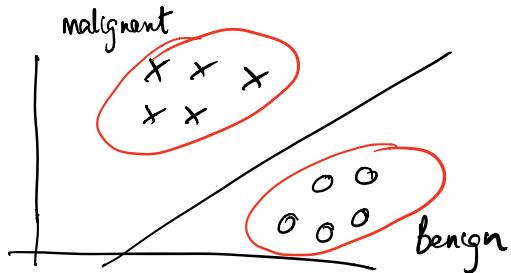


## Generative Learning Algorithms

- Gaussian Discriminant Analysis (GDA)
- Generative & Discriminative Comparison
- Naive Bayes



## Discriminative Learning Algorithm

learns  $p(y|x)$

or learns  $h_\theta(x) = \begin{cases} 0 & \text{directly} \\ 1 & \end{cases}$

$$x \rightarrow y$$

## Generative Learning Algorithm

learns  $p(x|y)$   
features

$p(y)$   
class prior

Bayes Rule:

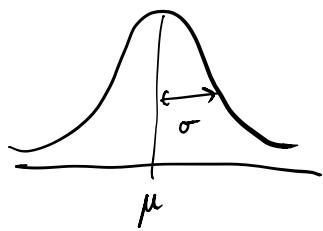
$$p(y=1|x) = \frac{p(x|y=1) \cdot p(y=1)}{p(x)}$$

$$p(x) = p(x|y=1)p(y=1) + p(x|y=0)p(y=0)$$

## Gaussian Discriminant Analysis (GDA)

Suppose  $x \in \mathbb{R}^d$  (drop  $x_0 = 1$  convention)

Assume  $p(x|y)$  is Gaussian



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

$\vec{\mu} \in \mathbb{R}^d$      $\Sigma \in \mathbb{R}^{d \times d}$

$$\mathbb{E}[Z] = \mu$$

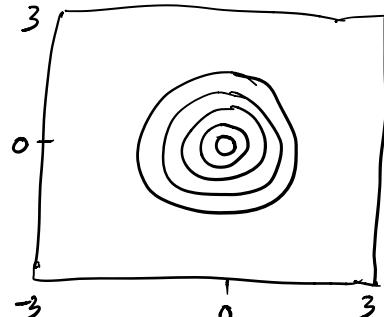
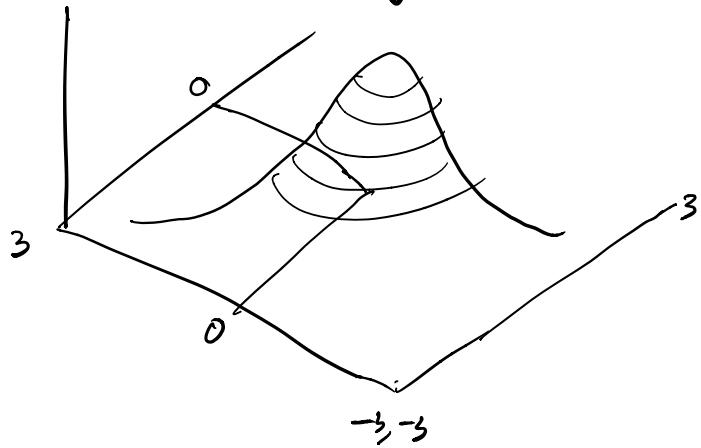
$$\Sigma_{ij} = \mathbb{E}[Z_i Z_j] - \mathbb{E}[Z_i] \mathbb{E}[Z_j]$$

$$\begin{aligned}\text{Cov}(Z) &= \mathbb{E}[(Z-\mu)(Z-\mu)^T] \\ &= \mathbb{E}[ZZ^T] - (\mathbb{E}Z)(\mathbb{E}Z)^T\end{aligned}$$

