

Day-3  
June 17, 2024

Last Session:

Linear Regression:

$$h_{\theta}(x^{(i)}) = \theta^T x^{(i)}$$

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (y^{(i)} - h_{\theta}(x^{(i)}))^2 \quad \text{MSE Loss}$$

$$\nabla J(\theta) = \frac{1}{m} \sum_{i=1}^m (y^{(i)} - h_{\theta}(x^{(i)})) x_j^{(i)}$$

Gradient Descent:

$$\theta_j^{(t+1)} \leftarrow \theta_j^{(t)} - \eta \cdot \nabla_{\theta_j} J(\theta)$$

$$\theta_j^{(t+1)} \leftarrow \theta_j^{(t)} - \eta \cdot \frac{1}{m} \sum_{i=1}^m (y^{(i)} - h_{\theta}(x^{(i)})) x_j^{(i)}$$

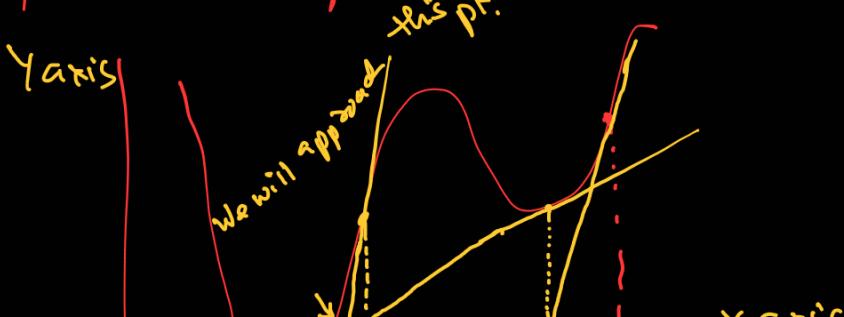
Normal Equations:

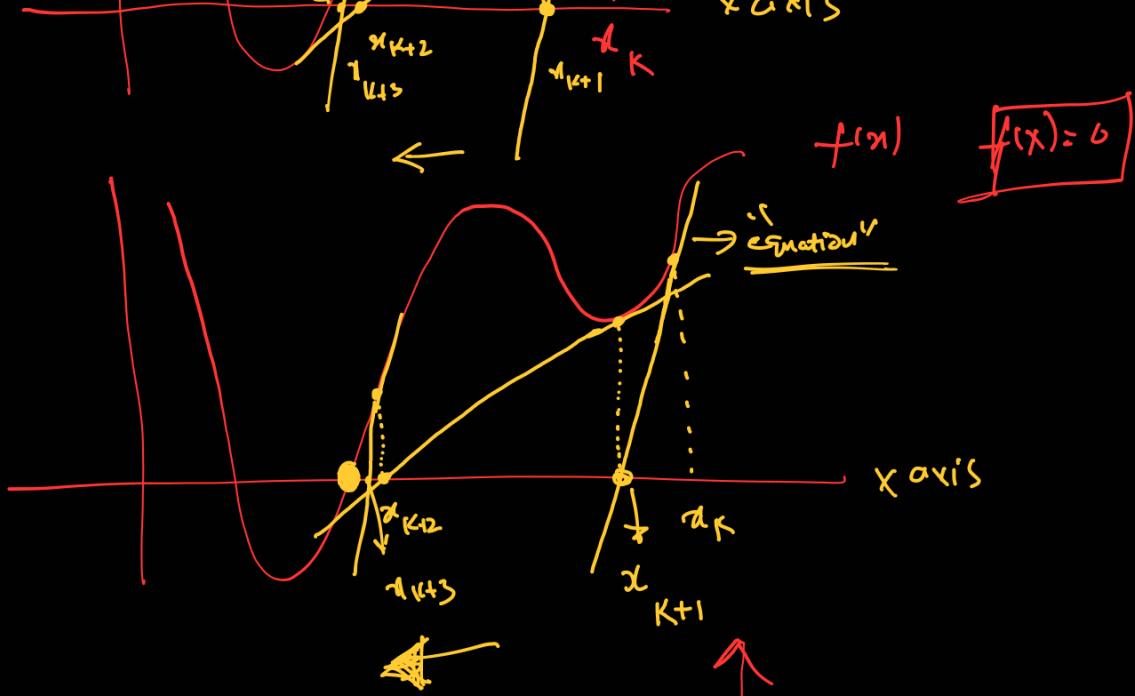
$$J(\theta) = \frac{1}{2m} (y - X\theta)^T (y - X\theta)$$

$$\boxed{\theta = (X^T X)^{-1} X^T y}$$

Newton's Method (Iterative method)

If we have find roots of an eqn  $f(x) = 0$





$$\boxed{x^{t+1} = x^t - \frac{f(x)}{f'(x)}} \quad \boxed{f'(x)=0} \quad \text{Newton's Update Rule}$$

## Optimization

We are not interested in  $f'(x)=0$ , rather we are interested in finding  $\boxed{f'(x)=0}$

$$f'(x) > 0 \rightarrow x^{t+1} = x^t - \frac{f(x)}{f'(x)}$$

$$\boxed{f'(x)=0 \rightarrow x^{t+1} = x^t - \frac{f'(x)}{f''(x)}} \quad \text{2nd order derivative term.}$$

In general, 2<sup>nd</sup> order methods are faster than 1<sup>st</sup> order methods. (X X X)

$$\text{Let } f: \mathbb{R}^n \rightarrow \mathbb{R}$$

+ the step mean a "less" using

Then given  $\alpha$ , we can approximate  $\min_{\mathbf{x}} f(\mathbf{x})$ .

