

# Unit - 1 : Introduction to Machine Learning.

## Machine Learning:

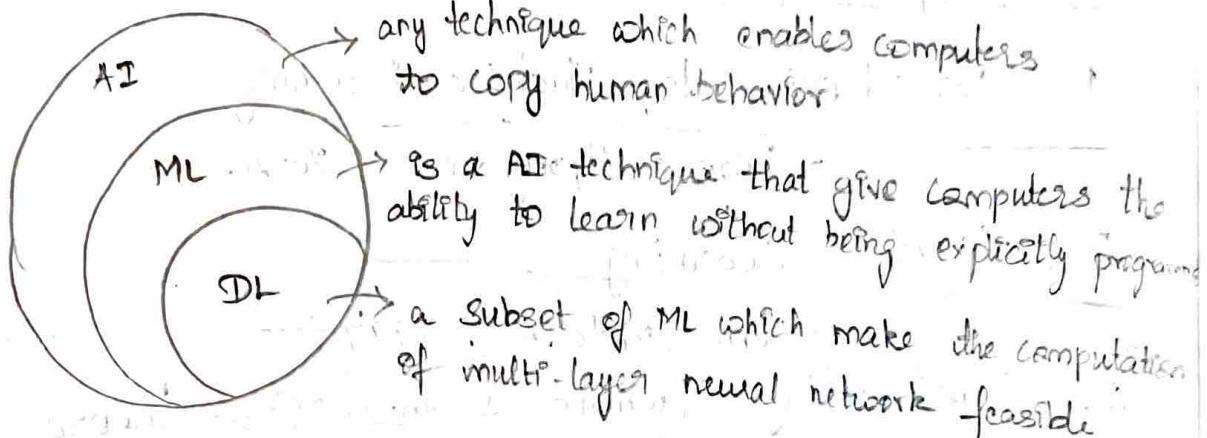
- \* "Field of study that gives computers the ability to learn without being explicitly programmed".

- Arthur Samuel.

## Learning:

- \* "Learning is any process by which a system improves performance from experience".

- Herbert A. Simon.



## What ML can do?

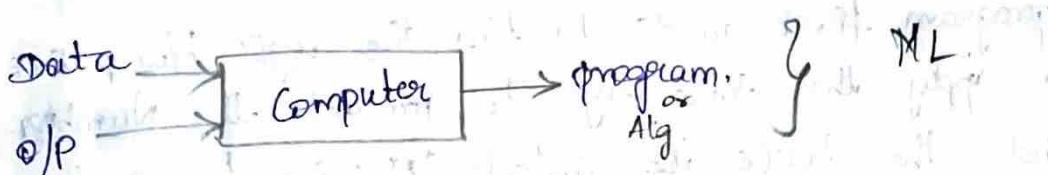
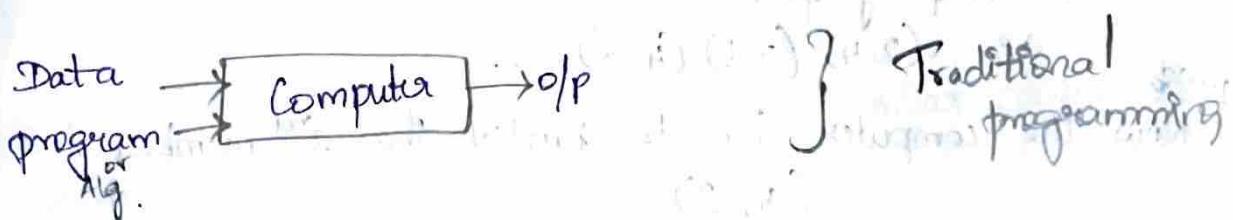
- \* find which category an object belongs to
  - by "Classification Algorithm".
- \* find what is strange
  - by "Anomaly Detection Algorithm".
- \* find how much and how many
  - by "Regression Algorithm".
- \* find how data is arranged.
  - by "Clustering Algorithm".
- \* What should I do next.
  - by "Reinforcement Algorithm".

## ML trends:

- self driving google car
- pattern recognition
- Sentiment analysis
- fraud detection

- online recommendation
- web search results

## Traditional programming Vs ML



## How developer creates a solution?

1. Developing an alg.
2. Implementing an alg. in code.
3. Usage [I/P parameters] → Implemented alg → Result

## How data engineer develops a solution using ML?

1. Data collection & preparation.
2. Experimenting with diff alg. to build a better model.
3. Usage [I/P parameters] → ML-model → Result

## How do you decide which Alg. to use?

- there are dozens of alg. used, so choosing the right alg. is trial and error method.
- even highly experienced data scientist can't tell whether an alg will work without trying it out.
- But alg selection, depends on size & type of data, insights you want to get from data & how those insights will be used.

## Example:

- \* Imagine you have some sets of pair of numbers.
- \* Put only 1 number of pair into machine to predict the other half of pair.  
i.e.,  $(2, 4)$   $(3, 6)$   $(4, 8)$  ...

## Question:

- \* Now the computer has to predict the 2<sup>nd</sup> number for  $(5, ?)$
- " The program first needs to find the logic b/w pairs & then apply the same logic to predict the Number. To find the logic it's called "Machine Learning".

Hence after finding the logic it can apply the same logic to predict each number.

## ML Workflow:

1. Get Data
2. Clean, prepare & Manipulate Data
3. Train Model
4. Test Data
5. Improve

## ML Applications:

- ✓ Recognizing patterns:
  - facial identities or facial expressions.
  - handwritten or spoken words
  - medical images
- ✓ Generating patterns:
  - generating images or motion sequences (demo)
- ✓ Recognizing anomalies:
  - unusual sequence of credit card transactions,
  - unusual sound in car engine.

