Last class:-

① Normalization ]

② Probabilistic interpretation of linear regression:-

$$y^{(i)} \mid x^{(i)}, \theta \sim N(\theta^T x^{(i)}; \sigma^2)$$

i.i.d $\quad \varepsilon \sim N(0, \sigma^2)$

ML estimate

$$\underset{\theta}{\text{argmax} \log} \left[ \prod_{i=1}^{m} P(y^{(i)} \mid x^{(i)}; \theta) \right] = \underset{\theta}{\text{argmax}} \; J(\theta)$$

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

## Newton's Method for Optimization -

Till Now:-

$$\underset{\theta}{\text{argmax}} \; f(\theta) \longrightarrow \underset{\theta}{\text{argmax}} \; J(\theta)$$

$$\theta \in \mathbb{R}^{n+1}$$

uses first order information $\left[ \hookrightarrow \text{gradient descent} \right]$

Converges to local minima.
convex:- global minima.

$\nabla_\theta f(\theta)$

$$\theta^{(t+1)} \longleftarrow \theta^{(t)} - \eta \cdot \nabla_\theta f(\theta)$$

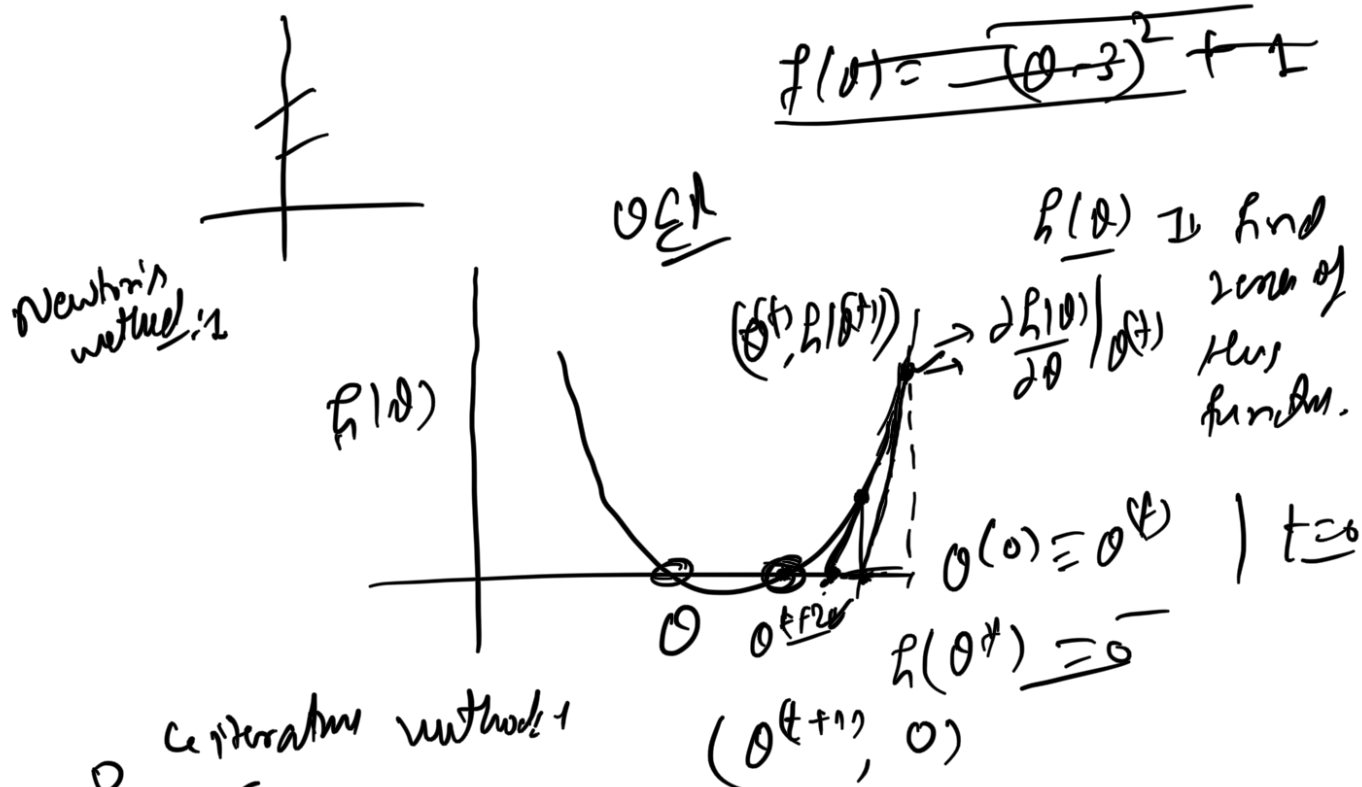↳ uses Second order information $\left[ \nabla^2_\theta f(\theta) \right) \checkmark$

(curvature) $\equiv$ (Hessian matrix) ]

Hopefully:- that you can make faster progress

Newton's Method for optimization :-

is for finding zeroes of a function.

$$f(\theta) = (\theta - 3)^2 + 1$$



Newton's
method : 1

$h(\theta)$ I find
zeros of
this
function.

