

Machine Learning

A computer program is said to learn from Experience E with respect to some class of task T and performance P , if its performance in task T , as measured by P , improves with experience E .

Type of machine learning

Based on the methods and way of learning, Machine learning is divided into four types

1. Supervised ML
2. Unsupervised ML
3. Semi-supervised ML
4. Reinforcement learning.

ML

Supervised ML

- House price prediction
- Medical Imaging
- Spam Detection
- Email filtering
- market trends
- weather prediction.
- Fraud detection
- speech recognition

Unsupervised ML

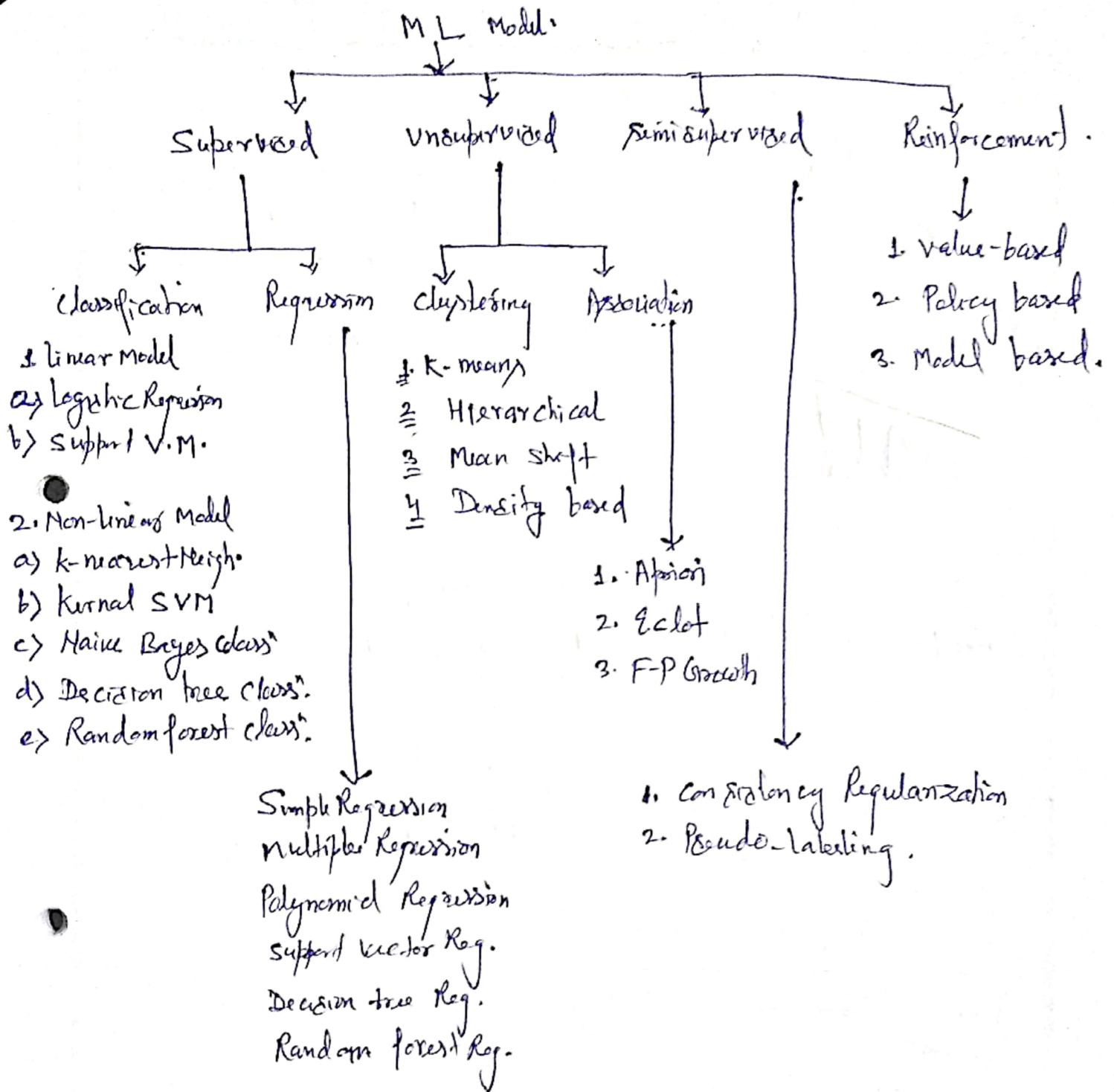
- Customer Segmentation
- Market basket analysis
- web usage mining
- continuous production
- Network analysis

Semi-supervised ML

- Text classification
- lane finding on GPS data.