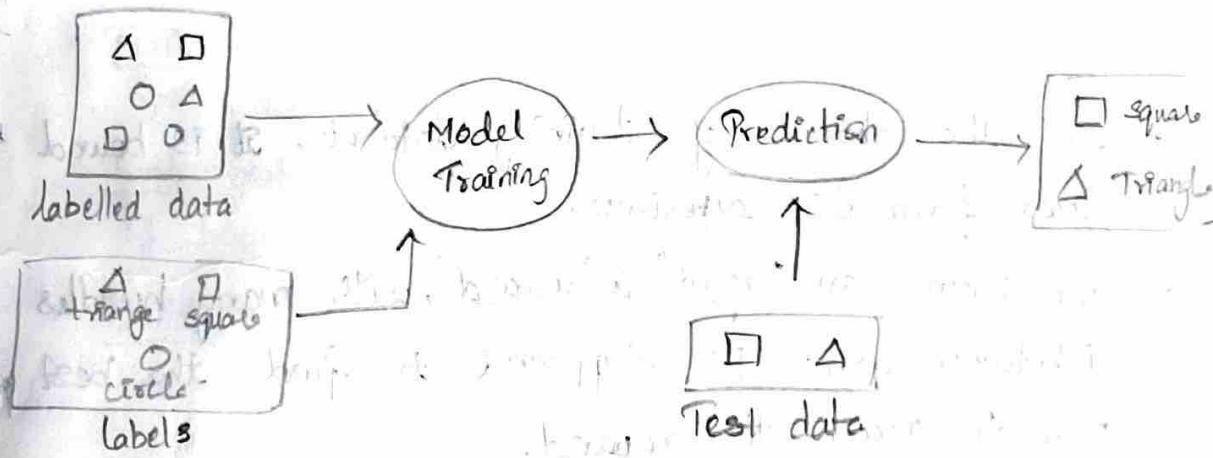
- ✓ Prediction:
  - future Stock prices.
- ✓ spam filtering, fraud detection.
- ✓ Recommendation Systems
- ✓ Information Retrieval.
  - find doc. or Images with similar content.
- ✓ Data Visualization:
  - display a huge DB in revealing way.

### Types of ML:

- Supervised Learning [Task driven i.e., predict next value]
- Unsupervised " [Data driven" i.e., identify clusters]
- Semisupervised "
- Reinforcement " [Learn from mistakes]

### ① Supervised Learning:

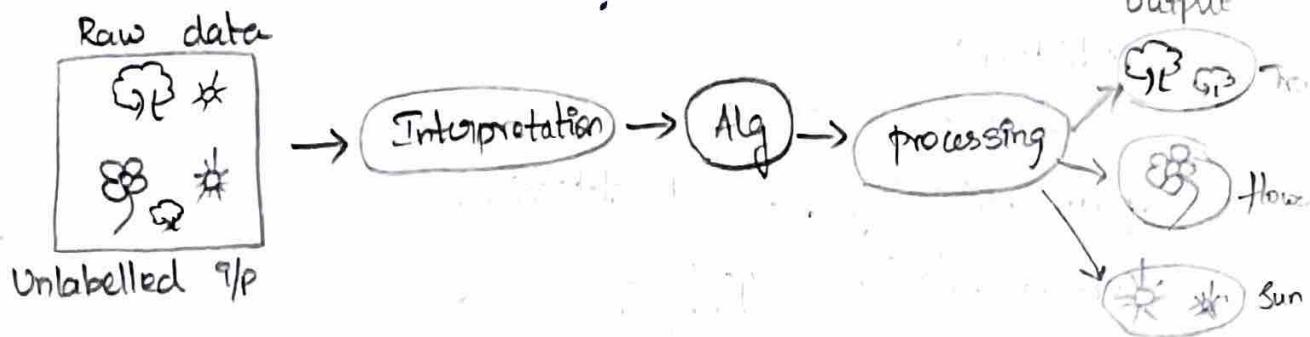
- It's a basic type of ML where ML alg is trained on labelled data.



Eg: Decision Tree, SVM, Naive Bayes, KNN ... etc.

### ② Unsupervised Learning:

- here labels are not known during training.
- since the labels are not given, hidden structure of data can also be founded.



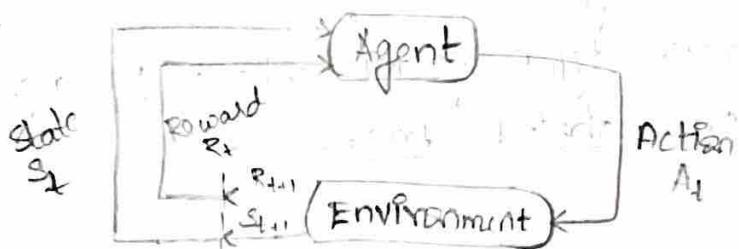
Eg: PCA, clustering, association learning... etc.

### ③ Semi supervised Learning

- it combines both labeled & unlabeled examples to generate an appropriate function or classifier.  
I.e. its a combination of supervised & unsupervised learning.
- goal is to classify some of unlabeled data using the labeled information set.

Eg: Speech analysis, Text document classifier etc...

### ④ Reinforcement Learning



- In the absence of training dataset, it is bound to learn from its experience.
- we have an "agent" & "reward", with many hurdles inbetween. agent is supposed to find the best possible path to reach the reward.
- favourable outputs are "encouraged" or "reinforced", and non-favourable o/p are "discouraged" or "punished".
- Interpreter will decide whether o/p is favourable or not, if its favourable, then it provides reward, otherwise it generates until it finds the better result.

Eg: Game playing, Robotics, Business ... etc.

Types: Positive & Negative reinforcement learning.

## Algorithms:

→ Supervised learning

- Prediction
- classification (discrete labels)
- Regression (real values)

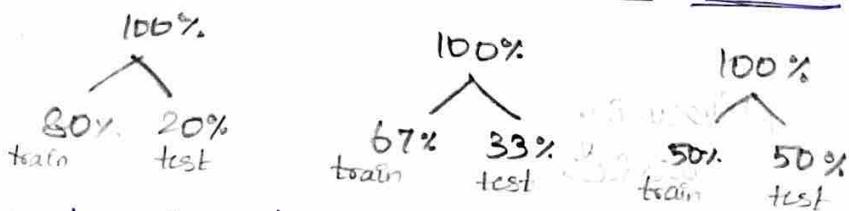
→ Unsupervised learning

- clustering
- probability distribution estimation.
- Finding estimation (in features)
- Dimension reduction.

→ Reinforcement learning.

