

# Lecture 5 GDA & Naive Bayes

## 1. Generative Learning Algorithms VS. Discriminative Comparison

### Discriminative Learning Algorithms:

- Learns  $p(y|x)$  or learns  $h(x) = \{ \cdot \}$ , directly,
- Learns a decision boundary
- Regressions, SVM

000,  
000,  
XXX  
XXX

### Generative Learning Algorithms:

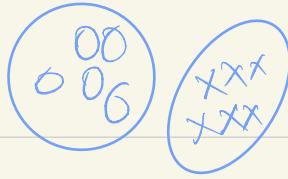
- Learns  $P(x|y)$  and  $P(y)$

feature class      class prior

### Using Bayesian Rules:

$$P(y=1|x) = \frac{P(x|y=1) \cdot P(y=1)}{P(x)}, P(x) = P(x|y=1)P(y=1) + P(x|y=0)P(y=0)$$

- Learns probability distribution of data
- GDA ; Naive Bayes



## 2. Gaussian Discriminant Analysis (GDA)

Assumption:

- $x \in \mathbb{R}^n$  , drop  $x_0 = 1$  (convention)
  - $P(x|y)$  is Gaussian
- $$\Rightarrow z \sim N(\vec{\mu}, \vec{\Sigma}), z \in \mathbb{R}^n$$
- $\mathbb{R}^n \leftarrow \mathbb{R}^{n \times n}$

$$E[z] = \mu; \text{Cov}(z) = E[(z-\mu)(z-\mu)^T] = E[z \cdot z^T] - (E[z])(E[z])^T$$

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

GDA Model:

$$P(x|y=0) = \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} \exp(-\frac{1}{2}(x-\mu_0)^\top \Sigma^{-1} (x-\mu_0))$$

$$P(x|y=1) = \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} \exp(-\frac{1}{2}(x-\mu_1)^\top \Sigma^{-1} (x-\mu_1))$$

$$P(y) = \phi^y (1-\phi)^{1-y} \quad (P(y=1) = \phi, \text{ Bernoulli})$$

Parameters:  $\mu_0, \mu_1, \Sigma, \phi$  ( $\text{RG}[0, 1]$ )

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

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

<compare generative>:  $L(\theta) = \prod_{i=1}^n P(y^{(i)}|x^{(i)}; \theta)$   $\stackrel{\text{conditional likelihood}}{\leftarrow}$   $= \prod_{i=1}^m P(x^{(i)}|y^{(i)}) P(y^{(i)})$

Maximum Likelihood Estimation:

$$\underset{\phi, \mu_0, \mu_1, \Sigma}{\text{Max}} l(\phi, \mu_0, \mu_1, \Sigma)$$

$$\Downarrow \log(L(\phi, \mu_0, \mu_1, \Sigma))$$

$$\Rightarrow \phi = \frac{\sum_{i=1}^m I\{y^{(i)} = 1\}}{m} \quad (\text{the proportion of } y=1)$$

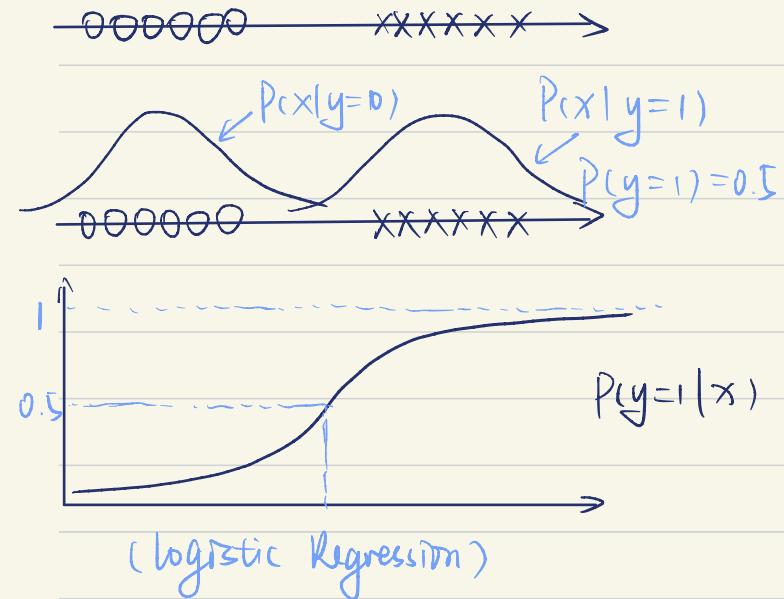
$$M_0 = \frac{\sum_{i=1}^m I\{y^{(i)} = 0\} x^{(i)}}{\sum_{i=1}^m I\{y^{(i)} = 0\}} \quad \leftarrow \begin{array}{l} \text{sum of feature vectors for data with } y=0 \\ \# \text{ data with } y=0 \end{array}$$

$$M_1 = \frac{\sum_{i=1}^m I\{y^{(i)} = 1\} x^{(i)}}{\sum_{i=1}^m I\{y^{(i)} = 1\}}$$

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

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

3. Compare GDA to Logistic Regression  
 for fixed  $\phi, M_0, M_1, \Sigma$ , plot  $P(y=1|x; \phi, M_0, M_1, \Sigma)$  as fn of  $x$



Generative

GDA assumes

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

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

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

stronger assumption

$$\begin{cases} x|y=1 \sim \text{Poisson}(\lambda_1) \\ x|y=0 \sim \text{Poisson}(\lambda_0) \end{cases}$$

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

Discriminative

Logistic Regression

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

" $x_0=1$ "

weaker assumption

How to choose which method to use?

- Large data set - Logistic Regression (okay to tell the algorithm less)
- Computational Efficiency Matters - GDA

#### 4. Naive Bayes (example: spam email classifier)

Feature Vector  $X$ ?

$$X = \begin{bmatrix} 1 \\ 0 \\ 1 \\ 0 \\ \vdots \end{bmatrix} \leftrightarrow \begin{bmatrix} a \\ \text{aardvark} \\ ; \\ ; \\ \text{zymurgy} \end{bmatrix}$$

↑  
10000  
↓

$X \in \{0, 1\}^n$   
 $X_i = \{\text{word } i \text{ appears in email}\}$

Assumption:  $X_i$ 's are conditionally independent given  $y$

$$\begin{aligned} P(X_1, \dots, X_{10000} | y) &= P(X_1 | y) P(X_2 | X_1, y) P(X_3 | X_1, X_2, y) \dots P(X_{10000} | \dots) \\ &\stackrel{\text{assume}}{=} P(X_1 | y) P(X_2 | y) P(X_3 | y) \dots P(X_{10000} | y) \\ &= \prod_{i=1}^n P(X_i | y) \end{aligned}$$

Parameters:  $\phi_j | y=1 = P(X_j = 1 | y=1)$  if it is a spam

$\phi_j | y=0 = P(X_j = 1 | y=0)$  if it is not a spam

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

Joint Likelihood:  $L(\phi_y, \phi_j | y) = \prod_{i=1}^m P(x^{(i)}, y^{(i)} | \phi_y, \phi_j | y)$

MLE:  $\hat{\phi}_y = \frac{\sum_{i=1}^m I\{y^{(i)} = 1\}}{m}$        $\hat{\phi}_j | y=1 = \frac{\sum_{i=1}^m I\{x_j^{(i)} = 1, y^{(i)} = 1\}}{\sum_{i=1}^m I\{y^{(i)} = 1\}}$