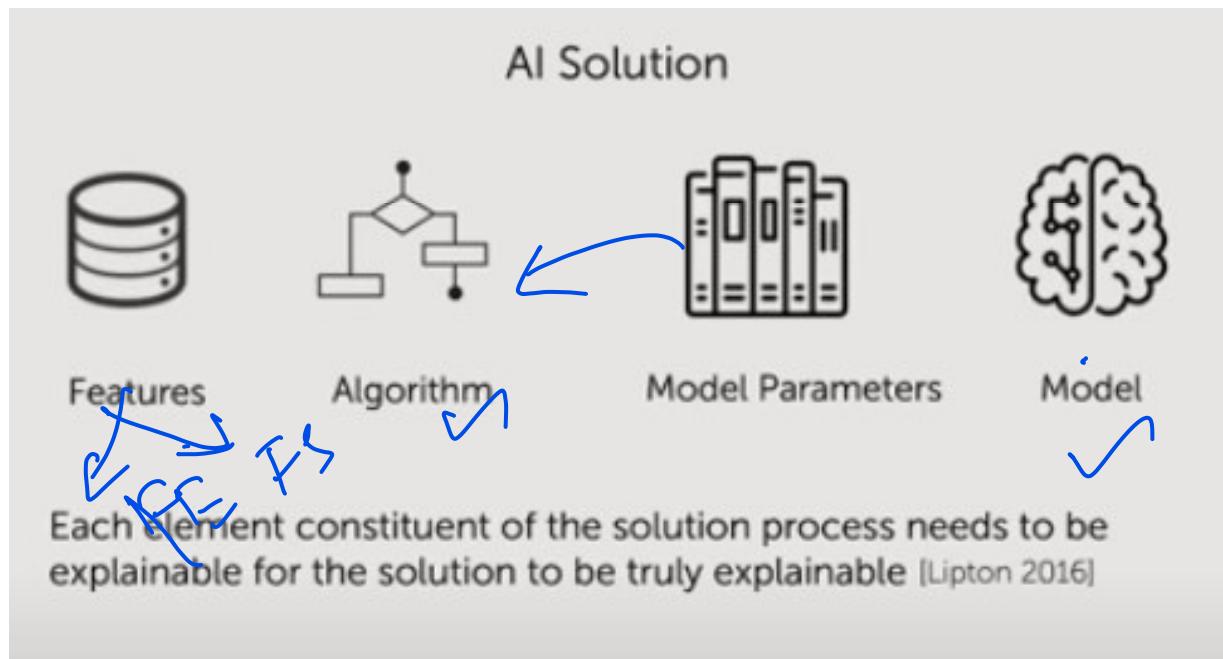
Decisions in Healthcare

- Heuristics
- Rule based systems → Expert system
- ML based systems – learning algorithms
 - Machine learning is concerned with using the right features to build the right models using algorithms that achieve the right task

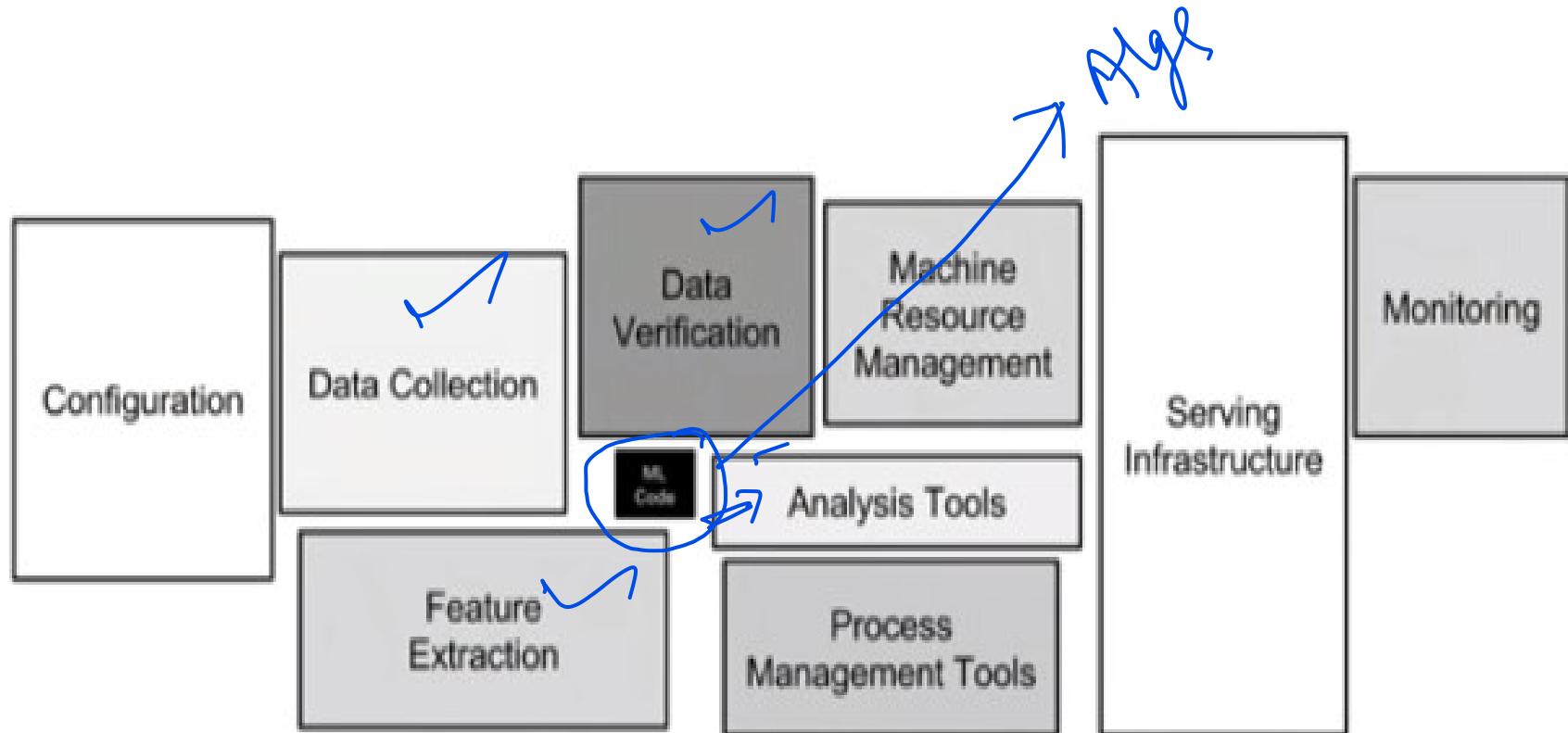
Age , history
symptoms
mycin

feature selection

AI Solution



Implementation of AI in Healthcare

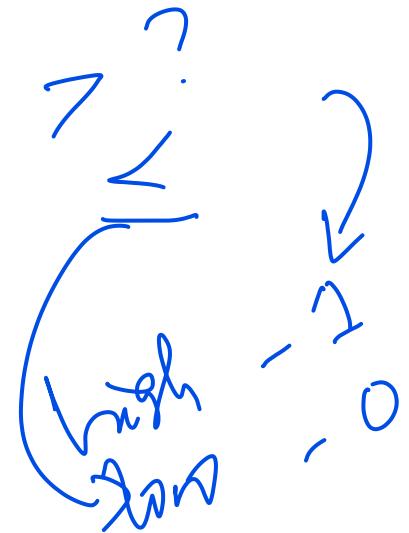


Only a small fraction of real-world machine learning systems actually constitutes machine learning code [Sculley 2015].

Complexity of Data

Data
Verification

- Syntactic correctness
 - Male, Female encoded as 1 and 0
- Morphological correctness
 - Blood pressure value is 500?
- Semantic correctness
 - Range of blood pressure value encoded as high or low, depends on Child or Adult(Age)



