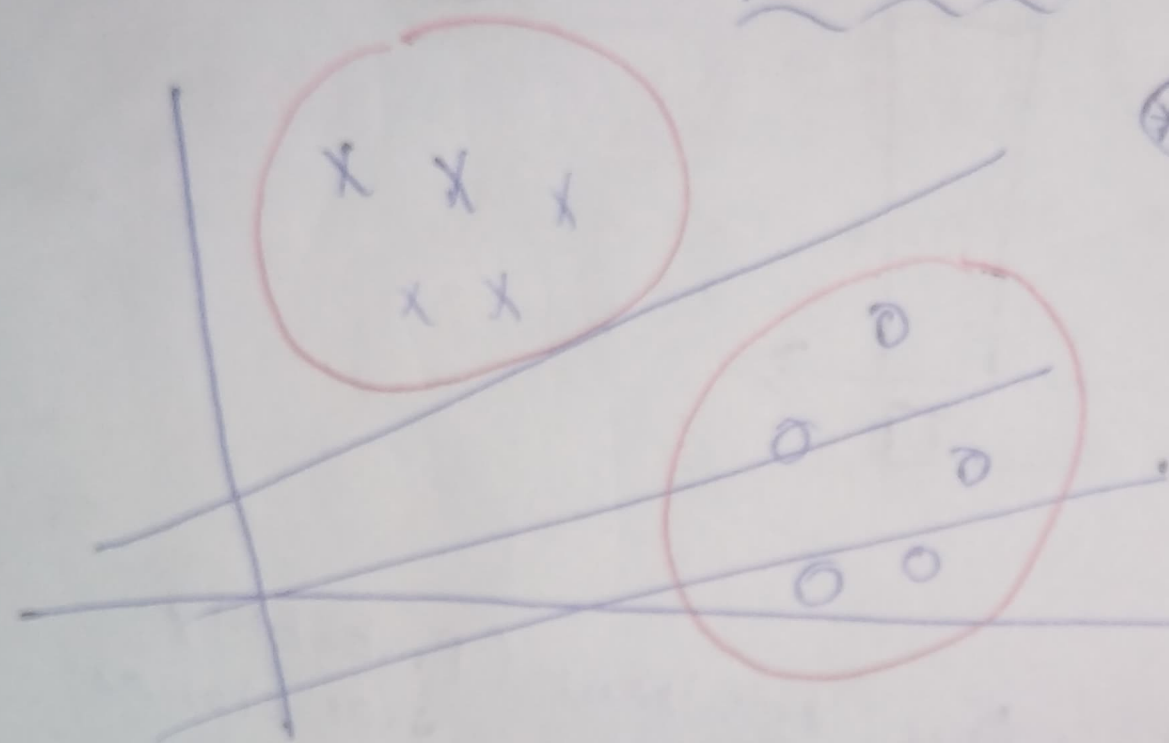


## Lecture-5



\* How would  
logistic Discrimination  
find a decision  
boundary?

↳ Discriminative  
model works  
by finding a  
decision boundary.

\* The read cluster shows  
how generative border  
works.

\* Rather than looking simultaneously at two  
classes to find a decision boundary  
(Discriminative model); generative model  
builds a model of what of the classes look like

# \* Discriminative

Learn  $P(Y|X)$

$$\text{Cor } \downarrow h_g(x) = \begin{cases} 0 \\ 1 \end{cases} \text{ directly}$$

## \* Generative learning algorithm:

learns  $P(X|Y)$

Features

class.

Also learns

$P(Y)$

class prior.

## \* Bayes Rules (GLA) $\rightarrow$ Testing Times

$$P(y=1|x) = \frac{P(X|y=1) \cdot P(y=1)}{P(X)}$$

$$P(X) = P(X|y=1) \cdot P(y=1) + P(X|y=0) \cdot P(y=0)$$

## \* Gaussian discriminant Analysis (GDA):

- Suppose  $x \in \mathbb{R}^n$  (drop  $x_0=1$  convention).
- Assume  $P(X|Y)$  is Gaussian.