Reinforcement Learning

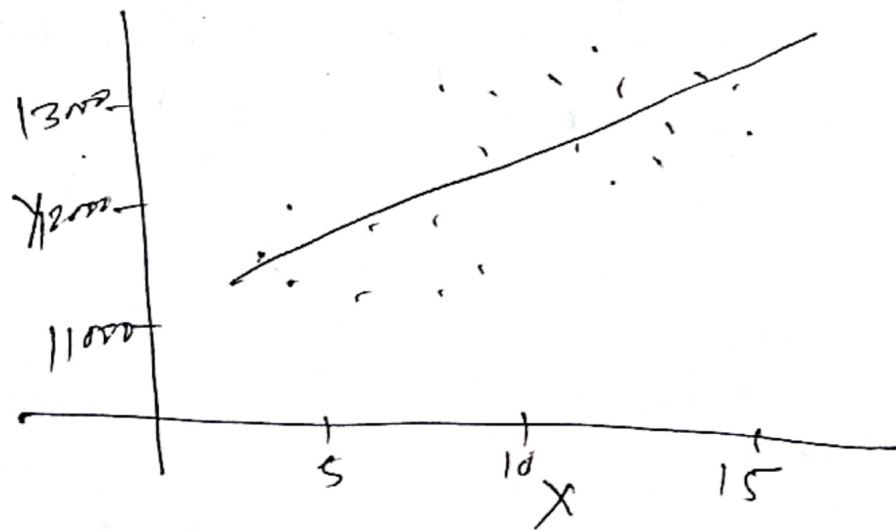
- optimized marketing
- Donorless cars
- Game theory
- OR
- Information theory
- Robotics
- Text mining



Linear models for regression: A Regression model is used for finding out the relationship b/w variables and forecasting. Different regression models differ based on

- the kind of relationship b/w dependent and independent variables they are considering
- The number of independent variables getting used.

The linear regression performs the task to predict dependent variable values (y) based on independent variable (x).



In figure X (input) is the work experience and (output) is the salary of a person. This regression line is the best fit line for our model.

Linear Regression: least square method

$$x = (x_1, \dots, x_n)$$

$$\bar{x} = \frac{\sum x}{n}$$

$$y = (y_1, \dots, y_n)$$

$$\bar{y} = \frac{\sum y}{n}$$

Fit a linear line of the form.

$$y = mx + c \quad \text{--- (1)}$$

$$\sum y = m \sum x + c \sum 1$$

$$\sum y = m \sum x + cn \quad \text{--- (2)}$$

$$xy = mx^2 + cx$$

$$\sum xy = m \sum x^2 + c \sum x \quad \text{--- (3)}$$

from (2) & (3)

$$m = \frac{n \sum xy - \sum x \sum y}{n \sum x^2 - (\sum x)^2} = \frac{\sum (x - \bar{x})(y - \bar{y})}{\sum (x - \bar{x})^2}$$

$$c = \bar{y} - m\bar{x}$$

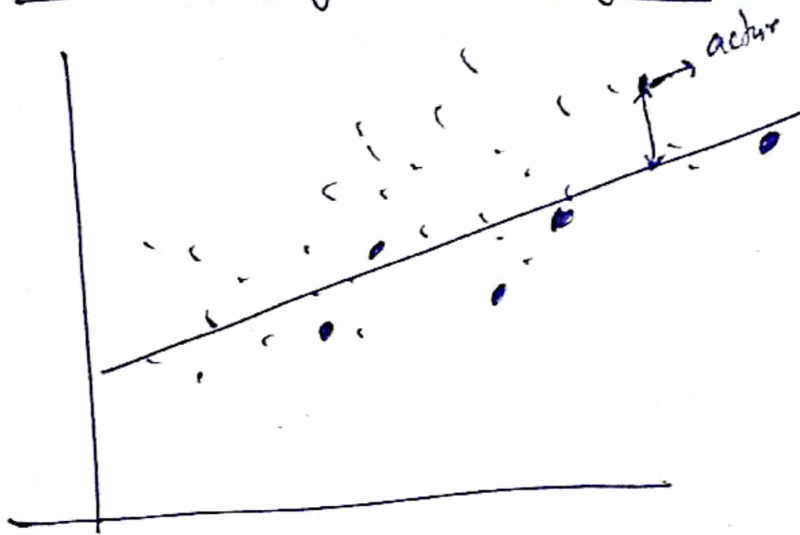
1)

