

# G.D.A.

$$\underbrace{P(x,y)}_{\text{Model}} = \underbrace{P(x|y) P(y)}_{\substack{\text{High} \\ \text{dimensional}}} \quad \begin{array}{l} \text{fraction of} \\ \text{age} \\ \text{belonging} \\ \text{to a} \\ \text{particular} \\ \text{class} \end{array}$$

$$\underbrace{P(y|x)}_{\text{Posterior}} = \frac{P(x|y) P(y)}{P(x)} = \frac{P(x|y) P(y)}{(P(x|y=0) P(y=0) + P(x|y=1) P(y=1))}$$

Prediction

$$\hat{y} = \arg \max_y p(y|x)$$

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

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

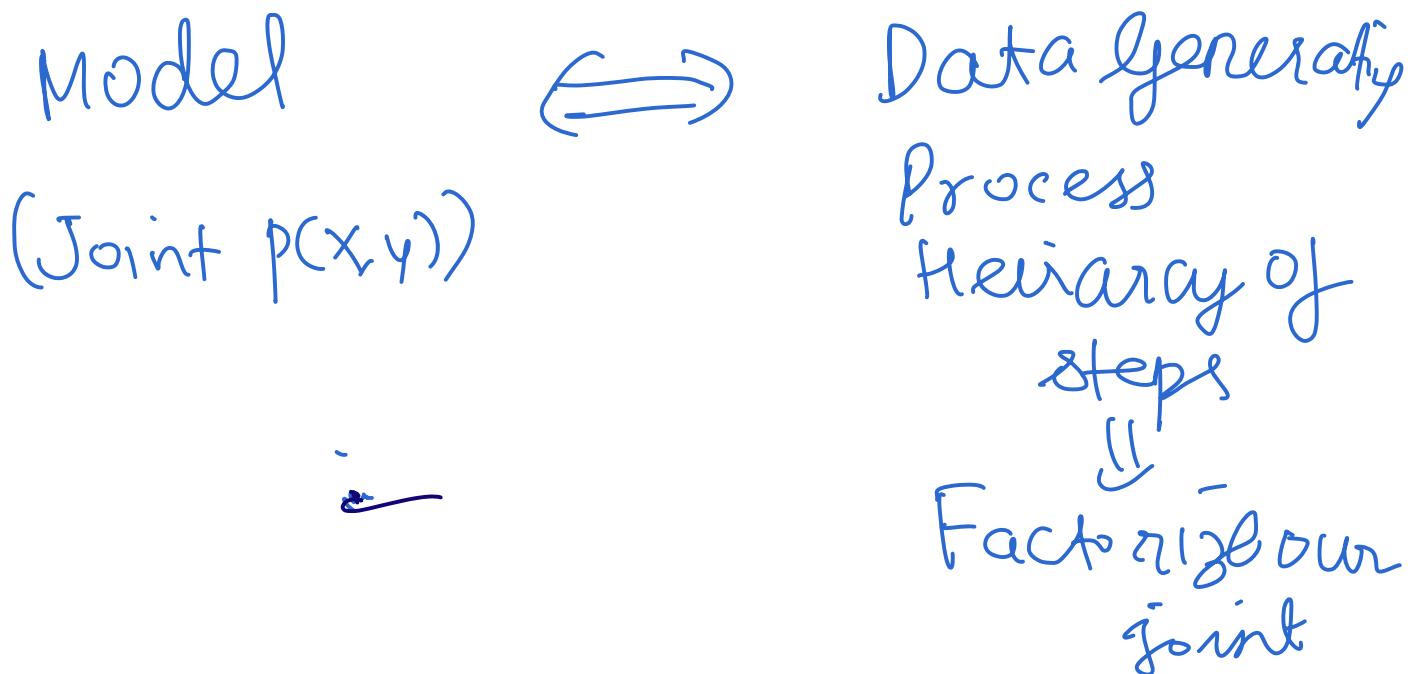
$p(x) \rightarrow$  prob. that you are going  
to encounter it

Two Algorithms

Both:  $y \in \{0, 1\}$

G.D.A.  $\rightarrow x \in \mathbb{R}^d$

Naive Bayes  $\rightarrow x$  is discrete (text classification)



G.D.A.

