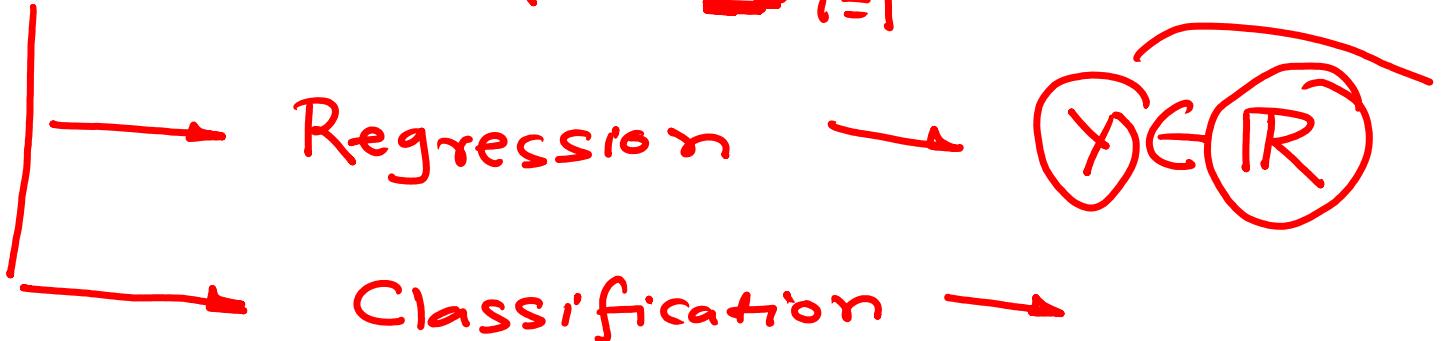


Supervised

$$\mathcal{D} = \{x_i, y_i\}_{i=1}^N$$



$$y \in \mathbb{R}$$

Fixed number of  
classes

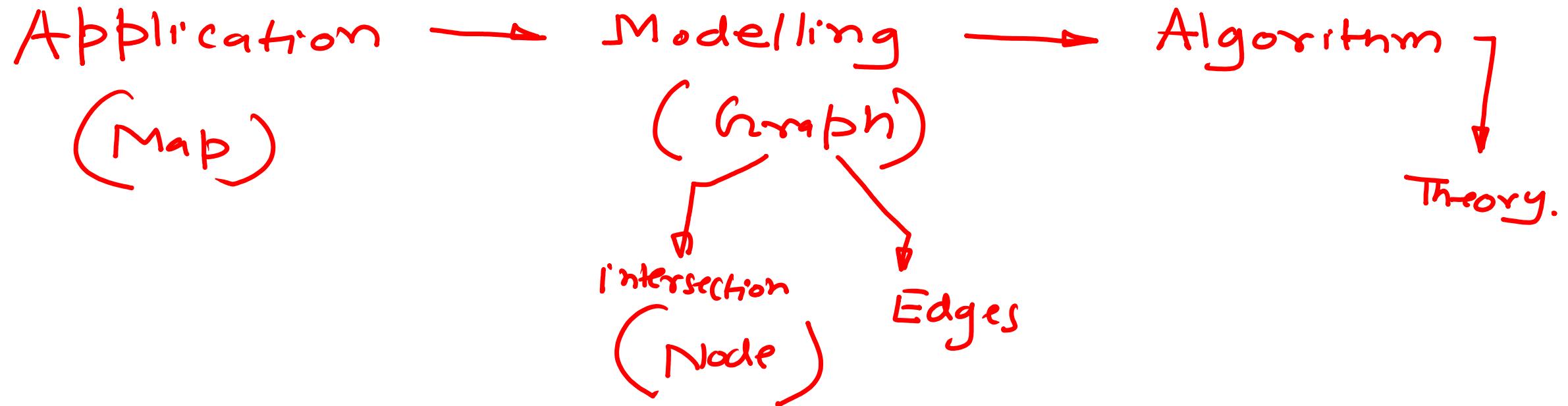
Obj :- to ~~pose~~ predict  
~~which~~  $x$  belongs to  
which class.