(39)

• What is multi-variate gaussian?

$$z \sim \mathcal{N}(\vec{\mu}, \Sigma)$$

$$z \in \mathbb{R}^n, \mu \in \mathbb{R}^n, \Sigma \in \mathbb{R}^{n \times n}$$

$$E[z] = \mu$$

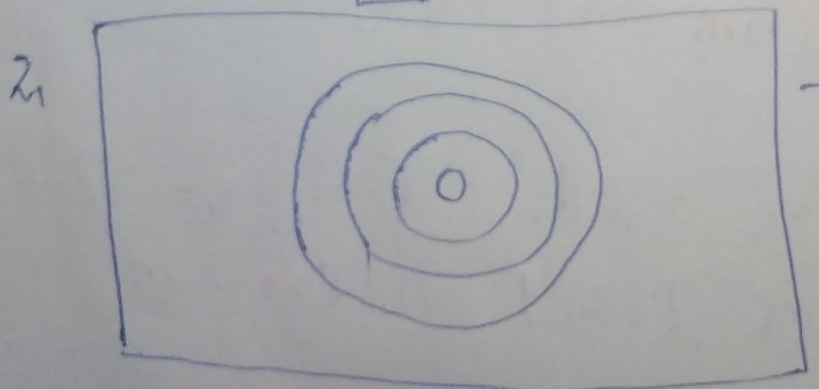
$$\text{cov}(z) = E[(z - \mu)(z - \mu)^T] \\ = E[zz^T - (Ez)(Ez)^T]$$

$$P(z) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(x - \mu)^T \Sigma^{-1}(x - \mu)\right)$$

PDF for multi-variate gaussian

Contour for Multivariate gaussian

$$\rightarrow \Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}; \mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$



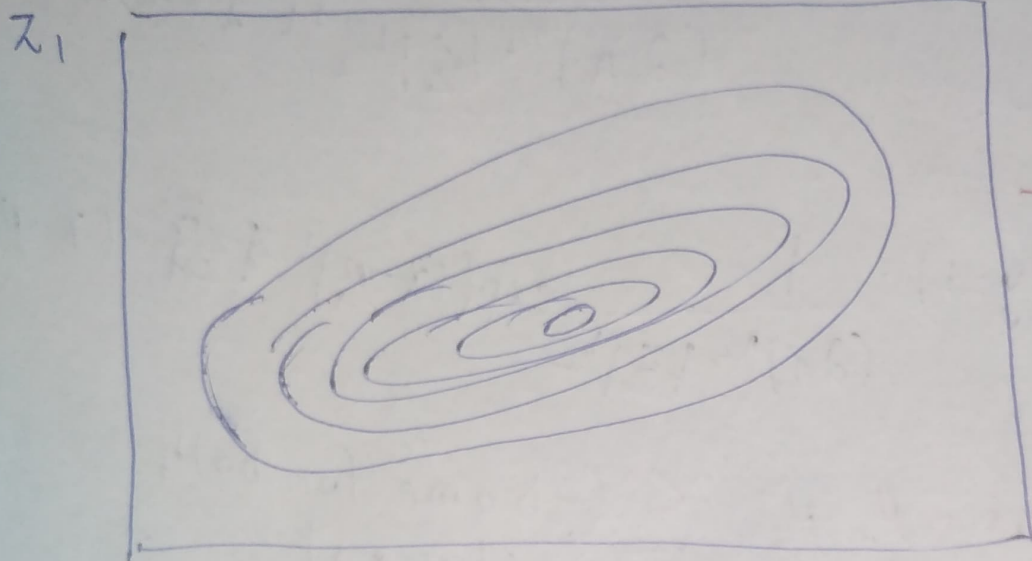
→ Perfectly Rounded Circles.

Fig.



