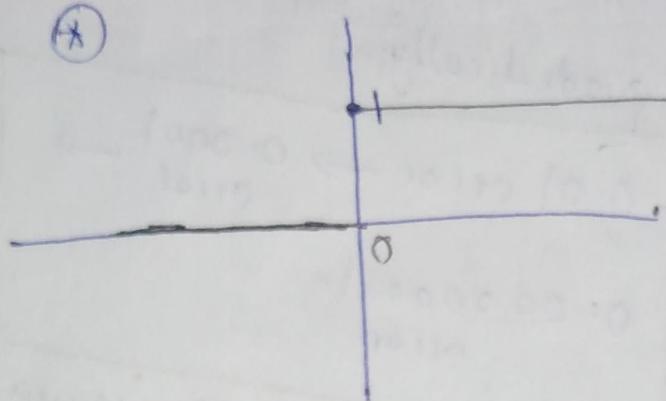


(24)

Lecture-4

* Perceptron Algorithm:-



$$g(z) = \begin{cases} + & z \geq 0 \\ - & z < 0 \end{cases}$$

→ Hard version of sigmoid.

$$\rightarrow h_s(x) = g(\theta^T \cdot x).$$

Update Rule for Perceptron algorithm.

④ $\boxed{\theta_j := \theta_j + \alpha (y^{(i)} - h_s(x^{(i)})) \cdot x_j^{(i)}}$

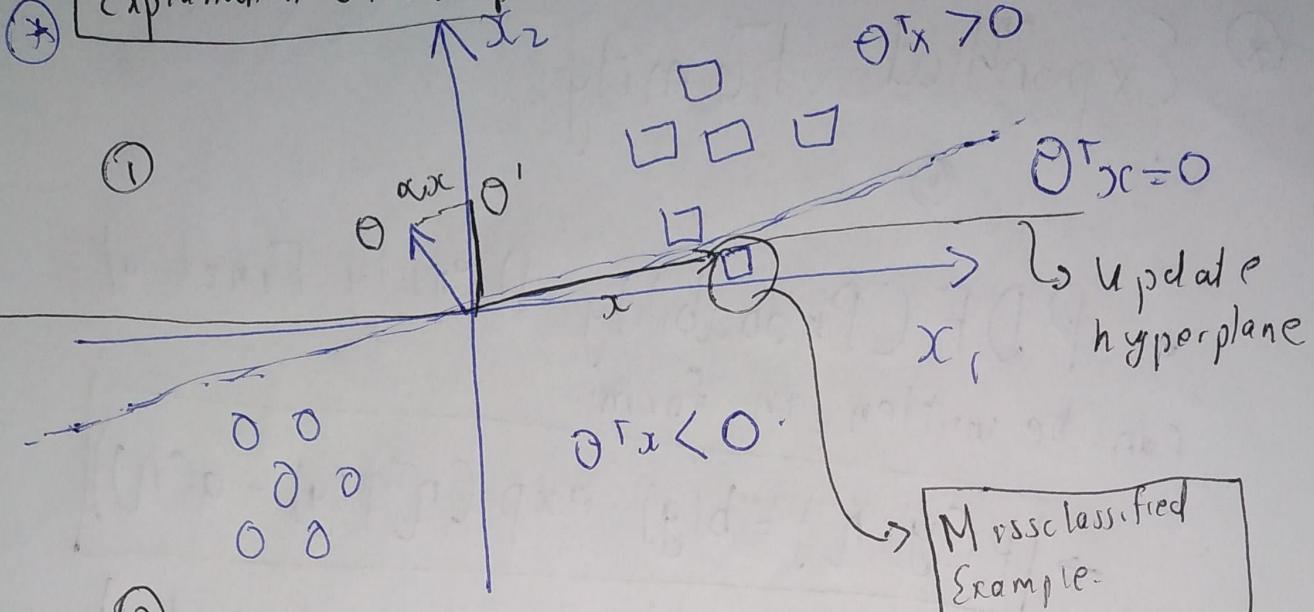
→ Same as logistic Regression.

• $y^{(i)} - h_s(x^{(i)})$ would be zero if correctly classify otherwise it would be negative.

