

## Supervised Learning

lets look we can be of two types.

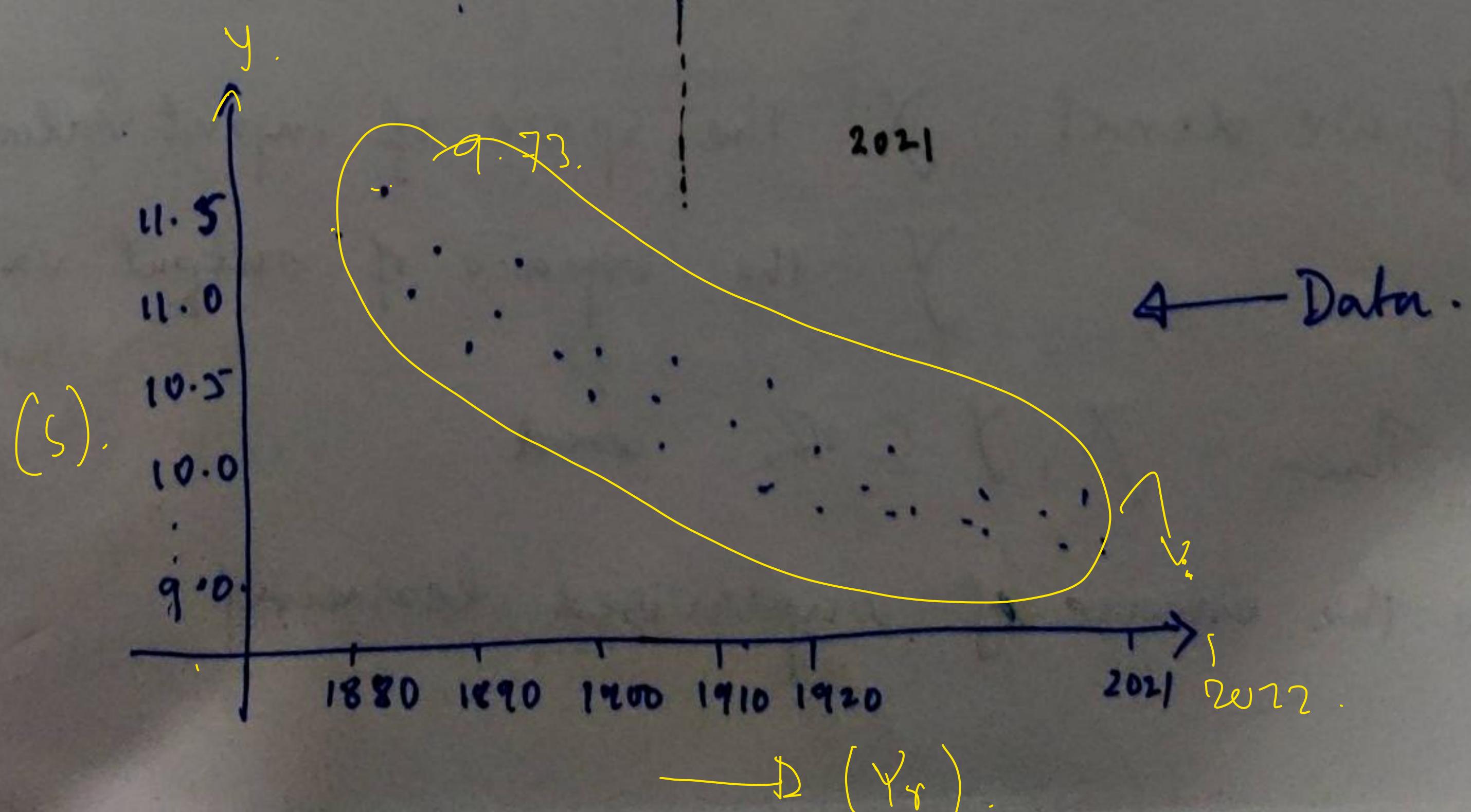
i) Regression

ii) Logistic regression / classification

Regression: Where target variables are continuous in nature.

lets say we have data for winning men's 100m times at summer olympics since, 1896. Two world wars interrupted in 1914, 1940

Time (s)	Years
11.3	1896
11.2	1897
11.25	1898
10.7	1899
10.85	:
:	



Given this data. we want to find the relation or predict the 100m winning time in the future.

Notation:  $\underline{x^i}$   $\rightarrow$  input variables (features)  
 $\underline{y^i}$   $\rightarrow$  Output / target variable.

in the above example

$x^i$  is the year.

$y^i$  is the time (winning) in seconds.

year ↓ time.

A pair  $(x^i, y^i)$   $\rightarrow$  training example.

$D = \{(x^i, y^i)\}_{i=1}^{40}$   $\rightarrow$  is the training set.  
 $= \{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(40)}, y^{(40)})\}$