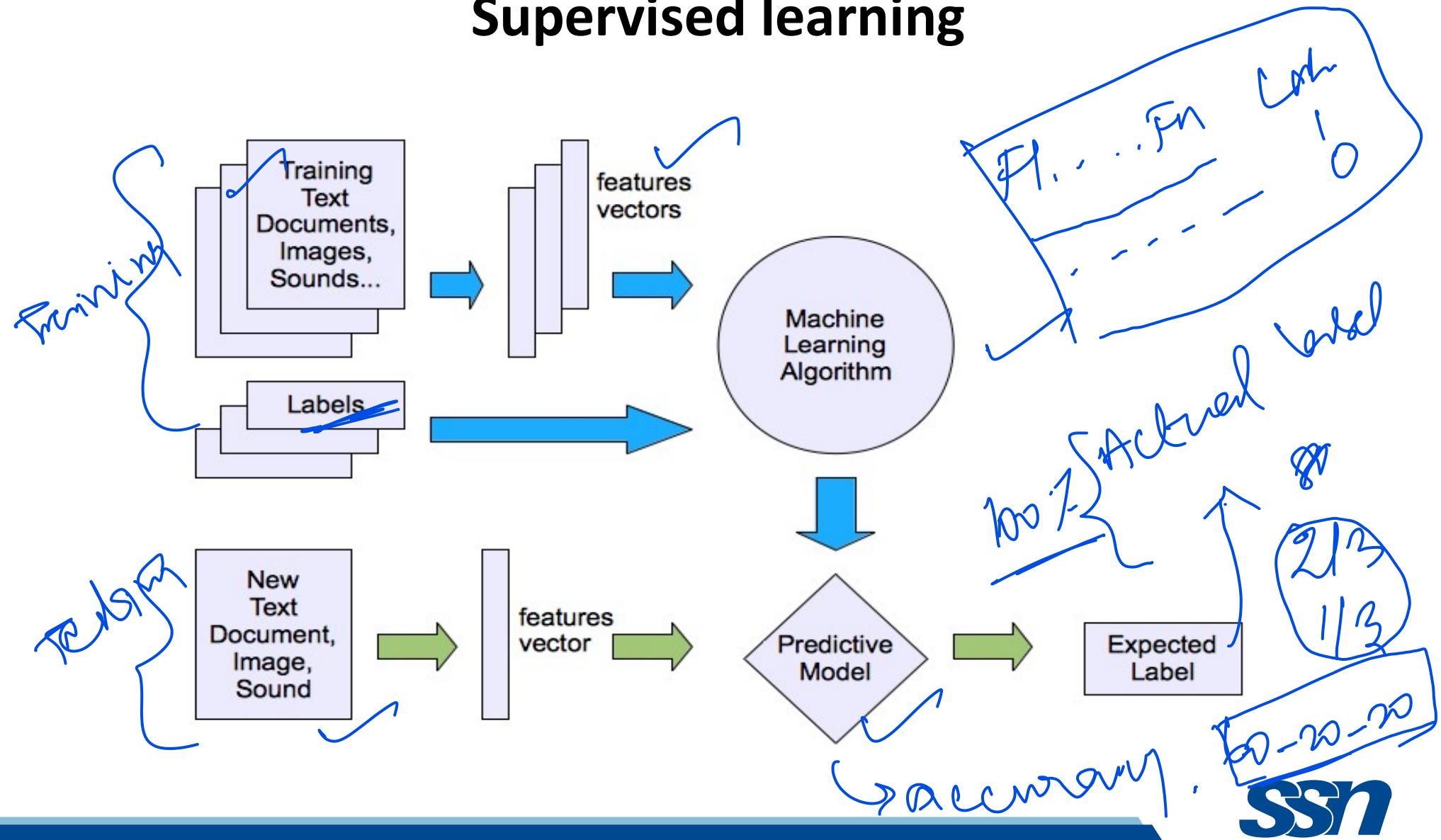
Learning Algorithms

- The success of machine learning system also depends on the algorithms.
- The algorithms control the search to find and build the knowledge structures.
- The learning algorithms should extract useful information from training examples.

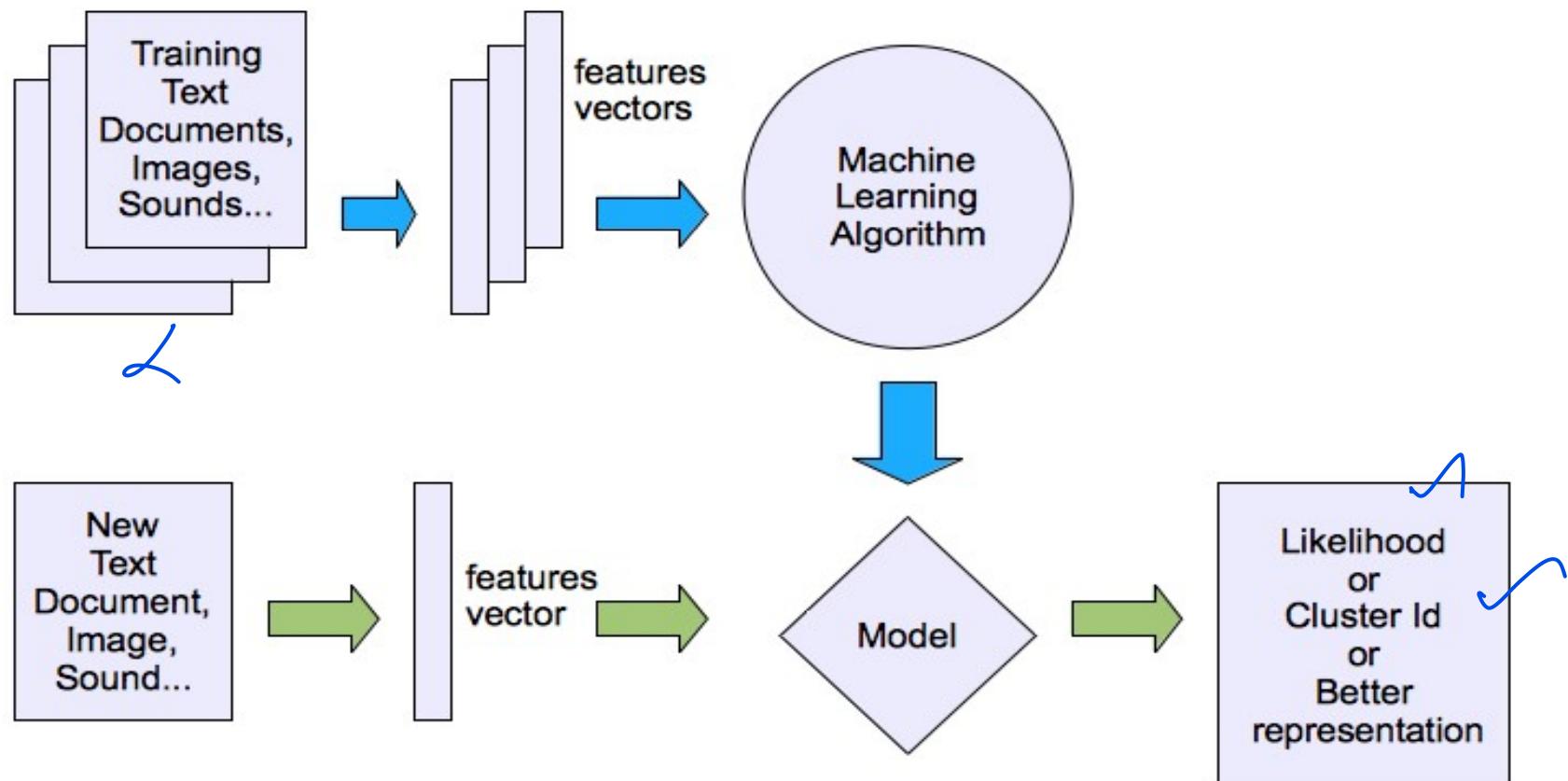
Classification of Algorithms

- **Supervised learning – task driven**
 - Prediction $\{x_n \in R^d, y_n \in R\}_{n=1}^N$
 - Classification (discrete labels), Regression (real values)
 - **Unsupervised learning – data driven**
 - Clustering $\{x_n \in R^d\}_{n=1}^N$
 - Probability distribution estimation
 - Finding association (in features)
 - Dimension reduction \rightarrow
 - **Semi-supervised learning**
 - **Reinforcement learning – learning from mistakes**
 - Decision making (robot, chess machine)
-
- Handwritten notes and diagrams include:
- A large curly brace grouping "Supervised learning – task driven" and "Unsupervised learning – data driven".
 - A bracket under "Semi-supervised learning" pointing to "GED" (General Equilibrium Dynamics) and "Birink".
 - A question mark next to "Subset > labeled".
 - Handwritten text "Evolutionary robotics" and "Learning" near the bottom right.
 - A logo in the bottom right corner consisting of the letters "ssn" in a stylized font.

Supervised learning



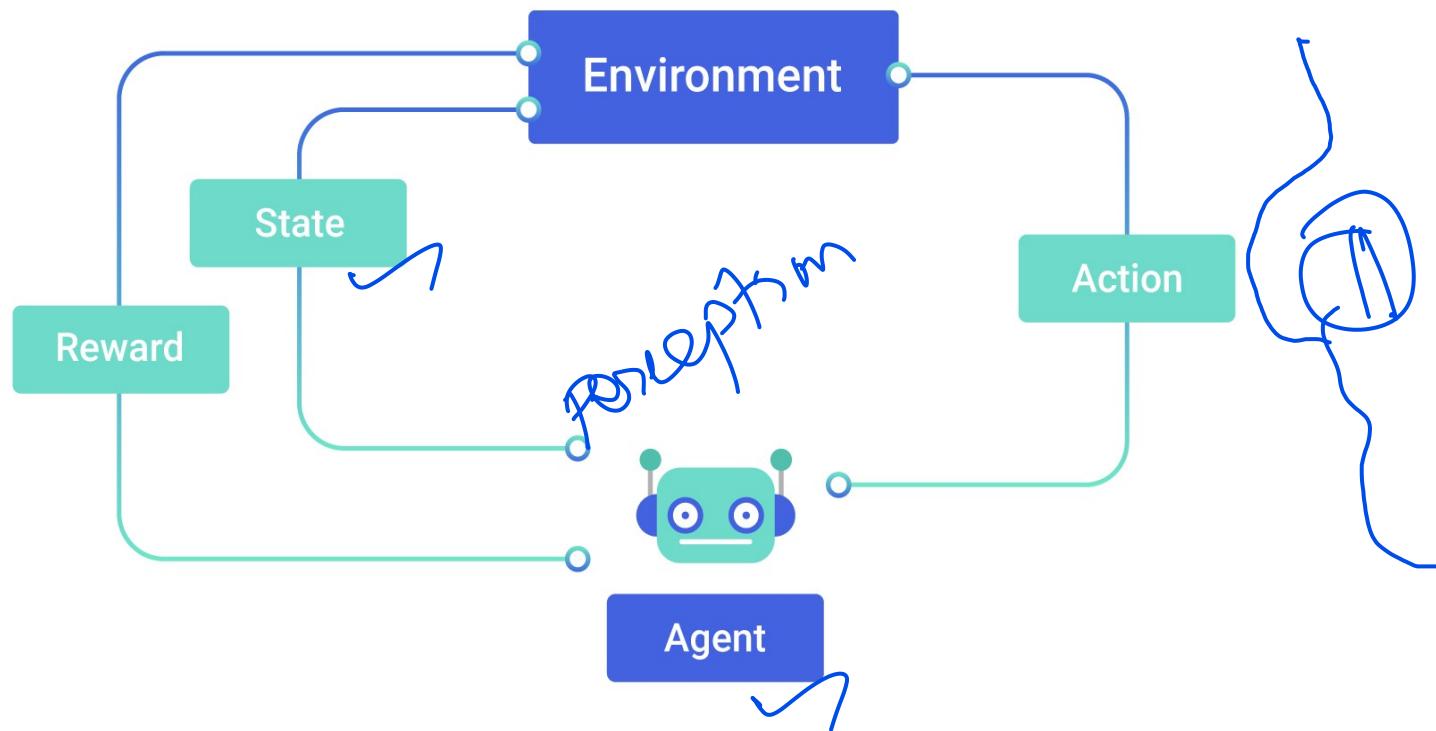
Unsupervised learning



Semi-supervised learning

- Training data includes a few desired outputs.
- Data is cheap, but labelled data is expensive.
- Use small labelled training set to build an initial model, which is then refined using the unlabelled data

Reinforcement learning



What are we seeking?