$$y \sim \text{Bern}(\phi)$$

$$x|y=0 \sim N(\mu_0, \Sigma)$$

$$x|y=1 \sim N(\mu_1, \Sigma)$$

Data Generating Process

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

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

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

Parameters  $\rightarrow \Sigma, \mu_1, \mu_0, \phi$

Likelihood over parameters

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

This part is p.o.f  
finding specific  $(x, y)$

$$= \log \prod_{i=1}^n p(x^{(i)} | y^{(i)}) p(y^{(i)})$$

$$\nabla L(\phi) = 0 \rightarrow \hat{\phi}, \hat{\mu}_0, \hat{\mu}_1, \hat{\Sigma}$$

$$\hat{\phi} = \frac{1}{n} \sum_{i=1}^n 1 \{ y^{(i)} = 1 \}$$

what fraction of cgs are 1

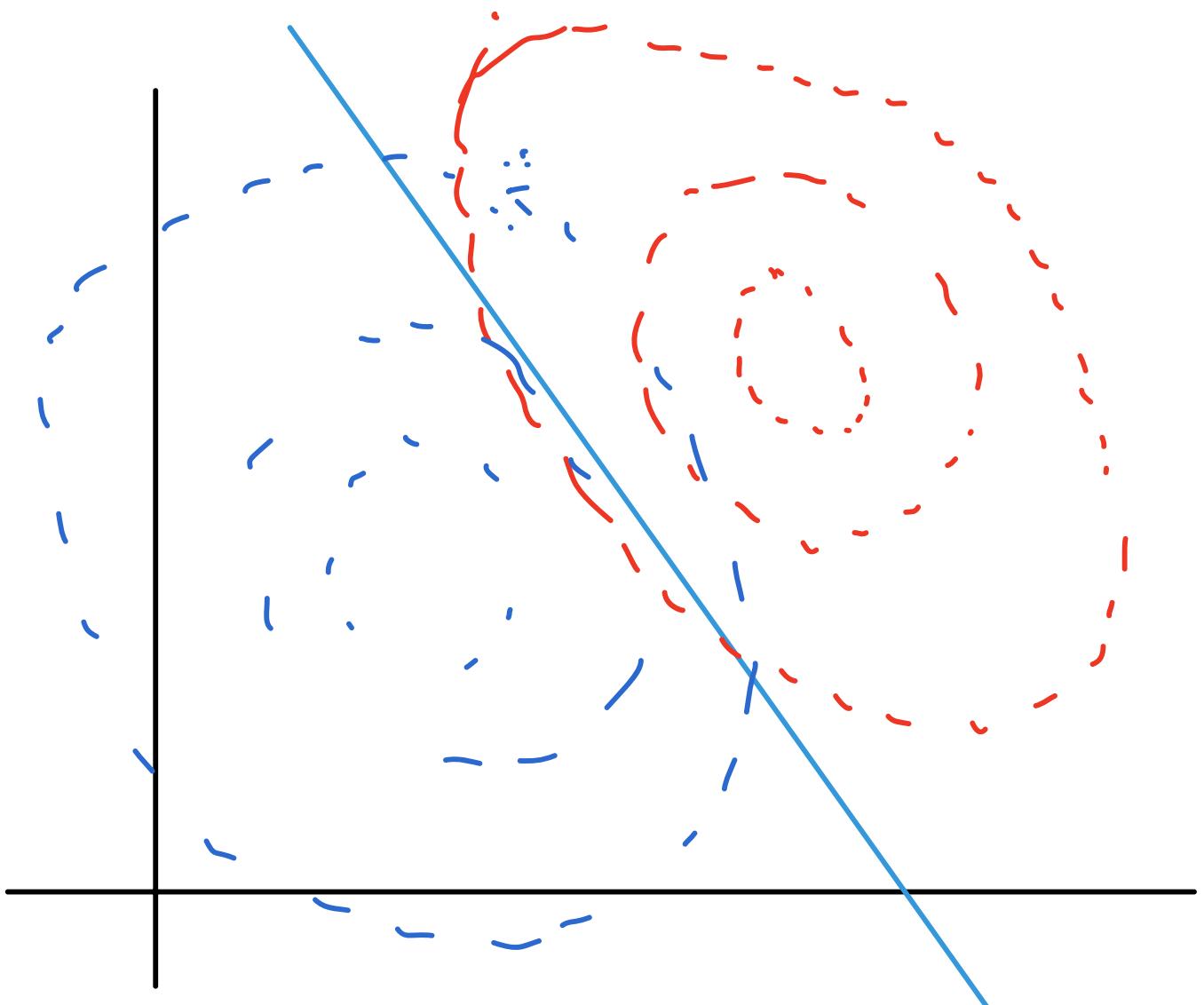
$$\hat{\mu}_0 = \frac{\sum_{i=1}^n 1 \{ y^{(i)} = 0 \} \cdot x^{(i)}}{\sum_{i=1}^n 1 \{ y^{(i)} = 0 \}}$$

$$\hat{\mu}_1 = \frac{\sum_{i=1}^n 1 \{ y^{(i)} = 1 \} \cdot x^{(i)}}{\sum_{i=1}^n 1 \{ y^{(i)} = 1 \}}$$

$$\Sigma = \frac{1}{n} \sum_{i=1}^n (x^{(i)} - \mu_x^{(i)}) (x^{(i)} - \mu_x^{(i)})^\top$$