If we denote  $\underline{X}$  the space of input values

$\underline{Y}$  the space of output values

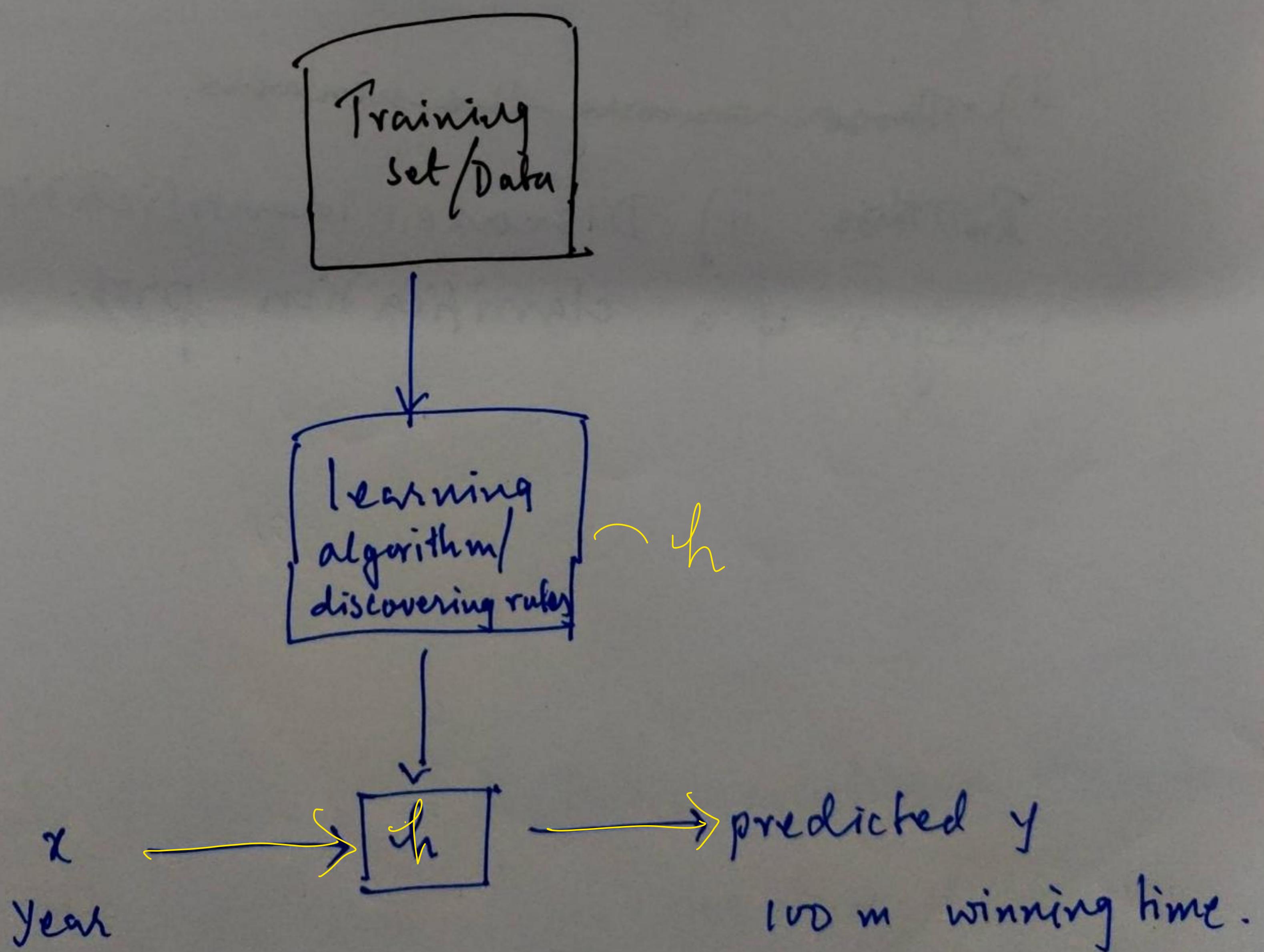
Then  $X, Y \subseteq \mathbb{R}$ . and.

Now the essence of supervised learning

$$D = \{(x^{(i)}, y^{(i)})\}_{i=1}^{n_0}$$

- ~ given a training set, to learn a function  $h: X \rightarrow Y$  so that  $h(x)$  is a good predictor for the corresponding value of  $y$ .
- ~  $h_w$  is called hypothesis function.

hence,



Classification: Logistic regression -

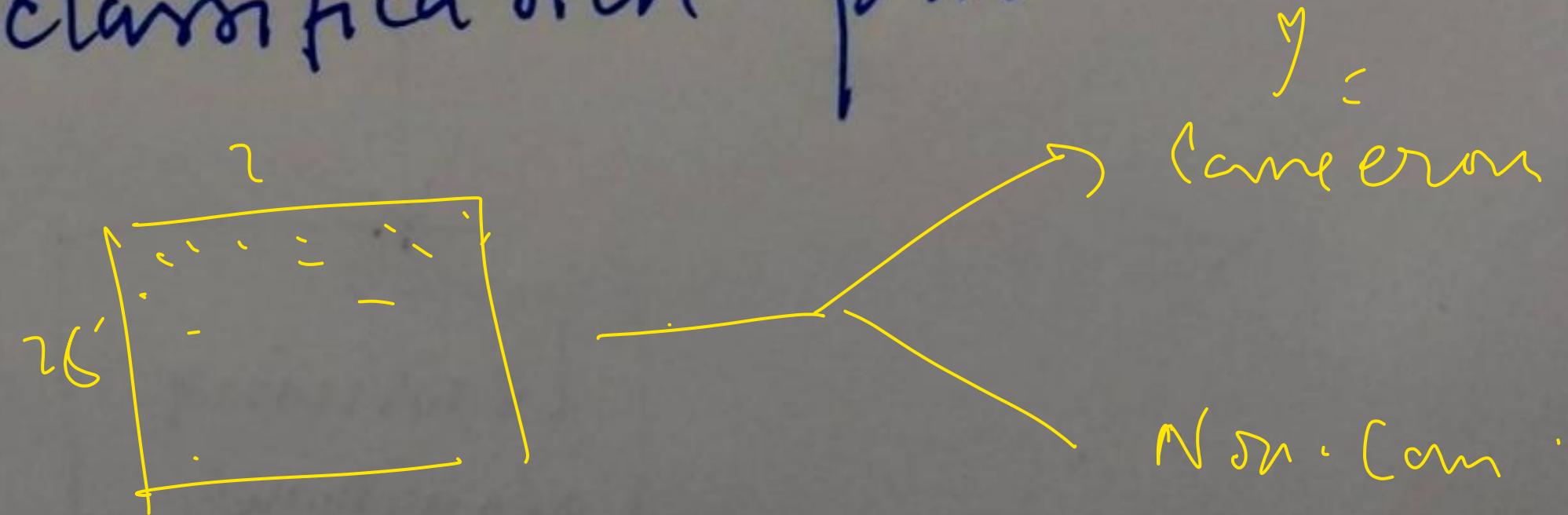
When the ~~not~~ target variable.  $y$  can take on only small number of discrete values.

example: i) If given the data, we want to know whether the winner is below.