- Supervised: Low E-out or maximize probabilistic terms

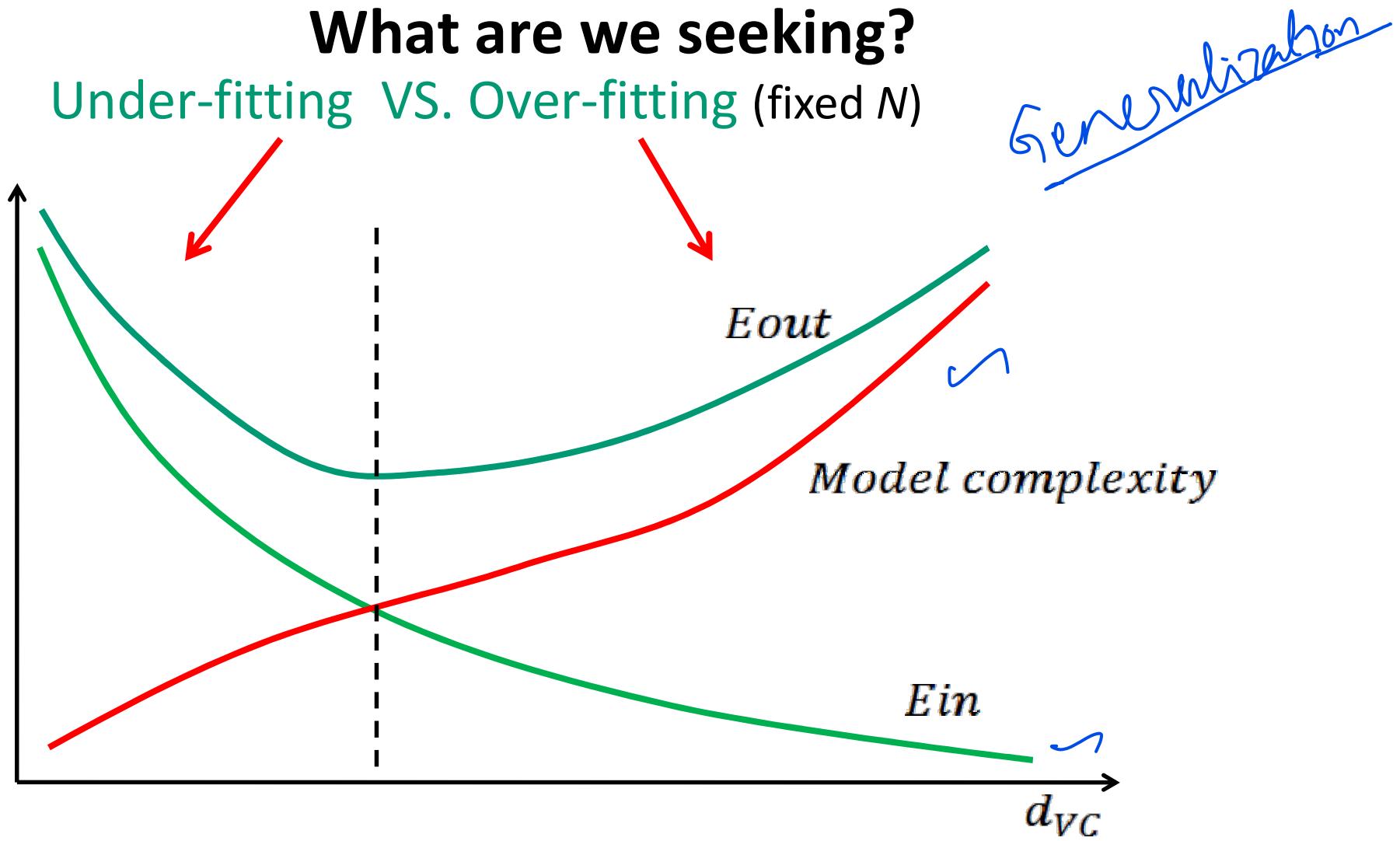
$$\text{error} = \frac{1}{N} \sum_{n=1}^N [y_n \neq g(x_n)]$$

E-in: for training set
E-out: for testing set
E-out(g) <= E-in(g)

- Unsupervised: Minimum quantization error, Minimum distance, MAP, MLE(maximum likelihood estimation)

What are we seeking?

Under-fitting VS. Over-fitting (fixed N)



Learning techniques

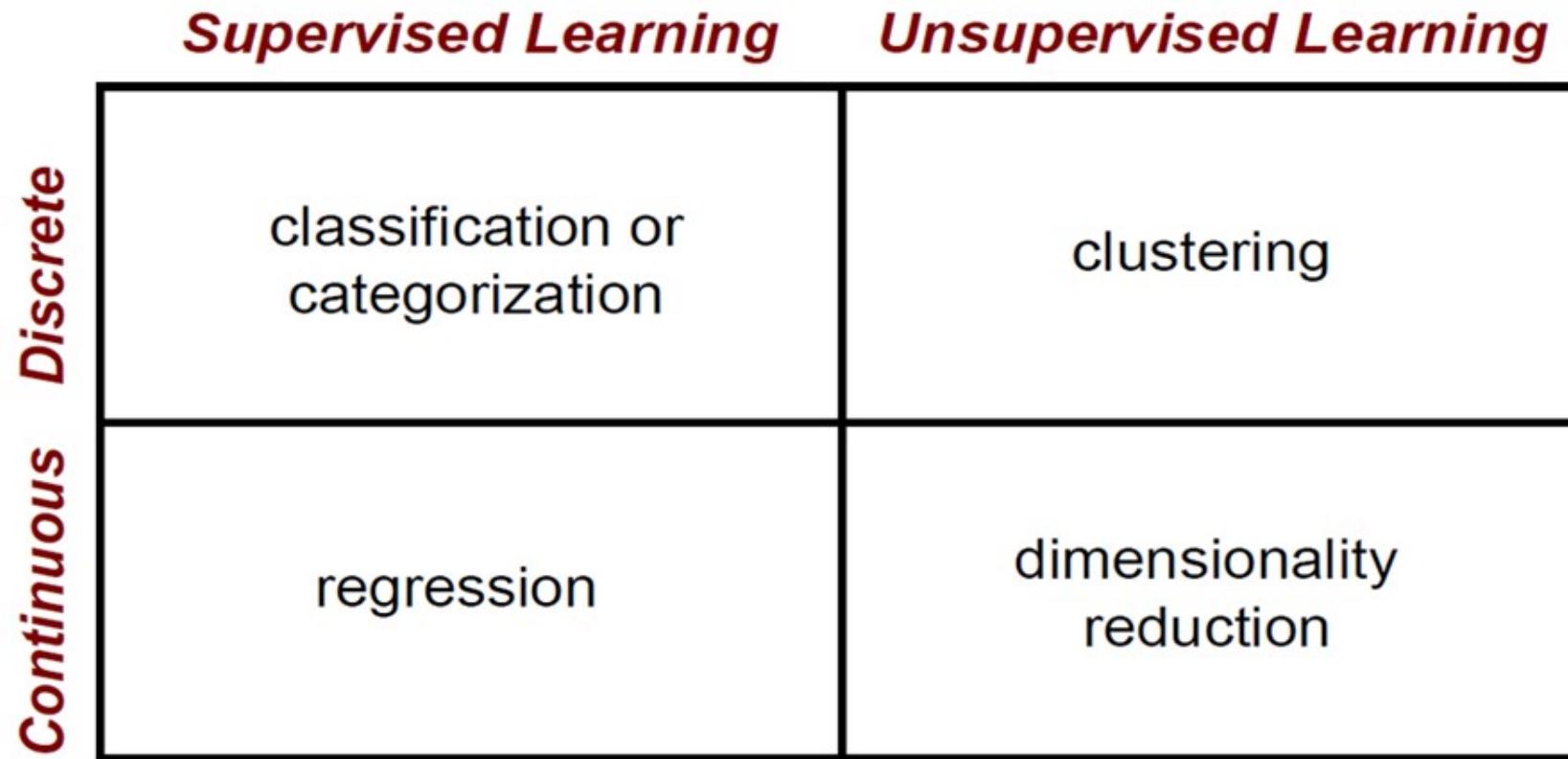
- Supervised learning categories and techniques
 - Linear classifier (numerical functions)
 - Perceptron, logistic regression, support vector machines (SVM), neural networks
 - Parametric (Probabilistic functions)
 - Naïve Bayes, Gaussian discriminant analysis (GDA), Hidden Markov models (HMM), Probabilistic graphical models
 - Non-parametric (Instance-based functions)
 - K-nearest neighbors, Kernel regression, Kernel density estimation, Local regression
 - Non-metric (Symbolic functions) ↗ *welwl*
 - Classification and regression tree (CART), decision tree
 - Aggregation → *hybrid model*
 - Bagging (bootstrap + aggregation), Adaboost, Random forest

Learning techniques

- Unsupervised learning categories and techniques
 - Clustering ✓
 - K-means clustering
 - Spectral clustering
 - Density estimation ✓
 - Gaussian mixture model (GMM)
 - Graphical models
 - Dimensionality reduction ✓
 - Principal component analysis (PCA)
 - Factor analysis

ICA

ML Problems

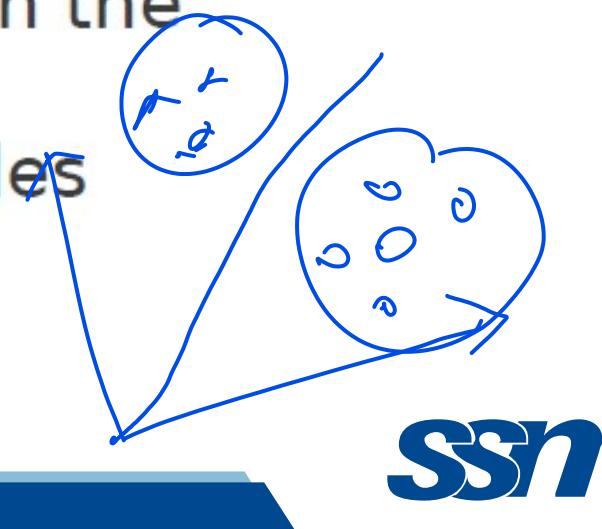


Linear vs Non-linear techniques

- Linear classification techniques
 - Perceptron
 - Logistic regression
 - Linear SVM
 - Naïve Bayes
- Non-linear classification techniques
 - k -nearest neighbors
 - Non-linear SVM
 - Neural networks (MLP)
 - Decision trees
 - Random forest