Obj :- we want to learn the relation  
between  $x$  and  $y$ .

$$y = f(x)$$

→ we want to find  $f$ .

\* 4 steps for applying AI for any  
real-life problem.



\* Linear Regression :-

Suppose  $\mathcal{D} = \{x_i, y_i\}_{i=1}^N$

$x_1$	$x_2$	$x_3$	$y$
$x_1^{(1)}$	$x_2^{(1)}$	$x_3^{(1)}$	$y^{(1)}$
$x_1^{(2)}$	$x_2^{(2)}$	$x_3^{(2)}$	$y^{(2)}$
:	:	:	:
$x_1^{(N)}$	$x_2^{(N)}$	$x_3^{(N)}$	$y^{(N)}$

we want to  
learn

$$Y = f(x_1, x_2, x_3)$$

\* In linear  
regression,

$$\begin{aligned}
 Y = & a_0 + a_1 x_1 + a_2 x_2 + a_3 x_3 + a_4 x_1 x_2 - (1) \\
 & + a_5 x_1 x_3 + a_6 x_2 x_3 + a_7 x_4^2 + a_8 x_2^2 + \\
 & a_9 x_3^2
 \end{aligned}$$

In a generic form,

$$Y = a_0 \phi_0(x) + a_1 \phi_1(x) + \cdots + a_k \phi_k(x) \quad -(2)$$

Comparing (1) and (2)

$$\phi_0(x) = 1 \qquad \phi_2(x) = x^2$$

$$\phi_1(x) = x_1$$

\* Eq. (2) is the generic expression of linear regression

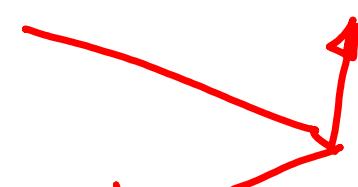
$$Y = \sum_{i=0}^K a_i \phi_i(X)$$

number of basis func.  
basis function  
unknown coefficients.

A For using linear regression in practice :-

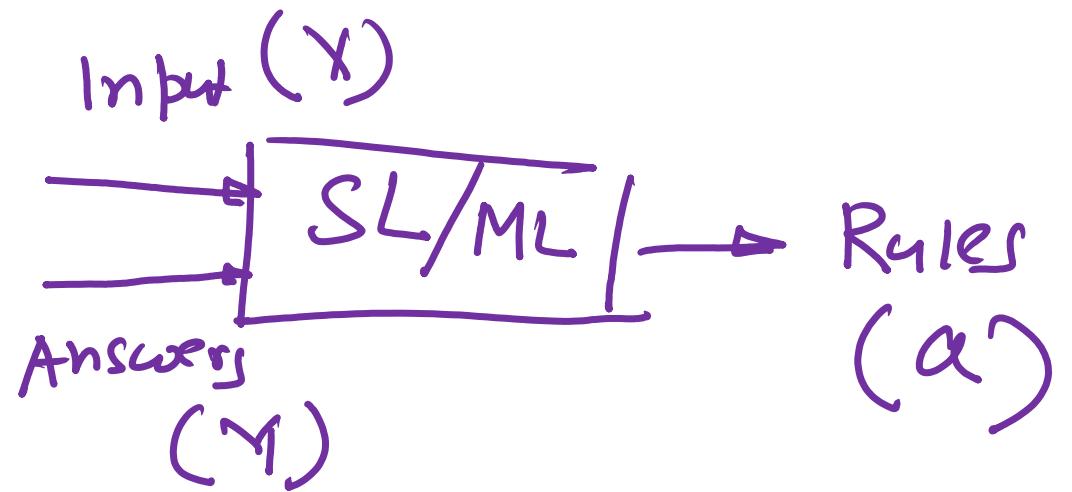
Hyperparameter

(1) we need to fix K



(2) we need to select  $\phi_i$  &  $a_i$

\* Learning in Linear regression means  
computing the unknown parameters from  
data.



\* Training :- (1)  $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$   
(2) Hyper parameters.