$\Rightarrow \hat{y} = mx + c$ 4 Error = $y - \hat{y}$

ex 1		ex 2	
x	2, 6, 5, 7	x	0, 1, 2, 3, 4, 5, 6, 7, 8, 9
y	3, 10, 4, 13	y	1, 3, 2, 5, 7, 8, 8, 9, 10, 12
\bar{x}		$\bar{x} - \bar{x}$	$y - \bar{y}$
\bar{y}			$(x - \bar{x})(y - \bar{y})$
2	3	-3	-4.5
6	10	1	2.5
5	4	0	-3.5
7	13	2	5.5
$\sum x = 20$	$\sum y = 30$	\sum	\sum
		0	1
			26.0

$$\bar{x} = \frac{20}{4} = 5, \quad \bar{y} = \frac{30}{4} = 7.5, \quad m = \frac{26}{0}$$

Efficiency checking.



Error: measure of how far the data is from the fitted regression line.

(i) mean absolute error = $\frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$
(MAE)

(ii) mean squared error = $\frac{1}{n} \sum (y_i - \hat{y}_i)^2$
(MSE)

(iii) Root mean squared error = $\sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$
(RMSE)

(iv) ~~R~~ relative absolute error. = $\frac{\sum_{i=1}^n |y_i - \hat{y}_i|}{\sum_{i=1}^n |y_i - \bar{y}|}$
(RAE)

(v) Relative squared error = $\frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$
(RSE)

(vi) ~~R-squared~~ error.

$$R^2 = 1 - RSE$$

(vii) ~~R-squared~~ adjusted = $1 - \frac{(1 - R^2)(n - 1)}{n - p - 1}$

$n \rightarrow$ sample size, $p \rightarrow$ No of features.

Gradient Descent Approach

Hypothesis funⁿ for Linear Regression.

$$Y = \theta_0 + \theta_1 \cdot X$$

where

X : input training data | θ_1 : intercept
 Y : labels to data. | θ_2 : coefficient of x .

when training the model - it fits the best line to predict the value of y for a given value of x . The model get best regression fit line by find the best θ_1 and θ_2 values.

Finding ~~best~~ ~~fit~~ values of θ_1 & θ_2 for best fit line

cost function: The model aims to predict y value & the error diff. b/w predicted value and true value is minimum. so it's important to update the θ_1 & θ_2 to reach the best value that minimize the error b/w predicted value (pred) and true value (y)

i.e minimize
$$\frac{1}{n} \sum_{i=1}^n (\text{pred}_i - y_i)^2$$

i.e
$$J = \frac{1}{n} \sum_{i=1}^n (\text{pred}_i - y_i)^2$$

Cost funⁿ of L.R is RMSE b/w predicted & true value.

* cost funⁿ of L.R. is RMSE b/w predicted & true value.

Gradient Descent: • cost funⁿ

$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m [h_{\theta}(x_i) - y_i]^2$$

Gradient Descent:

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

$$\stackrel{\text{Now}}{=} \frac{\partial}{\partial \theta_j} J_{\theta} = \frac{\partial}{\partial \theta_j} \frac{1}{2m} \sum_{i=1}^m [h_{\theta}(x_i) - y_i]^2$$

$$= \frac{1}{m} \sum_{i=1}^m [h_{\theta}(x_i) - y_i] \cdot \frac{\partial}{\partial \theta_j} (h_{\theta}(x_i))$$

$$\frac{\partial}{\partial \theta_j} J_{\theta} = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x_i) - y_i) x_j$$

$$\stackrel{\text{Now}}{=} \theta_j = \theta_j - \frac{\alpha}{m} \sum_{i=1}^m [(h_{\theta}(x_i) - y_i) x_j]$$

θ_j : weight of hypothesis.

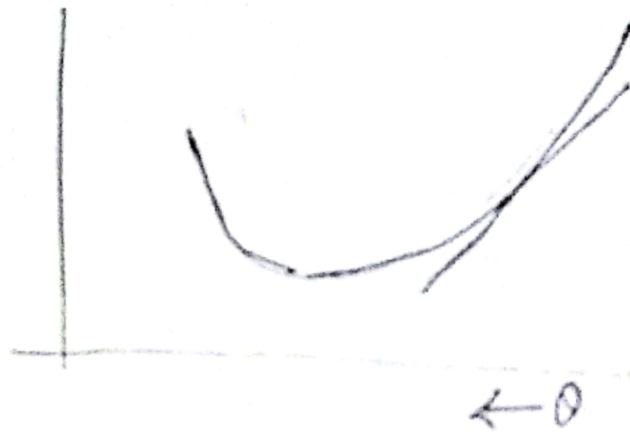
$h_{\theta}(x_i)$: predicted value of i^{th} input.

