

Doubt - Class-2

"2 degrees"

$$\rightarrow f(\mathbf{x}) \approx f(\mathbf{x}') + \mathbf{g}^\top (\mathbf{x} - \mathbf{x}') + \frac{1}{2} (\mathbf{x} - \mathbf{x}')^\top H (\mathbf{x} - \mathbf{x}')$$

Why H is
Symmetric

$$\mathbf{g} = \nabla_{\mathbf{x}} f(\mathbf{x}')$$

$$H = \nabla_{\mathbf{x}}^2 f(\mathbf{x}')$$

?

\hookrightarrow Symmetric

\mathbf{x} - vector



$$f(\mathbf{x}) = (1 \times 1)$$

$$\begin{bmatrix} n_1 & n_2 & \dots & n_m \end{bmatrix} \xrightarrow{\left[\begin{array}{cccc} \frac{\partial f(\mathbf{x})}{\partial x_1} & \frac{\partial f(\mathbf{x})}{\partial x_2} & \dots & \frac{\partial f(\mathbf{x})}{\partial x_n} \end{array} \right]}$$

$$\sum \nabla_{\mathbf{x}} f(\mathbf{x}) =$$

$$\begin{bmatrix} n_1 \\ n_2 \\ \vdots \\ n_m \end{bmatrix}$$

$$\left[\begin{array}{cccc} \frac{\partial f(\mathbf{x})}{\partial x_1} & \frac{\partial f(\mathbf{x})}{\partial x_2} & \dots & \frac{\partial f(\mathbf{x})}{\partial x_n} \end{array} \right]$$

$\underline{D \times 1}$

$\mathbf{I} \times \mathbf{D}$

$$\left[\begin{array}{cccc} \frac{\partial f(\mathbf{x})}{\partial x_1} & \frac{\partial f(\mathbf{x})}{\partial x_2} & \dots & \frac{\partial f(\mathbf{x})}{\partial x_D} \\ \frac{\partial f(\mathbf{x})}{\partial x_1} & \frac{\partial f(\mathbf{x})}{\partial x_2} & \dots & \frac{\partial f(\mathbf{x})}{\partial x_D} \\ \frac{\partial f(\mathbf{x})}{\partial x_1} & \frac{\partial f(\mathbf{x})}{\partial x_2} & \dots & \frac{\partial f(\mathbf{x})}{\partial x_D} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f(\mathbf{x})}{\partial x_1} & \frac{\partial f(\mathbf{x})}{\partial x_2} & \dots & \frac{\partial f(\mathbf{x})}{\partial x_D} \end{array} \right]$$

$\underline{D \times 1 \times 1 \times D}$

$$= \boxed{D \times D}$$

$$\boxed{1 \times D}$$

$\boxed{(1 \times 1)}$

$$\begin{aligned}
 & \left(\frac{\partial f(x)}{\partial x_1}, \frac{\partial f(x)}{\partial x_2}, \dots, \frac{\partial f(x)}{\partial x_D} \right) \\
 & \xrightarrow{\text{matrix form}} \frac{\partial f(x)}{\partial x} \\
 & \quad \left(\frac{\partial f(x)}{\partial x_1}, \frac{\partial f(x)}{\partial x_2}, \dots, \frac{\partial f(x)}{\partial x_D} \right) \\
 & \quad \xrightarrow{\text{matrix form}} \frac{\partial f(x)}{\partial x} = \frac{\partial f(x)}{\partial x_i} \frac{\partial x_i}{\partial x_j} = \frac{\partial f(x)}{\partial x_i \partial x_j}
 \end{aligned}$$

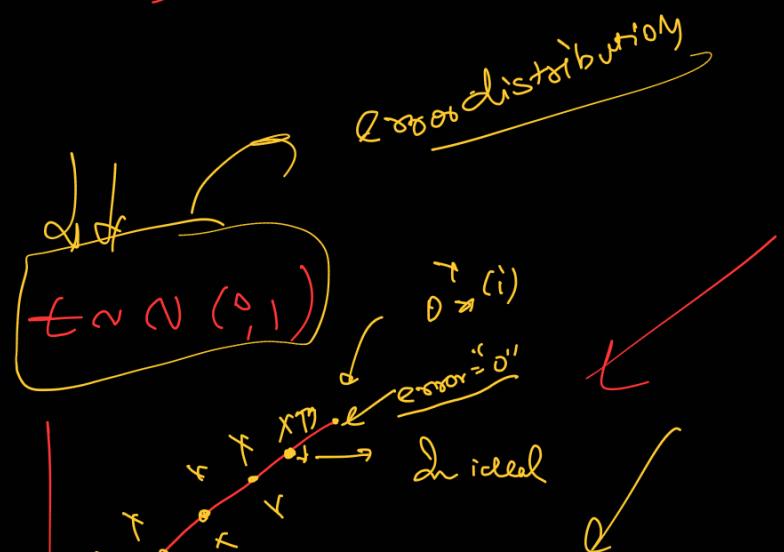
$$\begin{aligned}
 y^{(i)} &= \theta^T x^{(i)} + \epsilon \quad \rightarrow \epsilon \sim N(\theta, \sigma^2) \\
 &\quad \hookrightarrow \text{why normal dist?}
 \end{aligned}$$