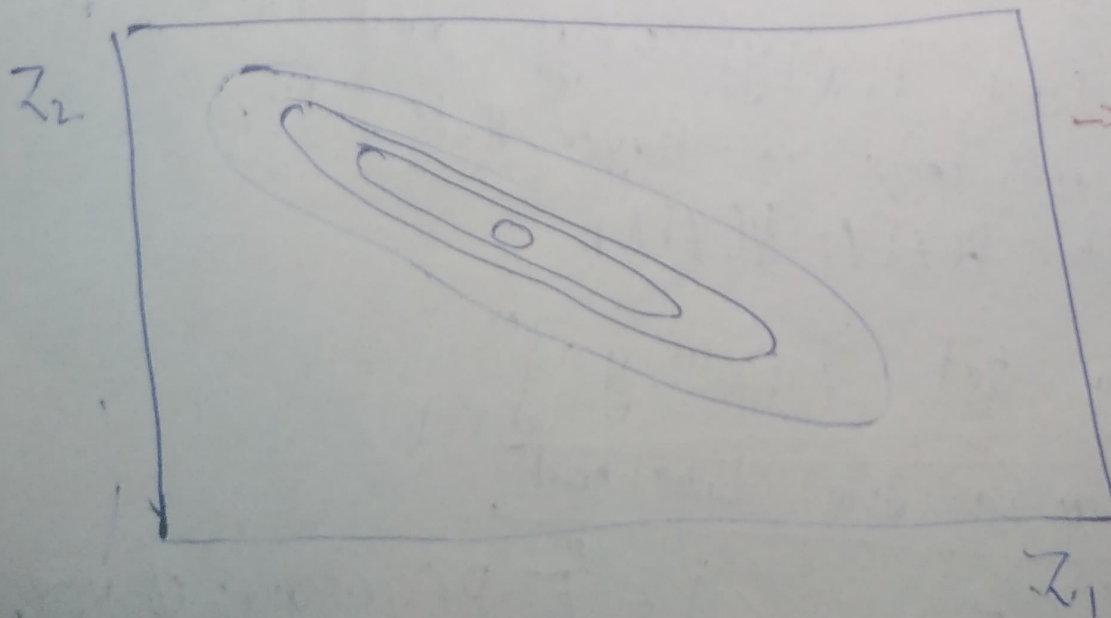
(40)

(\*)  $\Sigma = \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1 \end{bmatrix}$ ;  $\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$



→  $z_1$  and  $z_2$  being positively correlated.

(\*)  $\Sigma = \begin{bmatrix} 1 & -0.5 \\ -0.5 & 1 \end{bmatrix}$ ;  $\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$



→  $z_1$  and  $z_2$  being negatively correlated.

④ GDA model:-

$$④ P(x|y=0) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \cdot \exp\left(-\frac{1}{2}(x-\mu_0)^T \Sigma^{-1}(x-\mu_0)\right)$$

$$④ P(x|y=1) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \cdot \exp\left(-\frac{1}{2}(x-\mu_1)^T \Sigma^{-1}(x-\mu_1)\right)$$

Parameters:  $\mu_0, \mu_1, \Sigma$  (Same for both classes)

Need to find these.

~~1/0~~

$$④ P(y) = \phi^y (1-\phi)^{1-y} \text{ (Bernoulli variable)}$$

④ We find these parameter and then find  $P(x|y)$ ,  $P(y)$  and then will use these in Bayes Rule to predict  $P(x|y) \cdot P(y|x)$ .

④ Training set  $\{(x^{(i)}, y^{(i)})\}_{i=1}^m$

④ Maximize Joint Likelihood:-

$$\mathcal{L}(\phi, \mu_0, \mu_1, \Sigma) = \prod_{i=1}^m P(x^{(i)}, y^{(i)}; \phi, \mu_0, \mu_1, \Sigma)$$



(A2)

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

→ But Remember for Discriminative model we will maximizing Conditional likelihood:-

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

→ Maximum Likelihood Estimation for LDA:-

$$\begin{aligned} \max_{\phi, \mu_0, \mu_1, \Sigma} &= \mathcal{L}(\phi, \mu_0, \mu_1, \Sigma) \\ &= \log \left( \prod_{i=1}^m P(x^{(i)} | y^{(i)}) \cdot P(y^{(i)}) \right) \end{aligned}$$

$$\rightarrow \nabla_{\phi, \mu_0, \mu_1, \Sigma} \mathcal{L}(\phi, \mu_0, \mu_1, \Sigma) = 0$$

↳ You derivative for  $\mathcal{L}(\dots)$  w.r.t to each of the parameter and set to zero. If you do this then for each of the parameter, the following equation pops-up:-

$$\textcircled{1} \quad \phi = \frac{\sum_{i=1}^m (y^{(i)})}{m} = \frac{\sum_{i=1}^m \mathbb{I}\{y^{(i)} = 1\}}{m}$$

Indicator notation (true if condition holds else False).

(2)

$$\mu_0 = \frac{\sum_{i=1}^m \mathbb{I}\{y^{(i)} = 0\} \cdot X^{(i)}}{\sum_{i=1}^m \mathbb{I}\{y^{(i)} = 0\}}$$

$$\sum_{i=1}^m X^{(i)} \quad \text{or} \quad \sum_{i=1}^m \mathbb{I}\{y^{(i)} = 0\}$$

Same explanation.