J_{θ} : feature index

α : learning rate. where $x_j = \begin{cases} 1 & \text{for } j=0 \\ x_i & \text{for } j=1 \end{cases}$

(i) If slope is +ve

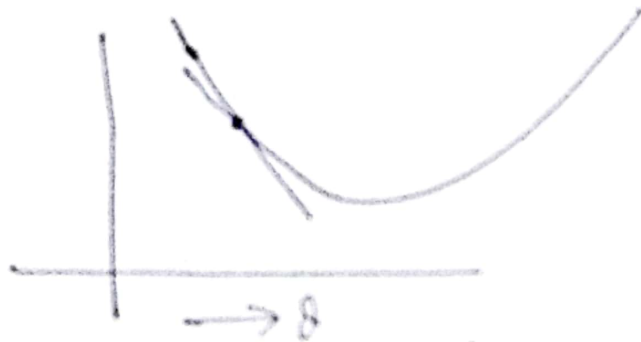
θ_0 decreases



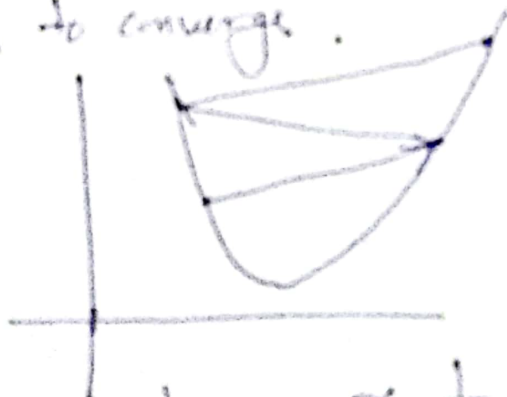
Note (i) for linear regression cost the function graph is always convex shaped.

(ii) is a feature index.

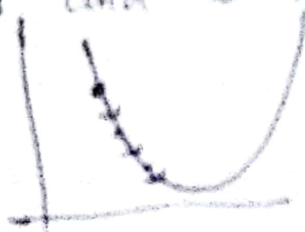
(ii) If slope is -ve: θ_0 increases.



(iii) If we choose α to be very large, Gradient Descent can overshoot the minimum. It may fail to converge.



(iv) If we choose α to be very small, Gradient descent will take small steps to reach local minimum and will take longer time to reach mini.



<u>ex</u>	x	y
	2	3
	6	10
	5	4
	7	13

$$h_{\theta} = \theta_0 + \theta_1 x$$

$$= \begin{bmatrix} \theta_0 & \theta_1 \end{bmatrix} \begin{bmatrix} 1 \\ x_1 \end{bmatrix}$$

$$= \begin{bmatrix} \theta_0 & \theta_1 \end{bmatrix} \begin{bmatrix} 1 & 1 & 1 & 1 \\ x_1 & x_2 & x_3 & x_4 \end{bmatrix}$$

Iteration 1

$$\theta_0 = 0, \quad \theta_1 = 0.$$

$$h_{\theta} = \begin{bmatrix} \theta_0 & \theta_1 \end{bmatrix} \begin{bmatrix} x_0 & x_0 & x_0 & x_0 \\ x_1 & x_2 & x_3 & x_4 \end{bmatrix}$$

$$= \begin{bmatrix} 0 & 0 \end{bmatrix} \begin{bmatrix} 1 & 1 & 1 & 1 \\ 2 & 6 & 5 & 7 \end{bmatrix}$$

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m [h_{\theta}(x_i) - y_i]^2$$

$$= 16.75$$

Here $\theta_j = 0$

$$\theta_j = \theta_j - \frac{\alpha}{m} \sum_{i=1}^m (h_{\theta}(x_i) - y_i) x_i$$

$$= 0 - \frac{0.01}{4} [(0-3) + (0-10) + (0-4) + (0-13)]$$

$$= 0.0075$$

$$\theta_j = \theta_j - \frac{1}{m} \sum_{i=1}^m (h(x_i) - y_i) x_i$$

$$= 0 - \frac{0.001}{4} [(0.3) \times 2 + (0.10) \times 6 + (0.4) \times 5 + (0.13) \times 7]$$

$$= -\frac{0.001}{4} [-6 - 60 - 20 - 91] = \cancel{0.02657}, 0.04425$$

⇒ we have

$$\theta_0 = \frac{\cancel{0.005}}{0.0075} \quad \& \quad \theta_1 = \frac{\cancel{0.02657}}{0.04425}$$