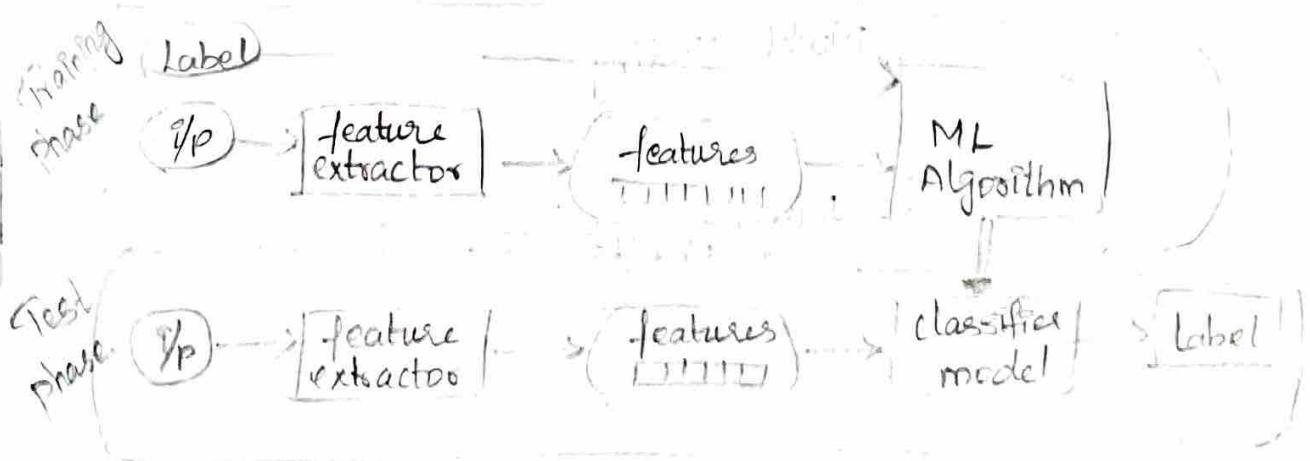
- Decision making (robot, chess machine)

## Training set, Test set and validation set:



Divide the total dataset into 3 subsets:

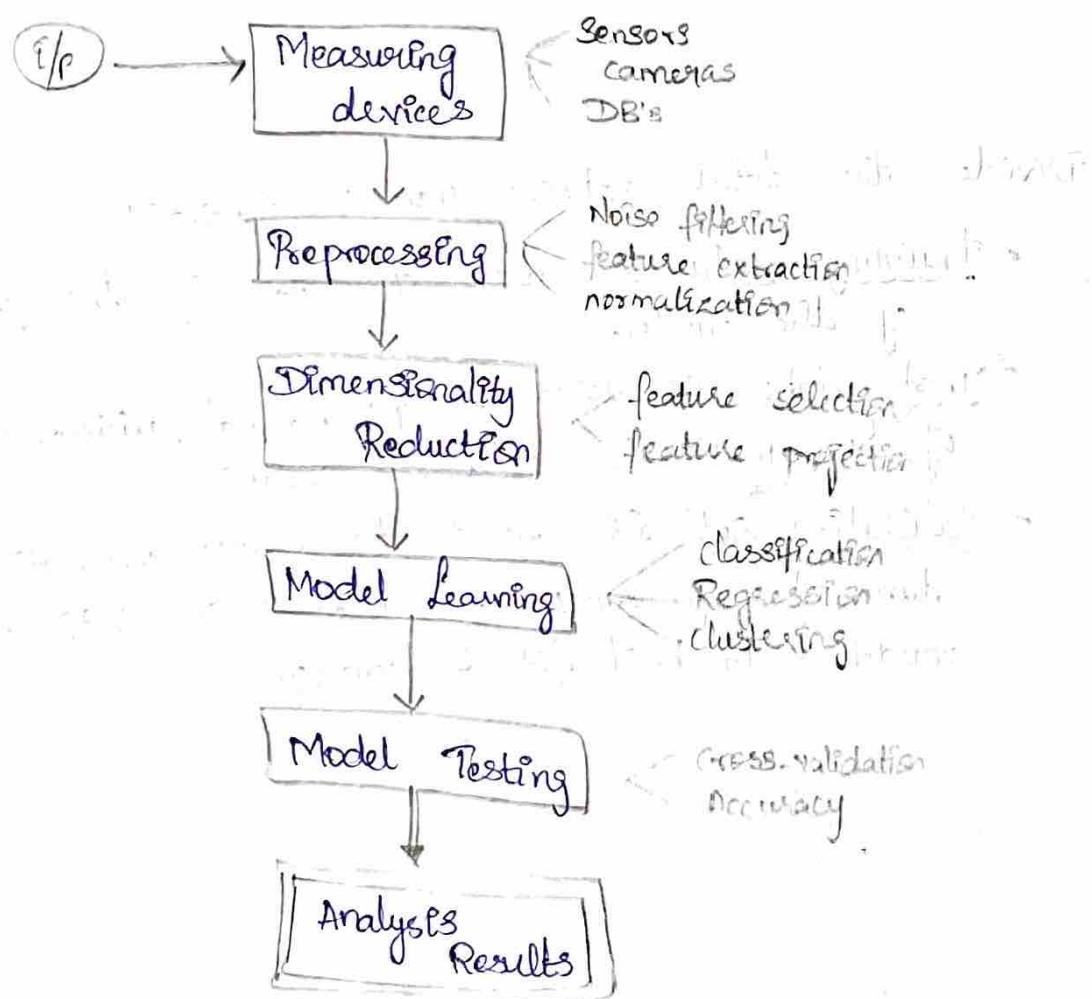
- Training data is used for learning the parameters of the model.
- Test data is used to get a final, unbiased estimate of how well the networks works.
- Validation set is used to reduce overfitting, tune the parameters - repeat this until we get model with best fit & accuracy.



## Hypothesis space:

- each setting of parameters in machine is a different hypothesis about function that maps I/p vectors to O/p vectors.
  - if data is noise free, each training example rules out a region of hypothesis space.
  - if data is noisy, each training example scales the posterior probability of each point in the hypothesis space in proportion to how likely the training example is given that hypothesis.
- Thus, the art of supervised ML is in:
- deciding how to represent the I/p & O/p.
  - selecting a hypothesis space that is powerful enough to represent the relationship b/w I/p's & O/p's.