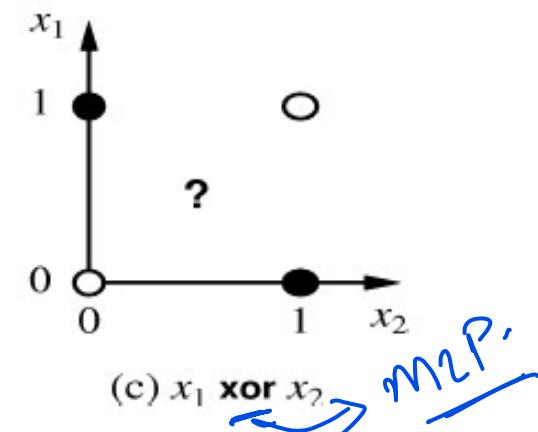
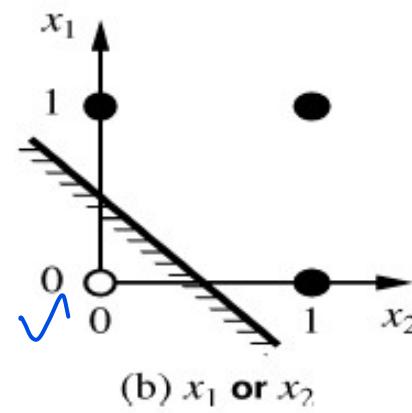
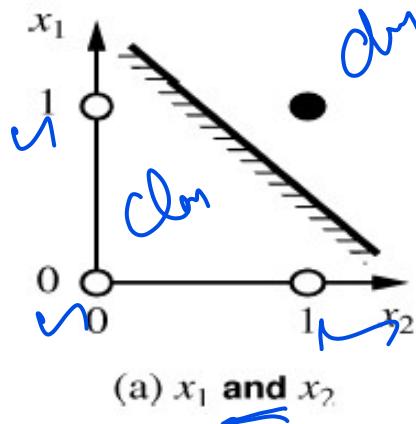
Linearly Separable

- Any pattern classification problem that can be solved by finding a linear decision boundary is called as **linearly separable**
- We need an algorithm to learn the weights, and thus a decision boundary, from given examples



Example – Logic gates

- Thus, single layer network represents a **linear separator** (line, plane, hyperplane) in the input space
- $\sum_j W_j x_j = 0$ or $\mathbf{W} \cdot \mathbf{x} = 0$
- This is not adequate for many pattern recognition problems!



Neural Networks

- Hebb's rule
- ✓ • McCulloch-Pitts neuron model
- Adaline → error
- Perceptron learning
- MLP- Backpropagation algorithm

ANN.

outpt

Σ

updating the
weight.

out Separation
problem

Adaline.

McCulloch-Pitts neuron model

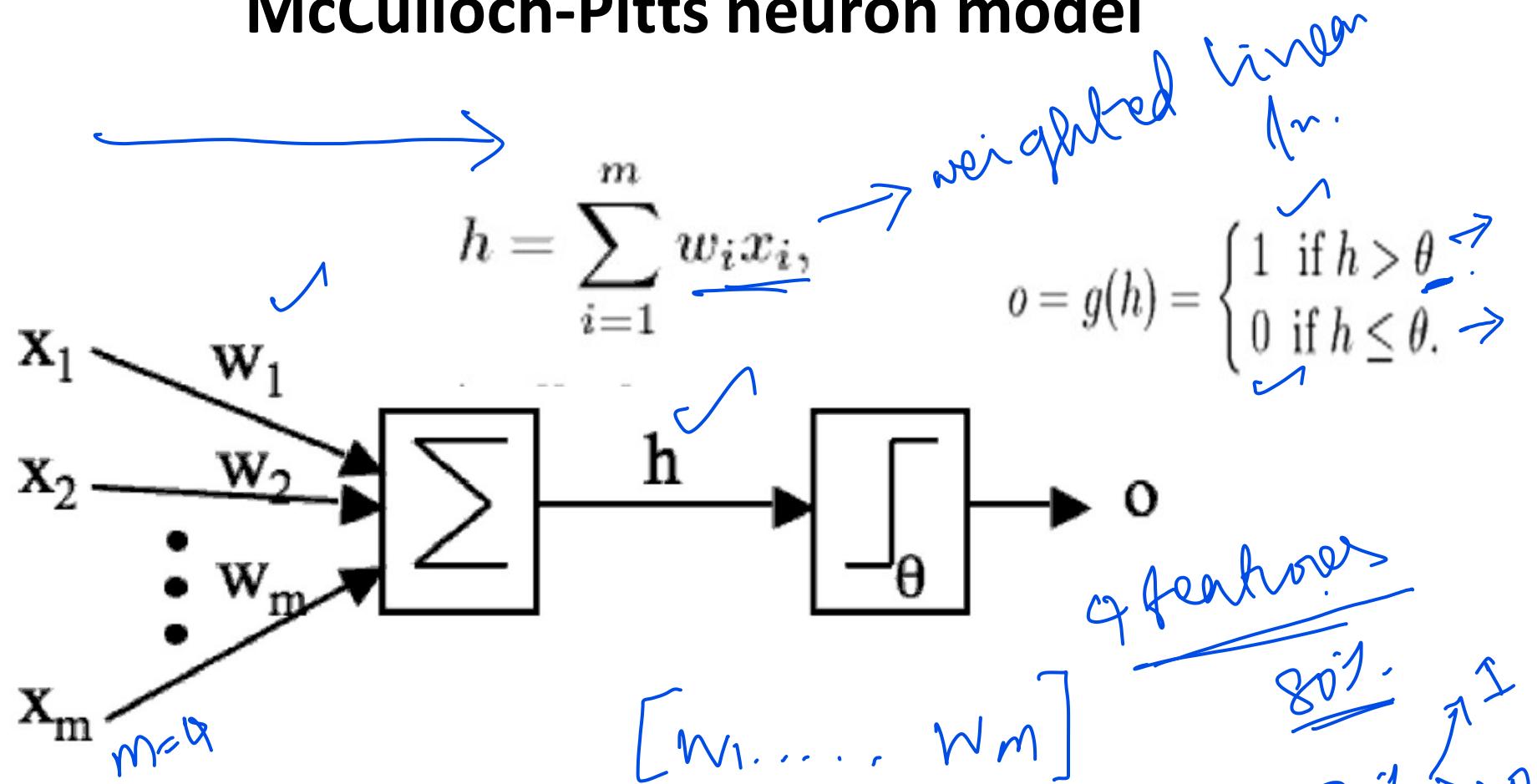


FIGURE 1.6: A picture of McCulloch and Pitt's mathematical model of a neuron. The inputs x_i are multiplied by the weights w_i , and the neurons sum their values. If this sum is greater than the threshold θ then the neuron fires, otherwise it does not.

Perceptron – learning algorithm

- The value of the learning rate decides how fast the network learns, typically $0.1 < \eta < 0.4$
- learning rate is a crucial parameter for convergence
- Weight updation: $w_{ij} \leftarrow w_{ij} + \eta(t_j - y_j) \cdot x_i$
- Bias: fixed (usually the value of $-\pm$ is chosen)



Perceptron – learning algorithm

- Initialisation



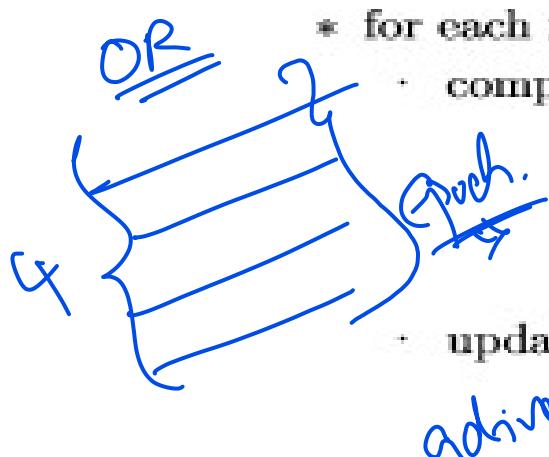
- set all of the weights w_{ij} to small (positive and negative) random numbers

- Training

- for T iterations or until all the outputs are correct:

* for each input vector:

- compute the activation of each neuron j using activation function g :



$$y_j = g \left(\sum_{i=0}^m w_{ij} x_i \right) = \begin{cases} 1 & \text{if } \sum_{i=0}^m w_{ij} x_i > 0 \\ 0 & \text{if } \sum_{i=0}^m w_{ij} x_i \leq 0 \end{cases} \quad (3.4)$$

