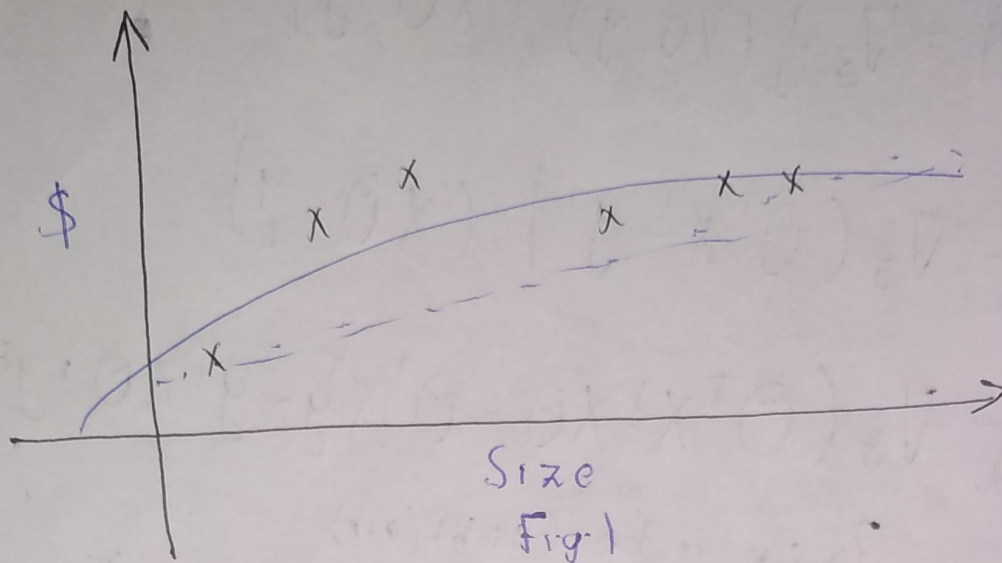


(12)

Lecture - 3

(*)



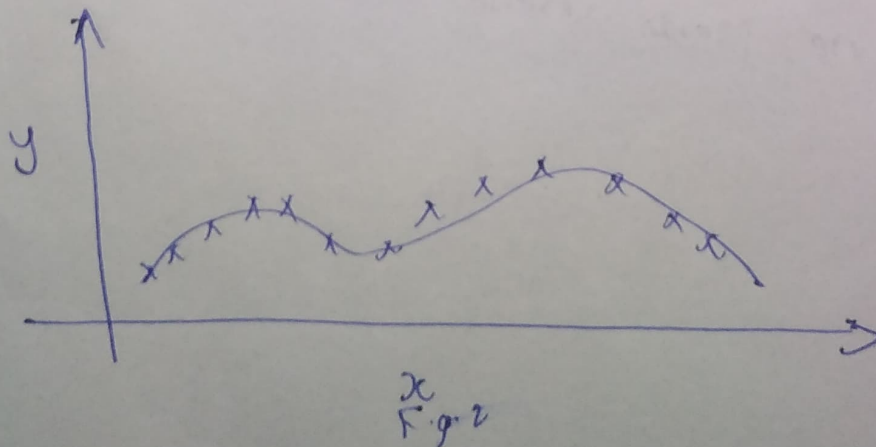
(*) You could define new feature to increase the complexity of your hypothesis

$$(*) \theta_0 + \theta_1 x_1$$

$$(*) \theta_0 + \theta_1 x_1 + \theta_2 x_1^2$$

$$(*) \theta_0 + \theta_1 x_1 + \theta_2 \sqrt{x_2}$$

(*) Locally Weighted Linear Regression:



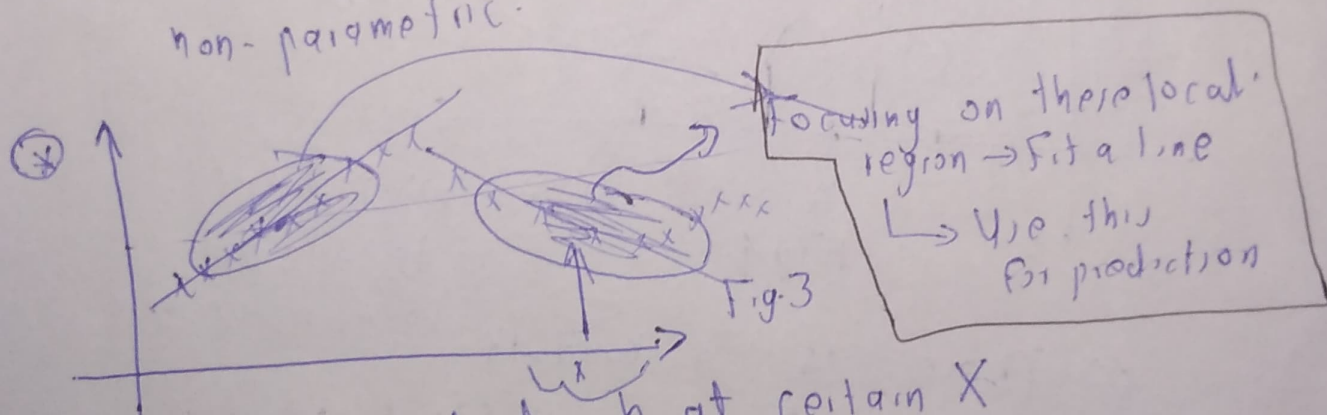
⑧ Difference between parametric and non-parametric method