20 years of age or, ~~above~~ <sup>above</sup> 20 to 25, & 25 >

ii) Those ~~among the~~ <sup>above</sup> ~~winner~~ <sup>winner</sup>  $D = \{(x^{(1)}, 1), (x^{(1)}, 0), (x^{(2)}, 1), (x^{(3)}, 1), \dots\}$

Both this. ii) Disease identification from images is a classification prob.



## Linear Regression

Andrew NG (Ref).

Let's consider a set of data of property price in Delhi.

Living area (ft <sup>2</sup> )	# bedrooms	Price (Lakhs) (INR)
2104	3	400
1600	3	330
2400	3	369
1416	2	232
3000	4	540
:	:	.
:	:	.
:	:	.

$$\underline{x} = \begin{pmatrix} x_1: \# \text{ of b. rooms} \\ x_2: \# \text{ area of the apt} \end{pmatrix} \in \mathbb{R}^2$$

let  $\underline{x} \in \mathbb{R}^2$  (a two dimensional vector)

and  $\underline{x}^i$  ( $i$  from  $m$  denotes the data set)

$i$  is the house identifier.

$$\underline{x}^{(i)} = \begin{pmatrix} x_1^{(i)} \\ x_2^{(i)} \end{pmatrix} \rightarrow \begin{array}{l} \text{feature 1 : } \underline{\text{living area size}} \\ \text{feature 2 : } \underline{\# \text{ of bedrooms}} \end{array}$$

Now before we do supervised learning on the given data, following steps \* need to be clarified

i) Defining the model

identifying the input  $x^i$  (with two attributes : living area, # of bed rooms.

Output  $y^i$  (the price in Lakhs)

ii) Modelling assumptions.  $y \approx h_{\omega}(x)$

This is an important step of the process. Whether the relationship bet output and input is linear OR higher order has to be decided. Based on this, parameter sets are identified and hypothesis is made.

if linear,  $\Rightarrow$   
say

$$h_{\omega}(x)$$

$w_0, w_i$  are the

$$\text{Hypothesis: } w_0 + w_1 x_1 + w_2 x_2$$

$$: w_0 + w_1 x_1 + w_2 x_1^2 + w_3 x_2 + w_4 x_2^2 + w_5 x_1 x_2$$

$w_0, w_i$  are the parameters. (or weights)

$$h_w(\underline{x}) : X \rightarrow Y$$

- ~ This is a linear mapping from  $X$  to  $Y$  space.
- ~ To give it a structure. lets say  $x_0 = 1$ .

the

$$h_w(\underline{x}) = w_0 x_0 + w_1 x_1 + w_2 x_2$$

$$= \underline{w^T \underline{x}} \quad \text{← Vectorization}$$

Both,  $\underline{w} = \begin{pmatrix} w_0 \\ w_1 \\ w_2 \end{pmatrix}$ ,  $\underline{x} = \begin{pmatrix} x_0 \\ x_1 \\ x_2 \end{pmatrix}$  are vector..

- ~  $x_0$  is the intercept term.

Ques iii) Defining a "goodness" model.

Given a data set, the problem now boils down to identifying the model parameters.

Ques: How do we define this "Goodness"

One assumption is to make the predicted output from  $h(\underline{x})$  close to  $y$  (actual output) for the known  $n$  data.

$$\text{Measure of } h_w(\underline{x}^{(i)}) - y^{(i)} = \underline{\Delta h}$$

→ And here we define a COST function, or loss loss function betw predicted and actual data  $y_i$

$$\sum_{i=1}^m h_f(x_i)^2$$

Cost function  $J(\theta_w)$  = Least Squares regression.

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

wf m: # of data points.

$x$ :  $x$  is the input feature vector.  $\in \mathbb{R}^n$ .

$y$ : output (price of the house)  $\in \mathbb{R}$

$h_w$ : Hypothesis function. (we have assumed it to be linear)

$\omega_0, \omega_1, \omega_L$ . we want to model.

The cost function  $J(\theta_w)$  is the least square error. in this case. But, there may be some other definition of cost function to ~~ratio~~ for defining the measure of goodness of the model.

$$\underline{\omega} = \begin{pmatrix} \omega_0 \\ \omega_1 \\ \omega_2 \end{pmatrix}$$

## No. Optimization (minimization) Algorithm for the cost function

- # We need to find out the weights <sup>w</sup> in defining the model.
- # Generally the weights are found so that they minimize the cost function. An iterative algorithm is adopted for the "optimization" prob. as defined below.