\* For a "good" value of  $a_{\alpha}$

$\tilde{y} = \sum_i a_i \phi_i(x)$  and the measured  
 $y$  should ~~not~~ be close.

↳ we need to define a  
error metric that measures  
the distance between  $y, \tilde{y}$   
 $E(a)$

\* We find the best  $a$  by minimizing the error

$$a^* = \arg \min_a E(a)$$

\* Some examples of error metric are :-

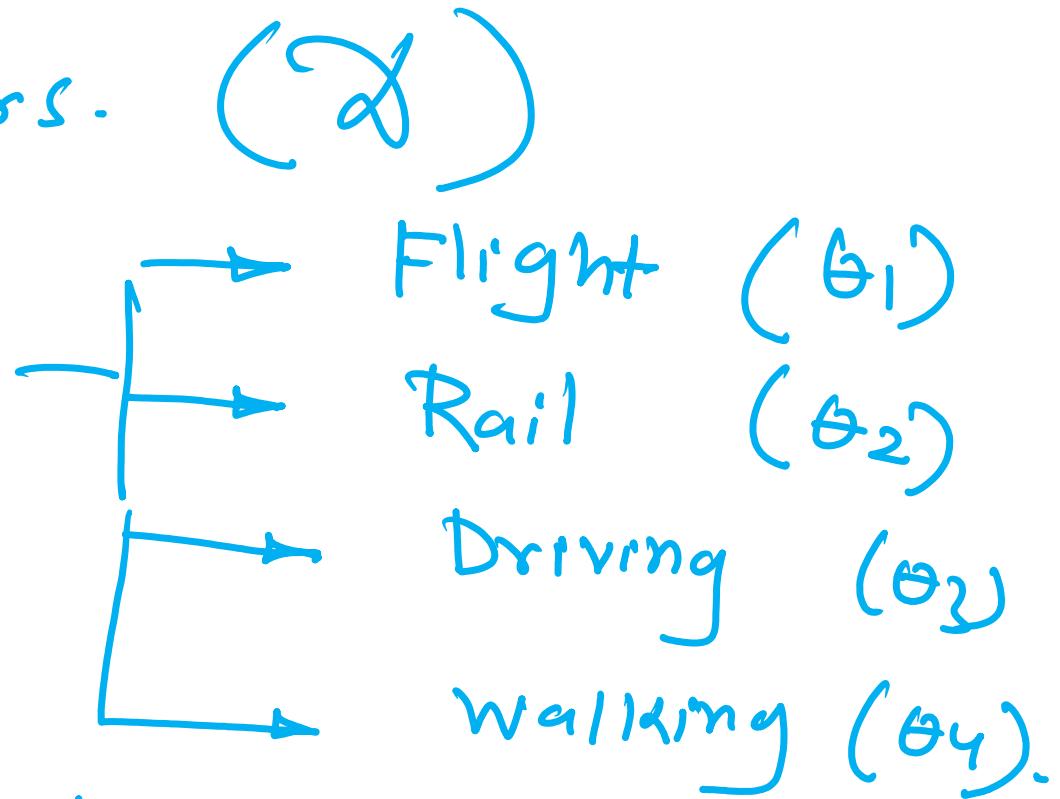
(1) Mean - squared error (MSE)

(2) Mean - absolute error (MAE)

\* Travelling from Delhi to Kolkata.

\* It took 24 hours. ( $\alpha$ )

\* Modes of travel



$$P(\alpha | \theta_1) = 0.1$$

$$P(\alpha | \theta_2) = 0.9$$

$$P(\alpha | \theta_3) = 0$$

$$P(\alpha | \theta_4) = 0$$

By generalizing,

$$\theta^* = \arg \max P(\mathbf{z} | \theta)$$

Maximum Likelihood  
Estimation.

- \* For linear regression, we have closed form solution.

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

$$\frac{dE(a)}{da} = 0$$

System of equations

\* But often closed form solutions are not available.