- update each of the weights individually using:

$$w_{ij} \leftarrow w_{ij} - \eta (y_j - t_j) \cdot x_i \quad (3.5)$$

- Recall \Rightarrow Test \rightarrow m

10×4

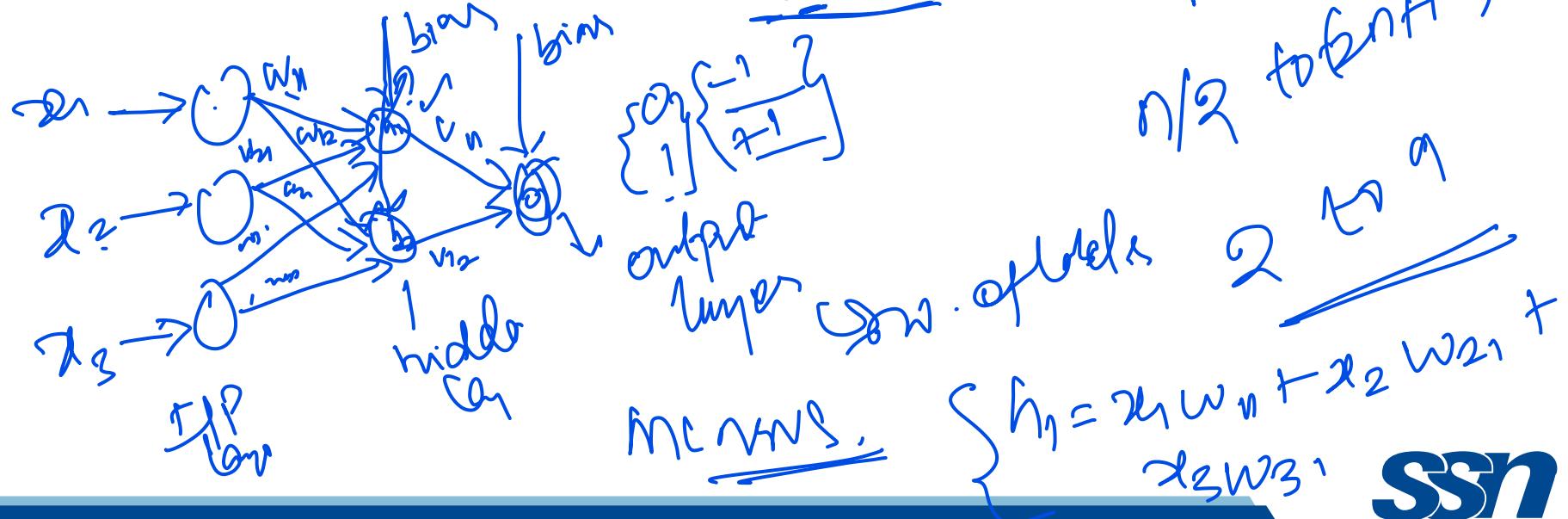
- compute the activation of each neuron j using:

$$\left\{ y_j = g \left(\sum_{i=0}^m w_{ij} x_i \right) = \begin{cases} 1 & \text{if } w_{ij} x_i > 0 \\ 0 & \text{if } w_{ij} x_i \leq 0 \end{cases} \right. \quad (3.6)$$

? - - -

Complexity

- The recall phase loops over the neurons, and within that loops over the inputs, so its complexity is $O(mn)$. ✓
- The training part does this same thing, but does it for T iterations, so costs $O(Tmn)$. ✓

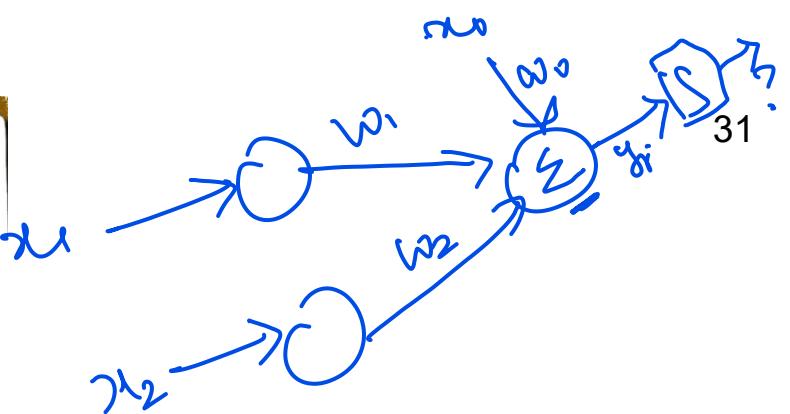


Epoch

Date: OR

x_1	x_2	Target	$(x_0) \text{ bias} = 0$
0	0	0	$w_0 - w_1 - w_2 = 0$
0	1	1	
1	0	1	
1	1	1	

Theta - 1
 $\eta = 1$
 $y_j = g(y_j) = \begin{cases} 1 & y_j \geq 1 \\ 0 & y_j < 1 \end{cases}$
 $\text{Actual } y_j = w_0 \times \text{bias} + x_1 \cdot w_1 + x_2 \cdot w_2$
 $w_i = w_i + \eta (t_j - y_j) \cdot x_i$



① $x_1 = 0 \quad x_2 = 0 \quad y_1 = 0 < 1 \quad w_0 = 0 + 1(0-0) \times 0 = 0$
 $\theta = 0 \quad w_1 = 0 \quad w_2 = 0$
 $x_1 = 0 \quad x_2 = 1 \quad y_2 = 0 < 1 \quad w_0 = 0 + 1(1-0) \times 0 = 0$
 $\theta = 0 \quad w_1 = 0 + 1(1-0) \times 0 = 0$
 $w_2 = 0 + 1(1-0) \times 1 = 1$

$\checkmark \quad x_1 = 1 \quad x_2 = 0 \quad y_3 = 0 < 1 \quad \therefore \text{output} = 0 \quad [0 \ 0 \ 1]$

$$w_0 = 0 + 1(1-0) \times 0 = 0$$

$$w_1 = 0 + 1(1-0) \times 1 = 1$$

$$w_2 = 0 + 1(1-0) \times 0 = 0 \quad [0 \ 1 \ 0]$$

$x_1 = 1 \quad x_2 = 1 \quad y_4 = 0 + 1 \times 1 + 1 \times 1 = 2 \geq 1$
 $\therefore \text{output} = 1$

$$w_0 = 0 + 1(1-1) \times 0 = 0$$

$$w_1 = 1 + 1(1-1) \times 1 = 1$$

$$w_2 = 1 + 1(1-1) \times 1 = 1 \quad [0 \ 1 \ 1]$$

Actual?
 $0 \quad 0 \quad 0 \quad 1$
 Target
 $0 \quad 0 \quad 0 \quad 1$

Epoch (II)

$$\begin{array}{l} \xrightarrow{\text{Input}} \\ \underline{x_1 = 0, x_2 = 0} \quad y_1 = 0 < 1 \quad \underline{\text{output} = 0} \end{array}$$

$$w_0 = 0 + 1(0-0) * 0 = 0$$

$$w_1 = 1 + 1(0-0) * 0 = 1$$

$$w_2 = 1 + 1(0-0) * 0 = 1$$

[0 1 1]

$$\begin{array}{l} \xrightarrow{\text{Input}} \\ \underline{x_1 = 0, x_2 = 1} \quad y_2 = 0 + 0 + 1 = 1 \quad \therefore \underline{\text{output} = 1} \end{array}$$

$$w_0 = 0$$

$$w_1 = 1 + 1(1-1) * 0 = 1$$

$$w_2 = 1 + 1(1-1) * 1 = 1$$

[0 1 1]

$$\begin{array}{l} \xrightarrow{\text{Input}} \\ \underline{x_1 = 1, x_2 = 0} \quad y_3 = 0 + 1 + 0 = 1 \quad \therefore \underline{\text{output} = 1} \end{array}$$

$$w_0 = 0 + 1(1-1) * 0 = 0$$

$$w_1 = 1 + 1(1-1) * 1 = 1$$

$$w_2 = 1 + 1(1-1) * 0 = 1$$

[0 1 1]

$$\begin{array}{l} \xrightarrow{\text{Input}} \\ \underline{x_1 = 1, x_2 = 1} \quad y_4 = 0 + 1 + 1 = 2 \quad \therefore \underline{\text{output} = 1} \end{array}$$

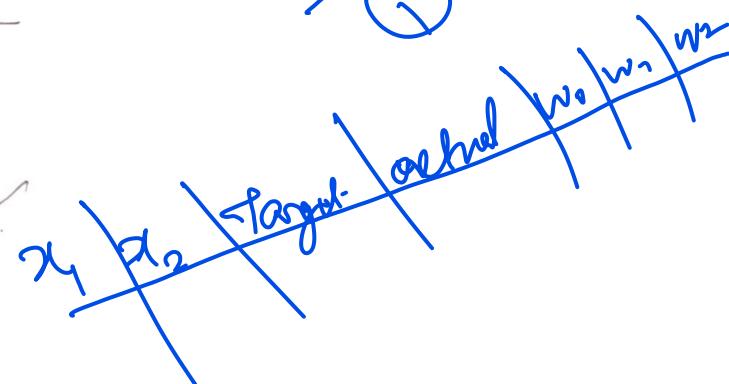
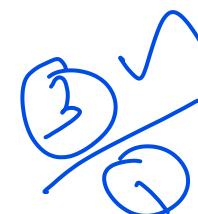
$$w_0 = 0 + 1(1-1) * 0 = 0$$

$$w_1 = 1 + 1(1-1) * 1 = 1$$

$$w_2 = 1 + 1(1-1) * 1 = 1$$

[0 1 1]

Epoch:



Pros and Cons of NN

- Neural networks are good to model with nonlinear data with large number of inputs.
- Once trained, the predictions are pretty fast.
- Neural networks are flexible for both regression and classification problems.
- Any data which can be made numeric can be used in the model, as neural network is a mathematical model with approximation functions.

Pros and Cons of NN

- It is computationally very expensive and time consuming to train with traditional CPUs.
- Neural networks are black boxes, that is how each independent variable is influencing the dependent variables.
- Neural networks depend a lot on training data. This leads to the problem of over-fitting and generalization.

