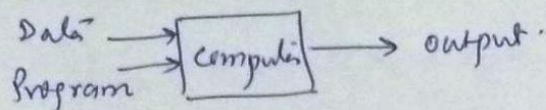


* Machine learning:-

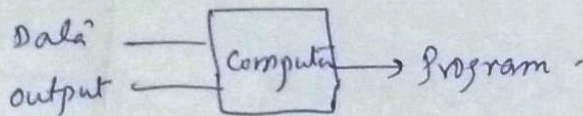
learning means = improving with experience at some task.

- Improve over task T ,
- With respect to performance measure, P
- Based on experience, E .
(data)

⇒ Traditional Programming:- Normal Computer vs ML



Machine learning:-



⇒ Types of Machine learning:-

└ Supervised learning.

└ providing right answers for every data.

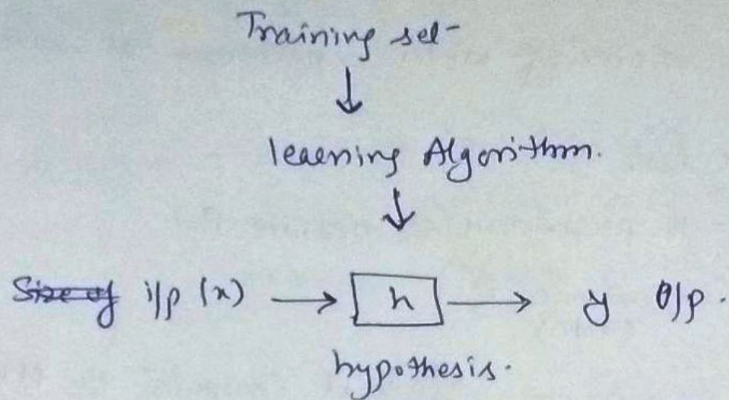
└ Classification → Discrete values

└ Regression = Continuous values. (infinite features).

└ Unsupervised learning - clustering the data with different arrangements.

└ Reinforcement learning:- trial and error method in finding the best outcome based on experience.

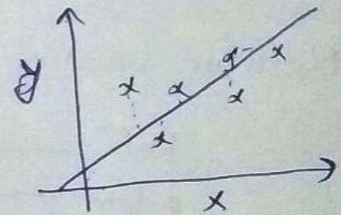
⇒ Regression :- It can be used in cases to determine casual relations between the independent and dependent variables.



So, how do we represent this hypothesis 'h', which will map the input to the output.

① Linear Regression :-

$$h_{\theta}(x) = \theta_0 + \theta_1 x.$$



So to minimize $h_0(x)$, we have to minimize θ_0, θ_1 .

for that- we can define a function such as cost function.

an :-

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

$$h_0(x'') = \text{predicted value}$$

$y^{(i)}$ = actual value

$$hox^{(1)} - y^{(1)} = \text{error of the given data.}$$

This cost function with two parameters can be very complex while computation, so we need some software to minimize the cost function -

Gradient- Descent-

Gradient Descent

③

for minimizing the cost function $J(\theta_0, \theta_1)$

$$\theta_j = \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$$

$\alpha \rightarrow$ learning parameter. (chosen smartly)

Gradient means \rightarrow slope of the line we draw on the graph.
 \hookrightarrow automatically take smaller steps, so no need to decrease α (learning rate) over & over again.

* Linear Regression with multiple variables - (multiple features)

$$h_\theta(x) = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \dots + \theta_n x_n.$$

$$\therefore h_\theta(x) = \theta^T x.$$

$\theta \in \mathbb{R}^{n+1} \rightarrow (n+1)$ variables.

Cost function:-

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$$

* Gradient Descent function \rightarrow

$$\theta_j = \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta) \quad \text{simultaneously update for every } j = 0, 1, 2, \dots, n.$$

for $n=1$

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_0^{(i)}$$

$$\theta_1 = \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_1^{(i)} \rightarrow x_1^{(i)} \text{ for variable } x_1.$$

for $n \geq 1$

$$\theta_n = \theta_n - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_n^{(i)}$$

⇒ Features scaling:-

↳ Reducing the number of features or range of features

↳ To reach the minima within less steps or to converge easily.

⇒ Scaling range $-1 \leq x \leq 1$ - Scaled using average methods.