## Learning Process:

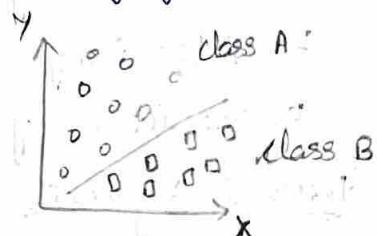


## Classification Algorithm:

- is a supervised learning technique that is used to identify the category of new observations on the basis of training data.
- It learns from given dataset & then classifies new observation into number of classes or groups, such as yes or no, 0 or 1, Spam or not spam etc...
- classes can be said as target or labels or categories
- unlike regression, o/p variable is category not a value.

$$y = f(x)$$

Categorical o/p      o/p variable



## Types of classification:

### \* Binary classifier

- has only two possible outcomes.
- Eg: male-female, Yes-No, 0-1, ... etc.

### \* Multi-class classifier

- has more than two possible outcomes
- Eg: classification of types of crops  
      " " music... etc.

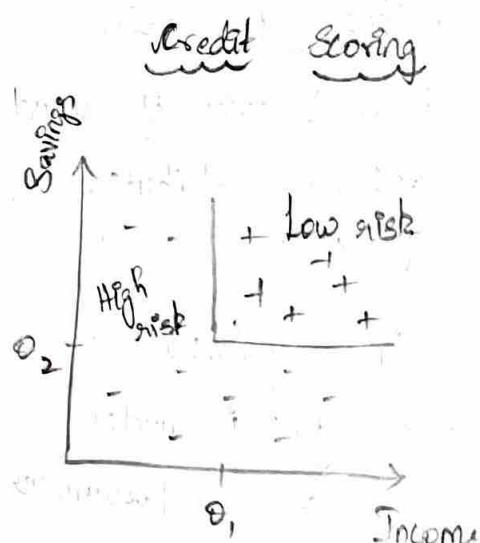
## Types of classification Alg:

### → Linear Models

- ↳ Logistic Regression
- ↳ SVM

### → Non-linear Models:

- ↳ KNN
- ↳ Naive Bayes
- ↳ Decision Tree
- ↳ Random forest



If  $\text{income} > \theta_1$  &  $\text{savings} > \theta_2$   
then low risk  
else  
high risk

## Regression Analysis:

- is used when o/p variable is continuous values such as 'weights', 'age', 'dollars' etc..
- is a supervised learning technique which helps in finding correlation b/w variables & enables to predict the continuous o/p variable based on one or more predictor variable.
- mainly used for prediction, forecasting, time series etc.

## Examples:

- ↳ prediction of rain using temperature & other factors
- ↳ determining market trends.
- ↳ prediction of road accidents due to rash driving

## Types of Regression:

- Linear regression
- Logistic "
- Polynomial "
- Support Vector "
- Decision Tree "
- Random forest "

## Prediction - Regression:

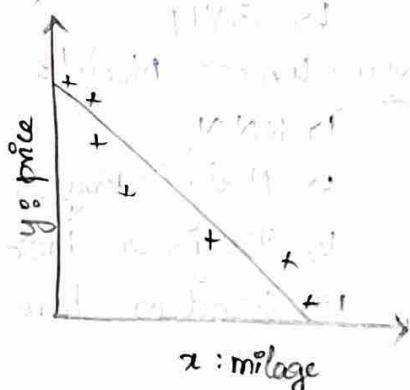
Predict price of used car

$x$ : car attributes

$y$ : price

$$y = g(x|\theta)$$

where  $g(\cdot)$  is model  
 $\theta$  is parameter



## Learning Associations:

### - Basket Analysis:

$P(Y|X)$  is a probability that somebody who buys  $X$  also buys  $Y$ ,

where  $X, Y$  are products, services...

$$\text{Eg: } P(\text{chips} | \text{beer}) = 0.7$$

i.e., 70% of customers who buy beer also buy chips.

### - Find frequent elements/items.

### - " the associations b/w them (o-s) relation b/w them.

### Eg: Single Dimensional Association rule:

Buy ( $X$ , 'Computer')  $\rightarrow$  Buy ( $X$ , 'S/W')

### Multi-Dimensional Association rule:

Age ( $X$ , '40-60') and Income ( $X$ , '50-lakhs')  $\rightarrow$  Buy ( $X$ , 'Computer')

i.e. Imagine that you are a sales manager at xyz company & you are talking to a customer who recently bought a PC & camera from store. What should you recommend to him/her next?

↳ frequent patterns & association rules are the knowledge that you want to mine in such scenario.

### Algorithms:

→ Apriori

→ Frequent pattern growth tree.

## Supervised Learning - Learning a class from Examples:

\* Lets take a set of cars. i.e., class - C: "family of cars"

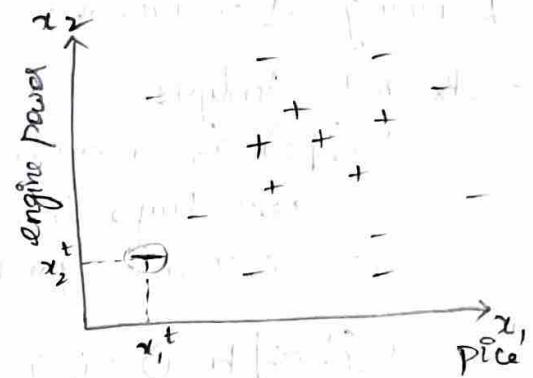
- A group of people look at the cars & label them as family car or not with two attributes as price & engine power.
- Ignore other attributes such as color, seating capacity etc.
- The cars that they believe as family cars are "positive examples" & other cars are "negative example".

- Single data point corresponds to one sample car.

- Coordinates?

