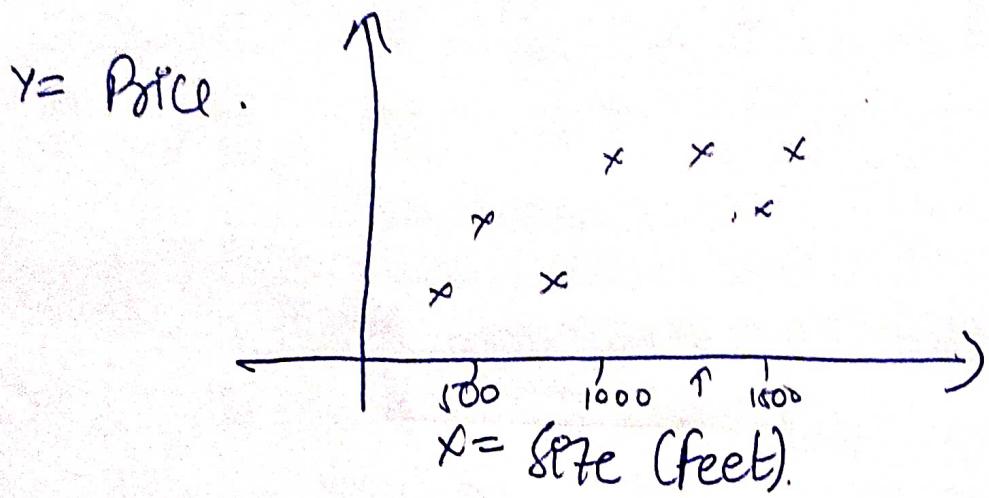


Machine Learning \rightarrow field of study that gives computers the ability to learn without being explicitly programmed.

Set Supervised learning

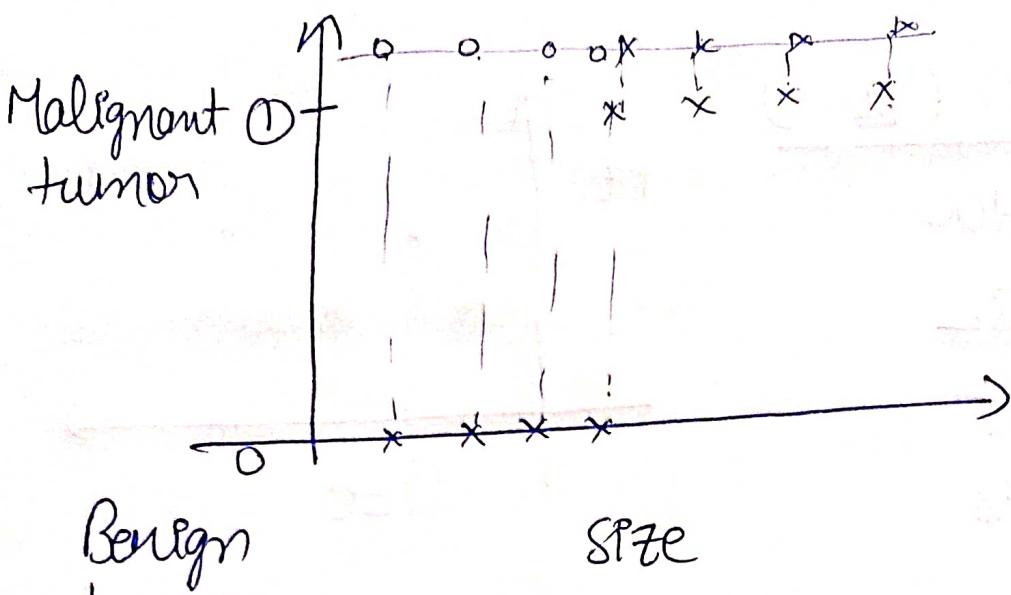


(x, y) , we need to learn mapping $x \rightarrow y$

This is example of Regression Problem.