Taylor's Theorem (Vector)

$$f(\mathbf{x}) \approx f(\alpha) + \mathbf{g}^T(\mathbf{x} - \alpha) + \frac{1}{2} (\mathbf{x} - \alpha)^T \mathbf{H} (\mathbf{x} - \alpha),$$

where  $\mathbf{g} = \nabla f(\alpha)$   
 $\mathbf{H} = \nabla^2 f(\alpha)$

Vector's

$$\rightarrow f(\mathbf{x}) \approx f(\alpha) + \nabla f(\alpha)(\mathbf{x} - \alpha) + \frac{\nabla^2 f(\alpha)}{2!} (\mathbf{x} - \alpha)^2 + \dots$$

$\downarrow$   
Scalar -

$$f(\mathbf{x}) = f(\alpha) + \mathbf{g}^T(\mathbf{x} - \alpha) + \frac{1}{2} (\mathbf{x} - \alpha)^T \mathbf{H} (\mathbf{x} - \alpha)$$

$$= f(\alpha) + \mathbf{g}^T(\mathbf{x} - \alpha) + \frac{1}{2} [(\mathbf{x}^T - \alpha^T) \mathbf{H} (\mathbf{x} - \alpha)]$$

$$= \frac{1}{2} [\mathbf{x}^T \mathbf{H} \mathbf{x} - 2\alpha^T \mathbf{H} \mathbf{x} + \alpha^T \mathbf{H} \alpha] + \mathbf{g}^T(\mathbf{x} - \alpha) + f(\alpha)$$

$$= \frac{1}{2} \mathbf{x}^T \mathbf{H} \mathbf{x} + b^T \mathbf{x} + c$$

where  $b = (\mathbf{g} - \mathbf{H} \alpha)$   
 $c = \text{left terms}$

I have to bring this form =  $\frac{1}{2} [\mathbf{x}^T \mathbf{H} \mathbf{x} - 2\alpha^T \mathbf{H} \mathbf{x} + \alpha^T \mathbf{H} \alpha] + \mathbf{g}^T(\mathbf{x} - \alpha) + f(\alpha)$

$$= \frac{\mathbf{x}^T \mathbf{H} \mathbf{x}}{2} - \frac{\alpha^T \mathbf{H} \mathbf{x}}{2} + \frac{\alpha^T \mathbf{H} \alpha}{2} + \mathbf{g}^T \mathbf{x} - \mathbf{g}^T \alpha + f(\alpha)$$

$$ax^2 + bx + c$$

$$= \frac{\mathbf{x}^T \mathbf{H} \mathbf{x}}{2} + \mathbf{g}^T \mathbf{x} - \frac{\alpha^T \mathbf{H} \mathbf{x}}{2} + \frac{\alpha^T \mathbf{H} \alpha}{2} - \mathbf{g}^T \alpha + f(\alpha)$$

$$= \frac{\mathbf{x}^T \mathbf{H} \mathbf{x}}{2} + \underbrace{(\mathbf{g}^T - \alpha^T \mathbf{H}) \mathbf{x}}_{\mathbf{T}} + \frac{\alpha^T \mathbf{H} \alpha}{2} - \mathbf{g}^T \alpha + f(\alpha) \quad (\alpha^T \mathbf{x})^T \\ (\mathbf{H} \alpha)$$

$$g(\mathbf{x}) = \frac{\mathbf{x}^T \mathbf{H} \mathbf{x}}{2} + (\mathbf{g} - \mathbf{H} \alpha)^T \mathbf{x} + \frac{\alpha^T \mathbf{H} \alpha}{2} - \mathbf{g}^T \alpha + f(\alpha)$$

2

$$\nabla g(\alpha) = 0$$

$$\frac{\nabla H\alpha + (g - H^T\alpha)}{2} = 0$$

$$H\alpha = H^T\alpha - g$$

$$\alpha = H^{-1}(H^T\alpha - g)$$

$$\alpha = H^{-1}(H\alpha - g)$$

$$\begin{cases} \alpha = \alpha - H^{-1}g \\ \tilde{\alpha} = \alpha - H^{-1}g \end{cases}$$

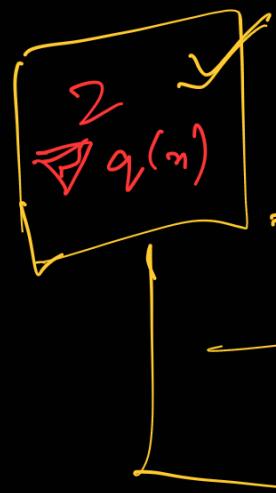
$$\theta^T A \theta = 2A\theta$$

$$\nabla_{\alpha} \alpha^T \alpha = \alpha$$

$$H = H^T$$

Symmetric matrix

$$\boxed{I \alpha = \alpha}$$



$$\nabla_{\alpha} g(\alpha) = H\alpha - (g - H^T\alpha)$$

$$\nabla_{\alpha} g(\alpha) = H$$

Not dependent on  $\alpha$   
Hessian Matrix

Hessian matrix  $\rightarrow$  That must be "Positive Semi" (PSD)  
Definite

$$\text{Vec} \geq \rightarrow \left( \begin{array}{c} Z^T H Z \geq 0 \\ \vdots \\ \alpha^T H \alpha \end{array} \right)$$