Decision tree

- A feature tree is a tree such that each internal node is labeled with a feature and each edge emanating from an internal node is labeled with a literal.
- The set of literals at a node is called a split.
- Each leaf of the tree represents a logical expression, which is the conjunction of literals on the path from root to the leaf.
- A leaf node represents a class label $h(x)$ or class label distribution

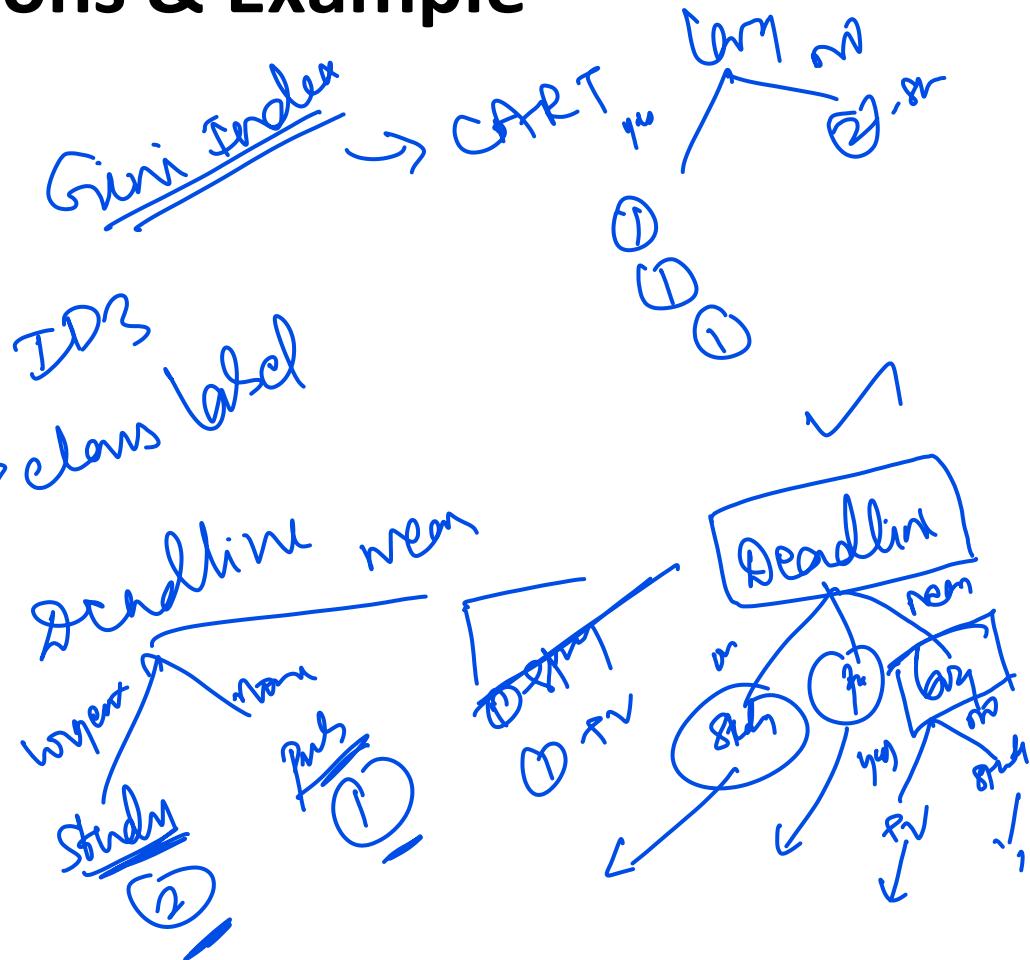
Impurity functions & Example

$$\checkmark \text{Entropy}(S) = \sum_{i=1}^c -p_i \log_2 p_i$$

$$\checkmark \text{Gain}(S, A) = \text{Entropy}(S) - \sum_{v \in \text{Values}(A)} \frac{|S_v|}{|S|} \text{Entropy}(S_v)$$

F1 *F2*

	Deadline	Lazy	Activity type
1	Urgent	Yes	Study
2	None	Yes	Pub
3	Near	No	Study
4	Near	Yes	TV
5	Urgent	No	Study



$$-\rho_1 \log_2 \rho_1$$

Date:

$$\text{Entropy } (S_1) = \frac{-3}{5} \log_2 \frac{3}{5} - \frac{1}{5} \log_2 \frac{1}{5} - \frac{1}{5} \log_2 \frac{1}{5}$$

$$= 1.371.$$

$$\text{Gaining } (S_1, \text{ sendline}) = \underline{\text{Entropy}}(S_1) - \frac{|S_{\text{urgent}}|}{|S_1|}$$

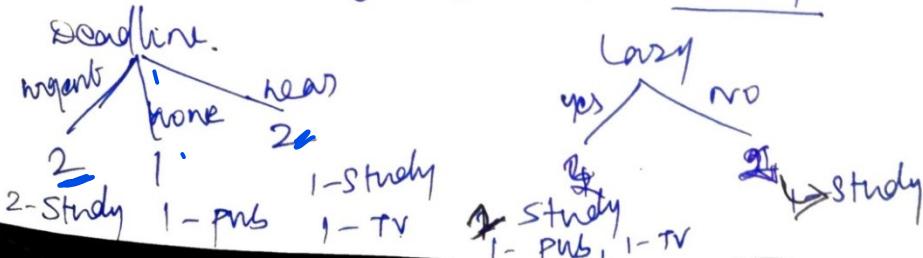
$$\text{Entropy}(S_{\text{irr}}) = \frac{|S_{\text{inone}}|}{|S_1|} \text{Entropy}(S_{\text{inone}})$$

$$-\frac{|S_{\text{near}}|}{|S_1|} \text{ Entropy } (S_{\text{near}})$$

$$= \cancel{1.87} - \frac{2}{5} \left(-\frac{2}{2} \log_2 \frac{2}{2} \right) - \frac{1}{5} \left(\cancel{1} \log_2 \frac{1}{1} \right)$$

$$= \frac{2}{5} \left(\frac{-1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} \right)$$

$$= 1 \cdot 371 - 0 - 0 - 0.4 = \underline{\underline{0.971}} \quad \checkmark$$



Date:

Gain (S_1 , lazy)

$$= \text{Entropy}(S_1) - \frac{|S_{1,\text{yes}}|}{|S_1|} \text{ entropy}(S_1)$$

$$- \frac{|S_{1,\text{no}}|}{|S_1|} \text{ entropy}(S_1)$$

$$= 1.371 - \frac{3}{5} \left(\frac{2}{3} \log_2 \frac{2}{3} - \frac{1}{3} \log_2 \frac{1}{3} \right)$$

$$\cancel{- \frac{1}{4} \log_2 \frac{1}{4}} - \frac{2}{5} \left(\frac{2}{2} \log_2 \frac{2}{2} \right)$$

$$= 1.371 - \left(\frac{3}{5} \right) + 0.933 - 0 = \underline{0.811} \quad \checkmark$$

$\text{ID3}(X, T, \text{Attrs})$

X : training examples; T : target attribute; Attrs : other attributes, initially all attributes

- Create Root node

If all X 's are +, return Root with class +

If all X 's are -, return Root with class -

If Attrs is empty return Root with class most common value of T in X

else

↙
 $A \leftarrow$ best attribute; decision attribute for Root $\leftarrow A$

For each possible value v_i of A :

- add a new branch below Root, for test $A = v_i$

- $X_i \leftarrow$ subset of X with $A = v_i$

- If X_i is empty then add a new leaf with class the most common value of T in X

- else add the subtree generated by $\text{ID3}(X_i, T, \text{Attrs} \leftarrow \{A\})$

return Root

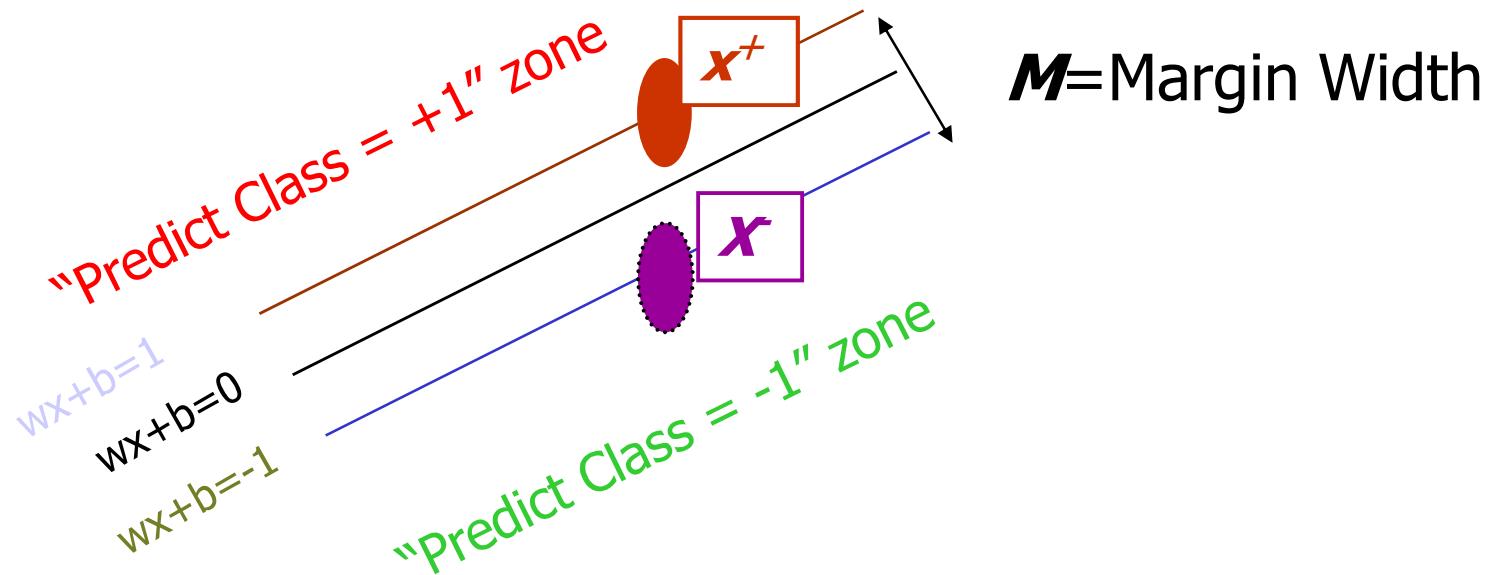
Strengths of decision tree

- Decision trees are able to generate understandable rules.
- Decision trees are able to handle both continuous and categorical variables.
- Decision trees provide a clear indication of which fields are most important for prediction or classification.

Weaknesses of decision tree

- Decision trees are less appropriate for estimation tasks where the goal is to predict the value of a continuous attribute.
- Decision trees are prone to errors in classification problems with many class and relatively small number of training examples.
- Pruning algorithms can also be expensive since many candidate sub-trees must be formed and compared.