28

Explanation of P.A.

①



②

$$\approx \theta \approx x \mid y=1$$

$$\approx \theta \not\approx x \mid y=0$$

θ and x dot product would be positive if they are similar otherwise

it would be negative therefore we want θ' to

be close to x as shown in the figure (in case of positive misclassification)

examples would be classified correctly).

③ Remember the update rule if example is correctly classified then we do nothing otherwise if example is positive and we misclassified it negative then we add to vector θ (See figure)

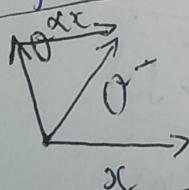
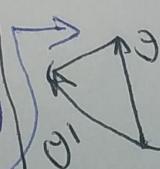


Fig.

Increasing similarity

④ Otherwise if example is negative meaning label is 0 and we classified it as positive (+) then we subtract vector



Decreasing similarity

* Exponential Family:

• PDF (Probability Density Function) can be written in form:-

$$\rightarrow P(y; \eta) = b(y) \cdot \exp(\eta^T T(y) - a(\eta))$$

* y - Data (output of data)

Can be a vector
 $T(y) = y$ scalar
 Assume

- η - natural parameters
- $T(y)$ - sufficient statistics
- $b(y)$ - base measure (function of y)
- $a(\eta)$ - log-partition (function of η)

↳ Reason why it's called

log-partition

$$P(y; \eta) = b(y) \cdot \frac{\exp(\eta^T T(y) - a(\eta))}{\exp(a(\eta))}$$

$$\text{II III} \quad \frac{\exp(n^T T(y))}{\exp(a(\eta))}$$

(Q7)

- ⑧ Partition function is a technical to indicate normalization constant of probability distribution. s.t whole equation $p(y; \theta)$ integrates to one.

- ⑨ For an choice of a, b, T a, long as expression integrate to one you have exponent exponent exponential family for different probability distribution

- ⑩ Algebraic massaging of PDF to bring to this form: if yes then it is exponential family.

~ Examples of showing PDFs belonging to exponential family:-

① Bernoulli Distribution [Binary outcome]:-

ϕ = probability of event]

$$P(y; \phi) = \phi^y \cdot (1-\phi)^{1-y} \quad (\text{take this to exponential family form}).$$

$$= \log \exp(\log(\phi^y \cdot (1-\phi)^{1-y}))$$

$$= \exp \left[y \log \left(\frac{\phi}{1-\phi} \right) + \log(1-\phi) \right]$$

28

$$\begin{aligned}
 &= \exp(y \log \phi + \log(1-\phi) - y \log(1-\phi)) \\
 &= \exp(y \log \phi - y \log(1-\phi) + \log(1-\phi)) \\
 &= \exp(y \log \frac{\phi}{1-\phi} + \log(1-\phi))
 \end{aligned}$$

Now comparing to the
original equation

- $b(y) = 1$
- $T(y) = y$ (as expected)
- $\gamma = \log\left(\frac{\phi}{1-\phi}\right) \Rightarrow \phi = \frac{1}{1+e^{-\gamma}}$???
- $a(n) = -\log(1-\phi)$
 $\Rightarrow -\log\left(1 - \frac{1}{1+e^{-\gamma}}\right)$
 $= \log(1+e^{-\gamma})$

→ This verifies Bernoulli is exponential family.

(29)

⑦ Example of Gaussian
(with fixed variance)

• Assume $\sigma^2 = 1$

$$P(y; \mu) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(y-\mu)^2}{2}\right)$$

$$= \underbrace{\frac{1}{\sqrt{2\pi}} e^{-\frac{y^2}{2}}}_{b(y)} \exp\left(\underbrace{\mu y - \frac{1}{2}\mu^2}_{\eta(\mu y) a(\mu)}\right)$$

$$\boxed{b(y) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{y^2}{2}\right)}$$

$$\eta(y) = y$$

$$\eta = \mu$$

$$a(\mu) = \frac{-\mu^2}{2} = \frac{\mu^2}{2}$$

→ Parameter matching
for Gaussian.