(3)

$$\mu_1 = \frac{\sum_{i=1}^m \mathbb{I}\{y^{(i)} = 1\} \cdot X^{(i)}}{\sum_{i=1}^m \mathbb{I}\{y^{(i)} = 1\}}$$

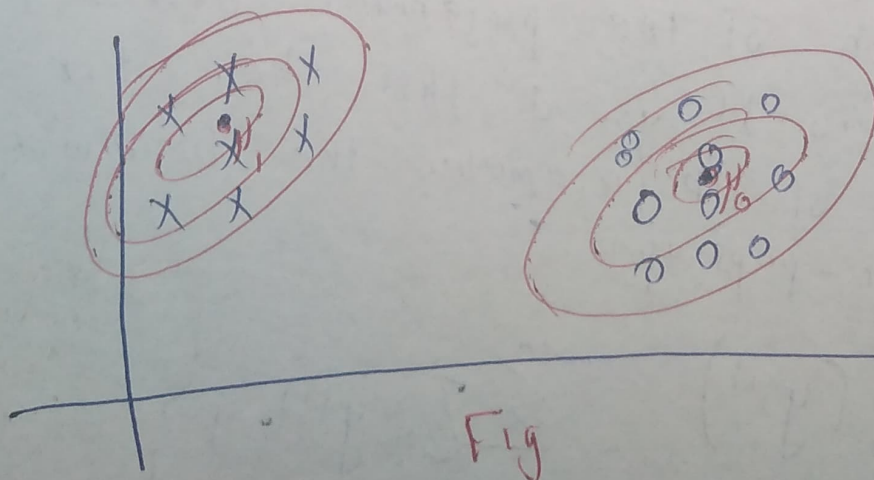
Sum of feature vector for example with  $y=1$ 

$$\sum_{i=1}^m \mathbb{I}\{y^{(i)} = 1\}$$

H of example with  $y=1$ 

(4)

$$\Sigma = \frac{1}{m} \sum_{i=1}^m (X^{(i)} - \mu_{y^{(i)}})(X^{(i)} - \mu_{y^{(i)}})^T$$





① Prediction Time :- (GDA) :-

$$\arg \max_y P(y|x) = \frac{\arg \max_y P(x|y) \cdot P(y)}{P(x)}$$

$$\min_z (z-5)^2 = 0, \arg \min_z (z-5)^2 = 5$$

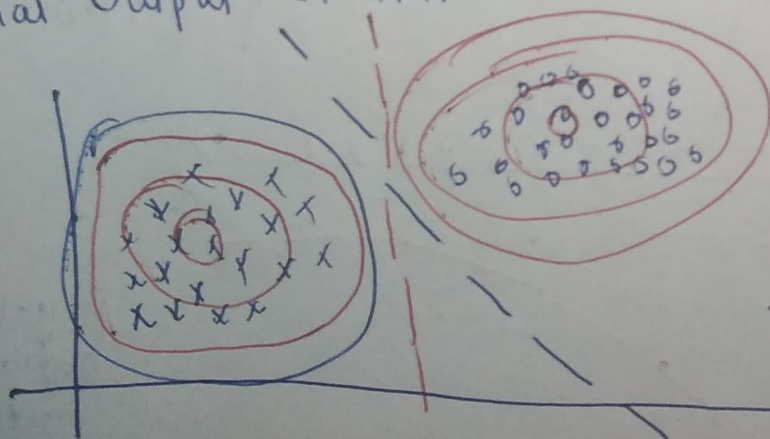
↳ Argmin defines the value you need to plug in to get the minimum.

→ What is argmin/argmax?

$$\arg \max_y P(y|x) = \arg \max_y (P(x|y) \cdot P(y))$$

↳ Denominator is not needed in case of argmax (It's a positive value).

Final Output of GDA:



→ GDA

implies the dashed line decision boundary.

Qu: P same but prior is different

→ Red dashed is the decision boundary of Logistic Regression.