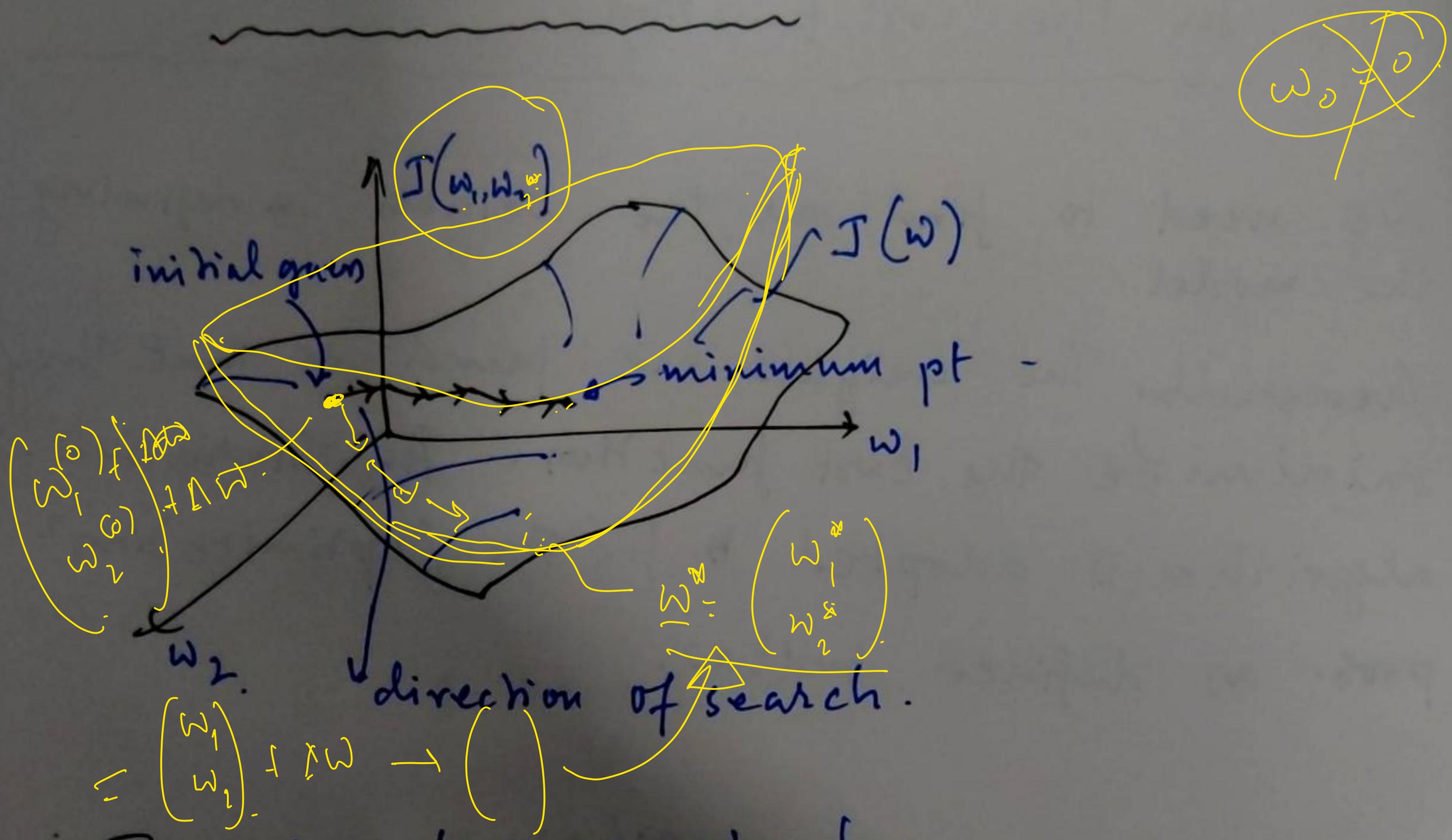
$$\begin{pmatrix} w_0 \\ w_1 \\ w_2 \end{pmatrix} \in \mathbb{R}^3 \leftarrow \underline{\underline{w}} = \underset{\text{for } J: \mathbb{R}^n \rightarrow \mathbb{R}, \underline{\underline{w}} \in \mathbb{R}^n}{\text{argmin } J(\underline{\underline{w}})} \quad \begin{aligned} &= \frac{1}{2m} \sum_{i=1}^m (h_w(x^{(i)}) - y^{(i)})^2 \\ &\text{Unconstrained optimization} \\ &J(w_0, w_1, w_2) \end{aligned}$$

Different algorithm exist:

- i) Gradient descent
- ii) Conjugate gradient descent
- iii) Newton's method.
- iv) Quasi-Newton method.
- v) Gauss-Jordan method

- # Least square minimization (LMS) is a convex optimization.

## Gradient descent



Towards to search for.

- # Gradient descent is an iterative procedure.
- # In order to search for the optimum (minimum) we need a search direction.

The "gradient" at a point of a real valued differentiable function gives the direction of max<sup>m</sup> increase. Hence, in order to find the max<sup>m</sup> decrease we need to search along -ive gradient dir.

$$\underline{d} = -\nabla f$$

$$= - \begin{pmatrix} \frac{\partial f}{\partial w_1} \\ \frac{\partial f}{\partial w_2} \end{pmatrix}$$

$$\underline{w}^k = \begin{pmatrix} w_1^k \\ w_2^k \end{pmatrix} - \alpha \cdot \left( \frac{\partial f}{\partial w_i} \right) \Big|_{w_1^k, w_2^k} = \underline{w}^{k+1}$$

Any iterative algorithm needs to have the following features.

$$\underline{w}^* = \arg \min J(w)$$

- i) An initial guess ( $\underline{w}^0$ )
- ii) A update rule. This depend on the search direction.

$$w^k \xrightarrow{w^k \rightarrow w^{k+1}} \underline{w}^{k+1} = \underline{w}^k + \underline{d}^k \xrightarrow{\text{FIR}} \text{search direction}$$

In. gradient descent  $\underline{d}^k = -\nabla f(\underline{x}^k)$ .

Widrow-Hoff L.R. iii) A step size or learning rate.

$\alpha$ . There are techniques to find this learning rate. using line search algorithm

iv) A stopping criteria.

The iterative algorithm needs to stop.

The magnitude of the update is proportional to the error. If the error is within a certain tolerance we stop.

$$\frac{|J(\underline{w}^{k+1}) - J(\underline{w}^k)|}{\max(1, |J(\underline{w}^k)|)} \leq \epsilon \quad \begin{matrix} 10^{-5} \\ J(\underline{w}^k) \approx J(\underline{w}^{k+1}) \end{matrix}$$

$$J(\omega) = \frac{1}{2} \sum_{i=1}^m \left( \sum_{j=1}^n w_j x_j^{(i)} - y_i \right)^2$$

$\underline{x} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix}$   $n=2$   
 division by  
m normalizes  
it.

$$\Rightarrow \frac{\partial J}{\partial w_j} = \frac{1}{2} \cancel{\sum_{i=1}^m} \left( \sum_{j=1}^n w_j x_j^{(i)} - y_i \right) x_j^{(i)}$$

i: # of data pts.

j: # of features.

$$\underline{\omega} = \begin{pmatrix} \omega_0 \\ \omega_1 \\ \omega_2 \\ \vdots \\ \omega_n \end{pmatrix} \rightarrow \begin{matrix} j=0 \\ j=1 \\ j=2 \\ \vdots \\ j=n \end{matrix}$$

Now the update rule.

$$\underline{w}_j^{k+1} = \underline{w}_j^k - \alpha \frac{\partial J}{\partial w_j} \quad j = 0, 1, 2, \dots, n$$

$$= \underline{w}_j^k - \alpha \cancel{\sum_{i=1}^m} \left( h_{\omega}(\underline{x}^i) - y_i \right) \underline{x}_j^{(i)}$$

repeat for  
k until  
convergence.