Fact

## Algorithm

$$\alpha_{t+1} = \alpha_t - \frac{f'(\alpha)}{f''(\alpha)} \rightarrow \nabla_{\theta} J(\theta)$$

$$f''(\alpha) \rightarrow \nabla_{\theta}^2 J(\theta)$$

Diff Gradient

Diff Gradient

Newton

$$\theta^{(t+1)} = \theta^{(t)} - (\nabla_{\theta}^2 J(\theta))^{-1} \nabla_{\theta} J(\theta)$$

GD:

$$\theta^{(t+1)} = \theta^{(t)} - \eta \cdot \nabla_{\theta} J(\theta)$$

## Probabilistic Interpretation of Linear Regression

In probabilistic interpretation, we are concerned how  $y^{(i)}$  are generated given by  $x^{(i)}$  parameterized by some  $\theta$ .

For Linear Regression,

$$y^{(i)} \approx \theta^T x^{(i)}$$



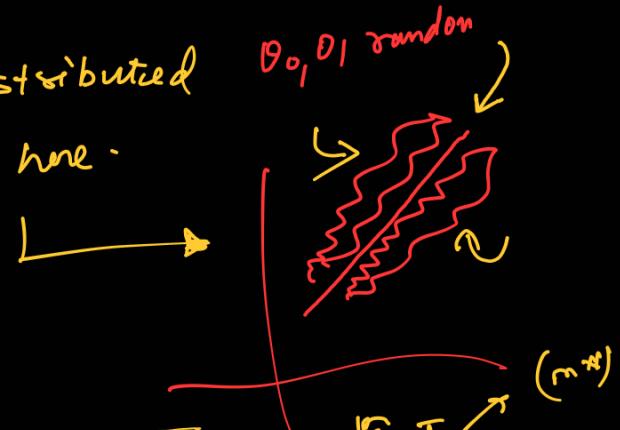
I am "expecting" that my data is distributed around (along) my model (i.e.  $\theta_0 + \theta_1 x$ ) here.

$$y^{(i)} \approx \theta^T x^{(i)}$$

approximation sign

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

"noise"



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

$\epsilon \sim N(\mu, \sigma^2)$

$$\theta: (n+1) \times 1$$

$$x: m \times (n+1)$$

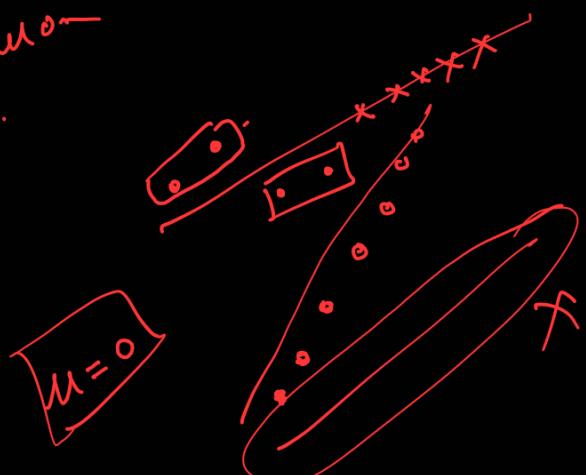
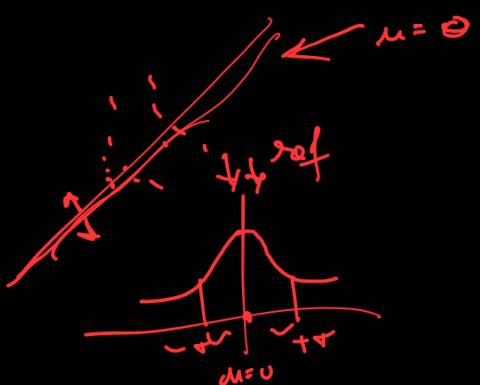
$$x\theta$$

$$[ ]^T \in \mathbb{R}^{(n+1) \times 1}$$

$$\underbrace{(m \times (n+1))}_{(m \times 1)} \leftarrow \text{Y dim}$$

$$y^{(i)} = \theta^T x^{(i)} + \epsilon \sim N(\mu, \sigma^2)$$

Now this becomes the eqn of how data is generated given some  $x^{(i)}$ . This is  $\theta$ , because we want on average the point must fall on the line.