(45)

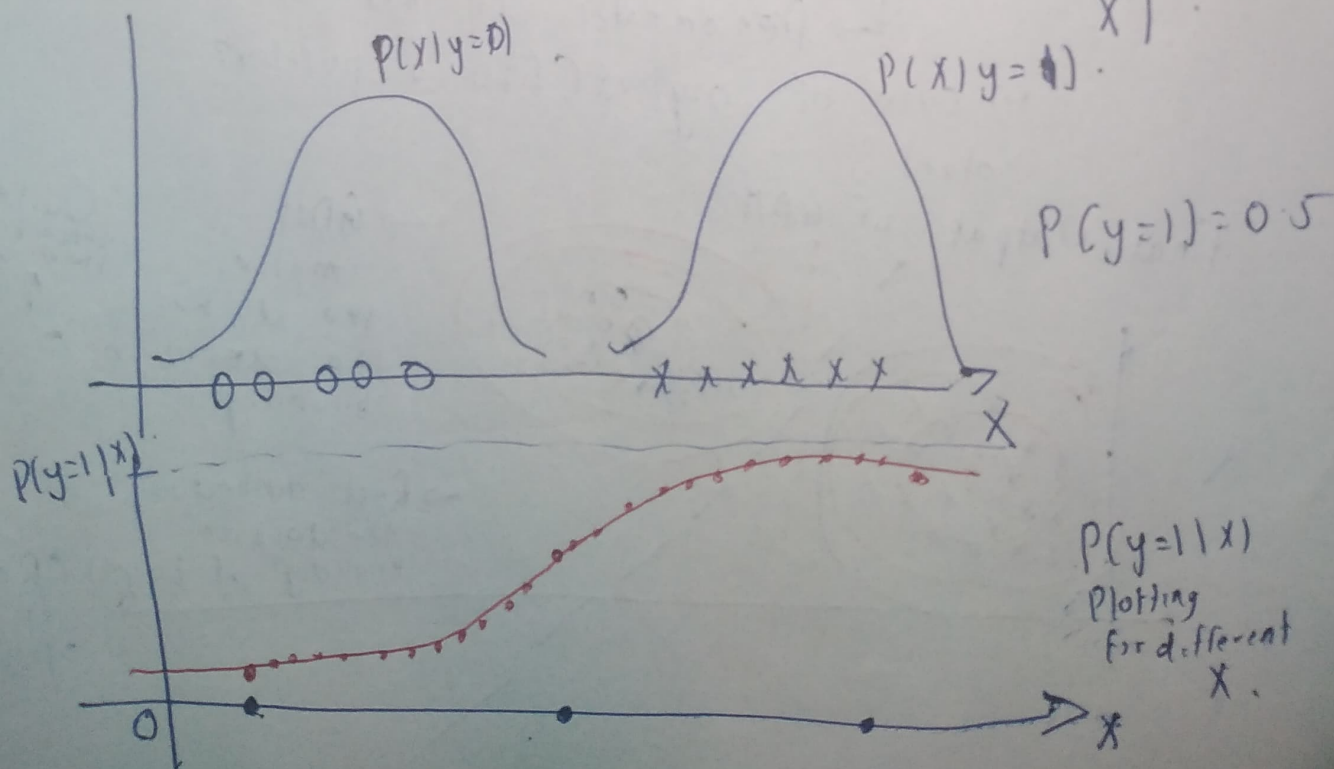
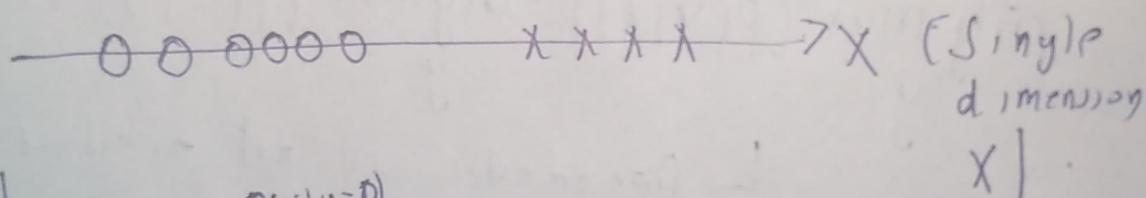
\* Compare GDA to logistic regression:

\* For fixed  $\phi, \mu_0, \mu, \Sigma$  lets

plot  $(P(y=1|x; \phi, \mu_0, \mu, \Sigma))$  as a function of  $x$

$$= \frac{P(x|y=1; \mu, \Sigma) \cdot P(y=1; \phi)}{P(x; \phi, \mu_0, \mu, \Sigma)}$$

\* For every value of  $x$  you can compute the above ratio (with fixed parameters).