⇒ Features can be scaled using Mean Normalization.

$$x_i = \frac{x_i - \mu_i}{s_i}$$

μ_i → avg value of s_i

s_i - range of values

* Normal Equations ::

↳ method to solve θ analytically by using $\theta = (X^T X)^{-1} X^T y$.

For 1D. ($\theta \in \mathbb{R}$) & $J(\theta) = a\theta^2 + b\theta + c$, to calculate the minima.

$$\frac{d}{d\theta} J(\theta) = 0 \text{ \& solve for } \theta.$$

For $\theta \in \mathbb{R}^{n+1}$ → multiple variables.

$$J(\theta_0, \theta_1, \theta_2, \dots, \theta_n) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

$$\frac{\partial}{\partial \theta_j} J(\theta) = \dots = 0 \text{ for every } j, \text{ solve for } \theta_0, \theta_1, \theta_2, \dots, \theta_n.$$

(5)

* Logistic Regression :- for classification problems

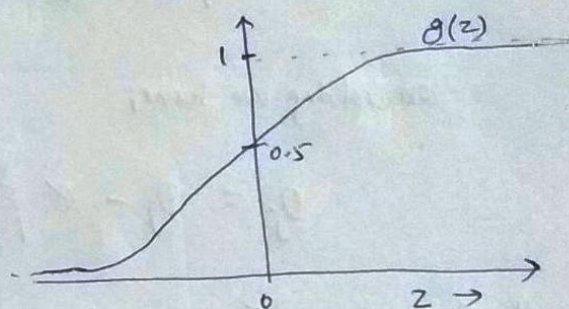
$$0 \leq h_{\theta}(x) \leq 1$$

So logistic regression function can be defined as :-

$$h_{\theta}(x) = g(z) \quad , \text{ where } g(z) = \frac{1}{1 + e^{-z}}$$

if we have $z = \theta^T x$ for linear regression, we conclude that.

$$h_{\theta}(x) = \frac{1}{1 + e^{-(\theta^T x)}}$$



$h_{\theta}(x)$ = estimated probability that $y=1$ on x .

$h_{\theta}(x) = P(y=1 | x; \theta)$ Probability that $y=1$, given x parameterized by θ .

$$P(y=0 | x; \theta) + P(y=1 | x; \theta) = 1.$$

* Decision Boundary :-

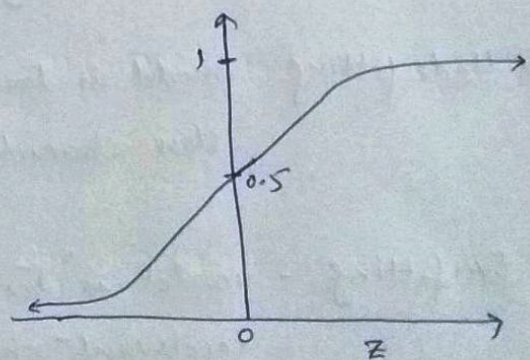
$$h_{\theta}(x) = g(\theta^T x) ; P(y=1 | x; \theta)$$

for $y=1$ if $h_{\theta}(x) \geq 0.5$

$$\Rightarrow \boxed{\theta^T x \geq 0}$$

for $y=0$ if $h_{\theta}(x) < 0.5$

$$\Rightarrow \boxed{\theta^T x < 0}$$



* Cost function of Logistic Regression :-

$$\text{Cost}(h_{\theta}(x), y) = -y \log(h_{\theta}(x)) - (1-y) \log(1-h_{\theta}(x))$$

for multiple values :-

$$J(\theta) = -\frac{1}{m} \left[\sum_{i=1}^m y^{(i)} \log h_{\theta}(x^{(i)}) + (1-y^{(i)}) \log (1-h_{\theta}(x^{(i)})) \right]$$

To fit parameter θ , minimize $J(\theta)$

\Rightarrow Gradient Descent

$$\theta_j = \theta_j - \alpha \frac{\partial J(\theta)}{\partial \theta_j}$$

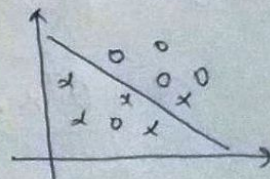
on solving we have;

$$\theta_j = \theta_j - \alpha \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

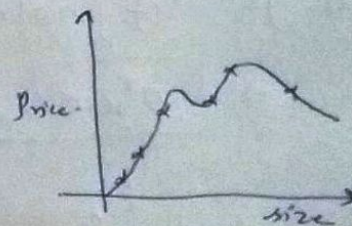
* Components of generalization error :-

- Bias :- It is a phenomenon that skews the result of an algorithm in favour or against an idea.
- Variance :- It refers to the changes in the model when using different portions of training data set.