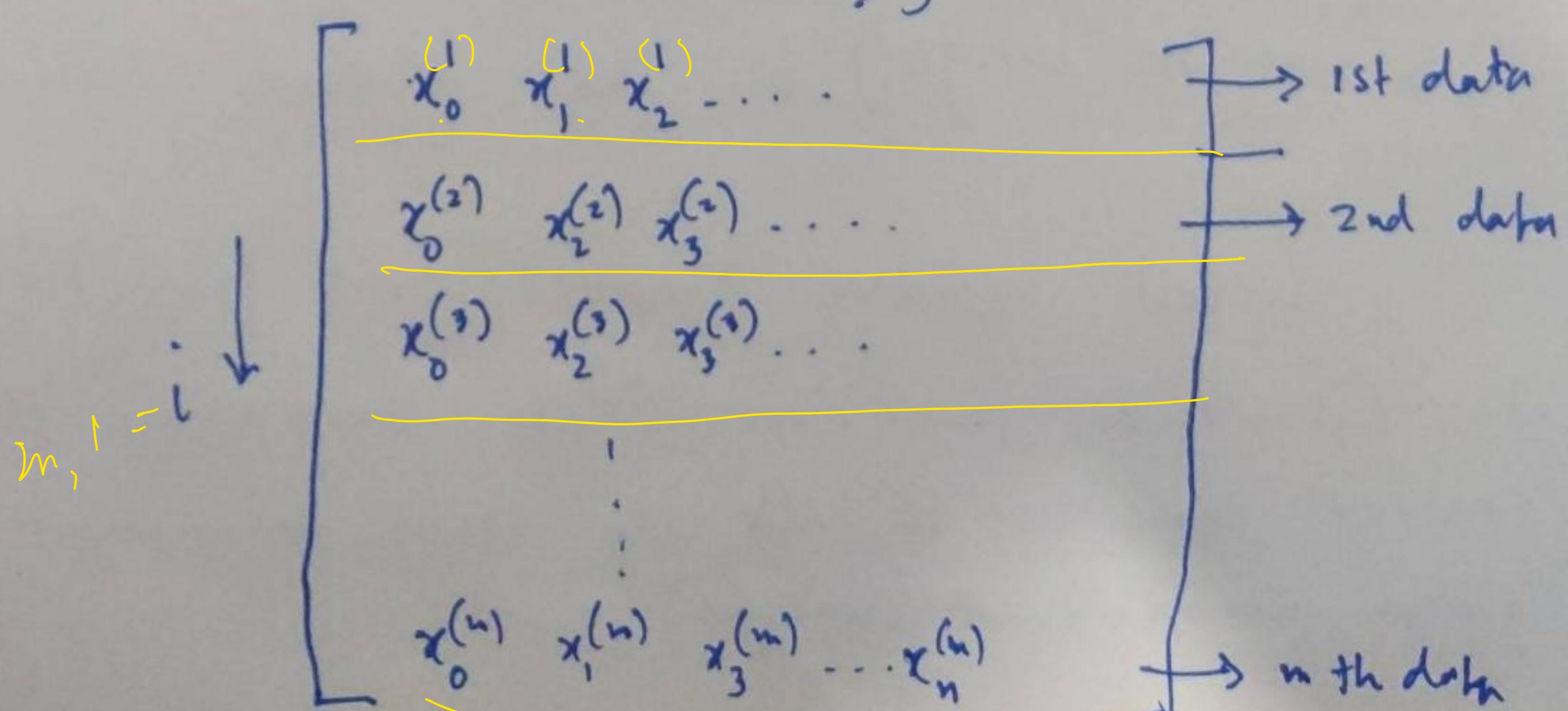
Fw  $\Leftrightarrow \forall j = 0 \text{ to } n$

$(\underline{x}^i, y^i)$  is a data pt

$\rightarrow j$

$$\underline{\omega}^{k+1} = \underline{\omega}^k - \alpha \nabla J$$

$$\nabla J = \begin{pmatrix} \frac{\partial J}{\partial \omega_0} \\ \frac{\partial J}{\partial \omega_1} \\ \vdots \\ \frac{\partial J}{\partial \omega_n} \end{pmatrix}$$



Batch gradient descent  $\Delta$ .

Stochastic gradient or (incremental)  
gradient descent

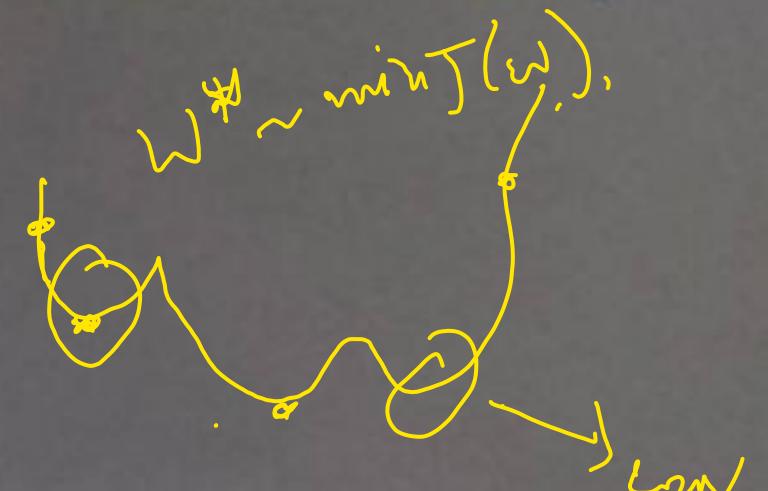
Batch Gradient descent:

⇒ i) The update rule in eq<sup>n</sup> ① is defined over the total "training set" i.e.,  $\{(x^i, y^i)\}_{i=1}^m$

ii) This method looks at every example or data pt in the training set at every step  $k$

⇒ Note: The gradient descent can be susceptible to local minimum in general, the optimization problem we have posed here for linear regression has only one global minima, and no other local optima; thus gradient descent always converges. (assuming the learning rate  $\alpha$  is not too large) to the global minimum.

