# Generative Learning Algorithms

Build models for what each class looks like
and classify new points using these models

Formally:

Learns $p(x|y)$ and $p(y)$

features class class prior

(given a tumor is malignant/benign, what do its features look like?)

(what is the probability of any tumor being malignant/benign?)

## Bayes Rule

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

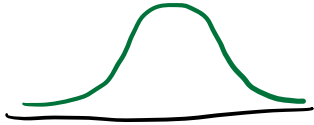$$= p(x|y=1)p(y=1) + p(x|y=0)\,p(y=0)$$

learn all these terms

# Gaussian Discriminant Analysis

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

Assume $p(x|y)$ is Gaussian

## Multivariate Gaussian

$$z \sim N(\vec{\mu}, \Sigma) \qquad z = (z_1, z_2, \dots z_d) \in \mathbb{R}^d$$

$\mathbb{R}^d$    $\mathbb{R}^{d \times d}$

$$E[z] = \vec{\mu}$$

$$Cov(z) = E\left[(z-\mu)(z-\mu)^T\right]$$

$$= E[zz^T] - (E[z])(E[z])^T$$

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

## GDA model

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

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

Parameters: $\underbrace{\mu_0, \mu_1}_{\mathbb{R}^d}, \underset{\mathbb{R}^{d \times d}}{\Sigma}, \phi \in [0,1]$

$$p(y) = \phi^y (1-\phi)^y$$

$$p(y=1) = \phi$$

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

Joint likelihood

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

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

# Maximum Likelihood Estimation

$$\underset{\phi, \mu_0, \mu_1, \Sigma}{\text{Max}} \quad \ell(\phi, \mu_0, \mu_1, \Sigma) = \log \mathcal{L}(\cdots)$$

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

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

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

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

# Prediction

$$\underset{y}{\arg\max}\ p(y|x) = \underset{y}{\arg\max}\ \frac{p(x|y)\, p(y)}{p(x)}$$

$$= \underset{y}{\arg\max}\ p(x|y)\, p(y)$$

# Comparison w/ LR

GDA assumes
$$x|y=0 \sim N(\mu_0, \Sigma)$$
$$x|y=1 \sim N(\mu_1, \Sigma)$$
$$y \sim Ber(\phi)$$

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

### Stronger assumption
$$x|y=1 \sim Poi(\lambda_1)$$
$$x|y=0 \sim Poi(\lambda_0)$$
$$y \sim Ber(\phi)$$

### Weaker assumption
$$p(y=1|x) \text{ vs logistic fn}$$

# Naive Bayes

Feature vector $X$      (ex. spam classification)

$$X = \begin{bmatrix} 1 \\ 0 \\ 0 \\ \vdots \\ \vdots \\ \end{bmatrix} \begin{matrix} a \\ aardvark \\ \vdots \\ \\ \\ zymurgy \end{matrix}$$

top-$d$ words

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

$x_i : \mathbb{1}\{\text{word } i \text{ in email}\}$

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

$2^d$ possible values of $x$

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

$\Rightarrow p(x_1 \dots x_d | y) = p(x_1|y)\, p(x_2|x_1,y) \cdots p(x_d|x_{d-1}\dots x_1, y)$

$\underset{\text{assume}}{=} \quad p(x_1|y)\, p(x_2|y) \cdots p(x_d|y)$

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

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 spam

$\phi_y = p(y)$      Pr(spam)