Example of Exponential family

Real - Gaussian
Binary - Bernoulli
Count - Poisson

R+ Gamma, exponential
Distrn - Beta, n. Distribution.

(30)

④ Properties of exponential family:-

- a):- MLE wrt $\eta \Rightarrow$ concave.
- b) \Rightarrow NLL wrt $\eta \Rightarrow$ convex

(*)

$$\begin{cases} \text{b)}: E[y|\eta] = \frac{\partial}{\partial \eta} a(\eta) \\ \text{c)}: \text{Var}[y|\eta] = \frac{\partial^2}{\partial \eta^2} a(\eta) \end{cases}$$

No Integrals.

⑤ Generalized Linear Model

(GLM):-

→ Assumption of Design choice:-

(i) $y|x; \theta \sim \text{Exponential Family}(n)$.(ii) $\eta = \theta^T x$ (Design choice) $\hookrightarrow \theta \in \mathbb{R}^n, x \in \mathbb{R}^n$ (iii) Test time: Output $E[y|x; \theta]$ \hookrightarrow Mean of the distribution $\Rightarrow h_\theta(x) = E[y|x; \theta]$.

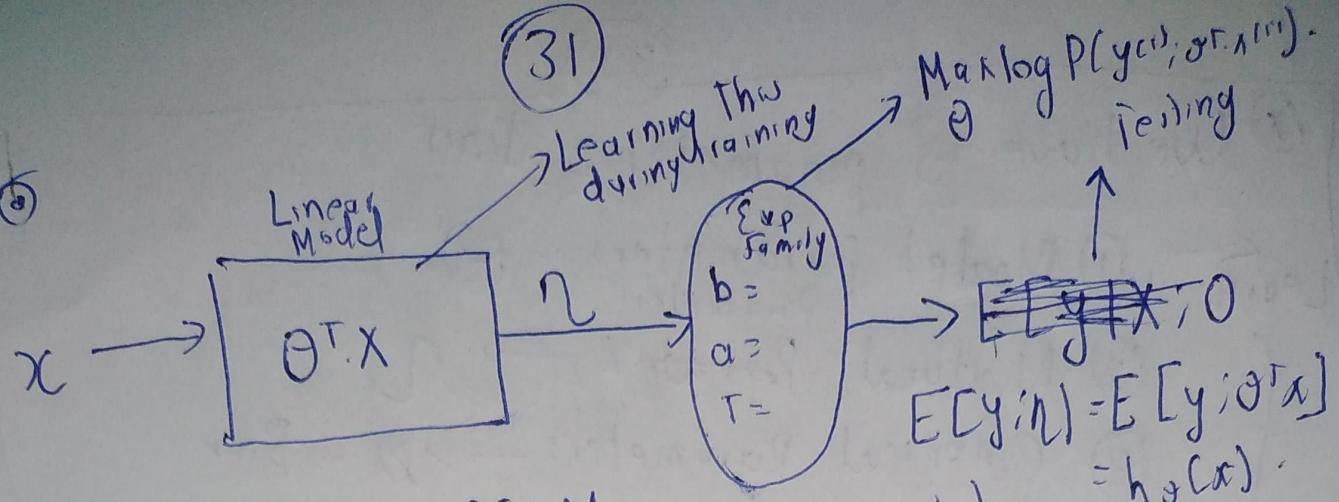


Fig GLM and Exponential family -
(Overall Summary)

④

GLM Training:-

④ Learning update Rule is
Same for all GLMs:-

$$\theta_j := \theta_j + \alpha (y^{(i)} - h_{\theta}(x^{(i)})) \cdot x_j^{(i)}$$

④ Terminology:

η - natural parameter

$$\mu = E[y; \eta] = g(\eta) \rightarrow \boxed{\text{Canonical Response function.}}$$

$$\eta = g^{-1}(\mu) \rightarrow \boxed{\text{Canonical Link function}}$$

Also note:

$$g'(n) = \frac{\partial}{\partial n} a(n)$$

* We have 3 parameterization

Learn \leftarrow ① Model parameters - Θ

② Natural Parameters $\rightarrow \gamma$

③ Canonical Parameters $\rightarrow \phi$ - Ber

$\phi \sigma^2 \rightarrow$ Gaussian

$\lambda \rightarrow$ Poisson

④ Design choice:-

$$\gamma = \theta^T x \rightarrow g(\gamma) \rightarrow \lambda$$

Fig

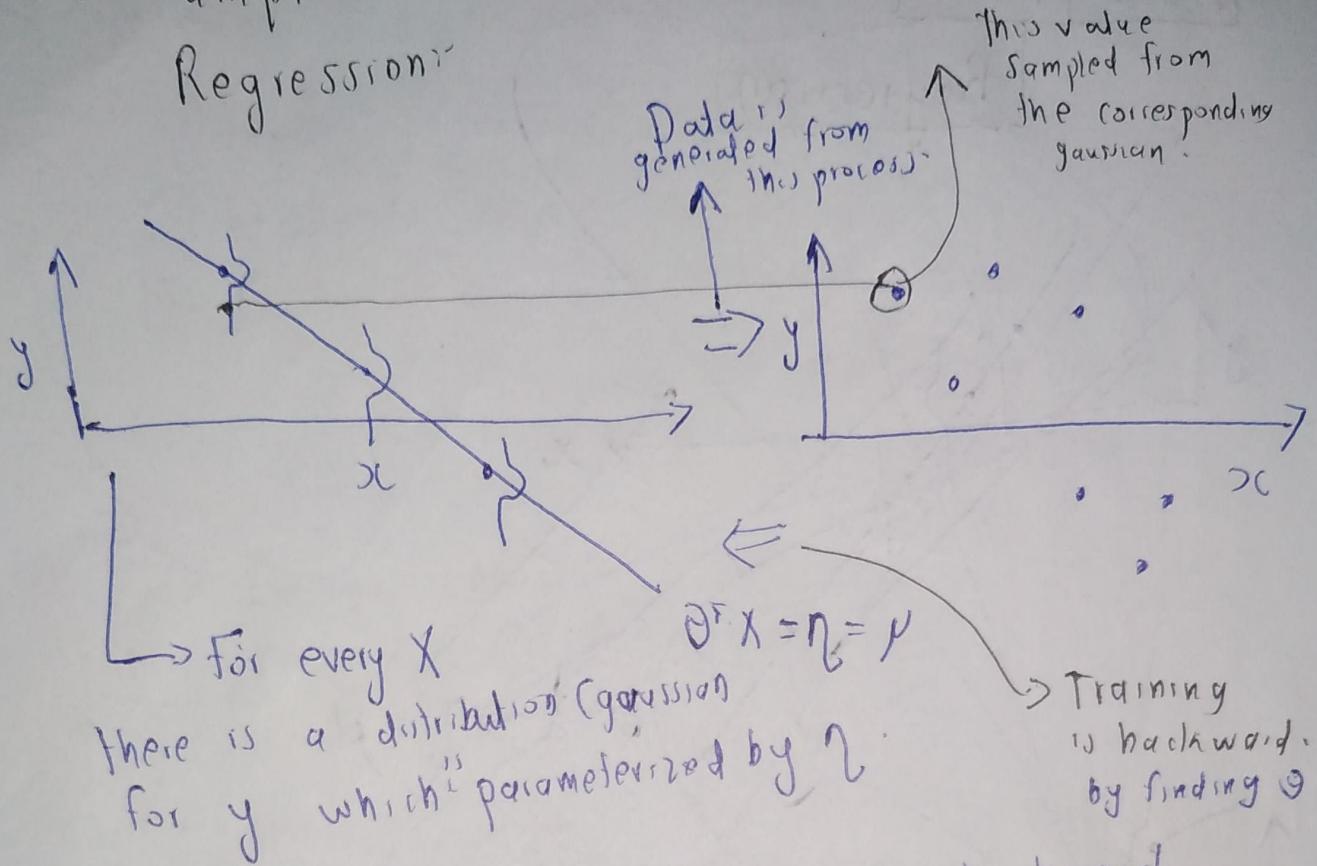
* For Logistic Regression:

$$h_{\theta}(x) = E[y|x; \theta] = \frac{1}{1+e^{-\theta^T x}}$$

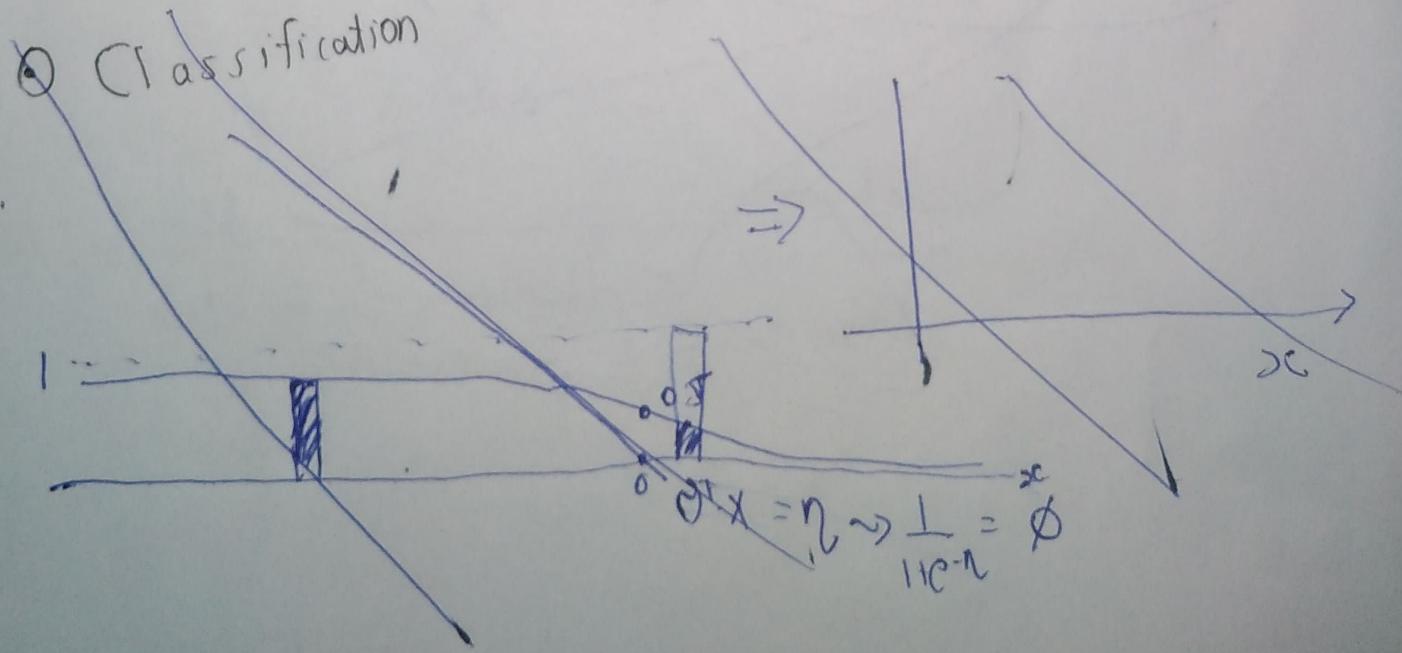
* HLM is a way to model
data \rightarrow Depending on the task
You will choose different models.

(33)