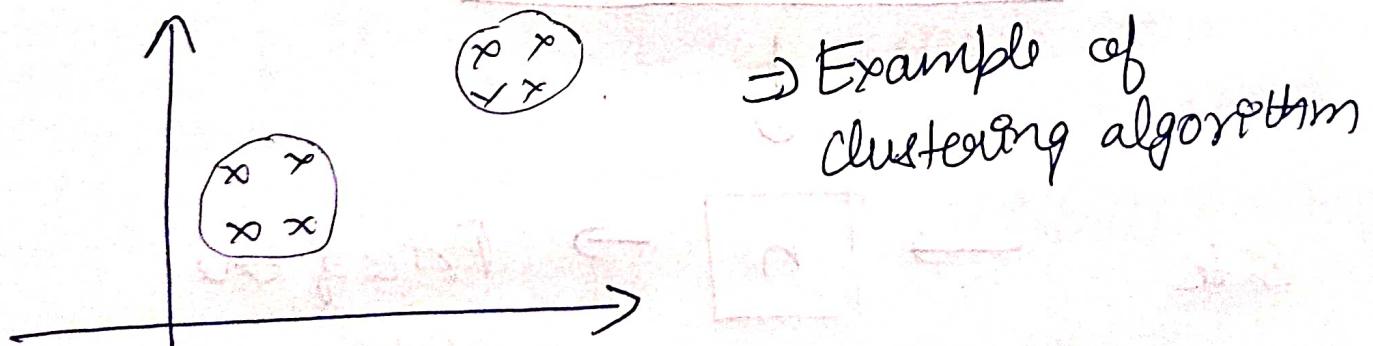
when there are discrete K number of outputs,
it is known as classification problem

Example of classification problem :-



In real world we would be having $n \rightarrow$ number of inputs to predict the feature.

Unsupervised learning

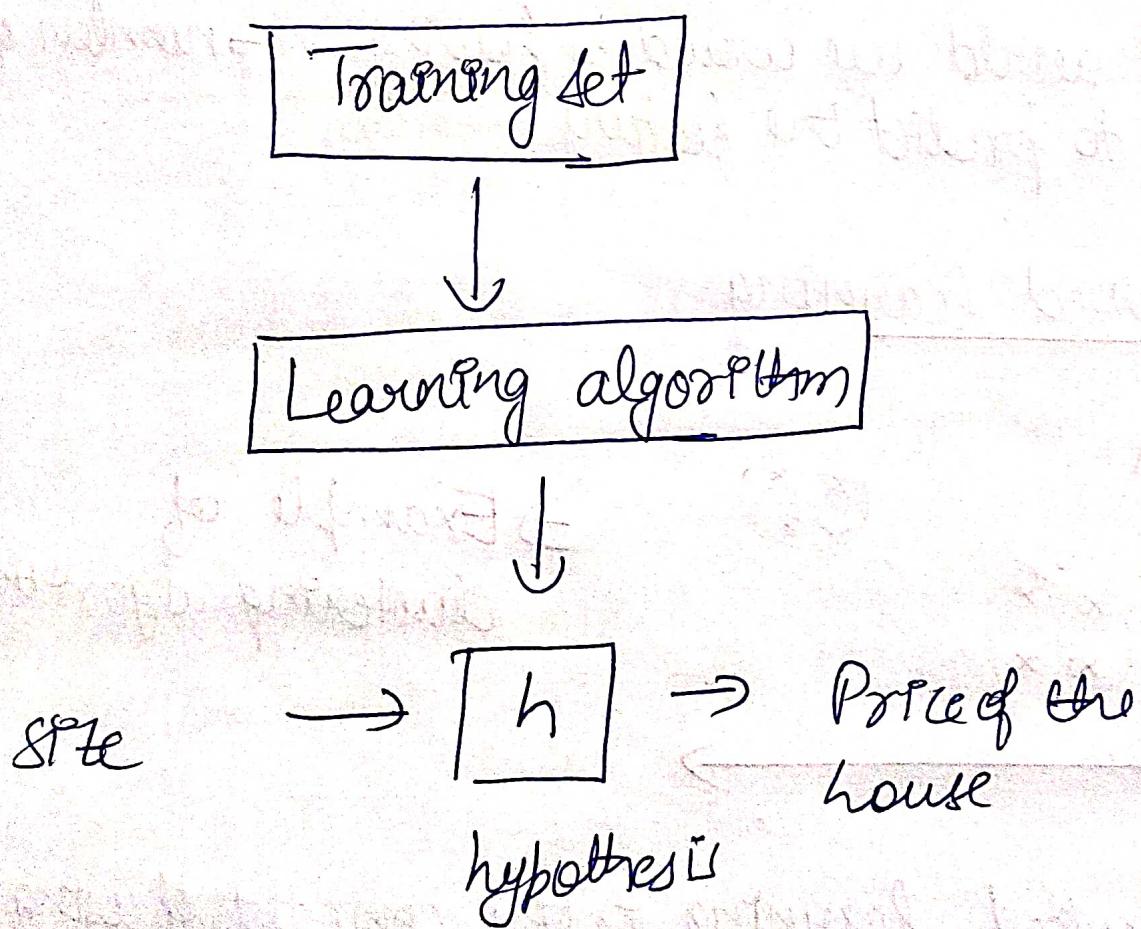
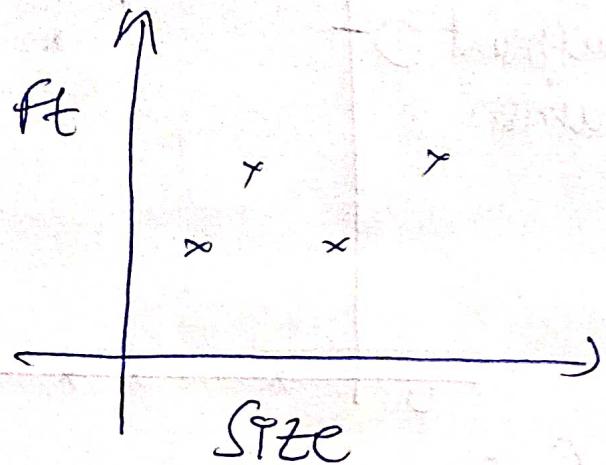