Linear SVM Mathematically



What we know:

- $\mathbf{w} \cdot \mathbf{x}^+ + b = +1$
- $\mathbf{w} \cdot \mathbf{x}^- + b = -1$
- $\mathbf{w} \cdot (\mathbf{x}^+ - \mathbf{x}^-) = 2$

$$M = \frac{(\mathbf{x}^+ - \mathbf{x}^-) \cdot \mathbf{w}}{\|\mathbf{w}\|} = \frac{2}{\|\mathbf{w}\|}$$

We can formulate a quadratic Optimization Problem and solve for w and b, for all i

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

Hard Margin v.s. Soft Margin

- **The old formulation:**

Find \mathbf{w} and b such that

$$\Phi(\mathbf{w}) = \frac{1}{2} \mathbf{w}^T \mathbf{w} \text{ is minimized and for all } \{(\mathbf{x}_i, y_i)\}$$

$$y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1$$

- **The new formulation incorporating slack variables:**

Find \mathbf{w} and b such that

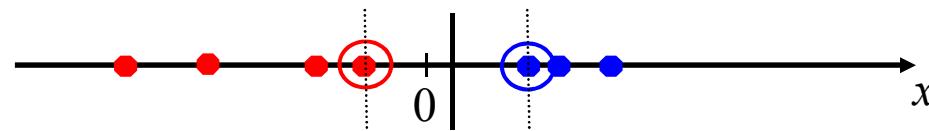
$$\Phi(\mathbf{w}) = \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum \xi_i \text{ is minimized and for all } \{(\mathbf{x}_i, y_i)\}$$

$$y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i \text{ and } \xi_i \geq 0 \text{ for all } i$$

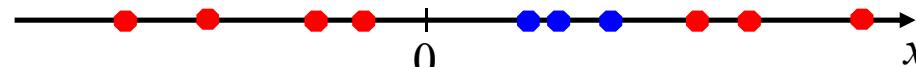
- **Parameter C can be viewed as a way to control overfitting.**

Non-linear SVMs

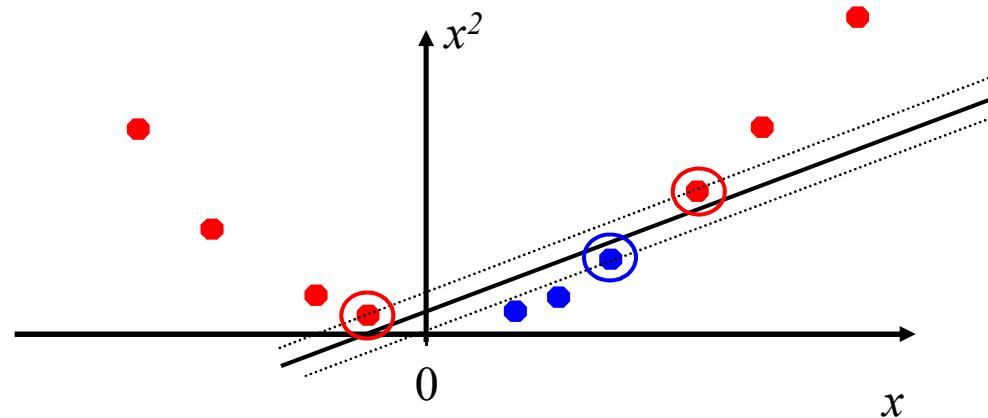
- Datasets that are linearly separable with some noise work out great:



- But what are we going to do if the dataset is just too hard?

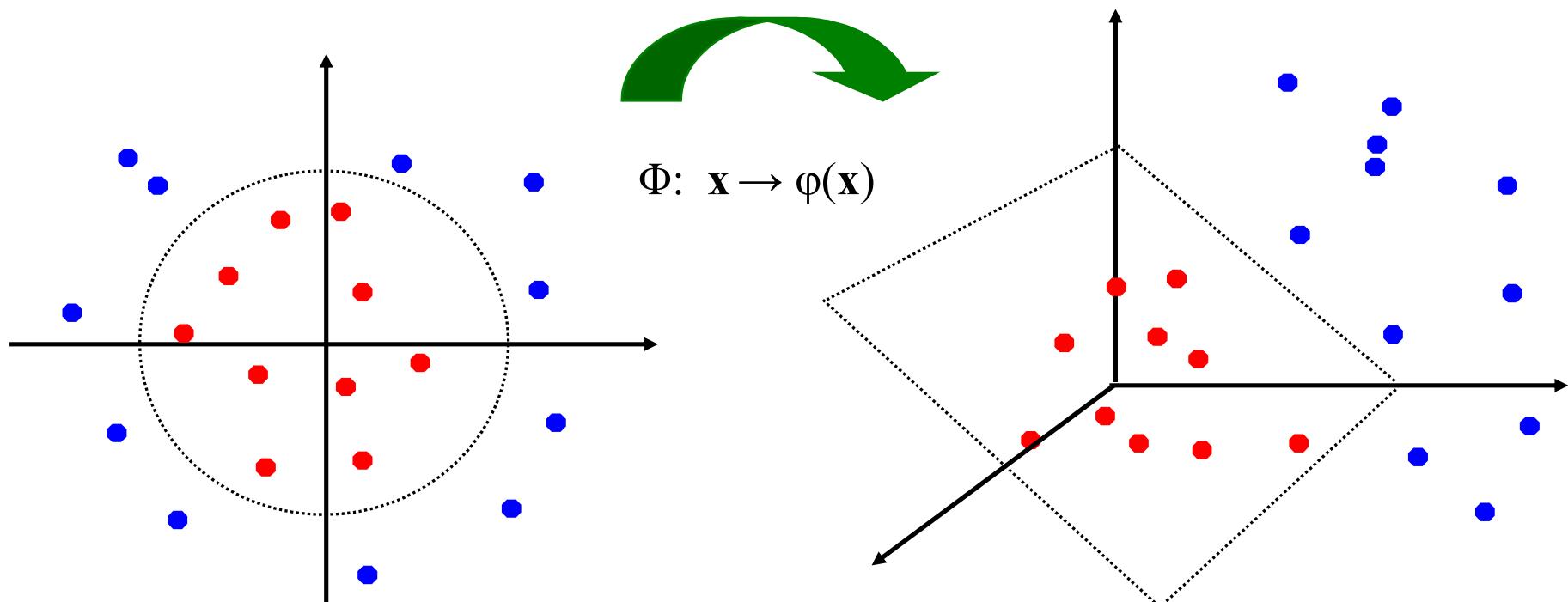


- How about... mapping data to a higher-dimensional space:



Non-linear SVMs: Feature spaces

- General idea: the original input space can always be mapped to some higher-dimensional feature space where the training set is separable:



Strengths and Weaknesses of SVM

Strengths

- Suitable for classification, regression and outliers detection.
- Effective in high dimensional spaces where the number of dimensions is greater than the number of samples.
- Uses support vectors

Weaknesses

- Sensitive to noise : small number of mislabeled examples can dramatically decrease the performance
- Considers two classes : how to do multi-class classification with SVM?
 - 1) with output classes m, learn m SVM's
 - SVM 1 learns "Output==1" vs "Output != 1"
 - SVM 2 learns "Output==2" vs "Output != 2"
 - :
 - SVM m learns "Output==m" vs "Output != m"
 - 2) To predict the output for a new input, just predict with each SVM and find out which one puts the prediction the furthest into the positive region.

Performance metrics – Overall

		True condition			
		Condition positive	Condition negative	Prevalence = $\frac{\sum \text{Condition positive}}{\sum \text{Total population}}$	Accuracy (ACC) = $\frac{\sum \text{True positive} + \sum \text{True negative}}{\sum \text{Total population}}$
Predicted condition	Predicted condition positive	True positive	False positive, Type I error	Positive predictive value (PPV), Precision = $\frac{\sum \text{True positive}}{\sum \text{Predicted condition positive}}$	False discovery rate (FDR) = $\frac{\sum \text{False positive}}{\sum \text{Predicted condition positive}}$
	Predicted condition negative	False negative, Type II error	True negative	False omission rate (FOR) = $\frac{\sum \text{False negative}}{\sum \text{Predicted condition negative}}$	Negative predictive value (NPV) = $\frac{\sum \text{True negative}}{\sum \text{Predicted condition negative}}$
		True positive rate (TPR), Recall, Sensitivity, probability of detection, Power = $\frac{\sum \text{True positive}}{\sum \text{Condition positive}}$	False positive rate (FPR), Fall-out, probability of false alarm = $\frac{\sum \text{False positive}}{\sum \text{Condition negative}}$	Positive likelihood ratio (LR+) = $\frac{\text{TPR}}{\text{FPR}}$	Diagnostic odds ratio (DOR) = $\frac{\text{LR}^+}{\text{LR}^-}$
		False negative rate (FNR), Miss rate = $\frac{\sum \text{False negative}}{\sum \text{Condition positive}}$	Specificity (SPC), Selectivity, True negative rate (TNR) = $\frac{\sum \text{True negative}}{\sum \text{Condition negative}}$	Negative likelihood ratio (LR-) = $\frac{\text{FNR}}{\text{TNR}}$	

https://en.wikipedia.org/wiki/Sensitivity_and_specificity

Performance metrics - Example

		Patients with bowel cancer (as confirmed on endoscopy)			
Fecal occult blood screen test outcome	Condition positive	Condition negative	Prevalence $= (TP + FN) / Total_Population$ $= (20 + 10) / 2030$ $\approx 1.48\%$	Accuracy (ACC) = $(TP + TN) / Total_Population$ $= (20 + 1820) / 2030$ $\approx 90.64\%$	
	Test outcome positive	True positive (TP) = 20 (2030 × 1.48% × 67%)	False positive (FP) = 180 (2030 × (100 – 1.48%) × (100 – 91%))	Positive predictive value (PPV), Precision $= TP / (TP + FP)$ $= 20 / (20 + 180)$ $= 10\%$	False discovery rate (FDR) $= FP / (TP + FP)$ $= 180 / (20 + 180)$ $= 90.0\%$
	Test outcome negative	False negative (FN) = 10 (2030 × 1.48% × (100 – 67%))	True negative (TN) = 1820 (2030 × (100 – 1.48%) × 91%)	False omission rate (FOR) $= FN / (FN + TN)$ $= 10 / (10 + 1820)$ $\approx 0.55\%$	Negative predictive value (NPV) $= TN / (FN + TN)$ $= 1820 / (10 + 1820)$ $\approx 99.45\%$
	TPR, Recall, Sensitivity $= TP / (TP + FN)$ $= 20 / (20 + 10)$ $\approx 66.7\%$	False positive rate (FPR), Fall-out, probability of false alarm $= FP / (FP + TN)$ $= 180 / (180 + 1820)$ $= 9.0\%$	Positive likelihood ratio (LR+) $= \frac{TPR}{FPR}$ $= (20 / 30) / (180 / 2000)$ ≈ 7.41	Diagnostic odds ratio (DOR) $= \frac{LR+}{LR-}$ ≈ 20.2	F_1 score = $2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$ ≈ 0.174
		False negative rate (FNR), Miss rate $= FN / (TP + FN)$ $= 10 / (20 + 10)$ $\approx 33.3\%$	Specificity, Selectivity, True negative rate (TNR) $= TN / (FP + TN)$ $= 1820 / (180 + 1820)$ $= 91\%$	Negative likelihood ratio (LR-) $= \frac{FNR}{TNR}$ $= (10 / 30) / (1820 / 2000)$ ≈ 0.366	

https://en.wikipedia.org/wiki/Sensitivity_and_specificity

Many classifiers to choose from

- Neural networks
- Naïve Bayes classifier
- Decision tree
- SVM
- Linear / Logistic regression
- K-nearest neighbor
- Ensemble approach
 - Boosting
 - Bagging
 - Random Forest

Which is the best one?