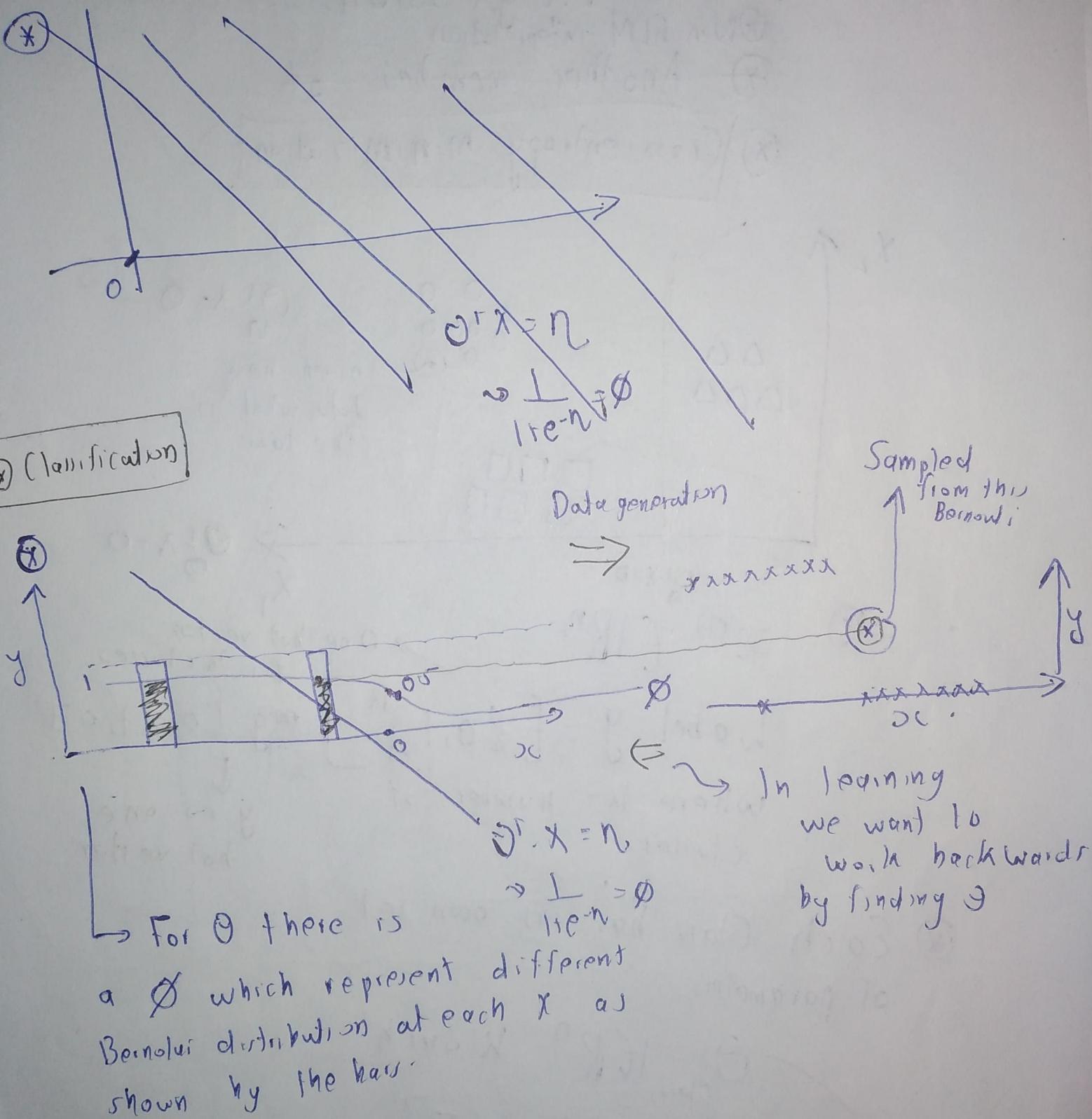
• Assumption:
Regression:



- With HLM we want to work backward the parameters from data to pre-modelling which generated the data.



(39)



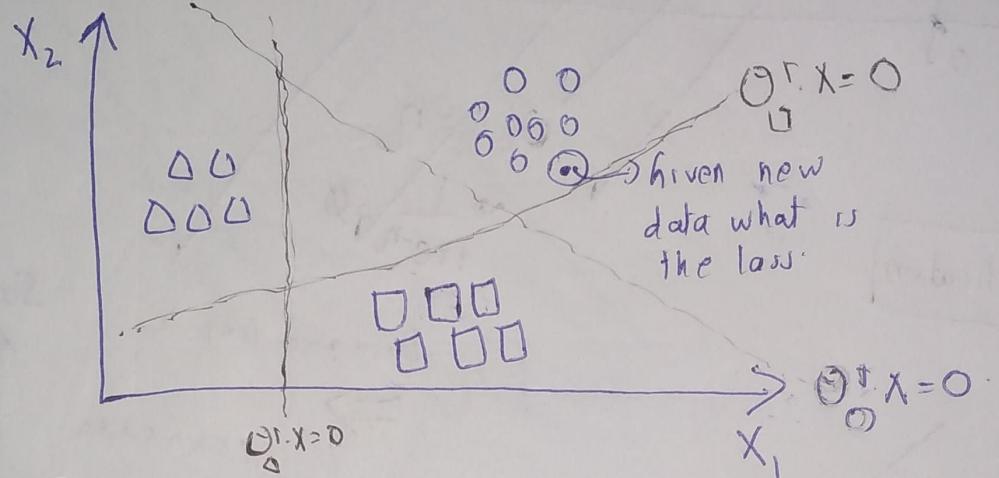
(35)