→ In parametric learning algorithm → You fit fixed set Θ to data

→ In non-parametric:-

- Amount of data/parameters you need to keep grow with size of data.

→ Locally Weighted Regression is non-parametric.



→ To evaluate h at certain x

LR: Fit Θ to minimize:-

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

Return $\Theta^T x$.

• For locally weighted regression:

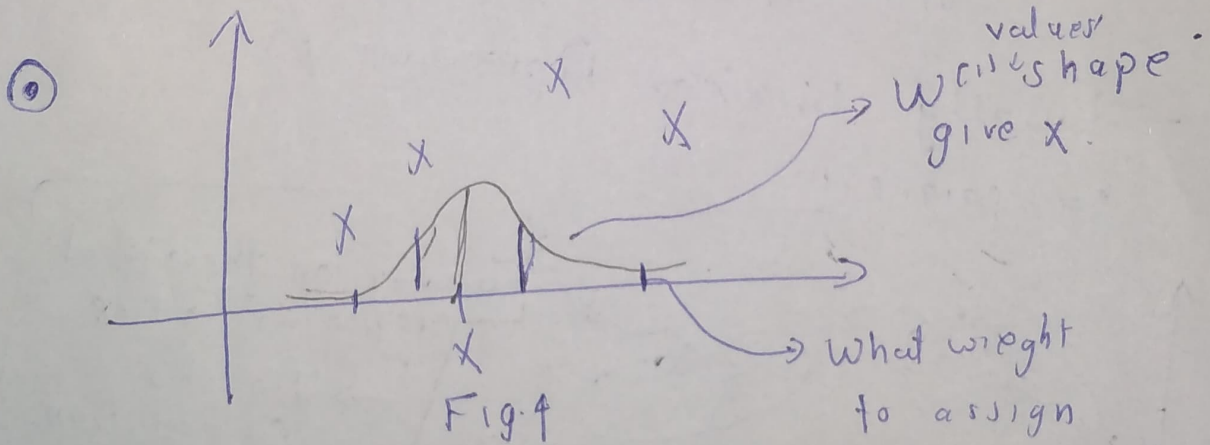
Fit Θ to minimize

$$\sum_{i=1}^m \boxed{w^{(i)}} (y^{(i)} - \Theta^T x^{(i)})^2 \quad \text{Eq. 1}$$

where $w^{(i)}$ is a weight function

$$w^{(i)} = \exp\left(-\frac{(x^{(i)} - x)^2}{2\tau^2}\right) \quad \text{Eq. 2}$$

⊗ If $|x^{(i)} - x|$ is small, $w^{(i)} \approx 1$
else it will be close to zero



⊗ There is a bandwidth parameter τ = bandwidth (Tau) to choose the width of bell-shape curve ($\uparrow \tau \rightarrow$ Larger width).

⊗ Use locally weighted when dimension of your dataset is small and you have a lot of data.

① Probabilistic Interpretation of Linear Regression:-

(Why Least

② Why Least Squares?

→ Assume ① $y^{(i)} = \theta^T x^{(i)} + \epsilon^{(i)}$

These independently identically distributed (i.i.d).

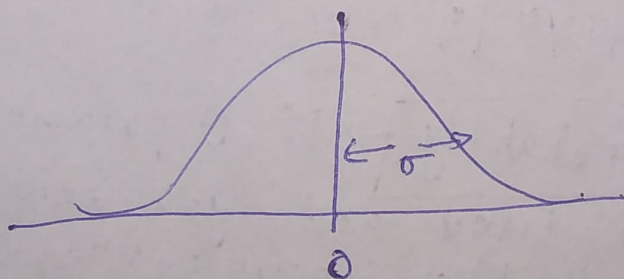
error term: unmodelled effects, random noise.

