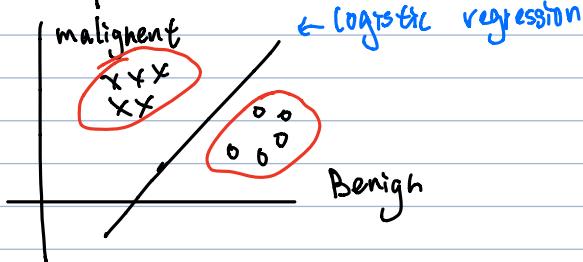


Generative Learning Algorithms

- Gaussian Discriminant Analysis (GDA)
- Generative X Discriminant Comparison
- Naive Bayes

Example: tumor detection



Discriminative Learning Algorithm

learns $p(y|z)$

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

$$x \rightarrow y$$

Generative Learning Algorithm

learns $P(x|y)$
↑
features ↓ class

$P(y)$ ← class prior.

$P(x|y), P(y)$ generative

$P(y|x)$. discriminative

Bayes Rule:

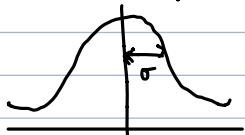
$$P(y=1|x) = \frac{P(x|y=1) P(y=1)}{P(x)}$$

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

Gaussian Discriminant Analysis (GDA)

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

Assume $p(x|y)$ is Gaussian



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

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

$$\vec{\mu} \in \mathbb{R}^d$$
$$\Sigma_{ii} = E[z_i z_i^T] - E[z_i] E[z_i^T]$$

μ

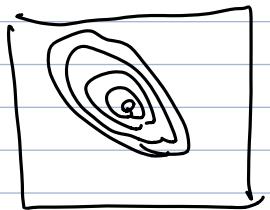
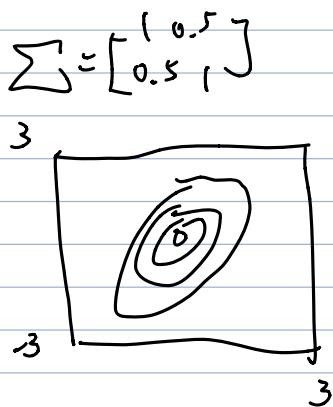
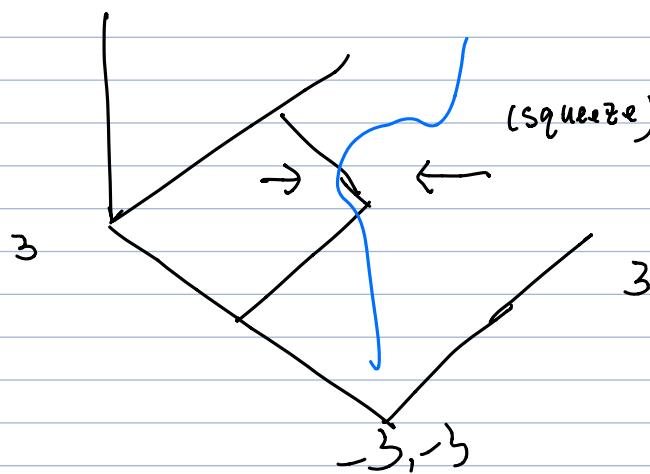
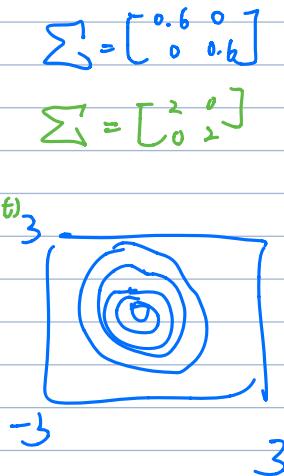
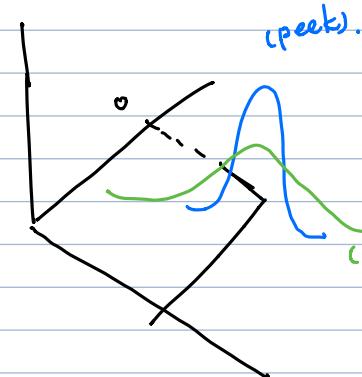
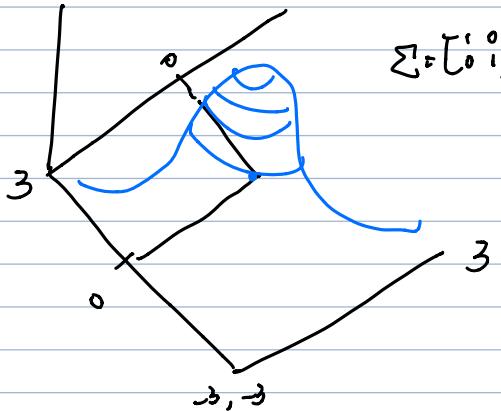
(covariance matrix)

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

$$\text{cov}(z) = \mathbb{E}[(z-\mu)(z-\mu)^T]$$

$$= \mathbb{E}[zz^T] - (\mathbb{E}[z])(\mathbb{E}[z])^T$$

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



GDA model

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

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

Parameters: $\mu_0, \mu_1, \Sigma, \phi$
(use same cov matrix)

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

$\mu_0, \mu_1 \in \mathbb{R}^d$ $\Sigma \in \mathbb{R}^{d \times d}$ $\phi \in [0, 1]$

Find $P(y=1|x)$, $P(y=0|x)$ by Bayes rule.