Iteration 2 :

predicted values after iteration 1 with linear regression hypothesis.

$$h\theta = [\theta_0 \quad \theta_1] \begin{bmatrix} x_0 & x_0 & x_0 & x_0 \\ x_1 & x_2 & x_3 & x_4 \end{bmatrix}$$

$$= \begin{bmatrix} \cancel{0.005} & \cancel{0.02657} \\ 0.0075 & 0.04425 \end{bmatrix} \begin{bmatrix} 1 & 1 & 1 & 1 \\ 2 & 6 & 5 & 7 \end{bmatrix}$$

$$= \cancel{[0.057 \quad 0.161 \quad 0.135 \quad 0.187]}$$

$$= [0.096 \quad 0.273 \quad 0.22875 \quad 0.31725]$$

Now, find θ_0 & θ_1 . (same as in previous step).

Q. 2

x	0	1	2	3	4	5	6	7	8	9
y	1	2	3	5	7	11	13	16	19	22

Multiple Linear Regression: Multiple Linear Regression

used to model two or more features and response by fitting a linear line to observed data.

Now, consider a data set with p features (Independent variables) and one response (dependent variable).

Also the dataset (training dataset) contains n observations. Then

X (feature matrix) = ~~Matrix~~ $[n \times p]$ where x_{ij} denotes the value of j th feature for i th observation.

y (response vector) = $[y_i]_n$

y_i is response of i th observation.

→ The regression line for p features is represented as

$$h(x_i) = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip}$$

$$h(x_i) = \beta_0 + \sum_{j=1}^p \beta_j x_{ij}$$

Therefore, we can write

$$y_i = h(x_i) + \epsilon_i \Rightarrow \epsilon_i = y_i - h(x_i)$$

n matrix for model can written as

$$y = X\beta + \epsilon$$

where

$$y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}, \quad \beta = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \vdots \\ \beta_p \end{bmatrix}, \quad \epsilon = \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix}$$

$$X = \begin{bmatrix} 1 & x_{11} & \dots & x_{1p} \\ 1 & x_{21} & \dots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & \dots & x_{np} \end{bmatrix}, \quad x_i = \begin{bmatrix} 1 \\ x_{i1} \\ x_{i2} \\ \vdots \\ x_{ip} \end{bmatrix}$$

Here we determine β for which residual error ϵ is minimum, using least square method.
And we get

$$\hat{\beta} = (X'X)^{-1} X'y$$

and Residual vector is estimated as

$$\hat{y} = X\hat{\beta}$$

Polynomial Regression:

General model for polynomial Regression.

$$y = \beta_0 + \beta_1 x + \beta_2 x^2 + \dots + \beta_m x^m$$

Classification :

Classification is a supervised learning approach which can be thought of as used for categorizing or classifying some unknown items into a discrete set of classes. Classification ~~at~~ attempts to learn the relationship between a set of feature variables and a target variable of interest.

The target attribute in classification is a categorical variable with discrete values.

Type of attributes:

1. Binary: Possesses only two values ex. True or False
 - * Symmetric: Both values are equally important in all aspects.
 - * Asymmetric: when both values may not be important.
2. Nominal: when more than two outcomes are possible.
ex classification for different colors like green, red, black
 - * Ordinal: values that must have some meaningful order.
ex grades: A, B, C, D.
 - * Continuous: may have infinite no of values.
ex weight
 - * Discrete: finite no of values:
ex colors:

Type of Classification algorithms:

Classification algorithms are categorized into two types:

1. Discriminative: It is a very basic classification algorithm used to determine one class row of data:

Ex Logistic Regression

2. Generative: It models the distribution of individual classes and tries to learn the model that generates the data behind the scenes by estimating assumptions and distribution of the model.

Ex Naive Bayes classifier.

Various classification models

1. Logistic Regression
2. Linear model
3. SVM
4. Nonlinear model
5. K-nearest Neighbour
6. Naive Bayes classifier.

Type of Classification algorithms:

Classification algorithms are categorized into two types:

1. Discriminative: It is a very basic classification algorithm used to determine one class row of data.
Ex Logistic Regression

2. Generative: It models the distribution of individual classes and tries to learn the model that generate the data behind the scenes by estimating assumptions and distribution of the model.
Ex Naive Bayes classifier.

Various classification models

1. Logistic Regression
2. Linear model
3. SVM
4. Nonlinear model
5. K-nearest Neighbour
6. Naive Bayes classifier.