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

Pf ϵ

→ $P(\epsilon^{(i)}) = \frac{1}{\sqrt{2\pi} \cdot \sigma} \exp\left(-\frac{(\epsilon^{(i)})^2}{2\sigma^2}\right)$

↳ PDF for $\epsilon^{(i)}$



→ This implies:-

→ $P(y^{(i)} | x^{(i)}; \theta) = \frac{1}{\sqrt{2\pi} \cdot \sigma} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right)$

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

↑
"parameterized as"
by.

*) $L(\theta) = P(y | x; \theta)$

↓
likelihood

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

$$= \prod_{i=1}^m \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right)$$

• What is the difference between probability & likelihood?

① → Likelihood → If you view $P(y^{(i)} | x^{(i)}; \theta)$ as a function of θ then $P(y^{(i)} | x^{(i)}; \theta)$ is likelihood. (Data is fixed)

② → Probability → If θ parameter is fixed and you vary the data then $P(y^{(i)} | x^{(i)}; \theta)$ is probability.

→ Log likelihood

$$\ell(\theta) = \log L(\theta)$$

$$= \log \prod_{i=1}^m \frac{1}{\sqrt{2\pi} \cdot \sigma} \exp(\dots)$$

$$= \sum_{i=1}^m \frac{1}{\sqrt{2\pi}} \log \left(\frac{1}{\sqrt{2\pi} \cdot \sigma} \cdot \exp \left(\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma} \right) \right)$$

$$= \sum_{i=1}^m \log \left(\frac{1}{\sqrt{2\pi} \cdot \sigma} \right) + \sum_{i=1}^m \log \exp \left(\frac{-(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma} \right)$$

$$= m \log \left(\frac{1}{\sqrt{2\pi} \cdot \sigma} \right) + \sum_{i=1}^m \log \left(\frac{-(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma} \right)$$

→ See this equation → the black part shows the MSE Formula so
 we maximize log-likelihood is equivalent
 to minimizing MSE i.e. $\frac{1}{2} \sum_{i=1}^m (y^{(i)} - \theta^T x^{(i)})^2$
 $\rightarrow J(\theta)$ (cost function).

→ Maximum Likelihood Estimator

→ Choose θ to maximize likelihood.

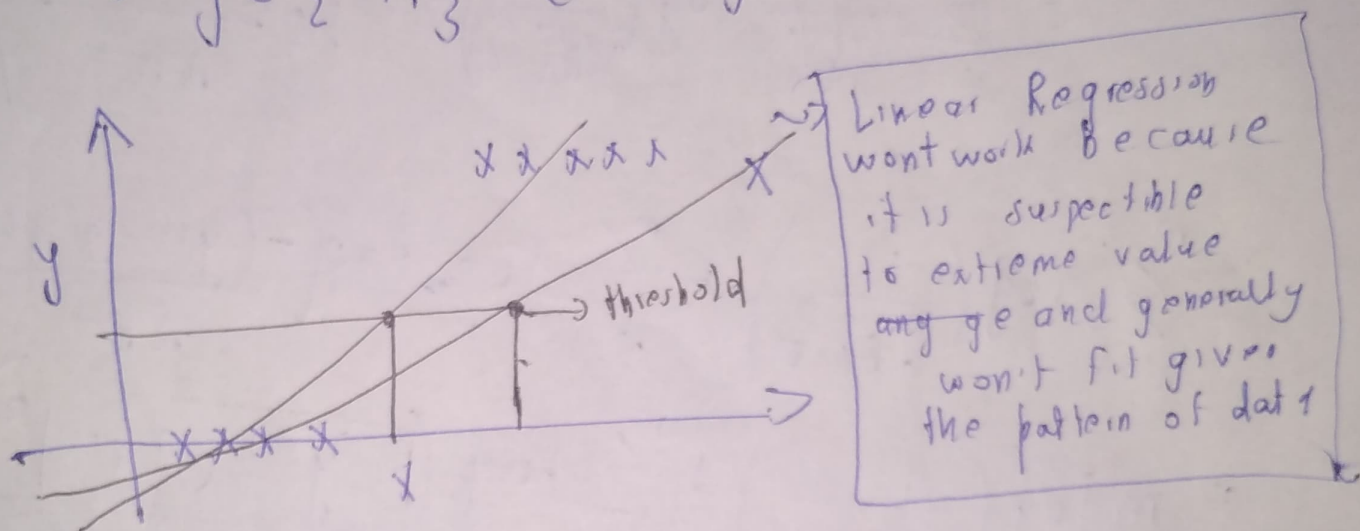
→ Easier to maximize log-likelihood $\ell(\theta)$

So in Linear Regression choose

② to minimize

⑧ Classification:-

- $y \in \{0, 1\}$ (binary classification).