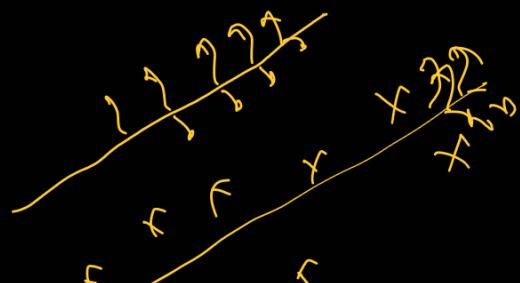
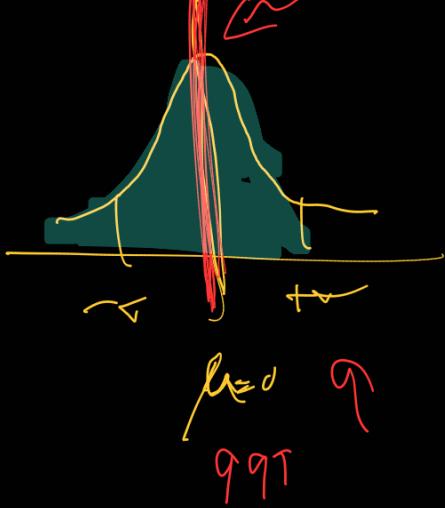
doubt

When we do not have any idea of underlying distribution over some variables, we take it as Normal distn because N.D is "universal" in nature and any phenomenon can be explained by using N.D as proxy.

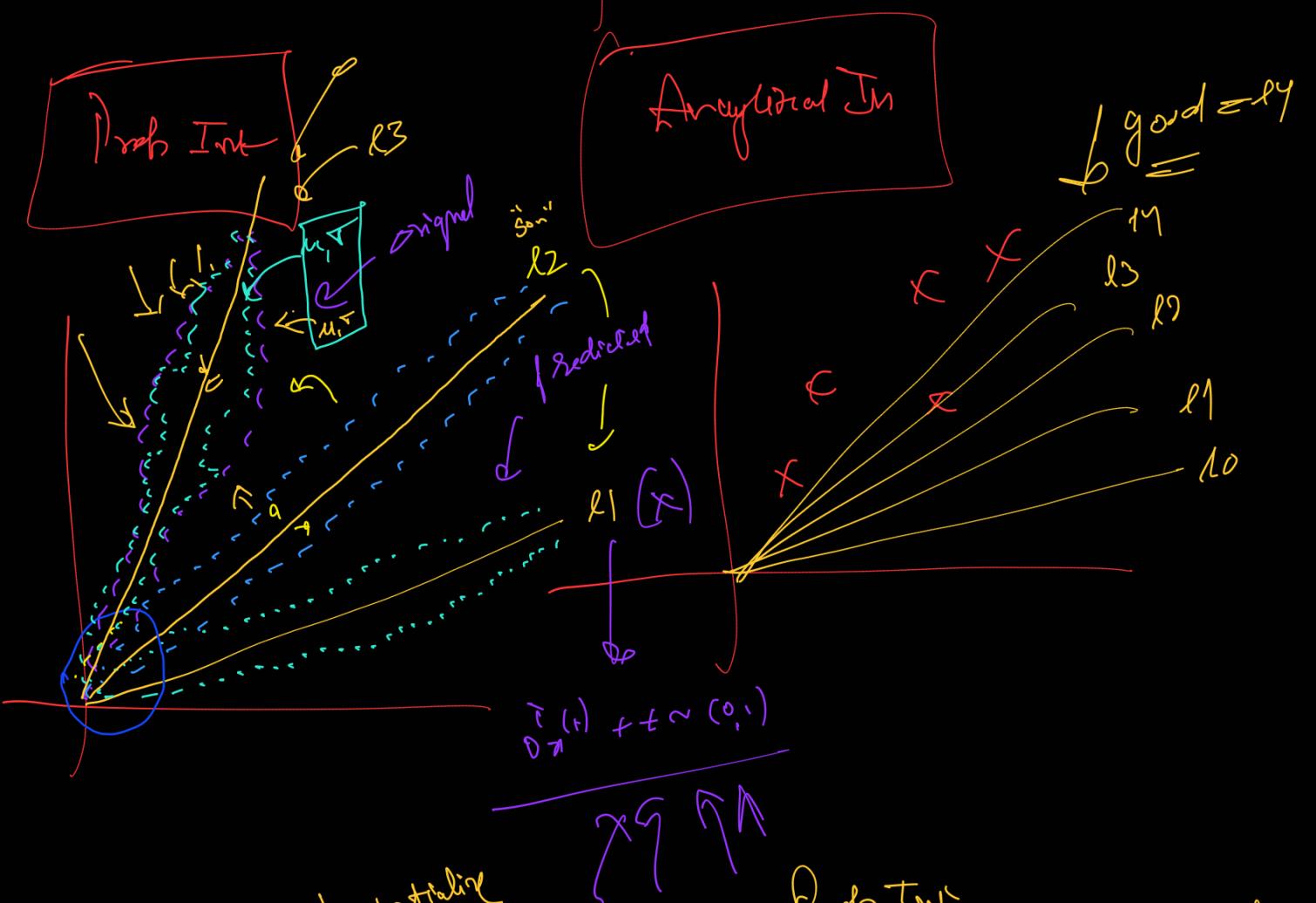
→ Why $\mu = 0$?