Least square "error" cost function  $J$  is CONVEX function. The property of a convex function is all its local minimum is global minimum. (unique)



## Stochastic gradient descent

The above results we obtained with 'batch gradient descent' which involve processing the entire training set in one go, can be computationally costly for large data set.

→ If the dataset is sufficiently large, it may be worthwhile to use "sequential learning" also known as "on-line learning" algorithms. In this data pts are considered one-at-a-time. model parameters are updated after each such presentation.

→ Sequential learning is also appropriate for real time applications in which the data observations are arriving in a continuous stream and predictions must be made before all the data points are seen.

→ Stochastic gradient descent also known as sequential gradient descent.

for LMS algo

$$J = \frac{1}{2} (y^1 - h_w(x^1))^2 + \frac{1}{2} (y^2 - h_w(x^2))^2 + \dots + J_m$$

$$\Rightarrow \nabla J = \underline{\nabla J_1} + \underline{\nabla J_2} + \dots + \underline{\nabla J_m}$$

The algorithm is.

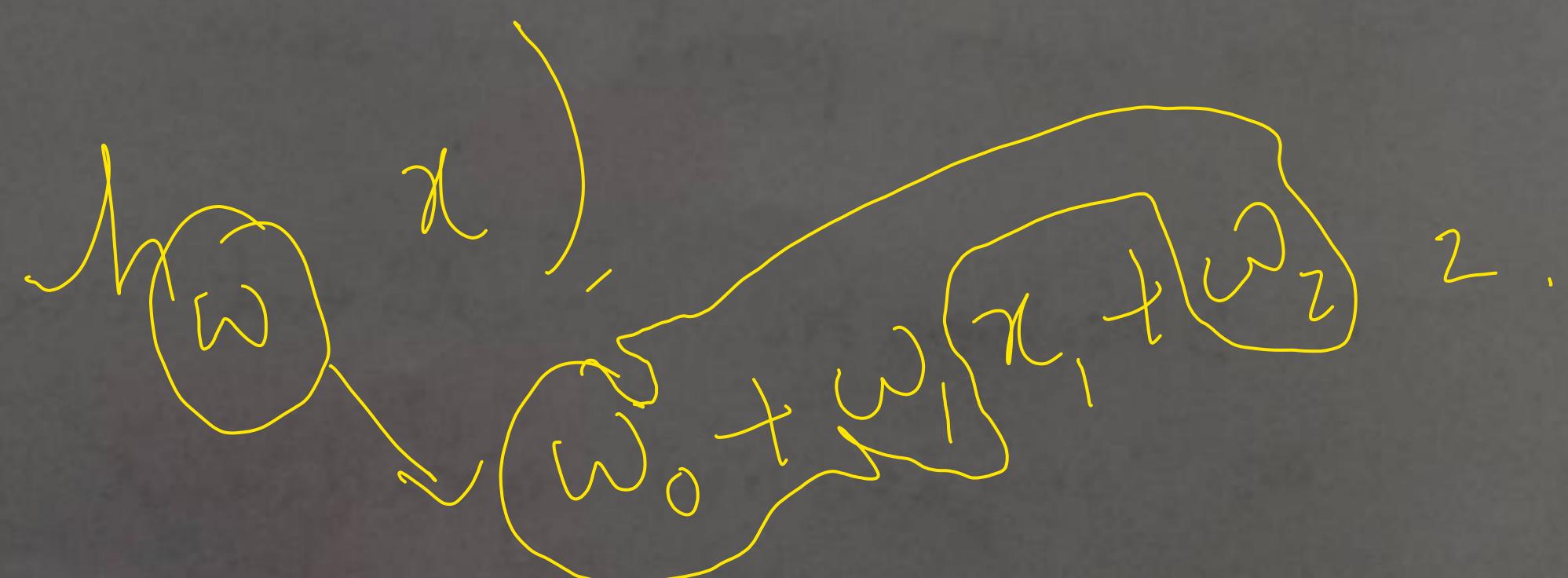
loop i=1 to m

$$\begin{aligned} \underline{w}^{k+1} &= \underline{w}^k - \alpha \nabla J_{\text{gm}} \quad \text{w/ } \nabla J_{\text{gm}} = (y^m - h_w(x^m))x_j^m \\ &= \underline{w}^k - \alpha (y^i - \underline{w}^{kT} \underline{x}_j^i) \underline{x}_j^i \\ &\quad \text{end} \\ &\quad \text{h}_w(\underline{x}^i) \end{aligned}$$

flash.



The algorithm is written in vector form.



## Normal Equations.

The property of a multivariate function like  $J(w_0, w_1, w_2, \dots) = J(\underline{w})$  to be minimum is.

$$J: \mathbb{R}^n \rightarrow \mathbb{R}$$

$\frac{\partial f}{\partial x} = 0, \frac{\partial^2 f}{\partial x^2} > 0$  local min

$\nabla_w J = 0 \quad \& \quad H_J(\underline{w}) > 0$

$\nabla_w J = \begin{pmatrix} \frac{\partial J}{\partial w_1} \\ \frac{\partial J}{\partial w_2} \\ \vdots \\ \frac{\partial J}{\partial w_n} \end{pmatrix} = 0$ 
 $H_J = \begin{bmatrix} \frac{\partial^2 J}{\partial w_1^2} & \frac{\partial^2 J}{\partial w_1 \partial w_2} & \cdots & \frac{\partial^2 J}{\partial w_1 \partial w_n} \\ \frac{\partial^2 J}{\partial w_2 \partial w_1} & \frac{\partial^2 J}{\partial w_2^2} & \cdots & \frac{\partial^2 J}{\partial w_2 \partial w_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 J}{\partial w_n \partial w_1} & \frac{\partial^2 J}{\partial w_n \partial w_2} & \cdots & \frac{\partial^2 J}{\partial w_n^2} \end{bmatrix} = (\nabla_w^T J)$

gradient
Hessian
 $a^T H_J a > 0$

: always a vector. =  $(\nabla_w J)^T$ 
derivative of  $J$

The Hessian of the function  $J(\underline{w})$  is.