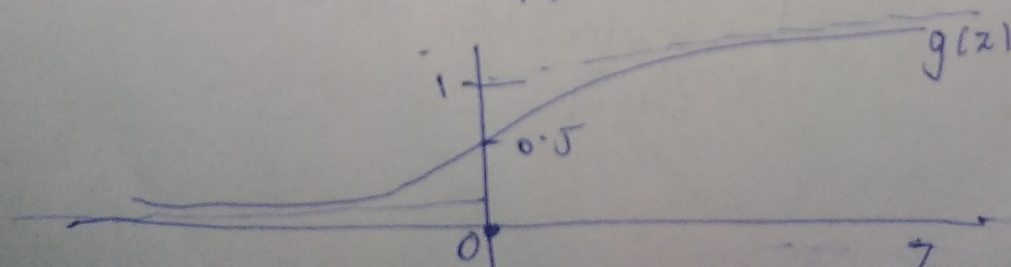
⑧ Logistic Regression:-

→ Want $h_0(x) \in [0, 1]$

$$\textcircled{a} h_0(x) = g(\theta^T x) = \frac{1}{1 + e^{-\theta^T x}}$$

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

→ Sigmoid or logistic function.



(19)

- ① Why $g(z) = \frac{1}{1+e^{-z}}$ because it then form generalized linear model (Explanation in future Lectures).

- ② let's make some assumption:

① $P(y=1 | x; \theta) = h_{\theta}(x)$

② $P(y=1 | x; \theta) \leq 1 - h_{\theta}(x)$

③ $y \in \{0, 1\}$

Putting / compressed in one equation $\rightarrow P(y | x; \theta) = h_{\theta}(x)^y \cdot (1 - h_{\theta}(x))^{1-y}$

If $y=1$:- $P(y | x; \theta) = h_{\theta}(x)$
if $y=0$:- $P(y | x; \theta) = 1 - h_{\theta}(x)$

③ $L(\theta) = P(\vec{y} | x; \theta)$

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

$$= \prod_{i=1}^m h_{\theta}(x)^{y^{(i)}} \cdot (1 - h_{\theta}(x))^{1-y^{(i)}}$$

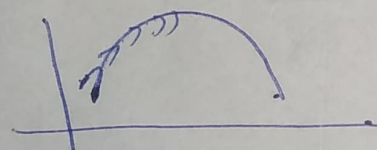
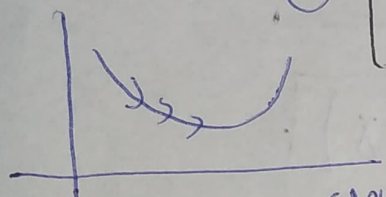
$$= \prod_{i=1}^m (h_{\theta}(x)^{y^{(i)}} \cdot (1 - h_{\theta}(x))^{1-y^{(i)}})$$

$$l(\theta) = \sum_{i=1}^m y^{(i)} \log h_{\theta}(x^{(i)}) + (1-y^{(i)}) \log (1-h_{\theta}(x^{(i)}))$$

Choose θ to maximize $l(\theta)$

① Batch gradient descent to choose θ :

$$\theta_j := \theta_j + \alpha \frac{\partial}{\partial \theta_j} l(\theta)$$



→ Plus sign because

We are trying to maximize log-likelihood

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

$$\frac{\partial l(\theta)}{\partial \theta_j}$$

② $l(\theta)$ is always concave
so there is global
optima.

→ Same update as linear regression
but $h_{\theta}(x^{(i)})$ is now different.

* Newton Method:-

* Gradient ascent takes a lot of baby steps to converge.

* N-M takes bigger step and need fewer iteration to converge.