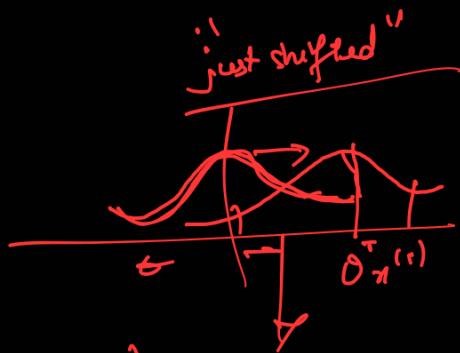


$$y^{(i)} = (\theta^T x^{(i)}) + N(\mu, \sigma^2)$$

This is interpreted as:

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

noise  $\mu = 0$



not a distribution of labels

Noise

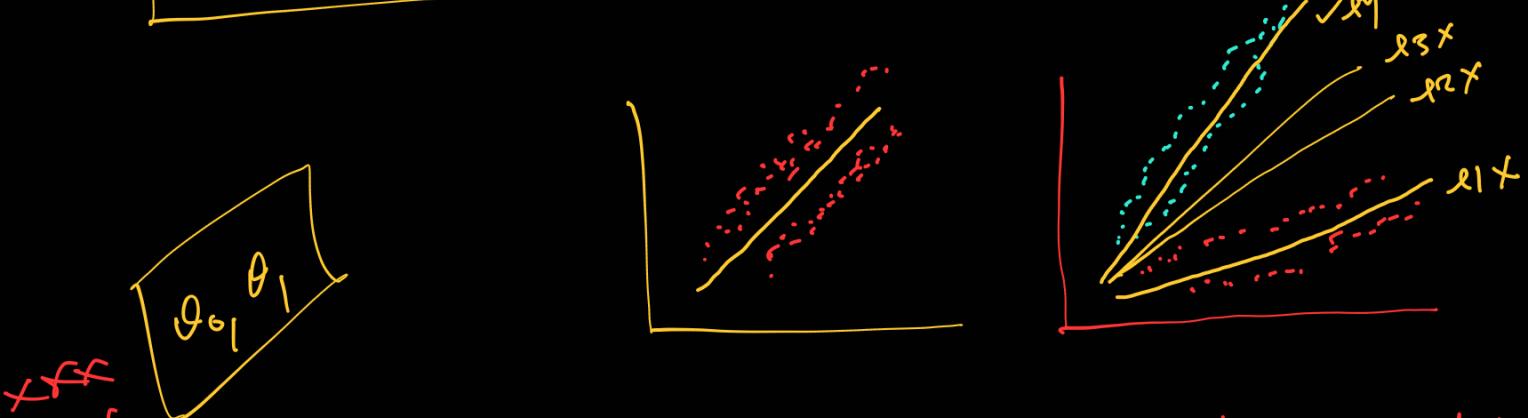
I am assuming that  $y^{(i)}$  is Normal with avg  $\theta_{\alpha}^{T(i)}$ , and some std dev. This also means that on average, my  $y^{(i)}$ 's should fall on the line  $\theta_{\alpha}^{T(i)}$

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

$$\cancel{P(y^{(i)} | x^{(i)}, \theta)} = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y^{(i)} - \theta_{\alpha}^{T(i)})^2}{2\sigma^2}\right)$$

Now we define the "LIKELIHOOD" of the data, given the model.

We are asking that "how likely" is this data generated by the given model.



We can perform "MLE" to find the model that best generates the data.  
(maximum likelihood estimation)

$$\theta_{MLE} = \underset{\theta}{\operatorname{argmax}} P(y^{(1)}, \dots, y^{(m)} | x^{(1)}, \dots, x^{(m)}, \theta)$$

from entire data

$$f(x) = 2x^2 + 3x + 4$$

max/min value of  $f(x)$

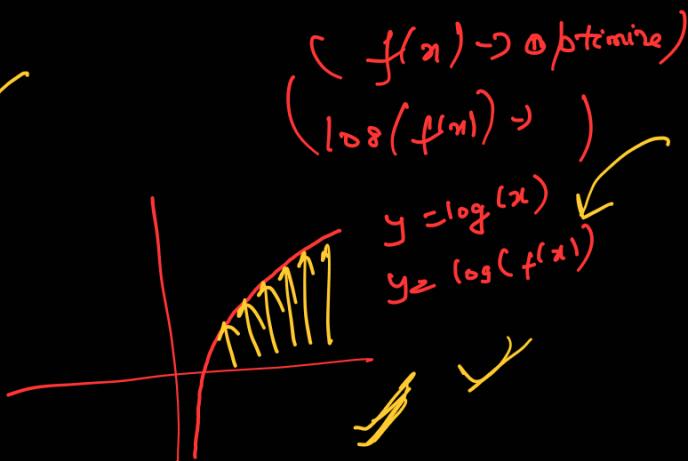
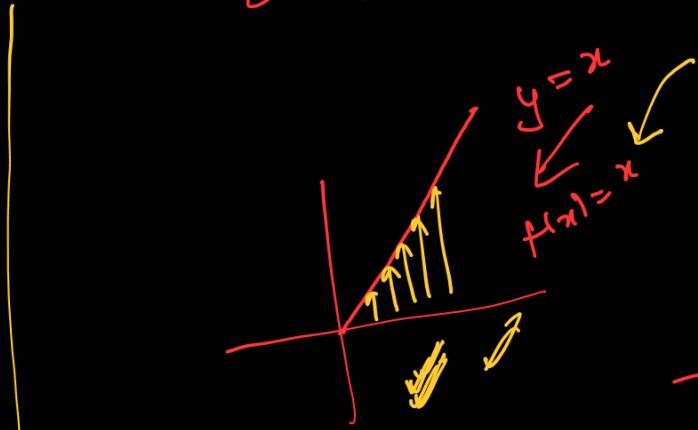
→ I am trying to find out that  $x_1$  which optimizes  $f(x)$  given  $-f'(x)$

$$\theta_{MLE} = \underset{\theta}{\operatorname{argmax}} \prod_{i=1}^m P(y^{(i)} | x^{(i)}, \theta)$$