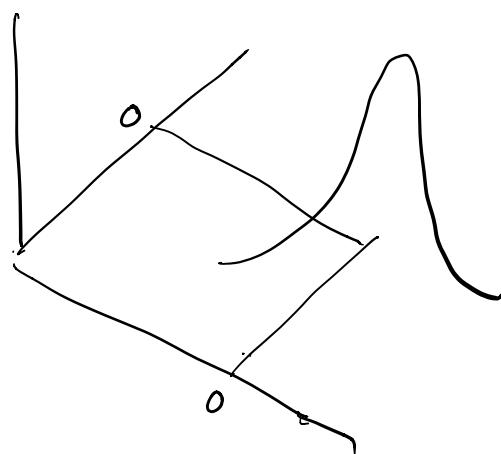
$$\mathbb{E}[Z] = \mathbb{E} Z$$

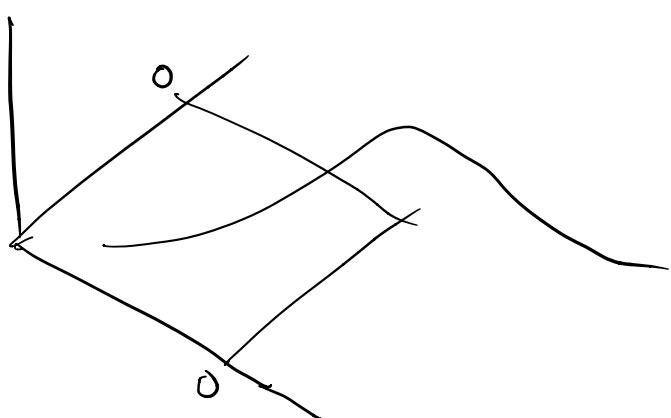
pdf  $p(Z) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu)\right)$

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

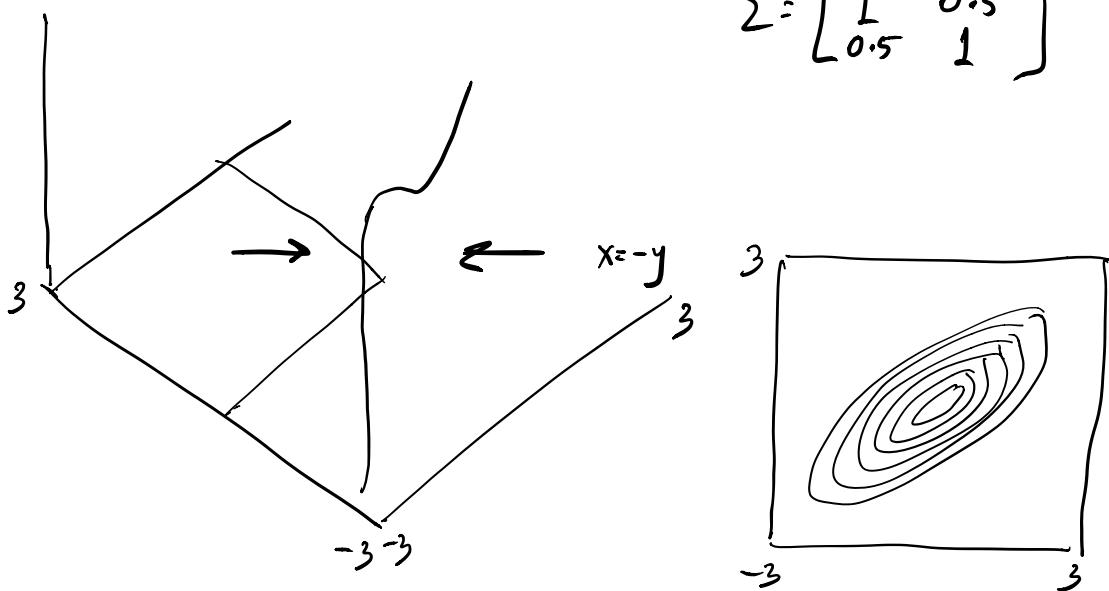


$$\Sigma = \begin{bmatrix} 0.6 & 0 \\ 0 & 0.6 \end{bmatrix}$$

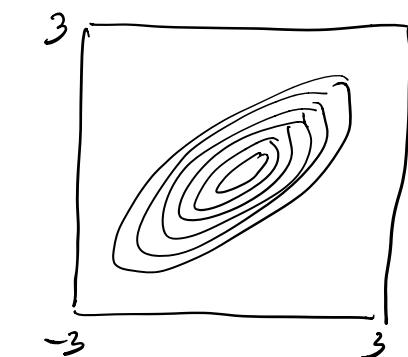




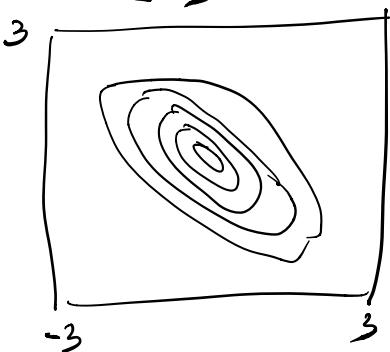
$$\Sigma = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}$$



$$\Sigma = \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1 \end{bmatrix}$$



$$\Sigma = \begin{bmatrix} 1 & -0.5 \\ -0.5 & 1 \end{bmatrix}$$



$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad \begin{bmatrix} 1 \\ 2 \end{bmatrix} \quad \begin{bmatrix} -1 \\ -3 \end{bmatrix}$$

