# MACHINE LEARNING

**Definition :** A computer program is said to learn from experience E with respect to some class of tasks T and performance measure P, if the performance on task T as measured by P increases with experience E.

→ Supervised learning & unsupervised learning

**Supervised learning**
Algorithm is fed the "right answers". Task of the algo is to predict more right ans

i> Regression - predict continuous real valued o/p.

ii> classification - Discrete value o/p

**unsupervised learning**
Algo is not fed with labels / instructions on performing a task.
used to organise clusters of data & identify them.
eg - Social n/w analysis
Market segmentation
Astronomical data analysis
The algo is responsible for finding structure in the data

**\* REGRESSION**

Notations

$m$ = number of training examples

$x$ = input variable / features

$y$ = output variable / target variable · to be predicted.

1 training example - $(x, y)$

$i^{th}$ " " $(x^{(i)}, y^{(i)})$

Supervised learning algorithm → Training example

Training set $(x^{(i)}, y^{(i)})$ $i = 1 \ldots m$
↳ Training set

↓

Learning algorithm

↓

i/p → hypothesi 'h' → o/p (estimated value)
(x)
$h: X \to Y$
(y)

$h_\theta(x) = \theta_0 + \theta_1 x$ → Linear regression with one variabl
univariate linear regression.

When the target variable that is to pro be predicted is
continuous, it is called a regression problem.
When y takes a small no: of discrete values → classificatio
problem.

* COST FUNCTION

$\theta_i$'s → parameters

With diff $\theta_i$'s we get different hypothesis.

Time series classification

→ Model based : Auto regression & HMM

→ distance based : Dynamic time warping

→ Feature based : extract meaningful features like
  DFT, STFT, DWT, PCA, SVD etc

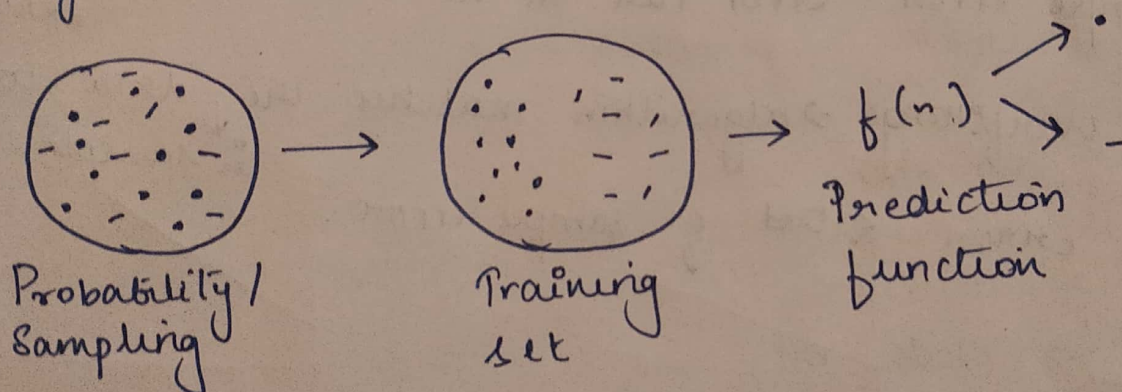Automatic feature based approaches using deep learning models → convolutional neural networks.

Training all parameters is done jointly using back propogation method.

# Central dogma of prediction

Any set of data requires probability / sampling to create a training set that differentiates the various samples.

We then extract characteristics from the data set by which we can model a function on to predict on new samples as to which class it belongs to.



Probability /
Sampling

Training
set

$f(n)$

Prediction
function

Components of a predictor

question → i/p data → features → algorithm → parameters → evaluat-
(what to                    (Characters)   (prediction    ion.
Predict ?)                                  function
                                            definition)

Properties of good features
→ Lead to data compression — We need to collect signi[ficant]
                                     features of lesser size.

→ Retains relevant info
→ Are created based on expert application knowledge.


Prediction is about accuracy tradeoffs
→ Interpretability vs accuracy
→ Speed vs accuracy
→ Simplicity vs accuracy
→ Scalability vs accuracy.