Joint dist  
=  $P(y^{(1)}, y^{(2)}, \dots, y^{(n)})$   
 $= p(y^{(1)}), p(y^{(2)}), p(y^{(3)}) \dots$

$$\theta_{MLE} = \underset{\theta}{\operatorname{argmax}} \prod_{i=1}^m \left[ P(y^{(i)} | x^{(i)}, \theta) \right]$$

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



$x \rightarrow \text{maximize}$   
 $\partial \rightarrow \text{maximize}$

$y \rightarrow \text{maximize}$   
 $\log(y) \rightarrow \text{maximize}$

$$\theta_{MLE} \rightarrow \log(\theta_{MLE})_{\max}$$

$$\theta_{MLE} = L(\theta) = \underset{\theta}{\operatorname{argmax}} \prod_{i=1}^m P(y^{(i)} | x^{(i)}, \theta)$$

$$LL(\theta) = \underset{\theta}{\operatorname{argmax}} \log \prod_{i=1}^m P(y^{(i)} | x^{(i)}, \theta)$$

The log likelihood

$$\begin{aligned}
 L(\theta) &= \underset{\theta}{\operatorname{argmax}} \left[ \log \prod_{i=1}^m p(y^{(i)} | x^{(i)}, \theta) \right] \\
 &= \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^m \log p(y^{(i)} | x^{(i)}, \theta) \\
 &\therefore \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^m \log \left( \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}} \right) \\
 &= \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^m \left[ \log \left( \frac{1}{\sqrt{2\pi\sigma^2}} \right) + \log \left( e^{-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}} \right) \right] \\
 &= \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^m \log e^{-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}} \quad f(n), \frac{f(n)}{4} \\
 &\underset{\theta}{\operatorname{argmax}} \sum_{i=1}^m -\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2} \quad \text{constant} \\
 &\underset{\theta}{\operatorname{argmax}} \frac{1}{2\sigma^2} \sum_{i=1}^m -\frac{(y^{(i)} - \theta^T x^{(i)})^2}{m} \\
 &\underset{\theta}{\operatorname{argmax}} \frac{1}{\sigma^2} \times \frac{1}{m} \sum_{i=1}^m -\frac{(y^{(i)} - \theta^T x^{(i)})^2}{m} \\
 &\therefore J(\theta) = \frac{1}{m} \sum_{i=1}^m (y^{(i)} - \theta^T x^{(i)})^2 \quad \text{In good descent} \\
 &\qquad \qquad \qquad f(n), f(n), \frac{f(n)}{4} \\
 &\theta \rightarrow \text{minimize} \\
 &\underset{\theta}{\operatorname{argmax}} \left[ \frac{m}{\sigma^2} \times (-1) \times \frac{1}{m} \sum_{i=1}^m (y^{(i)} - \theta^T x^{(i)})^2 \right]
 \end{aligned}$$

$$\text{argmax } (-1)^k \sum_{i=1}^k \sum_{r=1}^R \text{(})$$

$$\text{argmax } (-1)^k J(\theta)$$

$$\text{argmax } -J(\theta)$$

If I  $\max -f(x)$ , then what I am effectively doing  $\rightarrow \min f(x)$

$$\text{argmax } -J(\theta)$$

$$\checkmark \boxed{\text{argmin } J(\theta)}$$

$$\boxed{\text{argmax } L(\theta) = \text{argmin } J(\theta)}$$

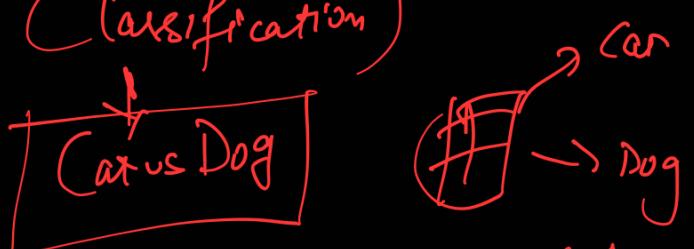
"We were right about Probabilistic Int of Lin Reg."

"& thus both approaches are interchangeable"

## Logistic Regression (Classification)

L-R  $\rightarrow$  continuous

[0.28, 79.67, 181821.229]



70.1.  $\rightarrow$  Cat

30.1.  $\rightarrow$  Dog

50.1.  $\rightarrow$  Dog

(100-50).1.  $\rightarrow$  Cat

Log Reg  $\rightarrow$  Labels depending how many categories I have.