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

$$\Theta \rightarrow \mu_0, \mu_1, \phi, \Sigma$$



$$P(X|Y=0) = P(X|Y=1)$$

→ separating  
hyperplane

GDA → Logistic Regression  
( $\Sigma$  same)

Principle axes of 2 ellipsoids  
are eigen vectors of  
covariance matrix

If unequal  $\Sigma$ , then nonlinear  
hyperplane (Logistic with  
non linear features)

When assumption of data

→ Gaussian, then G.D.A  
works very well.

Logistic is more robust,  
but when data is  
less and we know  
cond'n satisfied, GDA  
better.

# NAIVE BAYES

$x \rightarrow$  discrete  
Text classification

$$P(x_j | x_k) = P(x_j) \quad \text{Independence}$$

$$P(X_j | X_{K \setminus j}, y) = P(X_j | y) \left[ \begin{array}{l} \text{cond} \\ \text{indep} \\ \text{on } y \end{array} \right]$$

"Buy our lottery" = { 0 0 0 | a  
                          : ardwor  
                          |  
                          | buy  
                          |  
                          | lottery  
                          :  
                          | our } d

$$x \in \{0, 1\}^d$$

$$x_j \in \{0, 1\}$$

1  
T

Model :  $p(y=1) \Rightarrow \text{Bern}(\phi_y)$

$$p(x_j=1|y=0) \Rightarrow \text{Bern}(\phi_j|y=0)$$

$$p(x_j=1|y=1) \Rightarrow \text{Bern}(\phi_j|y=1)$$

↳ d

• → params

what fraction  
spam/nonspam

if email is

spammy,  
 $x_j$  appear  
prob.

$$l(\phi_y, \phi_j | y=0, \phi_j | y=1)$$

$$= \log \prod_{i=1}^n p(x^{(i)}, y^{(i)}; \phi)$$

$$[p(x, y) = p(y)p(x|y)]$$

$$\begin{aligned} &= p(y)p(x_1, x_2, x_3|y) \\ &= p(y)p(x_1|y) \cancel{p(x_2|y)} \\ &\quad p(x_3|y) \end{aligned}$$

$$= \log \prod_{i=1}^n p(y^{(i)}; \phi_y) \left( \prod_{j=1}^d p(x_j^{(i)}|y^{(i)}; \phi_j) \right)$$

$$\phi_j | y=1 = \frac{\sum_{i=1}^n 1 \{ x_j^{(i)} = 1 \wedge y^{(i)} = 1 \}}{\sum_{i=1}^n 1 \{ y^{(i)} = 1 \}}$$

↑

p that  $j^{th}$   
word in  
vocab will

Show up in  
spammy  
email

$$\phi_j | y=0 = \frac{\sum_{i=1}^n 1 \{ x_j^{(i)} = 1 \wedge y^{(i)} = 0 \}}{\sum_{i=1}^n 1 \{ y^{(i)} = 0 \}}$$

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

Prediction !!!

$$P(Y_2|X) = \frac{P(X|Y) P(Y)}{P(X)}$$

$$= \frac{P(X|y) P(Y)}{P(X|y=0) P(Y=0) + P(X|y=1) P(Y=1)}$$

$$= \frac{\prod_{j=1}^d P(x_j^{(i)} | y=1) P(y=1)}{\text{deno}}$$

$$P(Y=0|x) = \frac{\text{" } y=0 \text{ " } y=0}{\text{deno}}$$

What if we encounter word  
not shown in training set?

Then %,  $\Rightarrow$  indefinite

The method doesn't give  
prediction we use Laplace  
Smoothing

$\chi = 0, 13$

$$\phi = P(H) \approx \frac{1}{10}$$

Count H

1

1

1

1

Count T

0

$$\underline{\sum \text{Count H}}$$

$$\sum \text{Count H} + \sum \text{Count T}$$

$$P(H) = 1$$

$$P(T) = 0$$

Laplace  $\rightarrow$  add 1 to both initial

$$P(H)z \frac{10}{11} \quad P(T)z \frac{1}{11}$$

Assume seen every word,  
once is spammy, one in  
non-spammy

Suppose  $j=2$

$$x_2^{*} = 0 \forall i*$$

$$\phi_2 | y=1 = \frac{0}{\# \text{positive examples}} = 0$$

$$\phi_2 | y=0 = \frac{0}{\# \text{negative ex}} = 0$$

$$p(x) = p(x|y=0)p(y=0)$$

$$+ p(x|y=1)p(y=1)$$

$$= \prod_{i=1}^d p(x_i|y=0)p(y=0)$$

$$\begin{aligned}
 & + \prod_{j=1}^o p(x_j|y=1)p(y=1) \\
 & \quad \downarrow \\
 & P(x_2=0|y=0) \\
 & \quad \downarrow 0 \\
 P(x_2=1|y=0) & = \phi_2|y=1 = 0
 \end{aligned}$$

