$\rightarrow$ Focus on $X$

$\rightarrow$ Model $\underbrace{P(x,y)}_{\text{Model}} = \underbrace{P(x|y)}_{\text{High}} \cdot \underbrace{P(y)}_{\text{class}}$

High dimensional. Prior        class Prior

$\begin{bmatrix} \text{Till now we} \\ \text{have modeled} \\ P(y|x) \end{bmatrix}$

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

$\begin{bmatrix} \text{when } y \text{ is} \\ \text{binary} \end{bmatrix}$

$\hat{y} = \arg\max_{y} P(y|x)$

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

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

# Two algorithms

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

GDA : GAUSS Discriminant Analysis

NB : Naive Bayes

GDA — $x \in \mathbb{R}^d$ (continuous)

NB — $x$ is discrete (text classification)

Model
(Joint $p(x,y)$)  $\iff$ Data Generating Process

Hierarchy of steps

$\Downarrow$

Factorise our joint.

## GDA

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

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

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

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

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

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

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

$x \in \mathbb{R}^d$

$\underbrace{P(x,y)}_{\text{Generative}} = \underbrace{P(y/x) \cdot P(x)}_{\text{Discr.}}$

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

## Max. Likelihood to learn parameters

Log likelihood

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

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

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

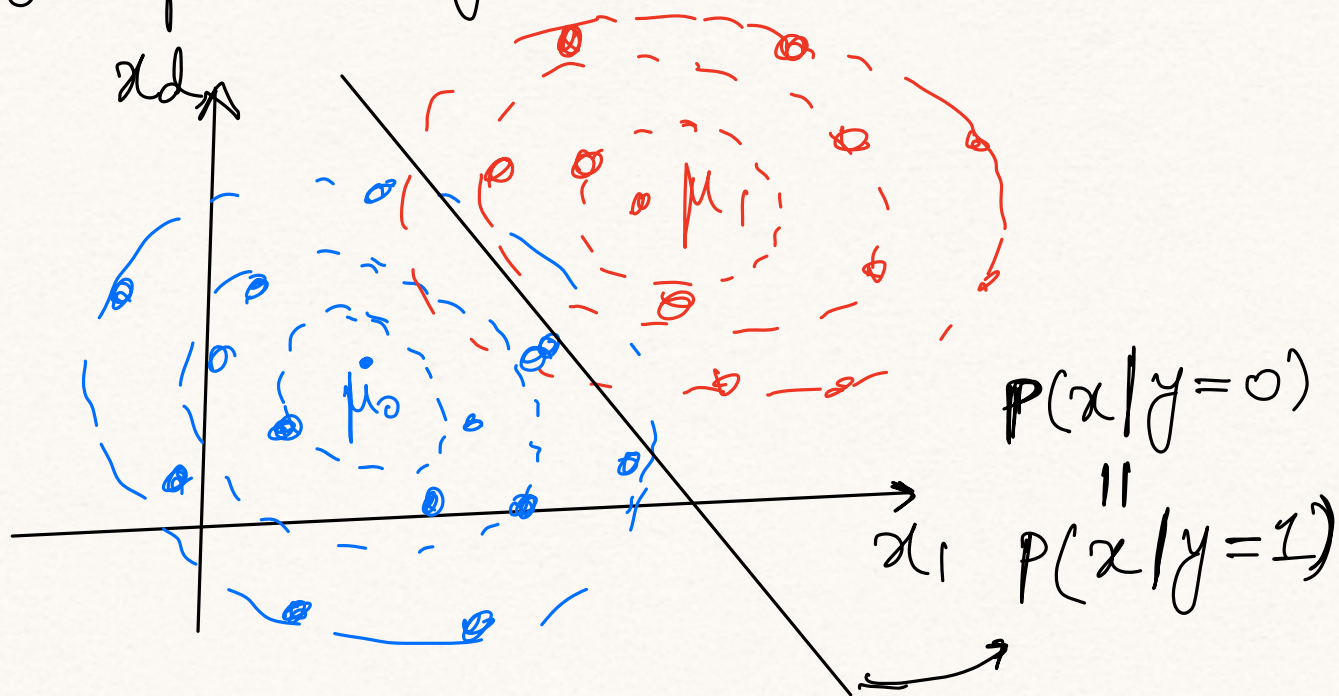$$\hat{\phi} = \frac{1}{n} \sum_{i=1}^{n} 1\{y^{(i)} = 1\}$$

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

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

$$\hat{\Sigma} = \frac{1}{n} \sum_{i=1}^{n} (x^{(i)} - \mu_{y^{(i)}})(x^{(i)} - \mu_{y^{(i)}})^T$$

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

$\theta$ depends only on $\mu_0, \mu_1, \phi, \Sigma$



$P(x|y=0)$
$=$
$P(x|y=1)$

GDA $\Rightarrow$ Logistic Regression.