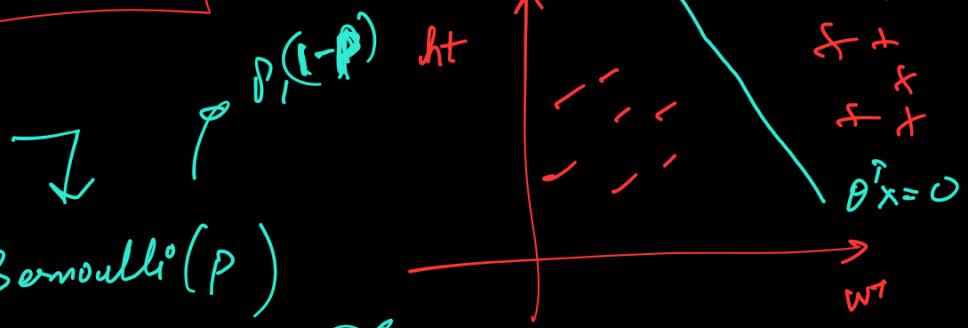
Binary Classification: 2 categories

$0.78 \rightarrow 1/0$

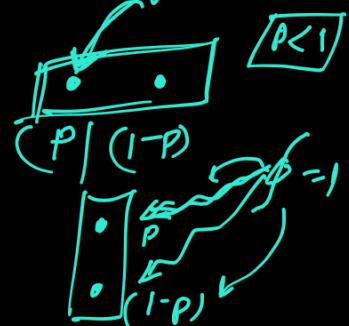
$$y = \{0, 1\}$$

$$y^{(i)} = \{0, 1\}$$

then  $y^{(i)} |_{x^{(i)}, \theta} \sim \text{Bernoulli}(\rho)$



Bernoulli  $\phi$



$$\begin{aligned} & nC_r p^r q^{n-r} \\ & (1-p) \end{aligned}$$

$$mC_m p^m (1-p)^{n-m}$$

w-l teams  $\rightarrow \text{Top } \ell$

$$T1, \dots, T8 \\ p_1, \dots, p_7, 1 - \sum_{i=1}^7 p_i$$

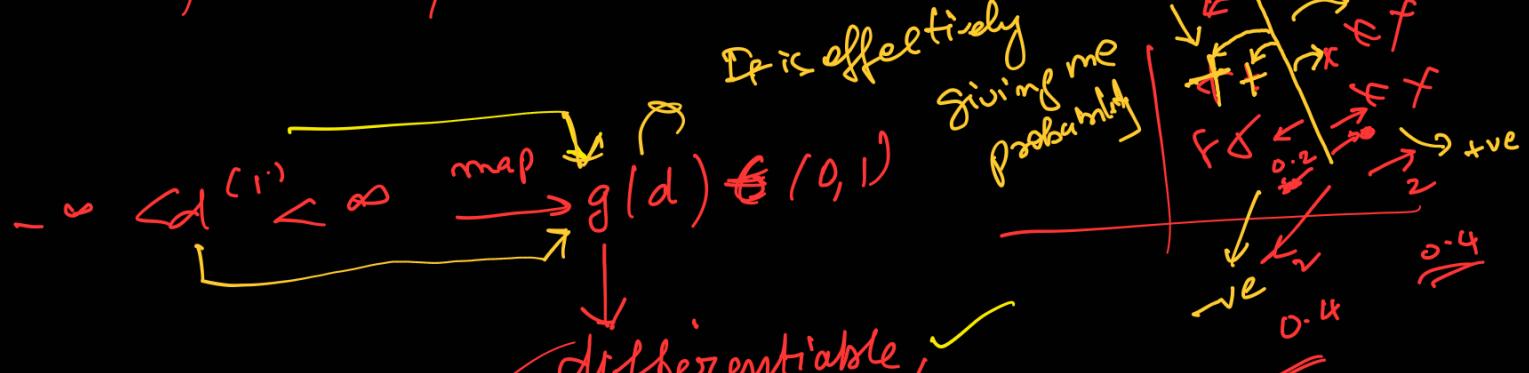
$\phi_m$

(1) 100 chocolates

$$\begin{array}{c} 40 \\ (0.4) \end{array} \quad \begin{array}{c} 20 \\ (0.1) \end{array} \quad \begin{array}{c} 50 \\ (0.5) \end{array}$$

$(y^{(i)} |_{x^{(i)}, \theta}) \sim \text{Bernoulli}(\rho)$

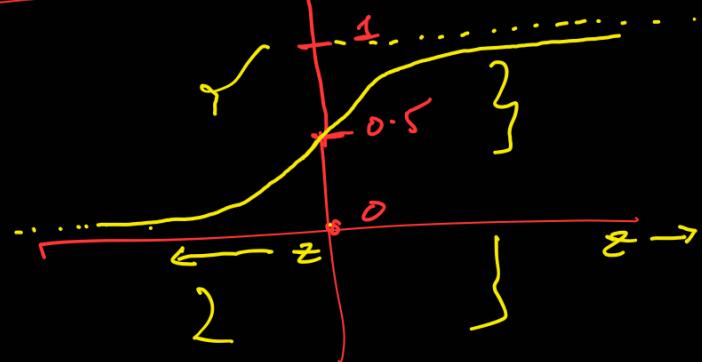
Let  $d^{(i)}$  be the dist b/w the line and the point  $(i)$ . Then ideally to have a high prob for a pt that is further from the line  $\Rightarrow$  vice-versa.