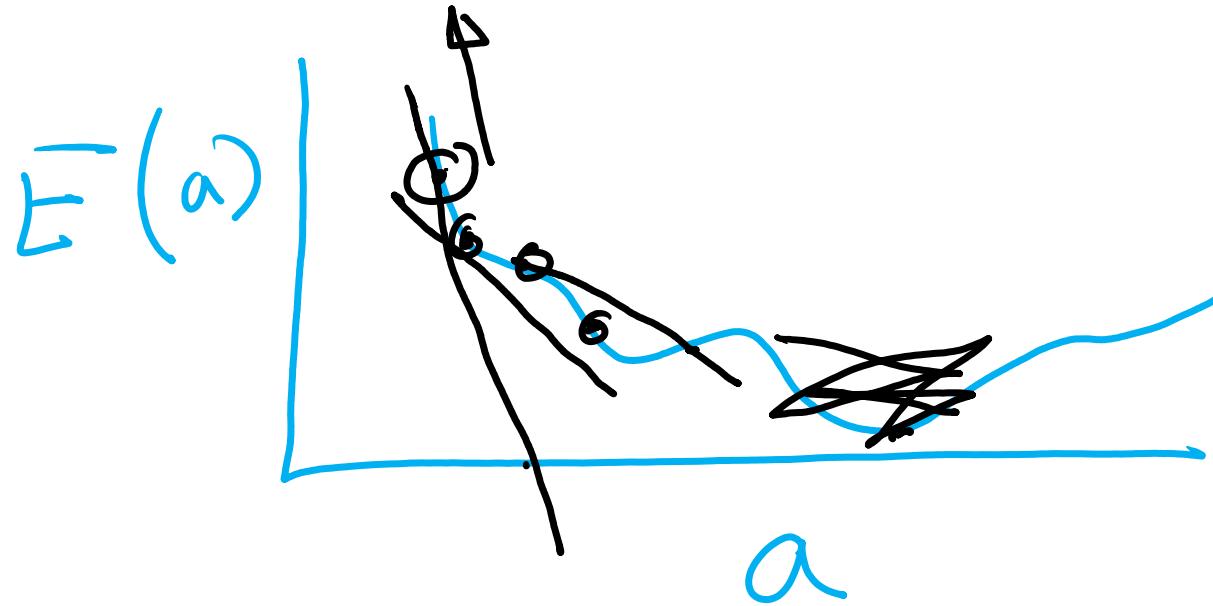
↳ we use numerical optimization

\* One such numerical optimization

algorithm is

$$a^{(t+1)} = a^{(t)} - \eta \underbrace{\nabla_a E(a)}_{\text{Gradient}} \quad \leftarrow \text{Gradient}$$

↳ Learning rate



Small  $\eta \leftarrow$  lots  
of  
iterations

\* Gradient  $\leftarrow$  direction  
of movement

large  $\eta \leftarrow$   
oscillation

η  $\leftarrow$  how much  
move

\* Special case :- Linear regression with least squared loss function has a closed form solution.

$$\hat{y} = \sum_{i=0}^K a_i \phi_i(x)$$

$$[y]_{N \times 1} = \phi [a]_{(K+1) \times 1} [x]_{(K+1) \times 1}$$

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix} = \begin{bmatrix} \phi_0(x') \\ \phi_1(x') \\ \vdots \\ \phi_K(x') \end{bmatrix}$$

$$\begin{bmatrix} \phi_0(x^{(1)}) & \phi_1(x^{(1)}) & \phi_2(x^{(1)}) & \dots & \phi_K(x^{(1)}) \\ \phi_0(x^{(2)}) & \phi_1(x^{(2)}) & \phi_2(x^{(2)}) & \dots & \phi_K(x^{(2)}) \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \phi_0(x^{(N)}) & \phi_1(x^{(N)}) & \phi_2(x^{(N)}) & \dots & \phi_K(x^{(N)}) \end{bmatrix} N \times (K+1)$$

$$Y = \underline{\phi} a$$

$$a = (\phi^T \phi)^{-1} \phi^T Y$$

Algorithm

Input:-  $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$ ,  $K, \phi \times 1$