Check your understanding

- Dataset without class labels are suitable for learning.
- Learning through reward function is called
- Machine learning is a sub domain of.....
- Learning without feature extraction is called.....
- Examples for probabilistic based learning algorithms areand.....
- SVM is more suitable for Dataset
- Pruning concept is related withalgorithm

Check your understanding

- Name any three impurity functions ?
- Misclassification is measured using Variable
- Name any three kernel functions of SVM?
- Name three important activation functions of NN?
- Number of nodes in input layer of NN is equivalent to
.....
- Widrow-Hoff rule is also called as and
- NN structure is derived from human.....
- The leaf nodes of decision tree are

References

- Machine Learning: An Algorithmic Perspective, by Stephen Marsland, Third Edition, CRC Press
- Machine Learning: The Art and Science of Algorithms that Make Sense of Data Paperback, by Peter Flach, Cambridge University Press



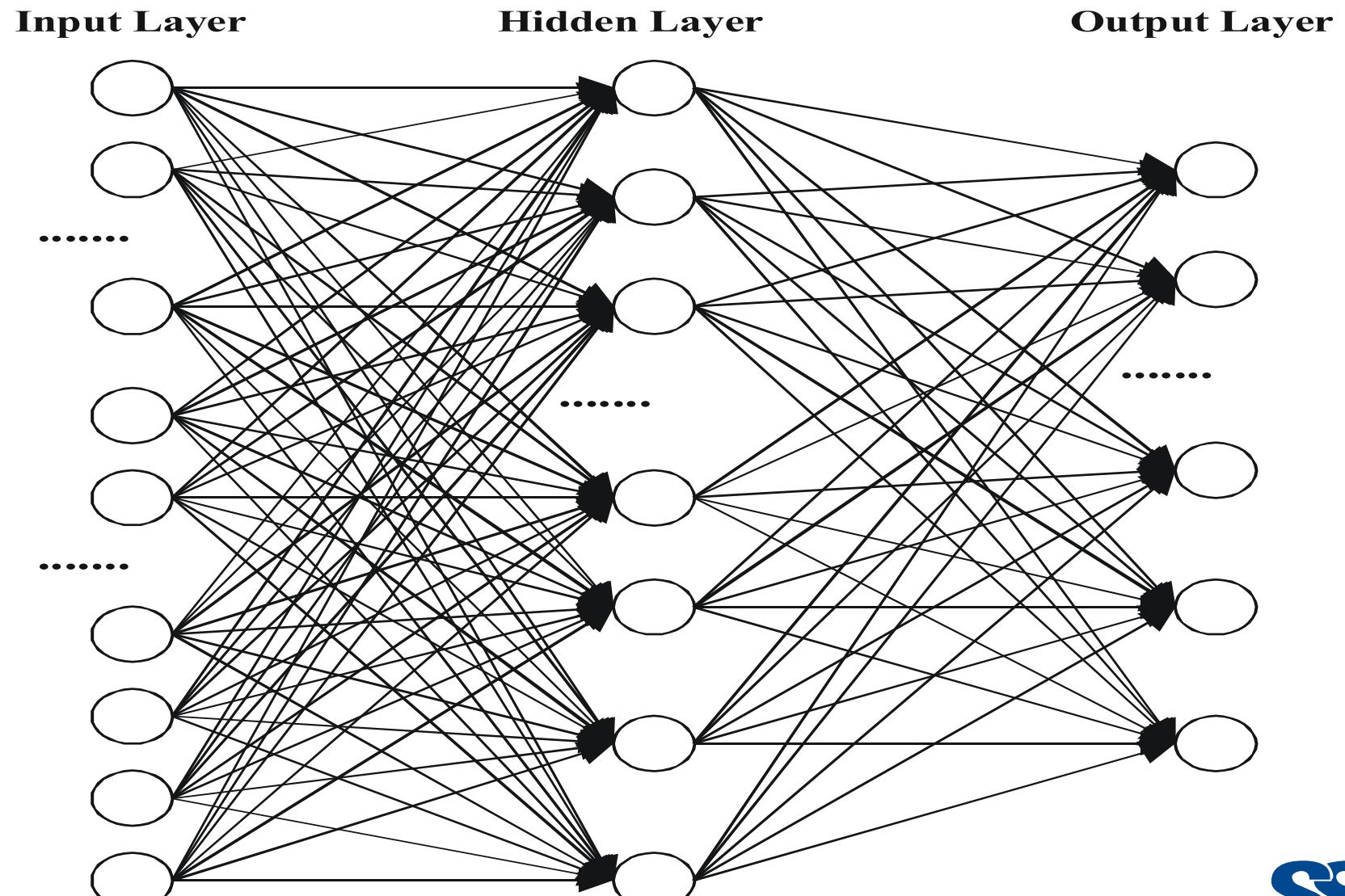
Thank You

Back propagation network

Training a network by back propagation involves
3 stages:

- Feed forward of the input training pattern
- Back propagation of the associated error
- Adjustment of weights

Multi Layer Network



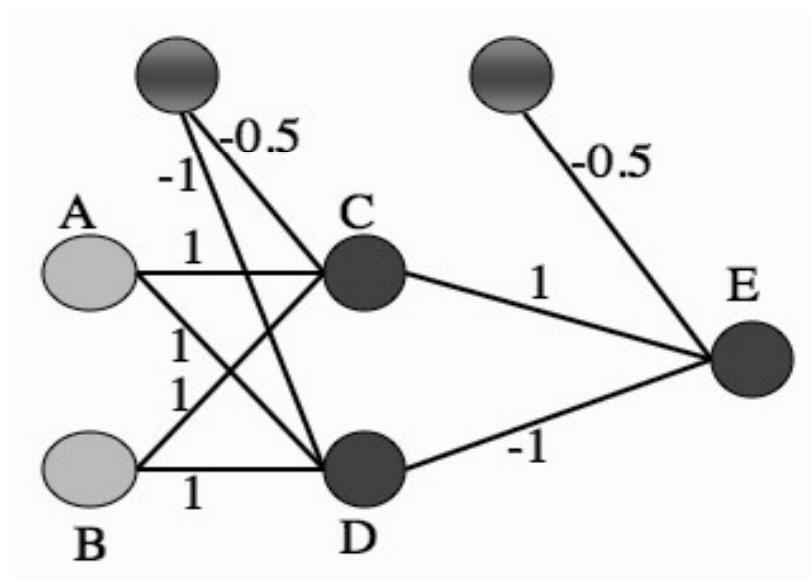
MLP training algorithm using back-propagation

- an input vector is put into the input nodes
- the inputs are fed *forward through the network*
- the inputs and the first-layer weights (here labelled as v) *are used to decide* whether the hidden nodes fire or not. The activation function $g(\cdot)$ *is the sigmoid* function
- the outputs of these neurons and the second-layer weights (labelled as w) *are used* to decide if the output neurons fire or not

Contd...

- the *error is computed as the sum-of-squares difference between the network outputs and the targets*
- this error is fed *backwards through the network in order to* first update the second-layer weights and then afterwards, the first-layer weights

MLP – XOR problem



- Multi-layer Perceptron network showing a set of weights that solve the XOR problem.

Multi-layer Perceptron Algorithm

- **Initialisation**
 - initialise all weights to small (positive and negative) random values
- **Training**
 - repeat:

* for each input vector:

Forwards phase:

- compute the activation of each neuron j in the hidden layer(s) using:

$$h_\zeta = \sum_{i=0}^L x_i v_{i\zeta} \quad (4.4)$$

$$a_\zeta = g(h_\zeta) = \frac{1}{1 + \exp(-\beta h_\zeta)} \quad (4.5)$$

- work through the network until you get to the output layer neurons, which have activations (although see also Section 4.2.3):

$$h_\kappa = \sum_j a_j w_{j\kappa} \quad (4.6)$$

$$y_\kappa = g(h_\kappa) = \frac{1}{1 + \exp(-\beta h_\kappa)} \quad (4.7)$$

Contd...

Backwards phase:

- compute the error at the output using:

$$\delta_o(\kappa) = (y_\kappa - t_\kappa) y_\kappa (1 - y_\kappa) \quad (4.8)$$

- compute the error in the hidden layer(s) using:

$$\delta_h(\zeta) = a_\zeta (1 - a_\zeta) \sum_{k=1}^N w_\zeta \delta_o(k) \quad (4.9)$$

- update the output layer weights using:

$$w_{\zeta\kappa} \leftarrow w_{\zeta\kappa} - \eta \delta_o(\kappa) a_\zeta^{\text{hidden}} \quad (4.10)$$

- update the hidden layer weights using:

$$v_\iota \leftarrow v_\iota - \eta \delta_h(\kappa) x_\iota \quad (4.11)$$

- * (if using sequential updating) randomise the order of the input vectors so that you don't train in exactly the same order each iteration

- until learning stops (see Section 4.3.3)

- Recall

- use the Forwards phase in the training section above

Practice Questions

4. b. Identify the suitable weights for the AND network for the following input and output.

X ₁	X ₂	T
1	1	1
-1	1	-1
1	-1	-1
-1	-1	-1

Where $w_0 = w_1 = w_2 = 0.5$, $\eta = 0.1$, bias = -1

Find the updated weight of each sample for one Epoch.

Practice Questions

- Construct a decision tree for “OR function” and write the rules for two inputs and one output