$$\begin{aligned}
 H_{ij} &= \left[ \frac{\partial^2 f}{\partial w_i \partial w_j} \right] \\
 &= \begin{bmatrix} \frac{\partial^2 f}{\partial w_1^2} & \frac{\partial^2 f}{\partial w_1 \partial w_2} & \cdots & \frac{\partial^2 f}{\partial w_1 \partial w_n} \\ \frac{\partial^2 f}{\partial w_2 \partial w_1} & \frac{\partial^2 f}{\partial w_2^2} & \cdots & \frac{\partial^2 f}{\partial w_2 \partial w_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial w_n \partial w_1} & \frac{\partial^2 f}{\partial w_n \partial w_2} & \cdots & \frac{\partial^2 f}{\partial w_n^2} \end{bmatrix}
 \end{aligned}$$

For the current prob in hand the training set data can be arranged in matrix form as

$$\underline{y} = \begin{pmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m)} \end{pmatrix} \xrightarrow{\longrightarrow} \underline{x} = \begin{bmatrix} x_0^{(1)} & x_1^{(1)} & \dots & x_n^{(1)} \\ x_0^{(2)} & x_1^{(2)} & \dots & x_n^{(2)} \\ \vdots & \vdots & \ddots & \vdots \\ x_0^{(m)} & x_1^{(m)} & \dots & x_n^{(m)} \end{bmatrix} \stackrel{n=2}{\downarrow} \underline{w} = \begin{pmatrix} w_0 \\ w_1 \\ w_2 \\ \vdots \\ w_m \end{pmatrix}$$

Hence  $J(w)$  can be written as

$$\begin{pmatrix} A.J_1 \\ A.J_2 \\ \vdots \\ A.J_m \end{pmatrix} = \begin{bmatrix} x_0^{(1)} & x_1^{(1)} & \dots & x_n^{(1)} \\ x_0^{(2)} & x_1^{(2)} & \dots & x_n^{(2)} \\ \vdots & \vdots & \ddots & \vdots \\ x_0^{(m)} & x_1^{(m)} & \dots & x_n^{(m)} \end{bmatrix} \begin{pmatrix} w_0 \\ w_1 \\ \vdots \\ w_m \end{pmatrix} - \begin{pmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m)} \end{pmatrix} =$$

$$\begin{aligned} \text{now } \underline{A.J(w)} &= \frac{1}{2^m} \sum_{i=1}^m A.J_i^2 \\ &= \frac{1}{2^m} \left[ (\underline{x} \underline{w} - \underline{y})^T (\underline{x} \underline{w} - \underline{y}) \right] \\ &= \frac{1}{2^m} \left( \underline{w}^T \underline{x}^T \underline{x} \underline{w} - 2 \underline{w}^T \underline{x}^T \underline{y} - \underline{y}^T \underline{x} \underline{w} + \underline{y}^T \underline{y} \right) \end{aligned}$$

$$\nabla_w J = \underline{x}^T \underline{x} \underline{w} - \underline{x}^T \underline{y} + 0$$

If can be shown that

$$\nabla_{\underline{w}} J = \underline{\underline{x}^T x} - \underline{\underline{x}^T y} = 0 \quad \textcircled{1}$$

Now, the condition for optimum  $\underline{\underline{w}}^*$  is  $\nabla_{\underline{w}} J = 0$ .

$$\underline{\underline{w}}^* = \arg \min J(\underline{w}).$$

From equ<sup>n</sup>  $\textcircled{1}$  we get.

$$\underline{\underline{x}^T x} \underline{\underline{w}} - \underline{\underline{x}^T y} = 0.$$
  
$$\Rightarrow \boxed{\underline{\underline{w}}^* = (\underline{\underline{x}^T x})^{-1} \underline{\underline{x}^T y}}$$

$$\underline{\underline{x}^T x} = \begin{bmatrix} x_0^{(1)} & x_1^{(1)} & x_2^{(1)} \\ x_0^{(2)} & x_1^{(2)} & x_2^{(2)} \\ x_0^{(3)} & x_1^{(3)} & x_2^{(3)} \end{bmatrix}$$

$\underline{\underline{w}}^*$  is  $\underline{\underline{w}}^k - \alpha \nabla f(\underline{\underline{w}}^k)$ . Closed form equ<sup>n</sup> for  $\underline{\underline{w}}^k$

Probabilistic viewpoint  $P(D|D) \prod_{i=1}^m P(y^{(i)} | x^{(i)}; w)$

Note: In machine learning literature, the -ive log of the likelihood function is called an error function. Because, the -ive logarithm is a monotonically decreasing function, maximizing the likelihood is equivalent to decreasing the error.

Why least square error function is a good estimator of a model's "goodness"

lets assume target variables and inputs are related through an equ<sup>n</sup>.

$$y^{(i)} = \underbrace{w^T x^{(i)}}_{\text{vector (multiple features)}} + \underbrace{\epsilon^{(i)}}_{\text{error in estimate.}}$$



What does  $\epsilon^{(i)}$  measure?

i)  $\epsilon^{(i)}$  is a random variable, which captures either unmodeled effects or features/attributes in the data. or random noise.

ii)  $\epsilon^{(i)}$  also follows a zero mean Gaussian distribution. i.e with.

$$\epsilon^{(i)} \sim \mathcal{N}(y^{(i)} | 0, \sigma^2)$$

↑ ↑ mean SD.