⇒ Example of clustering algorithm

Unsupervised learning is the concept of using unlabelled data

Linear Regression and Gradient descent

Size (ft)	Price (\$1000)
2104	400
1416	232
1534	315
852	718



① How do we represent hypothesis?

$$h(x) = \theta_0 + \theta_1 x$$

If we have more ~~BA~~ information, then

$$\cancel{h(x) = \theta_0 + \theta_1 x}$$

$$x_1 = \text{size}$$

$$x_2 = \text{no. of bedrooms}$$

$$h(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2$$

$$h(x) = \sum_{j=0}^n \theta_j x_j$$

$$x_0 = 1$$

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \end{bmatrix} \quad x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix}$$

θ = "parameters"
 m = training example
or rows

$$(x^{(i)}, y^{(i)}) = \text{training example}$$

x = input / features

y = output / target variable.

(x, y) = training example

We need to choose θ such that $h(x) \approx y$ for training examples.

$$h_{\theta}(x) = h(x)$$

$h_{\theta}(x) \Rightarrow$ hypothesis depends on both parameters and input feature x

In fact in linear regression we need to

$$\text{minimise}_{\theta} \sum_{i=1}^m (h_{\theta}(x) - y)^2$$

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

Gradient Decent

Start with θ (say $\theta = \vec{\theta}$).

Keep changing θ to reduce $J(\theta)$.

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta) \quad (j=0, 1, 2, \dots, n)$$

\rightarrow Learning rate

$$\boxed{\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta) \quad (j=0, 1, 2, \dots, n)}$$

~~$\frac{\partial}{\partial \theta_j} J(\theta)$~~

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

$$= 2 \frac{1}{m} (h_\theta(x^{(1)}) - y^{(1)}) \cdot \frac{\partial}{\partial \theta_j} (h_\theta(x^{(1)}) - y^{(1)})$$

$$= (h_\theta(x^{(1)}) - y^{(1)}) \cdot \frac{\partial}{\partial \theta_j} (h_\theta(x^{(1)}) - y^{(1)})$$

$$= [h_\theta(x^{(1)}) - y^{(1)}] \cdot x^{(1)}$$

$$\boxed{\frac{\partial}{\partial \theta_j} J(\theta) = (h_\theta(x) - y) \cdot (x_j)}$$

$$\theta_j := \theta_j - \alpha (h_\theta(x) - y) \cdot x_j$$

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

Repeat until convergence

$J(\theta)$ does not have local optima, the only local optima is the global optimum.