$$p(x) = 0$$

But example  
does occur

$$\begin{aligned}
 P(y=1|x) &= \frac{p(x|y=1)p(y=1)}{p(x)} \\
 0 &
 \end{aligned}$$

tail

$Z \sim \text{Bern}(\phi)$

n trials

$Z^{(1)} = \text{tail}$

$\vdots$  } tail = 0

$Z^{(n)}$

likelihood of  $\phi$

$$\begin{aligned}
 &= (1-\phi)(1-\phi) \dots \phi \\
 &= (1-\phi)^{\# \text{tails}} \phi^{\# \text{heads}}
 \end{aligned}$$

$$\begin{aligned}
 \arg \max_{\ell(\phi)} &= \frac{\# \text{heads}}{\# \text{heads} + \# \text{tails}}
 \end{aligned}$$

Suppose 3 draws, all tails

$$\hat{\phi} = \frac{0}{0+3} = 0$$

So according to this,  
coin will never give heads.

$$\phi_{j|y=1} = \frac{\sum_{i=1}^n 1(x_j^{(i)}=1, y^{(i)}=1) + 1}{\sum 1(y^{(i)}=1) + 2}$$

$$\phi_{j|y=0} = \frac{1 + \sum_{j=1}^n 1 \left\{ x_j^{(i)} = 1 | y^{(i)} = 0 \right\}}{2 + \sum 1 \{ y^{(i)} = 0 \}}$$

# Multinomial Event Model

$$y \sim \text{Bern}(\varphi) = 1$$

$$x|y=0 \sim \text{Categorical}(\varphi_k|y=0) \\ |\mathcal{V}| - 1$$

$$x_j \in \{1, \dots, |\mathcal{V}|\}$$

MLE

$$\hat{\varphi}_k|y=1 = \frac{\sum_{i=1}^n \sum_{j=1}^{d_i} \mathbb{I}\{x_j^{(i)} = k | y^{(i)}\}}{\sum_{i=1}^n \mathbb{I}\{y^{(i)} = 1\} \cdot d_i}$$

total words  
across  
messages  
in train  
sets

Sum over all words in  
all messages and.

count only those which  
are word k & belong to  
our class

$\phi_{k|y=0}$  = replace above  
with  $y=0$

Laplace Smooth

$$\frac{(\text{+ Num})}{(\text{N}) + \text{Deno}}$$

		Bernoulli	... a word	buy ... exam ...
"buy our lottery"	b	0		0 1
"buy this watch, this watch only"	c	0		0 1
Laplace	0	0	0	0 1 1 1
Smoothly	1	1	$\frac{1}{2}$	1 1 1 1
"When is your exam"	0	0	0	0 1 0 0
"When is the homework due"	0	0	0 1	1 0 0
				phi 1450

In Bernoulli, we normalize locally. We just see if it occurs / not.

Even if word occurs 10000 times or 1 time, it is the same.

In Multinomial, we normalize over the whole dataset & we count how many times each word appears

Multinomial				
a	aymm	buy	oam	watch
0	0	1		1 0
,	0 0 0	1	0 0 0 0 0	2
	4 2 3	2 3 2 2		
	1 5 4+2+3+2+1	1 1 1 1 1		
	Laplace Smooth			

We can also think of it as concatenating over as a long sequence

Training  $P(x, y)$

Prediction  $P(y|x)$

GD A (continuous  $x \in \mathbb{R}^d$ )

$$p(x|y) \sim N(\mu_y, \Sigma)$$
$$p(y|x) = \frac{1}{1 + \exp(-\theta^T x)}$$

NB (discrete  $x$ )

- \* Conditional Independence Assumption  $p(x_j|y, x_k) = p(x_j|y)$
- Bernoulli Event Model  $p(x_j|y) \sim \text{Bernoulli} [x_j \text{ is } j^{\text{th}} \text{ word in Vocabulary}]$
- Multinomial Event Model  $p(x_j|y) \sim \text{Multinomial} [x_j \text{ is } j^{\text{th}} \text{ word in message}]$

Binary  
If Multiclass  
classifi,  
then  
 $p(y|x) = \text{softmax}$   
for all  $p(x|y)$   
in exp-family

Laplace Smoothing  
+1 for each class upfront before looking at data

Stanford