* In sample error — the error rate you get on the same dat[a]
that was used to build the predictor.   RE~~DISTRI~~ SUBSTITUITION ERROR

* Out of sample error — error rate on new data set . GENERALIZ[ATION]
                                                        ERROR.

* Reason for overfitting → algorithm matches the data too well
                                                      ↓
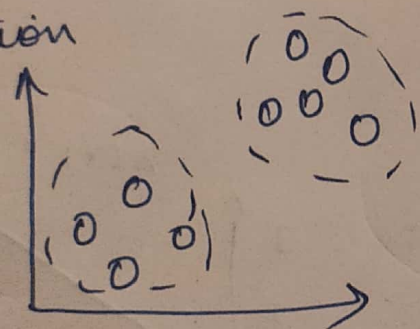                                                   in sample.

* In sample error < Out of sample error .

Supervised Learning — Right answers are given prior to building the model.

Regression — Prediction of continuous valued o/p.
Classification — Discrete valued output (0 or 1) (0,1,2.. )

Unsupervised learning — no information on the data is given. The task is to find some structure to the data which can later on, be used to extract some information

Eg :-

→ clustering. The algo might decide based on the data — like it could be seperated as shown
— clustering algo.

We can derive structure from the data where we dont necessarily know the effect of the variables. we can obtain this structure by clustering data based on relationships among the variables in the data.

✳ No feedback based on the prediction results.

## Model Representation

$x^{(i)}$ → i/p (features)

$y^{(i)}$ → o/p (target)

$m$ → number of features

$i = 1 \ldots m$ → training set.

$(x^{(i)}, y^{(i)})$ → training example.

Flow of supervised learning :

$$\boxed{\text{Training set}}$$

$$\downarrow$$

$$\boxed{\text{Learning algo}}$$

$$\downarrow$$

$x \to \boxed{h} \to$ predicted $y$

hypothesis

Goal is : given a training set, to learn a fn $h : x \to y$ so that $h(x)$ is a good predictor for the corresponding value of $y$.

$$h_\theta(x) = \theta_0 + \theta_1(x) \to \underset{\wedge}{\text{Linear regression}}$$
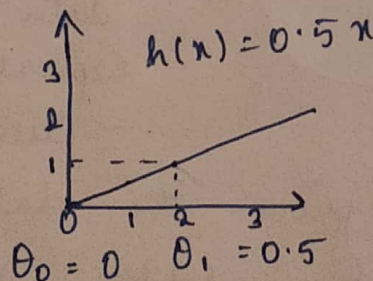
Univariate

$\theta_i$'s → parameters of the model.

Eg: –



$\theta_0 = 1.5 \quad \theta_1 = 0$
$h(x) = 1.5$

$h(x) = 0.5 x$

$\theta_0 = 0 \quad \theta_1 = 0.5$
$h(x) = 0.5 x$

$h(x) = 1 + 0.5 x$

$\theta_0 = 1 \quad \theta_1 = 0.5$
$h(x) = 1 + 0.5 x$

Idea : Choose $\theta_0, \theta_1$ such that $h(x)$ is close to $y$ for our training example $(x, y)$

∴ minimize $\underset{\theta_0, \theta_1}{\text{ }} \frac{1}{2M} \sum_{i=1}^{M} (h_\theta(x^{(i)}) - y^{(i)})^2$

// why are we averaging?

$\frac{1}{2M}$ ?

Square error func
or
↑ mean squared
error,

Cost function → $J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^{M} (h_\theta(x^{(i)}) - y^{(i)})^2$

$\underset{\theta_0 \theta_1}{minimize}$

$\underset{\theta_0 \theta_1}{minimize} \underbrace{J(\theta_0, \theta_1)}_{Cost\ func.}$

Basically the difference between the predicted value and the actual value.

// $\frac{1}{2}$ used as convenience for gradient descent because it can later be used to cancel out when derivative is taken.