$x_1$  = price

$x_2$  = engine power



- '+' → positive example (family car)

- → negative " (not a family car)

### Training Set:

→ Variables ' $x$ ' and ' $y$ '

where  $x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$

$y = \begin{cases} 1, & \text{if } x \text{ is +ve example (family car)} \\ 0, & \text{if } x \text{ is -ve example (not family car)} \end{cases}$

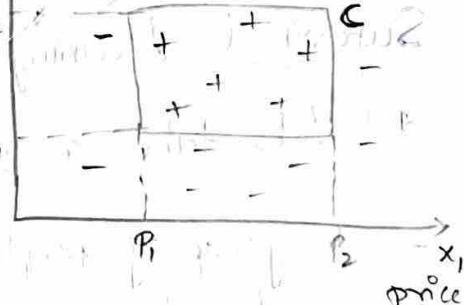
→ Each car is represented by ordered pair  $(x, y)$  and training set contains ' $N$ ' such examples.

i.e., 
$$X = \left\{ x^t, y^t \right\}_{t=1}^N$$

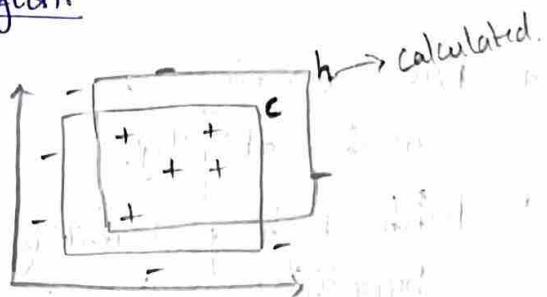
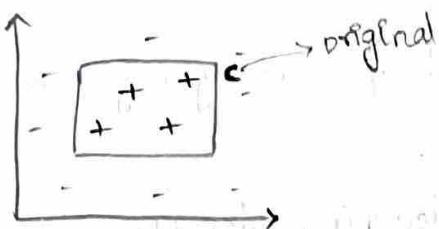
### Hypothesis class:

- If a car is family car, then its price & engine power should be in certain range.  
i.e.,  $(p_1 \leq \text{price} \leq p_2)$  and  $(e_1 \leq \text{engine} \leq e_2)$
- the class of family car is a rectangle in price-engine power space.
- Hypothesis  $h$  is specified by quadruple of  $(p_1, p_2, e_1, e_2)$  to approximate ' $c$ '.

$$h(x) = \begin{cases} 1, & \text{if } h \text{ classifies } x \text{ as +ve example (family car)} \\ 0, & \text{if } h \text{ " " " " -ve " (not a family car)} \end{cases}$$

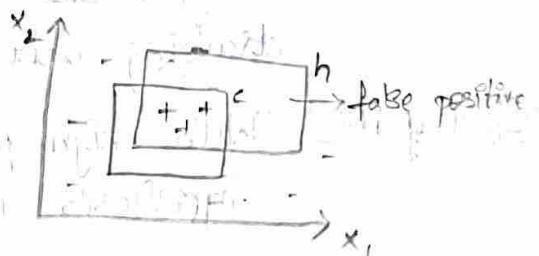
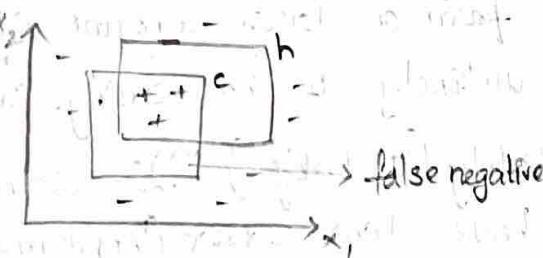


- 'c' is a target function.
- In real life we don't know  $c(x)$ , so we cannot evaluate how well  $h(x)$  matches  $c(x)$ .
- Instances within rectangle represents family car and outside are not family cars.
- There may be some error region.



False positive & false Negative:

- 'c' is actual class & 'h' is our produced hypothesis.
- point where c is 1 & h is 0 is false negative.
- " " " c is 0 but h is 1 is false positive.
- true positive & true negative - are correctly classified.

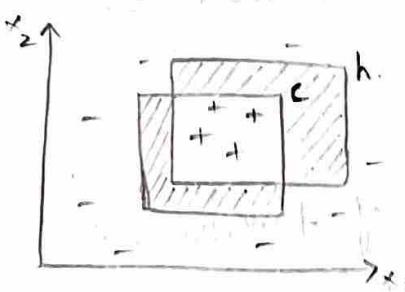


Error:

- the empirical error is the proportion of training instances where predictions 'h' do not match required values given in  $x$ .

$$E(h|x) = \sum_{t=1}^N I(h(x^t) \neq r^t)$$

$$\begin{aligned} h(x) &= \begin{cases} 1, & \text{true} \\ 0, & \text{false} \end{cases} \\ r &= \begin{cases} 1, & \text{true} \\ 0, & \text{false} \end{cases} \end{aligned}$$



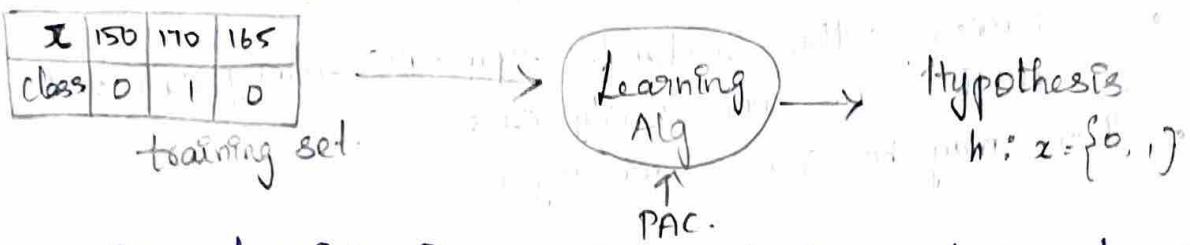
- probability of error should be small.

Error region :  $P(c \neq h) \leq \epsilon$

## Supervised Learning: PAC Learning:-

↳ Probably Approximately correct (PAC).

- Cannot expect the learner to learn a concept exactly.



\* PAC learning is a framework for analysing learning alg. mathematically.

\* With high probability, a PAC learning algorithm finds a hypothesis that is approximately identical to the hidden target concept.

✓ Given training set of labelled examples, learning alg generates a hypothesis. Run hypothesis on test set to check how good it is.

✓ Intuition: A hypothesis built from a large amount of training data is unlikely to be wrong. i.e. PAC

✓ Goal of PAC: With high probability ('probably') the selected hypothesis will have low error ('approximately correct')

\* PAC learning requires,

- Small parameters  $\epsilon$  and  $\delta$ .
- with probability atleast  $(1-\delta)$ , a system learns the concept with error at most  $\epsilon$ .

### Notations:

1) Instance Space  $x$ .

2) Concept class,  $C$  - family of functions, 'c'  
 $c : x \rightarrow \{0,1\}$

3) Hypothesis class,  $H$  - family of functions, 'h'  
 $h : x \rightarrow \{0,1\}$

4) Probability distribution,  $f$  - training samples are generated randomly from  $x$  according to  $f$ .