GDA model

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

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

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

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

$$\begin{array}{ccc} \mu_0, \mu_1 & \sum_{\mathbb{R}^d} & \phi \\ & \mathbb{R}^{d \times d} & [0, 1] \end{array}$$

$p(y=1|x), p(y=0|x)$  via Bayes rule

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

Joint likelihood fn:

$$\begin{aligned} L(\phi, \mu_0, \mu_1, \Sigma) &= \prod_{i=1}^n P(x^{(i)}, y^{(i)}; \phi, \mu_0, \mu_1, \Sigma) \\ &= \prod_{i=1}^n p(x^{(i)} | y^{(i)}) \cdot p(y^{(i)}) \end{aligned}$$

Cost fn & j<sup>t</sup> fn of  $x, y$

Discriminative:  $L(\theta) := \prod_{i=1}^n p(y^{(i)} | x^{(i)}; \theta)$

Maximum Likelihood Estimator

$$\max_{\phi, \mu_0, \mu_1, \Sigma} l(\phi, \mu_0, \mu_1, \Sigma) = \log L(\quad)$$

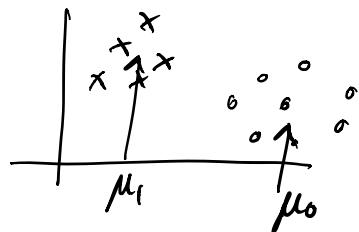
$$\phi = \frac{\sum_{i=1}^n y^{(i)}}{n}$$

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

$$\mathbb{1}_{\{\text{true}\}} = 1$$

$$\mathbb{1}_{\{\text{false}\}} = 0$$

$$\mu_0 = \frac{\sum_{i=1}^n \mathbb{1}_{\{y^{(i)} = 0\}} x^{(i)}}{\sum_{i=1}^n \mathbb{1}_{\{y^{(i)} = 0\}}}$$



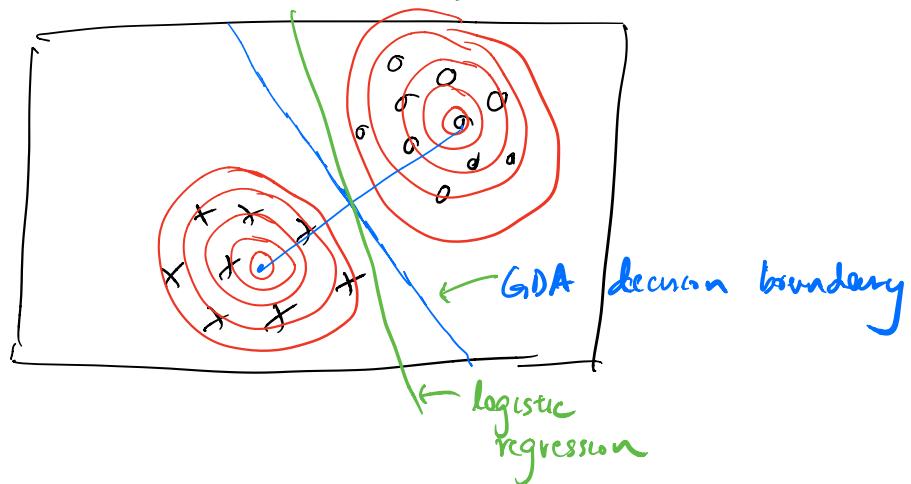
$$\mu_1 = \frac{\sum_{i=1}^n \mathbb{1}_{\{y^{(i)} = 1\}} x^{(i)}}{\sum_{i=1}^n \mathbb{1}_{\{y^{(i)} = 1\}}}$$

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

Prediction:

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

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

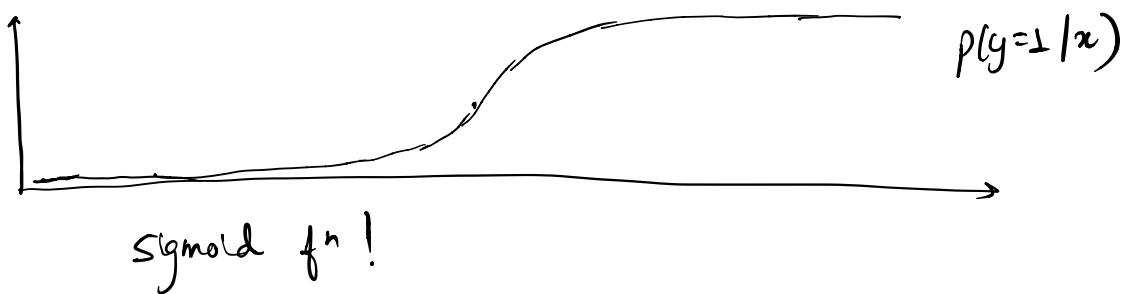
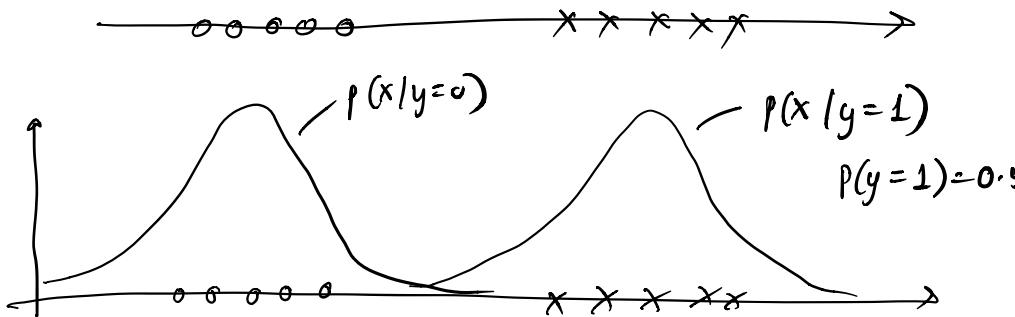


Comparison to logistic regression

for fixed  $\phi, \mu_0, \mu_1, \Sigma$

plot  $p(y=1|x; \phi, \mu_0, \mu_1, \Sigma)$  as a fn of  $X$

$$\Rightarrow \frac{p(x|y=1; \mu_1, \Sigma) p(y=1; \phi)}{p(x; \mu_0, \mu_1, \Sigma, \phi)}$$



Sigmoid fn!

generative  
GDA assumes

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

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

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

discriminative  
Logistic Regression

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

( $"x_b=1"$ )

Stronger  
assumption

Weaker  
assumption

$$\left. \begin{array}{l} x|y=1 \sim \text{Poisson } (\lambda_1) \\ x|y=0 \sim \text{Poisson } (\lambda_0) \\ y \sim \text{ber}(\phi) \end{array} \right\} \Rightarrow p(y=1|x) \text{ logistic}$$

Naive Bayes

Feature vector  $X$ ?

English dict

$x = \begin{bmatrix} 1 \\ 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 1 \end{bmatrix}$	aardvark	↑ top 10k
	aardwolf	
	:	
	buy	
	cs229	
	:	
	zygmurgy	

$$x \in \{0, 1\}^d \quad d = 10,000$$

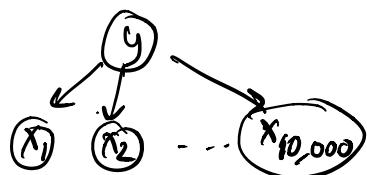
$$X_i = 1 \{ \text{word } i \text{ appears in email} \}$$

Want to model  $p(x|y)$ ,  $p(y)$

$2^{10,000}$  possible values of  $X$

Assume  $X_i$ 's are conditionally independent given  $y$

$$\begin{aligned} p(x_1, \dots, x_{10,000}|y) &= p(x_1|y) \cdot p(x_2|x_1, y) \cdot p(x_3|x_1, x_2, y) \\ &\quad \dots \quad p(x_{10,000}|\dots, y) \\ \text{assume} \quad &= p(x_1|y) \cdot p(x_2|y) \cdot p(x_3|y) \dots p(x_{10,000}|y) \end{aligned}$$



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

Parameters

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

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

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

Joint Likelihood

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

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

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