$h_\theta(x) = \theta_1 x$    // simplified version

$J(\theta_1) = \frac{1}{2m} \sum_{i=1}^{M} (h(x) - y)^2$

what is $J(0) = ?$

$h_\theta(x^i) = ?$   for every value of $x$, $h_\theta(x) = 0$

$J(\theta_1) = \frac{1}{2 \times 3} \sum_{i=1}^{3} [(0-1)^2 + (0-2)^2 + (0-3)^2]$

$= \frac{1}{6} [1 + 4 + 9] = \frac{14}{6}$

YAY!

*Gradient descent - Algorithm used for minimizing any fn.

Repeat until Convergence $\{$

$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$   for $j=0$ & $j=1$

↓ Assignment   ↳ learning rate    $\alpha \uparrow\uparrow$ v. large → aggressive gradient descent

$\}$

We want to min $J(\theta_0, \theta_1)$ by gradient descent
$\theta_0, \theta_1$

① start with an arbituary $\theta_0, \theta_1$
② keep varrying $\theta_0$ & $\theta_1$ until minimum is acheived.

\* Simultaneously update $\theta_0$ & $\theta_1$

Step 1 : $temp0 := \theta_0 = \theta_0 - \alpha \dfrac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$

$temp1 := \theta_1 = \theta_1 - \alpha \dfrac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$

step 2 : $\theta_0 = temp0$

$\theta_1 = temp1$

Eg :- $\theta_0 = 1$  $\theta_1 = 2$ , find the update.

$temp0 :- \theta_0 = 1 - \alpha \dfrac{\partial}{\partial \theta_0} J(\overset{1}{\theta_0}, \overset{2}{\theta_1})$

where $J(\theta_0, \theta_1) = \dfrac{1}{2M} \overset{M}{\underset{i=1}{\sum}} (h_\theta(x^i) - y^i)^2$  ? $\circ$ // Not finished.

// Vector = $n \times 1$ matrix. $n$ depicts the dimension.

Can add & sub matrices only if they are of the same dimensions.

Prediction = Data matrix $\times$ Parameters.

Learning Problems in ML :-

Classification - assign a category.

Regression - Predict a real value

Rank - order items : in

clustering - Seperation of region into homogeneous segme

Dimensionality reduction /Manifold learning - learn a

lower dimensional representation.

→ Supervised learning : Classification , Regression
→ Un " : Density estimation, Clustering,
demensionality reduction.
→ Semi - supervised "
→ Reinforcement " etc.


Decision Trees

Basic algorithm

1. Start with all variables in one group

2. Find the variable / split that best seperates
the outcomes

3. Divide data into two groups ("leaves") on
that split ("node")

4. Within each split find the best variable / sp
that seperates the outcomes.

5. Continue till the groups are too small / succintly

# Measure of Impurity.

$\longrightarrow$ signifies the number of times that particular class appears in that leaf

$$\hat{P}_{mk} = \frac{1}{N_m} \sum_{x_i \text{ in leaf } m} \mathbb{1}(y_i = k)$$

$\hat{P}_{mk}$ $\downarrow$

Probability to estimate

$m^{th}$ leaf    $k^{th}$ class

$N \rightarrow$ no. of objects in that class.

$N \rightarrow$ no. of objects in that class -

Misclassification Error $\rightarrow$ $1 - \hat{P}_{mk(m)}$ ; $K(m) =$ most common class.

$0 =$ perfect purity

$0.5 =$ no purity .    $\rightarrow$ when leaves are perfectly balanced; we don't get homogeneity. If 1, then perfect purity in another direction?

Gini Index:—

$$\sum_{k \neq k'} \hat{P}_{mk} \times \hat{P}_{mk'} = \sum_{k=1}^{K} \hat{P}_{mk}(1 - \hat{P}_{mk}) = 1 - \sum_{k=1}^{K} P_{mk}^{2}$$

$\hookrightarrow$ Same for this .

Deviance / Info gain

$$- \sum_{k=1}^{K} \hat{P}_{mk} \log_2 \hat{P}_{mk}$$

$1 =$ no purity    $0 =$ purity .