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





25000 la mean
20000 lignin
10000 hair



$D \rightarrow$ randomly initialize

$\theta > 0$

$\hat{\theta}_{MLE}$

D-T-M

[Signature]

→ Ex. 11.1.1.

why EM

→ I am trying to learn
from π of actual data
distribution

From \checkmark \checkmark , efficiency
 $\hat{\theta}$ same as $\hat{\theta}_{MLE}$ (can say)
 $\hat{\theta}_{MLE}$ \leftarrow $\text{maximum likelihood}$
 $\sum \text{on residuals}$ \rightarrow sum of
 sum \leftarrow sum

~~Part 3) estimate σ^2 data $\hat{\sigma}^2$ using~~
~~likelihood maximize $\sum \sqrt{y_i}$~~
 $\hookrightarrow \hat{\theta}_{MLE} //$

$\hat{\theta}_{MLE}$ \checkmark eqⁿ $\frac{\partial L}{\partial \theta}$ \rightarrow $\hat{\theta}_{MLE}$ |

$$L = \phi^y (1-\phi)^{(1-y)}$$

Bernoulli distribution

$\xrightarrow{\text{Binomial}}$ $\xrightarrow{\text{Repeated Bernoulli Events}}$

Distribution

$$n \choose r \times \phi^r (1-\phi)^{n-r} = 1$$

In total it should occur

$$\text{Bernoulli} = 1 = \phi$$

$$\phi^y (1-\phi)^{1-y}$$

$$\rightarrow \begin{bmatrix} 0 \\ 0 \end{bmatrix} \leftarrow y \quad \rightarrow \begin{bmatrix} 1 \\ 1 \end{bmatrix} \leftarrow (1-y)$$

$$\left(\phi^y (1-\phi)^{1-y} \right)^m$$

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

I am trying to maximize the likelihood of all the points inside my training data.