\Rightarrow Underfitting : model is too "simple" to represent all the relevant class characteristics.



\Rightarrow Overfitting :- model is "too complex" and fits irrelevant characteristics (or noise) in the data



* Regularization :-

- ↳ keep all the parameters or features but reduce the magnitude or values of parameters.
- ↳ works well when we have lot of features, each of which contributed a bit to predicting y .

* Regularized cost functions:-

- Linear Regression:-

$$J(\theta) = \frac{1}{2m} \left[\sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n \theta_j^2 \right]$$

- Gradient Descent:-

$$\theta_j = \theta_j - \alpha \left[\frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)} + \lambda \theta_j \right]$$

- Normal Equation:-

$$\theta = \left(X X^T + \lambda \begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \right)^{-1} X^T y$$

3×3

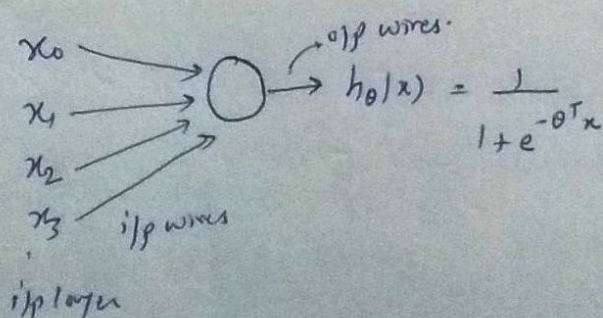
for $n=2$, matrix will be $(n+1)(n+1)$

- Logistic Regression:-

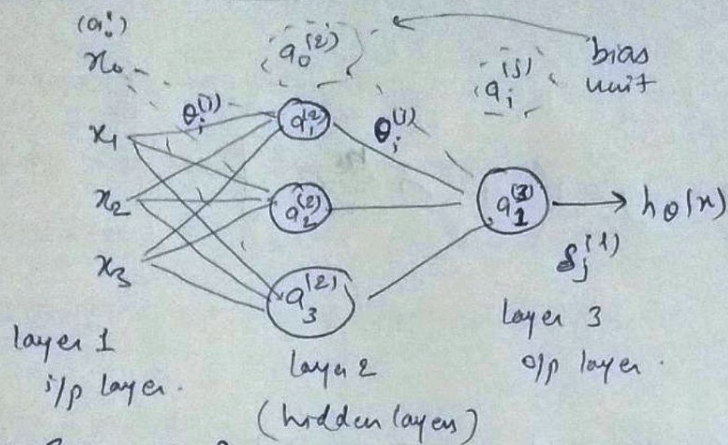
$$J(\theta) = \left[-\frac{1}{m} \sum_{i=1}^m y^{(i)} \log(h_{\theta}(x^{(i)})) + (1-y^{(i)}) \log(1-h_{\theta}(x^{(i)})) \right] + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$

* Neural Network Representation

↳ Model Representation : logistic unit



⇒ Neural networks :-



⇒ Forward Propagation :

$$a^{(1)} = x ; \quad z^2 = \theta^{(1)} x$$

$$z^2 = \theta^{(1)} a^{(1)} , \quad a^{(2)} = g(z^2) \Rightarrow a^{(2)} = \theta^{(1)} a^{(1)}$$

$$a_1^{(2)} = g(\theta_{10}^{(1)} x_0 + \theta_{11}^{(1)} x_1 + \theta_{12}^{(1)} x_2 + \theta_{13}^{(1)} x_3)$$

+ so on

$$h_\theta(x) = a_1^{(3)} = g(\theta_{10}^{(2)} a_0^{(2)} + \theta_{11}^{(2)} a_1^{(2)} + \theta_{12}^{(2)} a_2^{(2)} + \theta_{13}^{(2)} a_3^{(2)})$$

⇒ Back propagation : minimizing error or gradient computing :-

$$\delta_j^{(1)} = \text{error of node } j \text{ in layer 1}$$

$$\text{eg. } \delta_j^{(3)} = a_j^{(3)} - y_j \quad \text{or} \quad \delta^{(3)} = a^{(3)} - y$$

$$\delta^{(2)} = (\theta^{(2)})^T \delta^{(3)} + g'(z^{(2)})$$