(1) Formulate matrix  $\phi$

(2) Form  $Y$

$$(3) \quad a^* = (\phi^\top \phi)^{-1} \phi^\top y$$

Return  $a^*$

\* Logistic Regression

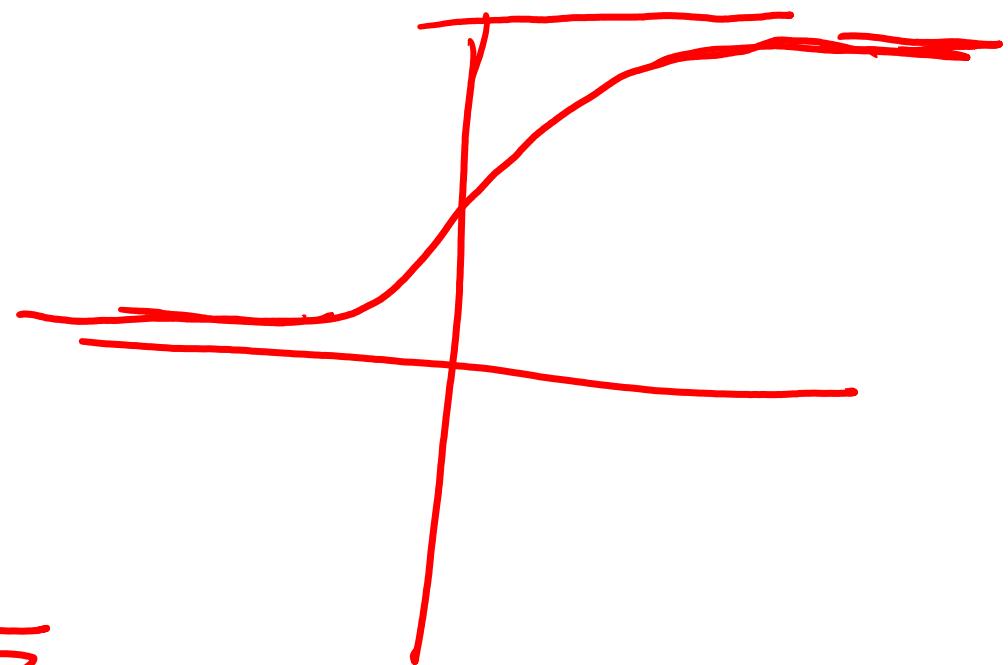
# Logistic regression is an algorithm  
for solving classification problems.  
(binary)

In logistic regression, the hypothesis is

$$f_{\theta}(x) = \sigma(ax)$$

“Sigmoid”

$$\sigma(z) = \frac{1}{1+e^{-z}}$$



Since we are dealing with binary classification  $y \in \{0, 1\}$

$$P(Y=1|x) = f_{\theta}(x) = \sigma(\underline{ax}) \rightarrow \sigma(\sum_i a_i x_i)$$

$$P(Y=0|x) = 1 - f_{\theta}(x)$$

$$P(y_i|x_i) = f_{\theta}(\cancel{x})(y_i) (1 - f_{\theta}(x))^{1-y_i}$$

$$P(x|\theta) = P(y|x; \theta) = \prod_{i=1}^N P(y_i|x_i; \theta)$$

Assumption  $\doteq$  IID assumption

I  $\rightarrow$  Independent

I  $\rightarrow$  Identical

D  $\rightarrow$  Dist

\* In logistic regression, gradient  
descent has to be used as there  
exist no closed form solution.

## Algorithm:- Logistic Regression

Input =  $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$

Initialize  $a^{(1)} = a_0$ .

for  $t = 1 : \text{epoch}$

$$a^{(t+1)} = a^{(t)} + \eta \nabla_a P(x/a)$$

end

Return  $a^{(t+1)}$

\* In practice, we work with the  
log likelihood  $\log P(\mathbf{z}|\mathbf{a})$ .

$$\mathbf{a}^{(t+1)} = \mathbf{a}^{(t)} + \gamma \nabla_{\mathbf{a}} \log P(\mathbf{z}|\mathbf{a}).$$





