In batch gradient descent:

Repeat {

for $i = 1$ to m {

~~θ_j~~ :

$$\theta_j := \theta_j - \alpha (h_\theta(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)}$$

}

}

$$\nabla_{\theta} J(\theta) = \begin{bmatrix} \frac{\partial J}{\partial \theta_0} \\ \frac{\partial J}{\partial \theta_1} \\ \frac{\partial J}{\partial \theta_2} \end{bmatrix}$$

↳ Derivative

$$A \in \mathbb{R}^{2 \times 2} \quad A = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix}$$

$$f: \mathbb{R}^{2 \times 2} \rightarrow \mathbb{R}$$

$$f(A) := A_{11} + A_{22}$$

$$f\left(\begin{bmatrix} 5 & 6 \\ 7 & 8 \end{bmatrix}\right) = 5 + 8$$

$$\nabla_A f(A) = \begin{bmatrix} \frac{\partial f}{\partial A_{11}} & \frac{\partial f}{\partial A_{12}} \\ \frac{\partial f}{\partial A_{21}} & \frac{\partial f}{\partial A_{22}} \end{bmatrix} = \begin{bmatrix} 1 & 2A_{12} \\ 0 & 0 \end{bmatrix}$$

if A is a square matrix

$\text{tr}(A) = \text{sum of diagonal entries}$



trace

$$\textcircled{1} \quad \text{tr}(A^T) = \text{tr}(A)$$

$$f(A) = \text{tr}(AB) \Rightarrow \nabla_A f(A) = B^T$$

$$\textcircled{2} \quad \text{tr}(AB) = \text{tr}(BA)$$

$$\textcircled{3} \quad \text{tr}(ABC) = \text{tr}(CAB)$$

$$\textcircled{4} \quad \nabla_A f(A) = \nabla_A \text{tr}(AA^T) = CA + C^T A$$

$$J(\theta) = \frac{1}{2} \sum_{q=1}^n (h(x^{(q)}) - y^{(q)})^2$$

$$\begin{aligned} x\theta &= \begin{bmatrix} -(x^{(1)})^T \\ -(x^{(2)})^T \\ \vdots \\ -(x^{(m)})^T \end{bmatrix} \begin{bmatrix} 0_0 \\ 0_1 \\ \vdots \\ 0_m \end{bmatrix} \\ &= \begin{bmatrix} x^{(1)\top} \theta \\ x^{(2)\top} \theta \\ \vdots \\ x^{(m)\top} \theta \end{bmatrix} \end{aligned}$$

$$\vec{y} = \begin{bmatrix} y^1 \\ y^2 \\ \vdots \\ y^m \end{bmatrix}$$

$$J(\theta) = \frac{1}{2} (x\theta - y)^T (x\theta - y)$$

$$\nabla_{\theta} J(\theta) = D \theta \frac{1}{2} (x\theta - y)^T (x\theta - y)$$

$$= \frac{1}{2} \nabla_{\theta} ((x\theta)^T x^T - y^T) (x\theta - y)$$

$$= \cancel{\frac{1}{2}} \cancel{x^T} \cancel{x\theta} [= \frac{1}{2} [x^T x\theta + x^T x\theta - x^T y - x^T y]$$

$$= x^T x\theta - x^T y \text{ setting } \cancel{\theta} = 0$$

$$x^T x\theta = x^T y$$

$$\theta = (x^T x)^{-1} x^T y$$

Multiple Linear Regression

$$\hat{y} = b_0 + b_1 x_1 + b_2 x_2 + \dots + b_n x_n$$

Assumptions of Linear Regression

- ① Linearity (Linear relationship b/w Y and each x_i)
- ② Homoscedasticity (Equal Variance)
- ③ Multivariate Normality (Normality of error dist.)
- ④ Independence of observation. Includes no autocorrelation
- ⑤ Lack of Multicollinearity :- $x_1 \sim x_2$
- ⑥ Outlier check.

5 methods of building models

- ① All-in
 - ② Backward elimination. 2
 - ③ Forward Selection. 3
 - ④ Bidirectional elimination. 4
 - ⑤ Scope Comparison.
- Stepwise Regression.