## 5) Learning Alg. A. or L

Defn:

A concept class  $C'$  is PAC learnable if, there is an alg. 'A' which for samples drawn with any probability distribution 'F' and any concept  $c \in C$ , with high probability produce a hypothesis  $h \in C$  where error is small.

Example:

Same from last topic Hypothesis space to error.

Error region:  $P(c \text{ XOR } h) \leq \epsilon$

Approximately Correct:

- the hypothesis  $h$ , that approximately correct and error is less than or equal to  $\epsilon$ .  
where  $0 \leq \epsilon \leq \frac{1}{2}$  i.e.  $P(c \text{ XOR } h) \leq \epsilon$

Probably Approximately Correct:

- less error with high probability  
 $[P(\text{Error}(h) \leq \epsilon)] \leq 1 - \delta$   
 $[P(P(c \text{ XOR } h) \leq \epsilon)] \leq 1 - \delta$

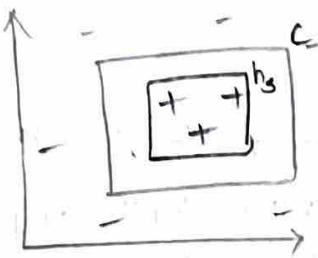
Formal defn:

A concept class,  $C$  is said to be PAC learnable by alg.  $L$  if for all  $c \in C$ , distribution  $F$  over  $X$ ,  $\epsilon$  such that  $0 < \epsilon < \frac{1}{2}$  and  $\delta$  such that  $0 < \delta < \frac{1}{2}$ , the learner  $L$  will output a hypothesis ' $h$ ', with probability atleast  $(1-\delta)$  and  $\text{error}_F(h) \leq \epsilon$ .

PAC learnability for axis-aligned rectangle:

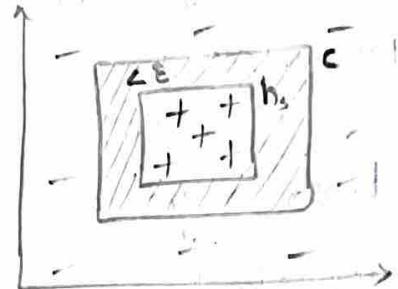
↳ Specialization:

- $h_s$  is the tightest possible rectangle around a set of positive training examples.
- $h_s$  is subset of  $C$ , hence  $\boxed{\text{error region} = C - h}$

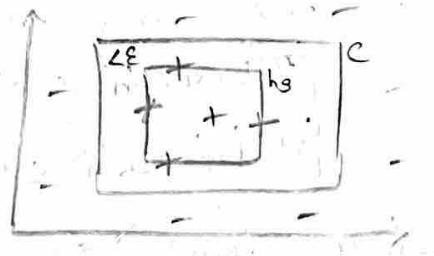
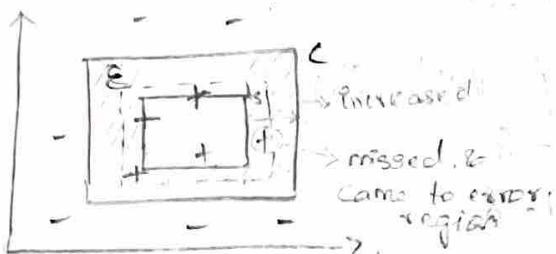


Approximately Correct:

- If an hypothesis lies b/w  $h$  &  $c$  (shaded region) then it is approximately correct.



- If the generated hypothesis does not touch any of these regions, then error region is greater than  $\epsilon$ . and not approximately correct, coz error region got increased.
- Should have atleast one +ve example at each side of rectangle.



Error Region:

Error region = sum of four rectangular strips  $< \epsilon$

Each strip is atmost  $\epsilon/4$

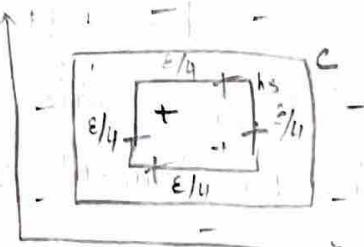
prob. of +ve example falling in any one of strip =  $\epsilon/A$

prob. of randomly drawn +ve example misses a strip =  $1 - \epsilon/4$

$P(m \text{ instance miss a strip}) = (1 - \epsilon/4)^m$

$P(m \text{ instance miss any strip}) < 4(1 - \epsilon/4)^m$

$$\therefore [m > 4/\epsilon \log 4/8]$$



### Ex. problem 1:

Hypothesis  $h_1$  generated the errors with respect to price & engine power of given 10 samples,

Given:  $\epsilon = 0.05$ ,  $\delta = 0.20$

$$\text{Now, } P(h_1) \geq 1 - \delta$$

$$P(h_1) = \frac{8}{10} \rightarrow \text{total 10 values where } 3^{\text{rd}}, 4^{\text{th}}, 8^{\text{th}} \text{ are greater than } \epsilon \\ = 0.80 \quad \text{so } [10 - 2 = 8]$$

$$\therefore 0.80 \geq 1 - 0.20$$

$$0.80 \geq 0.80$$

Hence  $h_1$  is probably approximately correct.

S.No	Error( $h_1$ )
1	0.001
2	0.025
3	0.07
4	0.003
5	0.035
6	0.045
7	0.027
8	0.065
9	0.012
10	0.036

### Ex. problem 2:

Hypothesis  $h_2$  generated the errors with respect to price & engine power of given 10 samples,

Given:  $\epsilon = 0.05$ ,  $\delta = 0.20$

$$P(h_2) \geq 1 - \delta$$

$$P(h_2) = \frac{7}{10} \rightarrow 3^{\text{rd}}, 4^{\text{th}}, 9^{\text{th}} \text{ are greater than } \epsilon \\ = 0.70 \quad [10 - 3 = 7]$$

$$0.70 \geq 1 - 0.20$$

$$0.70 < 0.80$$

Hence  $h_2$  is not probably approximately correct.

### Noise and Model Complexity:

→ 'Noise' is any unwanted anomaly in the data & due to noise, the class may be more difficult to learn and zero error may be infeasible with simple hypothesis class.