✓ differentiable, ✓  
 $(0, 1)$

✓  $z \rightarrow \infty \rightarrow g(z) = 1$  ✓  
 $z \rightarrow -\infty \rightarrow g(z) = 0$  ✓

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



$$d^{(i)} = \omega^T (\alpha^{(i)} - \hat{\alpha})$$

$$d^{(i)} = \omega^T \alpha^{(i)} - \boxed{\omega^T \hat{\alpha}}$$

$\boxed{d^{(i)} = \omega^T \alpha^{(i)} + b}$

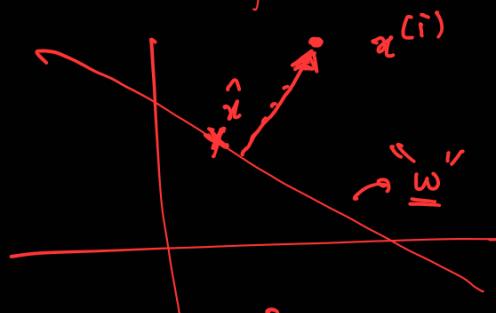
$$\boxed{\partial \omega^T \alpha^{(i)}}$$

$$\rightarrow \theta_0 x_0 + \theta_1 x_1 + \dots + \theta_n x_n$$

$$\boxed{\theta_0} [\theta_1, \theta_2, \dots, \theta_n]$$

$$\begin{bmatrix} \alpha_0 \\ \alpha_1 \\ \vdots \\ \alpha_n \end{bmatrix}$$

$$\theta_0 x_0 + \dots + \theta_n x_n + \theta_0 (\eta_0 = 1)$$



$$d^{(i)} = \omega^T (\alpha^{(i)} - \hat{\alpha})$$

fact:  $\nabla$

$$\nabla_{\omega} d^{(i)}$$

$$P(y^{(i)} = 1 | \alpha^{(i)}, \theta) = g(d^{(i)})$$

$$g(\nabla_{\omega} d^{(i)})$$

$$\nabla_z \sigma(z) = \sigma(z)(1 - \sigma(z))$$

Derive ↑↑↑↑↑↑  
Fact.

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

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

$$P(y^{(i)} = 0 | \alpha^{(i)}, \theta) = 1 - P(y^{(i)} = 1 | \alpha^{(i)}, \theta) = 1 - \frac{1}{1 + e^{-\nabla_{\omega} d^{(i)}}}$$

$$P(y^{(i)} \geq 1 | x^{(i)}, \theta) = \frac{1}{1 + e^{-\theta^T x^{(i)}}}$$

Bernoulli<sup>o</sup>:  $P(y=0) = \phi^{y} (1-\phi)^{1-y}$

$V_{\text{as}}(\phi) = \phi - \phi^y (1-\phi)^{1-y}$

Fact:  $\int \phi^y \cdot y = \phi$   
 $\int (1-\phi)^{1-y} \cdot y = 0$

log Likelihood:  $\log \text{Prob. You sum}$

$$\text{LL}(\theta) = \sum_{i=1}^m \log P(y^{(i)} | x^{(i)}, \theta)$$

$$\text{LL}(\theta) = \sum_{i=1}^m y^{(i)} \log P(y^{(i)} = 1 | x^{(i)}, \theta) + (1-y^{(i)}) \log (1 - P(y^{(i)} = 1 | x^{(i)}, \theta))$$

$$\begin{aligned} \partial_2 \text{LL}(\theta) &= \phi^y (1-\phi)^{1-y} \\ \log(L) &= \log \phi^y (1-\phi)^{1-y} \\ &= y \log \phi + (1-y) \log (1-\phi) \end{aligned}$$

$$\text{LL}(\theta) = \sum_{i=1}^m y^{(i)} \log \underbrace{P(y^{(i)} = 1 | x^{(i)}, \theta)}_{\sigma(\theta^T x^{(i)})} + (1-y^{(i)}) \log \underbrace{(1 - P(y^{(i)} = 1 | x^{(i)}, \theta))}_{1 - \sigma(\theta^T x^{(i)})}$$

$$\text{LL}(\theta) = \sum_{i=1}^m y^{(i)} \log \sigma(\theta^T x^{(i)}) + (1-y^{(i)}) \log (1 - \sigma(\theta^T x^{(i)}))$$

$$\begin{aligned} \nabla \text{LL}(\theta) &= \sum_{i=1}^m \left( \frac{y^{(i)}}{\sigma(\theta^T x^{(i)})} \cdot \sigma(\theta^T x^{(i)}) \cdot (1 - \sigma(\theta^T x^{(i)})) \cdot x^{(i)} \right) \\ &\quad + (1-y^{(i)}) \cdot [D - \sigma(\theta^T x^{(i)}) / (1 - \sigma(\theta^T x^{(i)})) \cdot x^{(i)}] \end{aligned}$$

$$(1 - \tau(\theta^T x^{(i)}))$$

$$\begin{cases} g_2 = \log \frac{1}{1 - \tau(x)} \\ g'(x) = \frac{1}{\tau(x)} \cdot \tau'(x) \end{cases}$$

$$\nabla_{\theta} LL(\theta) = \sum_{i=1}^m y^{(i)} (1 - \tau(\theta^T x^{(i)})) x^{(i)} + (-1)(1 - \tau(\theta^T x^{(i)})) \tau(\theta^T x^{(i)}) x^{(i)}$$