equivalent

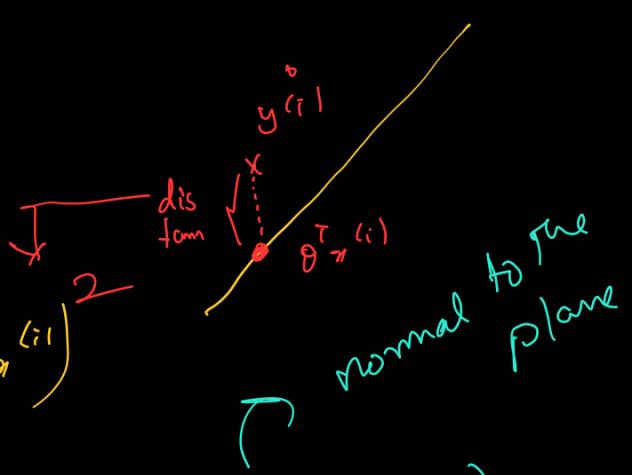
Minimizing $J(\theta)$

Gradient Descent

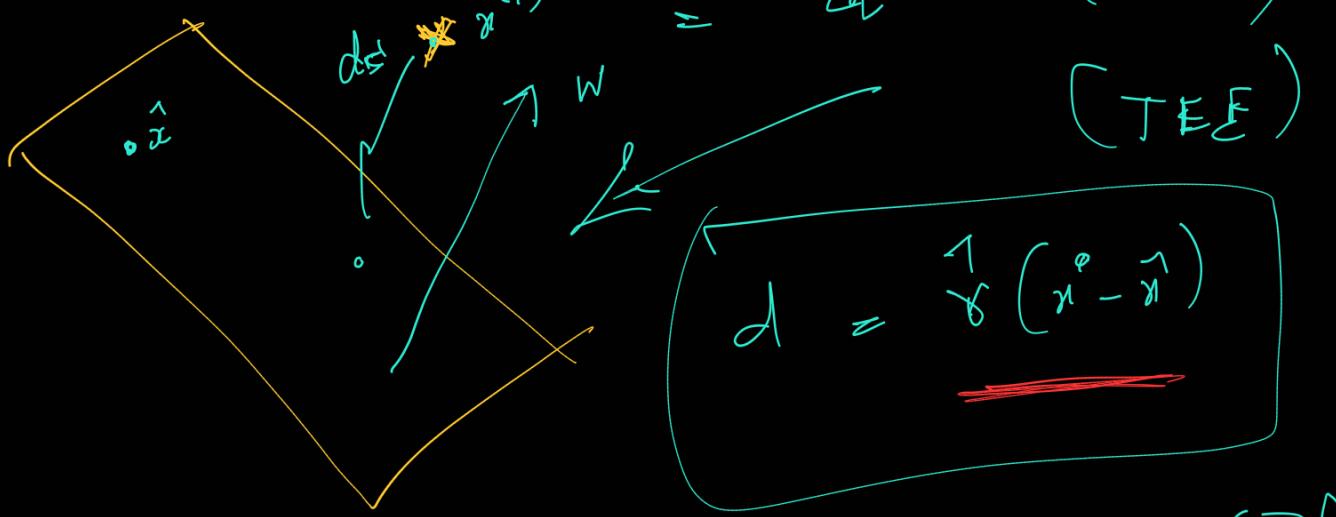
5 datapoints

doubt. Why distance?

$$(y^{(i)} - \theta^T x^{(i)})^2$$



$$d^{(i)} = w^T (x^{(i)} - \bar{x}) e$$



Vector 3D \rightarrow \mathbb{R}^1 (JFF)

$$d^{(i)} = w^T (x^{(i)} - x)$$

$$x^{(i)} - x$$

In will tel jom

$$J(\theta) = (y - x^\theta)^T (y - x^\theta)$$

$$y$$

$$= (y^T - \theta^T x^T) (y - x^\theta)$$

$$J(\theta) = 2\theta^T x^T = y^T (y - x^\theta) - \theta^T x^T (y - x^\theta)$$

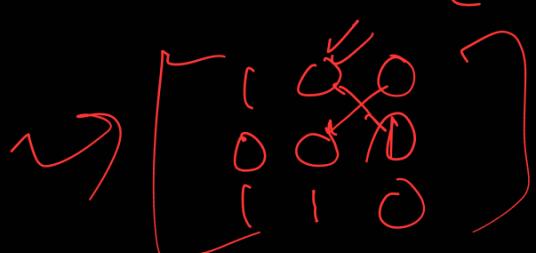
$$\left. \begin{aligned} \theta^T x^T \theta &= \theta \\ \theta^T x^T x^T \theta &= \theta \end{aligned} \right\} \quad y^T y - \underbrace{y^T x^\theta}_{\text{cancel}} - \underbrace{\theta^T x^T y}_{\text{cancel}} + \theta^T x^T x^\theta \theta$$

$$y^T y - \underbrace{y^T x^\theta}_{\text{cancel}} - \underbrace{y^T x^\theta}_{\text{cancel}} + \theta^T x^T x^\theta \theta$$

$$y^T y - \underbrace{y^T x^\theta}_{\text{cancel}} - \underbrace{y^T x^\theta}_{\text{cancel}} + \theta^T x^T x^\theta \theta$$

X : feature matrix

$$\begin{matrix} (m \times n) & (n \times m) \\ \downarrow & \downarrow \\ (m \times m) \end{matrix}$$



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