#### Reason for noise:

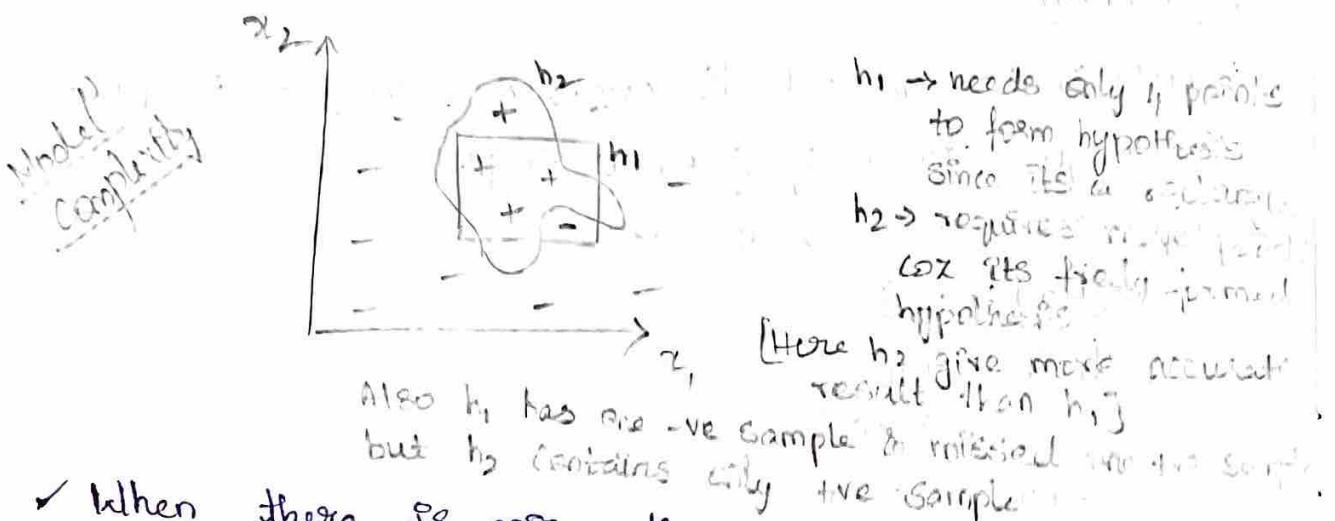
- There may be imprecision in recording the I/P attributes, which may shift the data points in I/P space.
- There may be errors in labeling the data points, which may relabel positive instance as negative & vice versa. This is called as 'teacher noise'.

- There may be additional attributes, which we have not taken into account, that affect the label of instance. Such attributes may be hidden or latent which may be unobservable.

↳ the effect of these neglected attributes is thus modeled as random component & included in 'noise'.

### Effect of noise:

- Noise distorts data.
- due to noise, the learning problem will not give any accurate results



- ✓ When there is noise, there is not a simple boundary between the +ve & -ve instances, & zero misclassification error may not be possible with simple hypothesis.
- ✓ A rectangle ( $h_1$ ) is a simple hypothesis with 4 parameters defining the corners.
- ✓ An arbitrary closed form can be drawn by piecewise functions with larger number of control points.

### Occam's Razor:

Using simple rectangle makes more sense coz of following.

- It is simple model to use.
- ↳ easy to check whether point is inside or outside rectangle.

2. It is simple model to train. & has fewer parameters.  
↳ easy to find corner values of rectangle than the control points of arbitrary shape.
3. It is simple model to explain.  
↳ a rectangle simply corresponds to defining intervals on two attributes. By learning a simple model, we can extract info. from raw data given in training set.
4. Given the empirical error, we say that a simple model would generalize better than complex model.  
↳ this principle is known as 'Occam's Razor', which states that simpler explanations are more plausible and any unnecessary complexity should be shaved off.

### Learning Multiple Classes:

→ Single class problem:

- Considering only single class while learning the examples are called as single class problem.

Eg: learning a family car, we have +ve example belonging to class family car & never consider other samples

→ Two class problem:

- Considering two classes while learning the examples are called two class problem.

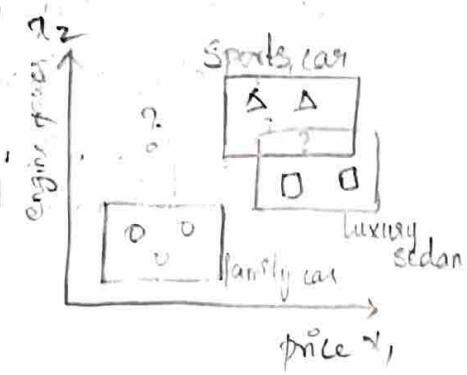
Eg: Learning a family car, we have +ve examples belonging to class family car & -ve examples belonging to all other cars.

→ Multiple classes:

- More than 2 classes are said to be multiple classes

Eg: Consider three classes:

- family car
- sports car
- luxury sedan.



- Here three hypotheses induced, each one covering instance of one class & leaving outside the instance of other two classes. (like when  $p_1, p_2, p_3$  are  $\frac{1}{3}$  some value then it is family etc)
- And '?' are reject regions where 'no' or more than one class is chosen.
- \* Here  $K$  classes denoted as, ' $C_p$ ' where  $p = 1, 2, \dots, K$ .  
an  $\text{I/P}$ , instance belongs to one & exactly one of them.

$$x = [x^t, r^t]_{t=1}^N$$

where  $r$  has  $K$  dimensions &

$$r_i^t = \begin{cases} 1, & \text{if } x^t \in C_p \text{ (belong to any one class)} \\ 0, & \text{if } x^t \in C_j, j \neq p \text{ (doesn't belong to any class)} \end{cases}$$

i.e., reject region

- \* The boundary separating the instance of one class from the instances of all other classes. (form hypothesis)
- \*  $K=2$  if it has two classes like +ve & -ve classes  
So it is  $K$ -class classification problem.
- \* The training set belonging to  $C_p$  are +ve instance of hypothesis ' $h_i$ ' & the examples of all other classes are -ve instance of ' $h_i$ '.

- \* So, in  $K$ -class problem, we have  $K$ -hypotheses to learn,

$$h_i(x) = \begin{cases} 1, & \text{if } x^t \in C_i \\ 0, & \text{if } x^t \in C_j, j \neq i \end{cases}$$

Error:

The total error takes the sum over prediction for all classes over all instances.

$$E\left(\{h_i\}_{i=1}^K | x\right) = \sum_{t=1}^N \sum_{i=1}^K I(h_i(x^t) \neq r_i^t)$$

$K$ -classes,  $t$  = No. of training samples.