⑧ If you  $P(Y=1|x)$  with different  $x$  then a Sigmoid Function appears

⑨ Mechanics of GDA and LR are different but they both end up choosing sigmoid function for  $P(Y|x)$

⑩ When is  $\Gamma$

⑪ GDA assumes <sup>Generative (Discriminative model)</sup>

$$\begin{aligned} x|y=0 &\sim \mathcal{N}(\mu_0, \Sigma) \\ x|y=1 &\sim \mathcal{N}(\mu_1, \Sigma) \\ y &\sim \text{Ber}(\theta) \end{aligned}$$

Logistic regression <sup>(Discriminative)</sup>

$$P(y=1|x) = \frac{1}{1 + e^{-\theta^T x}} \quad ("x_0 = 1")$$

⑫ As illustrative above that these set of assumes implies

↳ However, the opposite is not true so if you know  $P(Y|x)$  is governed by sigmoid function you cannot imply  $P(x|y)$  is Gaussian ( $\neq$ )

⑧ This shows that GDA takes stronger assumption while LR do not since  $P(X|Y)$  under (Logistic Regression) LR can be of any distribution.

→ Which works better when?

⑧  $P(X|Y)$  is gaussian then

GDA would work better because of its assumption or else logistic regression

GDA is more computationally efficient.

⑧ Let's say:-

- ①  $X|Y=1 \sim \text{Poisson}(\lambda_1)$
- ②  $X|Y=0 \sim \text{Poisson}(\lambda_2)$
- ③  $Y \sim \text{Bi}(\phi)$

→  $P(Y=1|X)$  is still logistic regression.

↓  
GDA would perform worst

Logistic Regression would still work fine

⑧ Key takeaways:

- Weak assumption make our model more robust
- Strong assumption would make our model restricted.