* Softmax Regression:-

* Non-HLM interpretation

* Another member of

* Cross-entropy minimization



Given new
data what is
the loss.

* $x^{(i)} \in \mathbb{R}^n$

One-hot vector
for K classes.

Label $y = [\{0, 1\}^K]$ e.g. $[0, 0, 1, 0]$

where $K = \text{number of}$
classes

$y \rightarrow$ one
hot vector

* Each class has its own set
of parameters.

$\theta_{\text{class}} \in \mathbb{R}^n$ K such

class $\in \{0, 0, 1\}$.

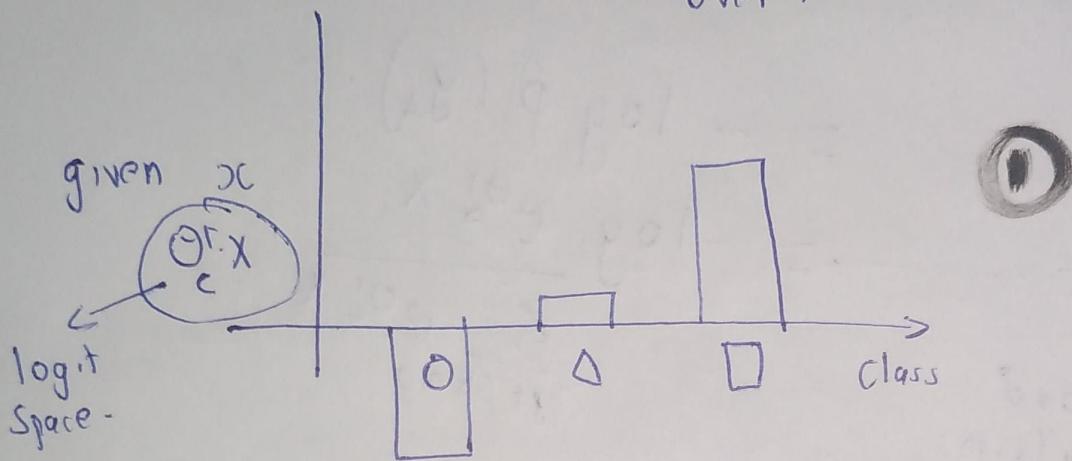
Can Represent as $\xrightarrow{\text{matrix}} K \times n$ matrix

$$\begin{bmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_K \end{bmatrix}$$

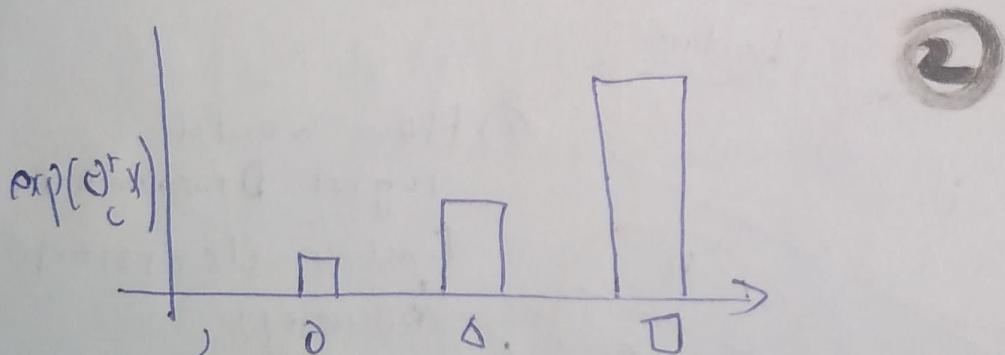
* Generalization
of logistic
Regression.