✱ GDA more efficient than logistic reg.

If $\Sigma$ was not same for both the distri⁀.
Then instead of st. line it could be a
curve of degree 2/3/4 -- , etc.

# Naive Bayes

$x$ — discrete

Text classification (spam filters)

## Conditional Independence

$$P(x_j | x_k) = P(x_j) \quad [\text{indep.}]$$

$$P(x_j | x_k, y) = P(x_j | y) \quad \begin{bmatrix} \text{conditional} \\ \text{indep.} \\ \text{on } y \end{bmatrix}$$

## Bernoulli Event Model

"Buy our lottery" $= \begin{bmatrix} 0 \\ 0 \\ 1 \\ \vdots \\ 1 \\ \vdots \\ 1 \\ \vdots \\ 1 \\ 0 \end{bmatrix} \left. \begin{matrix} a \\ \text{aardvark} \\ \\ \vdots \\ buy \\ \\ lottery \\ \\ our \\ zygmorgy \end{matrix} \right\} d$

$\underbrace{\qquad\qquad}_{\text{Vocabulary}}$
— $d$ dim.

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

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

Model: $P(y=1) = \text{Bernoulli}(\phi_y)$ — 1

$P(x_j | y=0) = \text{Bernoulli}(\phi_{j|y=0})$ — $d$

$$P(x_j | y = 1) = \text{Bernoulli}(\phi_{j|y=1}) - d$$

$$\underset{1}{=} 1$$

$$\ell(\underbrace{\phi_y}_{1}, \underbrace{\phi_{j|y=0}}_{d}, \underbrace{\phi_{j|y=1}}_{d}) = \log \prod_{i=1}^{n} P(x^{(i)}, y^{(i)}; \phi)$$

$$= \log \prod_{i=1}^{n} P(y^{(i)}; \phi(y)) \left[ \prod_{j=1}^{d} P(x_j^{(i)} | y^{(i)}; \phi) \right]$$

$$P(x_1, x_2, \cdots, x_d | y) = P(x_1 | y) \cdot P(x_2 | x_1, y) \cdot$$
$$P(x_3 | x_1, x_2, y) - \cdots$$

If conditionally independent $\Rightarrow P(x_1 | y) \cdot P(x_2 | y) \cdots$

## MLE

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

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

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

$$P(y|x) = \frac{P(x|y) \cdot P(y)}{P(x)} = \frac{P(x|y) \cdot P(y)}{\begin{array}{c} P(x|y=0) \cdot P(y=0) \\ + \\ P(x|y=1) \cdot P(y=1) \end{array}}$$

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

## PROBLEM !

For new words $\rightarrow$ $P(y|x) = \frac{0}{0+0}$ $\rightarrow$ Thus, fails!!

$\rightarrow$ Can be corrected by Laplace smoothing.

Assume that each word has been seen once in each spam email as well as non-spam email.

## Laplace Smoothing Ver.

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

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

## Multinomial Event Model

$$y \sim Bernoulli(\phi_y)$$

$$x|y=0 \sim Categorical(\phi_{k|y=0}) \quad -1 \quad -|v|-1$$

$$x_j \in \{1, - -, |v|\}$$

$$x^{(i)} \in \{1, - -, |v|\}^{d_i}$$

## MLE

$$\phi_{k|y=1} = \frac{\sum_{i=1}^{n} \sum_{j=1}^{d_i} 1\{x_j^{(i)} = k \wedge y^{(i)} = 1\}}{\sum_{i=1}^{n} 1\{y^{(i)} = 1\} d_i}$$

$$\phi_{k|y=0} = \underline{\phantom{xxxxxxxx}}$$

## L.S ver

$$\phi_{k|y=1} = \frac{\sum_{i=1}^{n} \sum_{j=1}^{d_i} 1\{x_j^{(i)} = k \wedge y^{(i)} = 1\} + 1}{\sum_{i=1}^{n} 1\{y^{(i)} = 1\} d_i + |V|}$$