* One-dimensional Justification for NM :-

1. Have f

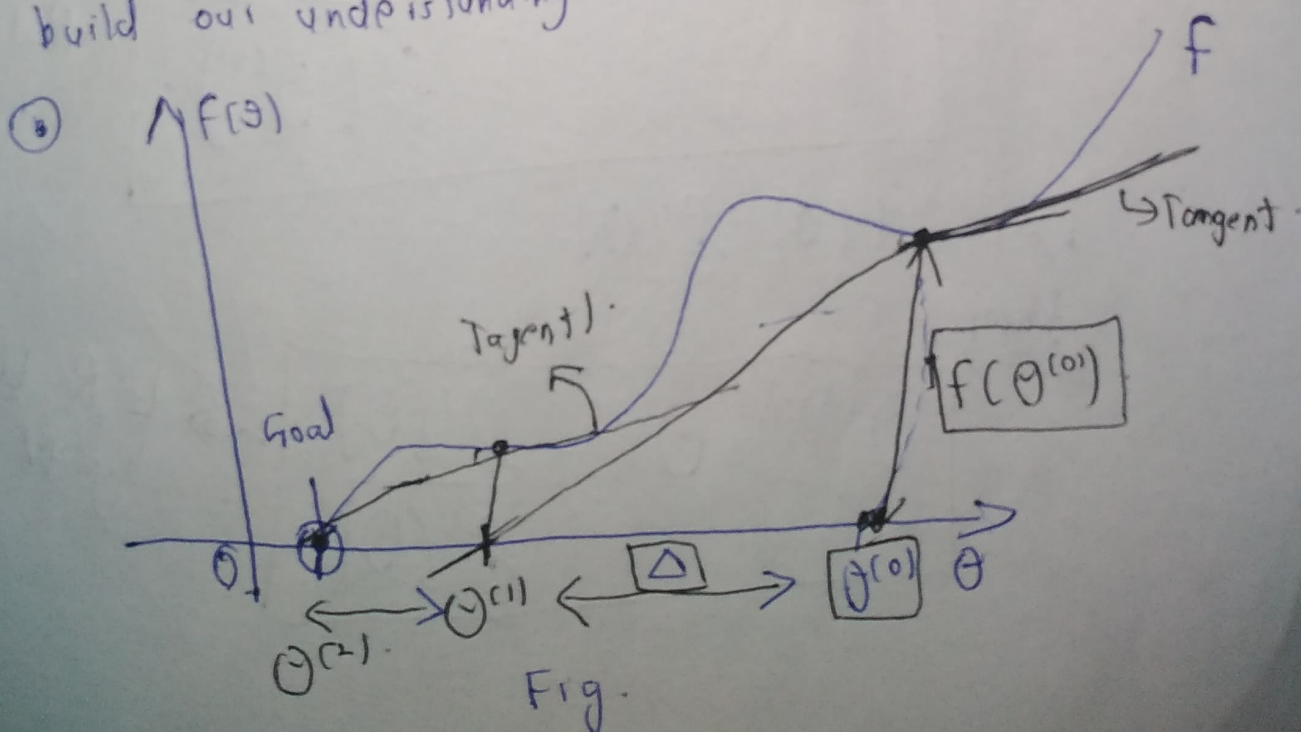
2. Want to find θ s.t $f(\theta) = 0$.

* [Want to maximize $\ell(\theta)$]

i.e want $\ell'(\theta) = 0$

↳ we will $f(\theta) = \ell'(\theta)$ in L.R

But We will work with this example to build our understanding.



Maths for NM

(*) From above fig:-

$$\theta^{(1)} = \theta^{(0)} - \Delta$$

$$\otimes \quad f'(\theta^{(0)}) = \frac{f(\theta)}{\Delta}$$

$$\Delta = \frac{f(\theta^{(0)})}{f'(\theta^{(0)})}$$

$$\odot \quad \theta^{(1)} = \theta^{(0)} - \frac{f(\theta^{(0)})}{f'(\theta^{(0)})}$$

Generally

$$\theta^{(t+1)} = \theta^{(t)} - \frac{f(\theta^{(t)})}{f'(\theta^{(t)})}$$

→ Update Rule for NM.

Let $f(\theta) = l'(\theta)$

then:-

$$\theta^{(t+1)} = \theta^{(t)} - \frac{l'(\theta)}{l''(\theta)}$$

→ L.R Update by N.M.

⑧ Newton Methods enjoys the property of Quadratic Convergence:

→ After one iteration the error of margin from optimal value i.e. $f(\theta) = 0$ decrease quadratically:-

⑨

0.01 error	→	0.0001 error	→
0.00000001 error			

→ Reason why N.M requires fewer iteration as compare to L.R.

⑩ NM update Rule when θ is a Vector:

$$\theta^{(t+1)} = \theta^{(t)} + \underbrace{H^{-1}}_{\text{Matrix of Inverse}} \underbrace{\nabla_{\theta} \ell}_{\text{Vector of partial derivative}}$$

• where H is the Hessian matrix

• Hessian matrix is matrix of partial derivatives

$$H_{ij} = \frac{\partial^2 \ell}{\partial \theta_i \partial \theta_j}$$

In high-dimension Newton Method is however computationally expensive because each iteration requires inverting a large matrix

Rule of Thumb

→ Less parameters → N.M

→ More parameters → G.D