(36)

- ① Our goal is to take these parameters and find ~~chart~~ of new example probability distribution over the classes.

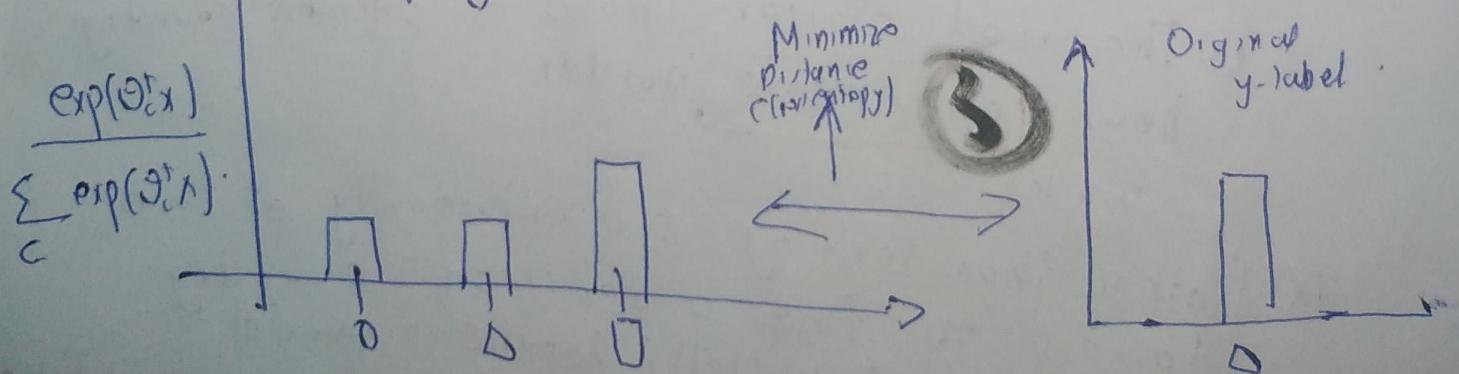


→ exp logits



Normalize

$\hat{P}(y|x) \rightarrow$ Our hypothesis. (Probability over classes)



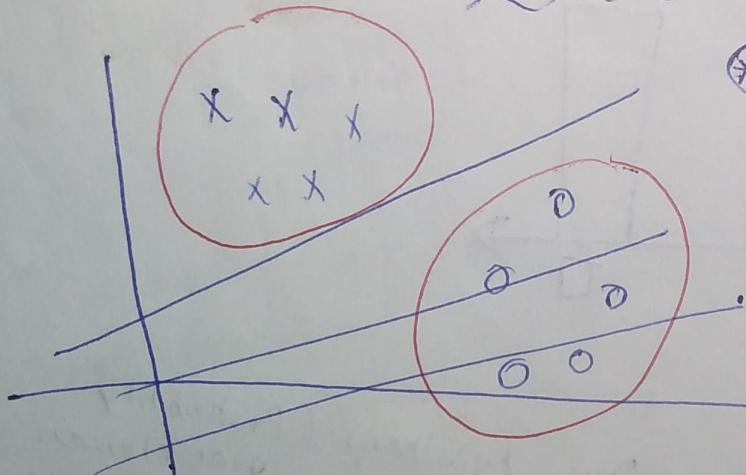
(37)

$$\text{Cross-entropy}(p, \hat{p}) = -\sum_{y \in C} p(y) \log \hat{p}(y)$$

$$= -\log \hat{p}(y_A) \\ = -\log \frac{e^{\theta_A^\top \cdot X}}{\sum_{y \in C} e^{\theta_y^\top \cdot X}}$$

Treat this
as a loss and
you perform
gradient descent
w.r.t parameters.

Lecture-5



- ④ How would logistic Discriminative find a decision boundary?

↳ **Discriminative**

model works by finding a decision boundary.

- ⑤ The red cluster shows how **generative** border works.

- ⑥ Rather than looking simultaneously at two models, find a decision boundary