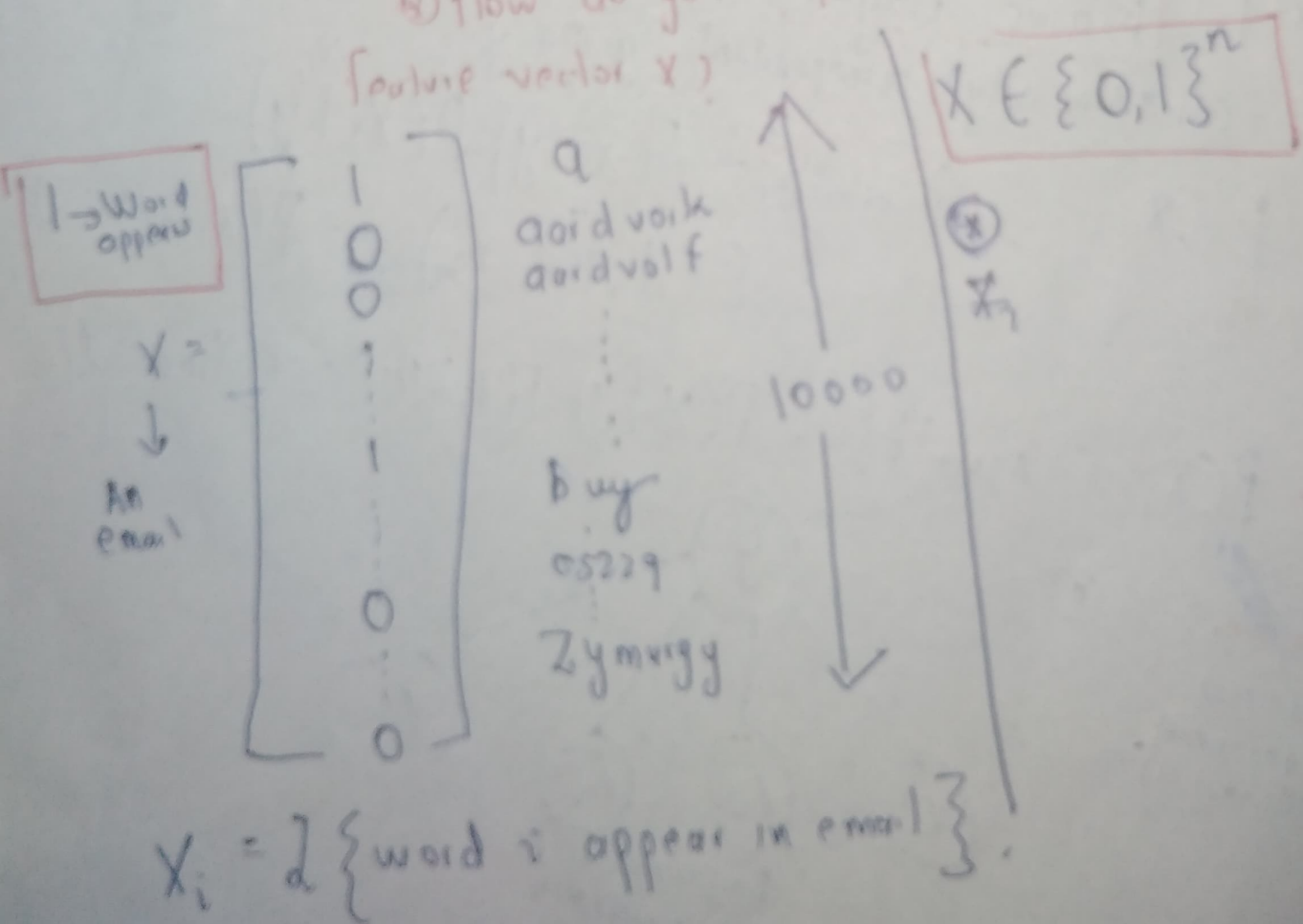


\* If  $P(X|Y)$  is from exponential family with their natural parameter then  $P(Y|X)$  would be logistic (sigmoid function).

## (\*) Naive Bayes

• Email classification problem:-

How do you represent a feature vector  $X$ ?



→ In naive bayes want to model  $P(X|y), P(y)$

→  $2^{10000}$  possible outcome of  $X$

⊛ → if we model of  $P(X)$  in straightforward way as a multinomial distribution over  $2^{10000}$  possible outcome then you need  $2^{10000}-1$  parameters.

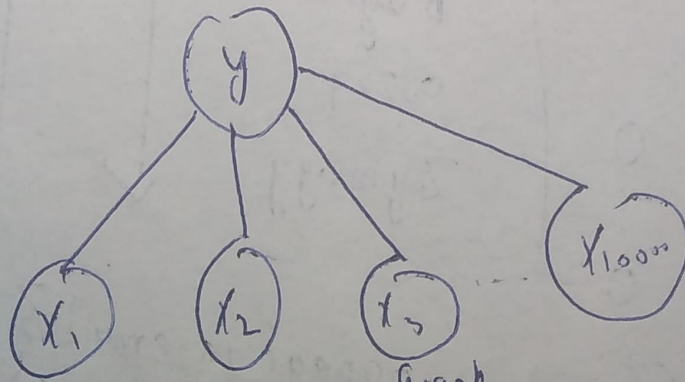
⊛ In Naive Bayes we assume  $X$ 's are conditionally independent given  $y$ .

←  $P(X_1, \dots, X_{10000}|y) = P(X_1|y) P(X_2|y) \dots P(X_{10000}|y)$

NO Assumption

←  $P(X_1, \dots, X_{10000}|y) = P(X_1|y) \cdot P(X_2|y) \cdot \dots \cdot P(X_{10000}|y)$

← Naive Bayes Assumption.  
(Conditional Independent Assumption).



Bayesian Graph Representation.

(56)

$$P(X_1, \dots, X_{10000}) = \prod_{i=1}^n P(X_i | Y)$$

\* Parameter of this model:-

~~$$\phi_{j|y=1} = P(X_{j=1} | y=1)$$~~

$$\phi_{j|y=1} = P(X_{j=1} | y=1)$$

$$\phi_{j|y=0} = P(X_{j=1} | y=0)$$

$$\phi_y = P(y=1)$$

→ Joint Likelihood

$$\mathcal{L}(\phi_y, \phi_{j|y}) = \prod_{i=1}^m P(X^{(i)}, y^{(i)}; \phi_y, \phi_{j|y})$$

$$\text{MLE: } (\nabla_{\phi_y, \phi_{j|y}} \mathcal{L}(\phi_y, \phi_{j|y}) = 0 \rightarrow \text{You}$$

find out

$$\phi_y = \frac{\sum_{i=1}^m \mathbb{1}\{y^{(i)} = 1\}}{m}$$

m

~~0~~



(5)

$$P(j | y=1) = \frac{\sum_{i=1}^m \mathbb{1}\{X_j = 1\}}{\sum_{i=1}^m \mathbb{1}\{y^{(i)} = 1\}}$$

What fraction of all spam email had word  $j$

All spam email

→ Testing time: - (Naive Bayes) :-

$$\arg \max_y P(y | X) = \arg \max_y \prod_{j=1}^p (X_j | y) \cdot P(y)$$