$$\nabla_{\theta} LL(\theta) = \sum_{i=1}^m x^{(i)} \left( y^{(i)} - y^{(i)} \tau(\theta^T x^{(i)}) + (y^{(i)} - \tau(\theta^T x^{(i)})) \tau'(\theta^T x^{(i)}) \right)$$

$$\nabla_{\theta} LL(\theta) = \sum_{i=1}^m x^{(i)} \left( y^{(i)} - \tau(\theta^T x^{(i)}) \right) \overbrace{\theta_0(x^{(i)})}^{h_0(x^{(i)})}$$

$$\nabla_{\theta} LL(\theta) = \sum_{i=1}^m (y^{(i)} - h_0(x^{(i)})) x^{(i)}$$

MSE loss derivative

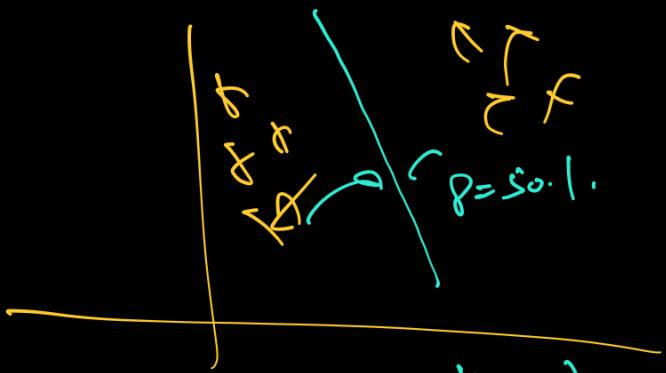
Logistic Regression and Linear Regression

have similar properties

$$\theta_j \leftarrow \theta_j^{(t+1)} - \eta \frac{\nabla J(\theta)}{\partial \theta}$$

$$\max_{\theta} LL(\theta) = \min_{\theta} J(\theta)$$

Generalized Class of Linear Models (GLM)



$$P(y^{(i)}=1) = \frac{1}{2} = P(y^{(i)}=0)$$

$$\frac{1}{1+e^{-\theta^T x^{(i)}}} = 1 - \frac{1}{1+e^{-\theta^T x^{(i)}}}$$

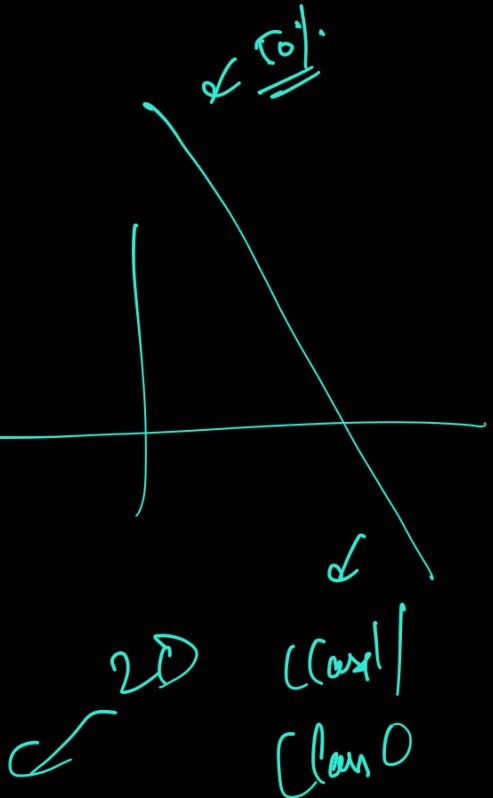
$$\frac{1}{1+e^{-\theta^T x^{(i)}}} = 1$$

$$1+e^{-\theta^T x^{(i)}} = 2$$

$$e^{-\theta^T x^{(i)}} = 1 \approx e^0$$

$$\boxed{\theta^T x^{(i)} = 0}$$

$$\theta_0 + \theta_1 x$$

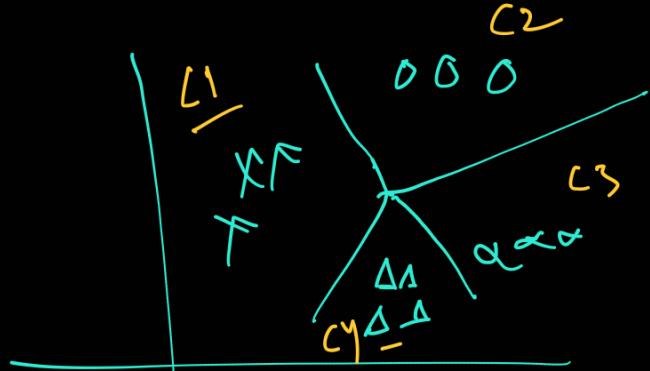


And this is the eqn of sep line.

$$\text{None kin} \quad p_{\text{soft}} = f_{\text{soft}}$$

1, 2, 3, 4  
1, 2, 3, 4  
1, 2, 3, 4

1, 2, 3, 4  
 $(z, r^u)$



2 vs 3

2 vs

Multi Class

WS (regression)

$$P(y=1) = \frac{\overset{100}{\cancel{4}}}{100} = P(y=4)$$

$$P(y=1) = 0.25 = P(y=4)$$

$$P(y=2) = \frac{\overset{100}{\cancel{4}}}{(100)} = P(y=3)$$