Eg: Handwritten digit recognition dataset,  
all digits has similar distribution & belong to same class.

$$\begin{array}{l} 0 \ 0 \ 0 \rightarrow C_1 \\ 1 \ 1 \ 1 \rightarrow C_2 \\ 2 \ 2 \ 2 \rightarrow C_3 \end{array}$$

## Model Selection and Generalization:

### Model :



- When you train an alg. with data it will become a model.

### Model Selection:

- is a process of selecting an optimal model from a set of candidate models, given data.

Data used in learning are:

- { ↳ Training data
- ↳ Validation data
- ↳ Test data

- finding an optimal model is minimizing both the error of train data & error of test data.

### Bias:

- ↳ can be defined as diff. b/w the expected prediction of our model & correct value which we try to predict
- ↳ High bias can cause our model to miss significant relations b/w our features ( $x$ ) & o/p's ( $y$ ) so it cannot learn training data.
- ↳ This is also known as 'under-fitting'.

make lot of assumptions which cause inaccurate predictions

## Variance:

- ↳ is the error from sensitivity to small fluctuation in the training data.
- ↳ high variance can cause random noise to be introduced into training data rather than intended o/p's.
- ↳ This is also known as "over-fitting".  
learns training data too well to cannot generalize to new data.

## Inductive bias:

- Set of assumptions we make to have learning possible is called Inductive bias of learning alg.
- Model Selection is also about choosing the right inductive bias.

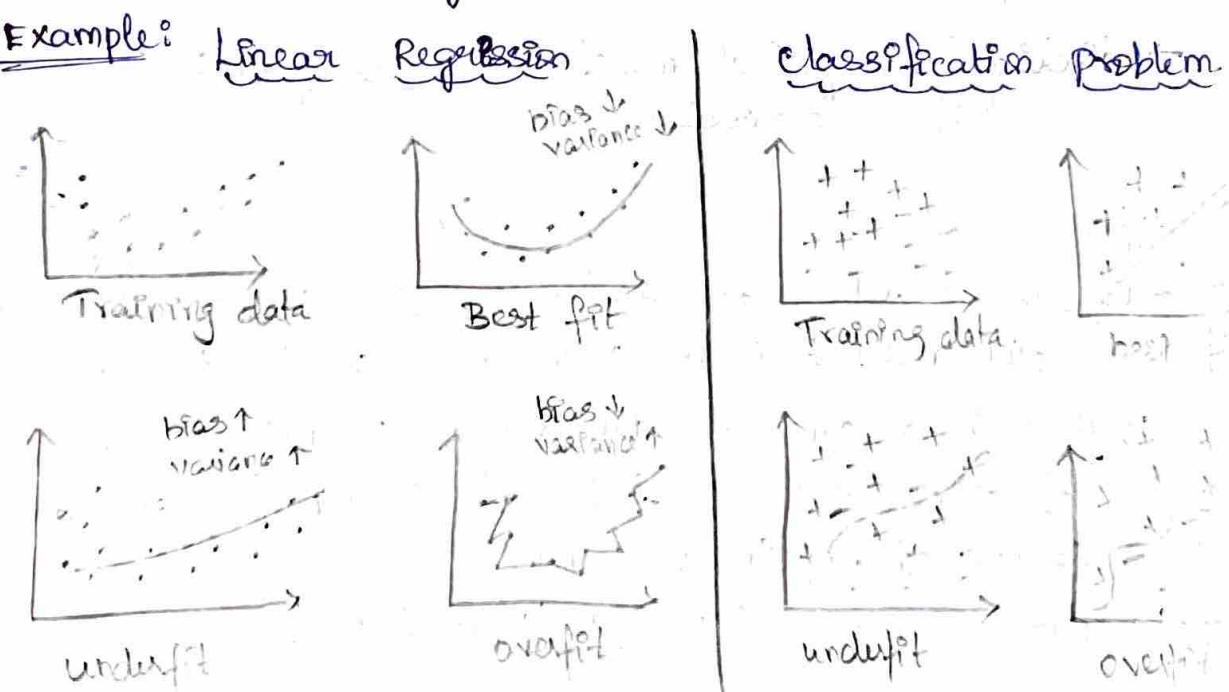
## Generalization:

- is how well a model trained on training set predicts the right o/p for new instances.
- the model should be selected having the best generalization.

## Need to avoid two problems:

- ↳ underfitting.
- ↳ overfitting.

## Example:



## underfitting:

- production of ML model that is not complex enough to accurately capture relationships b/w dataset features & target variables.

## overfitting:

- production of model which corresponds exactly to dataset (which may contain noise) & therefore fail to predict new data

## Dimensions of Supervised ML Algorithms:

\* We know that,

$$\mathcal{X} = [x^t, y^t]_{t=1}^N$$

where,  $t$  is training set with 1 to  $N$  samples.

$x^t$  is arbitrary I/P

$y^t$  is desired o/p. (Is a real value in regression)

for two-class learning  
for ( $K > 2$ ) multi-class, exactly 2 dimension if  $t$  is 1 & other is 0

\* The aim is to build a good & useful approximation to  $y^t$  using the model  $[g(x^t|\theta)]$  parameters.

\* For doing this, three decisions must be made:

1. Model we use in learning
2. Loss function.
3. Optimization procedure.

→ Model used in learning:

$$[g(x|\theta)]$$

where  $g(\cdot)$  is model

$x$  is I/P (sample)

$\theta$  are parameters

→ Loss function:

→ If there is a error, there will be a difference b/w desired o/p ' $y^t$ ' & can be found using the loss function  $L(\cdot)$ .

$$E(\theta|x) = \sum_t L(y^t, g(x^t|\theta))$$

• approximation error (or) loss is the sum of losses over the individual instances.

where  $E(\cdot)$  is error

$\theta$  is parameters

$x$  is sample

$t$  is training set from 1 to N.

$L(\cdot)$  is loss function.

$y^t$  is d/p.

$g(x^t|\theta)$  is model.

→ optimization procedure:

- is used to find  $\theta^*$  that minimizes total error.

$$\theta^* = \arg \min_{\theta} E(\theta|x)$$

where,

$\arg \min$  → returns argument that minimizes.

$E(\theta|x)$  → error.

parameters → sample.