$(\theta^{(t)}, h(\theta^{(t)})) \Rightarrow \dfrac{\partial h(\theta)}{\partial \theta}\Big|_{\theta^{(t)}}$

$\theta^{(0)} = \theta^{(t)}$ } $t = 0$

$h(\theta^*) = 0$

$(\theta^{(t+1)}, 0)$

o iterative method 1

$$\dfrac{h(\theta^{(t+1)}) - h(\theta^{(t)})}{\theta^{(t+1)} - \theta^{(t)}} = h'(\theta)\Big|_{\theta^{(t)}}$$

$$\boxed{\theta^{(t+1)} - \theta^{(t)} = -\dfrac{h(\theta^{(t)})}{h'(\theta^{(t)})}}$$

$$\theta^{(t+1)} = \theta^{(t)} - \dfrac{h(\theta^{(t)})}{h'(\theta^{(t)})}$$

Newton's update 1 for finding
zeros of a
function

↳ ⑪ how does it converge

with our problem?

$\checkmark \left[\underset{\theta}{\text{argmin}} \; f(\theta)\right]$ [ Convex ]

Finding zeros of $\overline{f'(\theta)}$ ]

$\Longleftarrow \nabla_\theta f(\theta) \Big|$  $\left[\dfrac{\partial f(\theta)}{\partial \theta}\right]$

$\Rightarrow$ Even if function is not conve, we will still find (local optima)

Newton's method for optimization:-

$\underset{\theta}{\text{argmin}} \; f(\theta)$

$\left[\; \theta^{(t+1)} \longleftarrow \theta^{(t)} \; - \; \left[\dfrac{f'(\theta^{(t)})}{f'(\theta^{(t)})}\right] \;\right]$

newton update

$f(\theta) = (\theta - 3)^2 + 1$

Newton's update:-  $\overline{\theta^{(t)}}$

$\diagup \; \theta^{(t+1)} \longleftarrow \theta^{(t)} \; - \; \dfrac{2(\theta - 3)\big|\theta^{(t)}}{2}$

$\boxed{\theta^{(t+1)} \Longleftarrow 3}$  $\therefore$ Newton's method converges in single iteration.

"Faster".   minimal Mascances:-

$f'(\theta) \geq 0$

$< 0$  maxima

**Higher Dimension:**

unividuate.   Scalar update

$$\theta^{(t+1)} \leftarrow \theta^{(t)} - \left[f''(\theta^{(t)})\right]^{-1} \left[f'(\theta^{(t)})\right]$$

et multi-variate   Vector update

$$\left[\theta^{(t+1)} \leftarrow \theta^{(t)} - H^{-1}\big|_{\theta^{(t)}} \nabla_\theta f(\theta)\big|_{\theta^{(t)}}\right]$$

$$\frac{H}{O(n^3)}$$



$\theta_2$   $\theta^{(0)}$

$\theta_1$

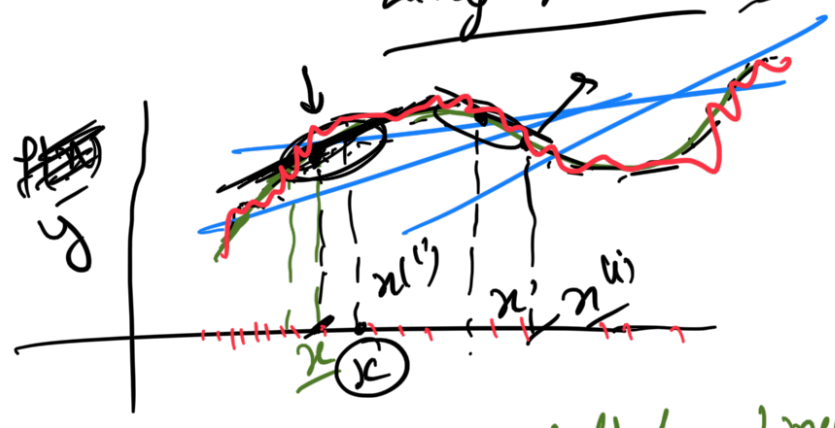$\theta_2$   $\theta^{(0)}$

$\theta_1$

[ Why would one even prefer gradient descent over Newton's method? ]

**Locally weighted Linear Regression :-**

**Non parametric methods :-**

"Lazy learning"



$y$

$x^{(j)}$   $x,$ $x^{(i)}$

$x$

$h_\theta(x) \simeq \boxed{\quad} y$

Multiple Linear Approximations to the fn.

$x^{(t)}$   $m$

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

inversely proportional to distance of $x^{(i)}$ from $x$

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

$\hookrightarrow$ ①

$$w^{(i)} = e^{\frac{-(x-x^{(i)})^2}{2\tau^2}} \longrightarrow \tau : \text{bandwidth}$$

locally weighted linear regression.

underfit

$\tau \rightarrow \infty$ $\qquad w^{(i)} \rightarrow e^{-0} = 1$

$\hookrightarrow$ Linear Regression

overfit

$\tau \rightarrow 0$ $\qquad$ You pay more & more attention to closer by example

"overfit"

$$\Rightarrow J(\theta) = \sum_{i=1}^{m} w^{(i)} [y^{(i)} - h_\theta(x^{(i)})]^2 / 2m$$

$n \times n$ matrix

$$w^{(i)} \longleftarrow e^{\frac{(x-x^{(i)})^T \hat{\Sigma}^{-1} (x-x^{(i)})}{2}}$$