$$\sigma^2 \approx \frac{1}{\beta}$$

precision  $\beta$  (inverse of variance)

$$\epsilon^{(i)} \sim N(0, \beta^{-1})$$

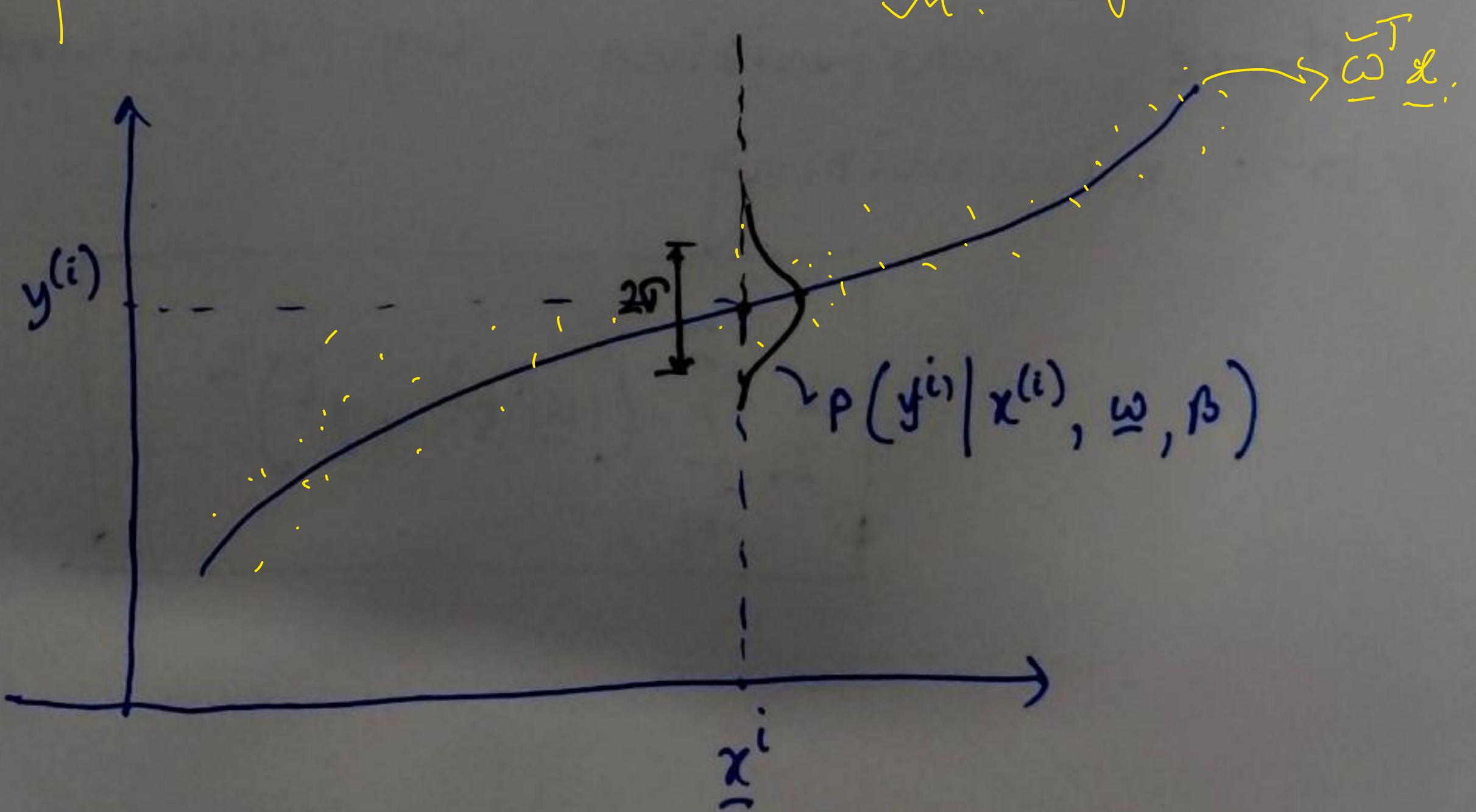
$$p(\epsilon^i) = \frac{1}{\sqrt{2\pi\beta}} \exp\left(-\frac{(\epsilon^i)^2}{2\beta}\right)$$

Now, the probability distribution of the target variable  $y^{(i)}$  is given by.

$$p(y^{(i)} | x^{(i)}, \underline{\omega}, \beta) = N(y^{(i)} | \underline{\omega}^T \underline{x}^{(i)}, \beta^{-1})$$

$$\text{I. } \mu = \underline{\omega}^T \underline{x}^{(i)}$$

$$\text{II. } \sigma^2 = \beta^{-1}$$



We now use the training dataset  $\{(x^i, y^i)\}_{i=1}^m$  to determine the values of the unknown parameters  $\underline{\omega}$  and  $\beta$  by "maximum likelihood".

As we assume every data pt is independent of each other, i.e.  $\epsilon^i, y^i \& x^i$  are independent random variables. So, the likelihood function is defined as

$$L = \left( \begin{array}{c} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m)} \end{array} \right) = P(y^t | \underline{x}^t, \underline{\omega}, \beta) = \prod_{i=1}^m N(y^{(i)} | \underline{\omega}^T \underline{x}^i, \beta^{-1}) = e^{-\frac{1}{2\beta}}$$

Likelihood:  $\prod_{i=1}^n p(y^i | x^i; \underline{w})$ , is the given.

$\rightarrow$  probabilistic model relating  $y^i$  with  $x^i$ .

Now what is ~~the~~ a reasonable way of choosing our best guess of the parameters  $\underline{w}$ ?

$$\log(p(y|x, \underline{w}, \beta)) = -\frac{\beta}{2} \sum_{i=1}^m (\underline{w}^T \underline{x}^i - y^i)^2 + \frac{m \ln \beta - \frac{m}{2} \ln(2\pi)}{2}$$

Hence, maximizing log (likelihood) is equivalent to minimizing

$$\boxed{\frac{1}{2m} \sum_{i=1}^m (\underline{w}^T \underline{x}^i - y^i)^2}$$

normalize by a constant

$f(n)$

$-f(n^*)$

least sq error

Note: const  $\beta, m$  &  $2\pi$  does not change the optimum pt. of the objective function.