Backward elimination

- ① Select significance level.
- ② Fit full model with all possible predictors
- ③ If $P > S.L$ go to 4 else FIN
- ④ Remove predictor.
- ⑤ Fit to model without this variable

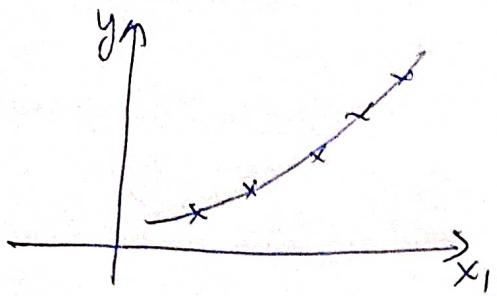
Forward selection

- ① Select significance level.
- ② Fit all simple regression ~~your~~. Selection with lowest p-value.
- ③ Keep this variable and fit all possible models with one extra predictor added to one already have
- ④ Consider the predictor with lowest P value
 - If PLSL go to Step 3 else FIN
↓
Keep the previous model.

Bi-directional Elimination

- ① Select a significance level to enter and to stay in the model. e.g. SLEnter = 0.05 and SLSay = 0.05
- ② Perform next step of forward selection (PLSL Enter).
- ③ Perform all steps of Backward elimination (old values must have PL(Say))
- ④ No variable can enter or exit.

Polynomial Regression



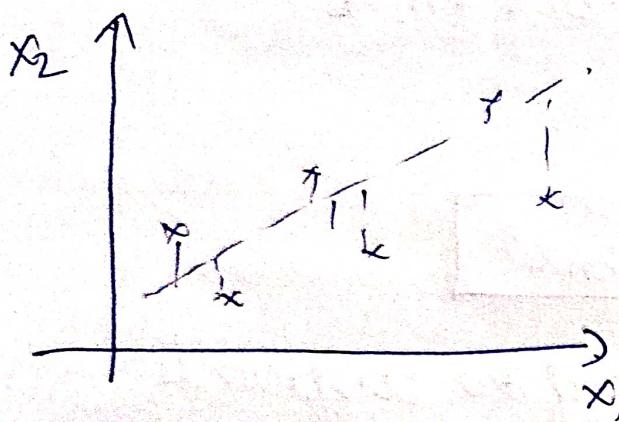
$$y = b_0 + b_1 x_1 + b_2 x_1^2$$

$$y = b_0 + b_1 x_1 + b_2 x_1^2 + \dots + b_{n-1} x_1^{n-1} b_n x_1^n$$

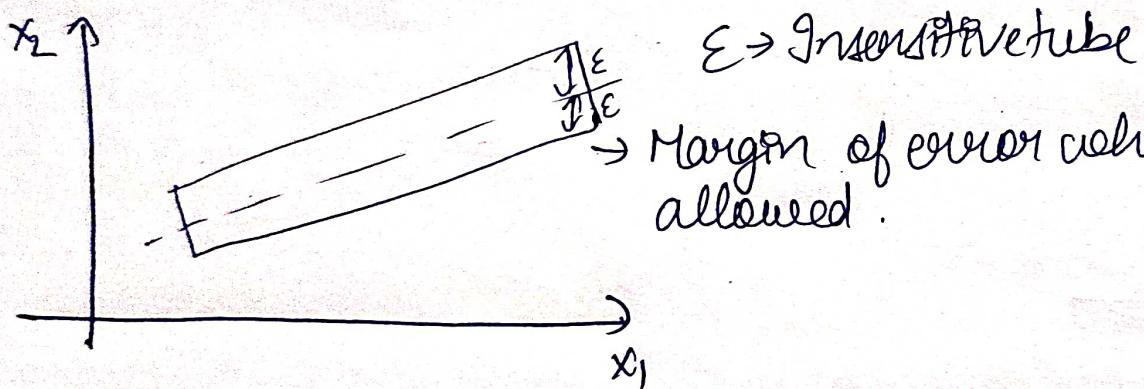
It is still called linear \Rightarrow becoz of the coefficients are linear.

~~SVR~~ Intro

SVR \rightarrow Support Vector Regression



Ordinary Least squares $\sum (y - \hat{y})^2 \rightarrow \min.$



$\epsilon \rightarrow$ Insensitive tube

\rightarrow Margin of error which is allowed.

If error is above insensitive tube then its E_i and if it is below insensitive tube then its E_i^*

$$\boxed{\frac{1}{2} \|\omega\|^2 + C \sum_{i=1}^m (E_i + E_i^*) \rightarrow \min}$$