How to fit the parameters?

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

• Joint likelihood func:

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

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

Cost fn. is joint fn. of x, y

Discriminative. $L(\theta) = \prod_{i=1}^n P(y^{(i)}|x^{(i)}; \theta)$ (Diff. from Generative)

Maximum Likelihood Estimator

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

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

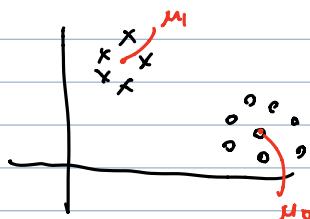
Indicator fn.

$\begin{cases} 1 & \text{if true} \\ 0 & \text{if false} \end{cases}$

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

(mean of all examples)
labeled 0

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



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

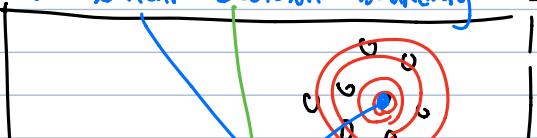
Prediction

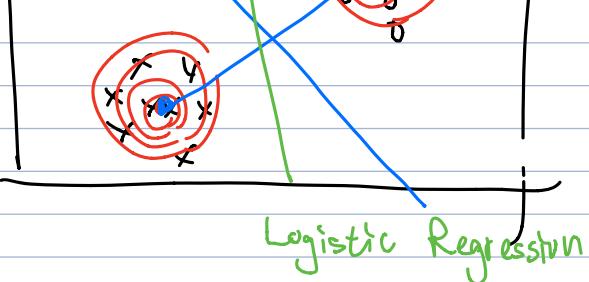
$$\arg \max_y P(y|x) = \arg \max_y \frac{P(x|y) P(y)}{P(x)}$$

(drop)

$$= \arg \max_y P(x|y) P(y)$$

GPA: linear decision boundary \Rightarrow Reason: Same Σ





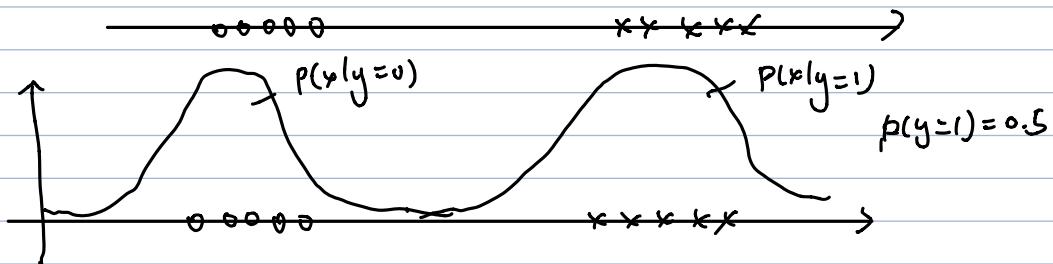
Comparison (w/ logistic regression)

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

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

$$\left[\frac{p(x|y=1; \mu_1, \Sigma)}{p(x; \mu_0, \mu_1, \Sigma, \phi)} \right]^\phi$$

Suppose a 1-D example



Prob.



Sigmoid fn.!

generative

(Does better if data IS multi-Gaussian)
(Efficient)

GDA assumes

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

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

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

discriminative

(More data, use it!)

Logistic Regression

. More robust)

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

($x_0 = 1$)

Stronger

assumption (Performs less strong)

Weaker

assumption. (Do well!)

$$x|y=1 \sim \text{Poisson}(\lambda_1)$$

$$x|y=1 \sim \text{Poisson}(\lambda_2)$$

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

$$\left. \begin{array}{c} \\ \\ \end{array} \right\} \Rightarrow$$

$P(y=1|x)$ is logistic

works for all generalized linear methods

Naive Bayes

Feature vector x ?

English dict

$$x = \begin{bmatrix} \text{apple} \\ \vdots \\ \text{zoo} \end{bmatrix} = \begin{bmatrix} 1 \\ \vdots \\ 0 \end{bmatrix} \quad d = 10,000$$

cone-hot-vec

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

$$x_i = \prod \{\text{word } i \text{ appears in email}\}$$

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

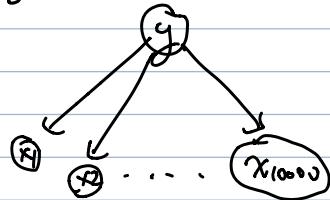
2^{10000} possible value of x .

Key assumption: x_i 's are conditionally independent given y

$$p(x_1, \dots, x_d) = p(x_1|y) p(x_2|x_1, y) \dots p(x_d|x_1, x_2, \dots, y)$$

$$\stackrel{\text{assume}}{=} p(x_1|y) p(x_2|y) \dots p(x_d|y)$$

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



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

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

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

^{ny}

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

Prove the GDA has linear decision boundary

Let $P \triangleq (\phi, \mu_0, \mu_1, \Sigma)$

$$\begin{aligned} P(y=1|x; \phi, \mu_0, \mu_1, \Sigma) &= \frac{p(x|y=1; \phi, \mu_1, \Sigma) p(y=1; \phi)}{p(x|y=0; \mu_0, \mu_1, \Sigma, \phi)} \xleftarrow{\phi} \\ &= \frac{p(x|y=1; \mu_1, \Sigma) \phi}{p(x|y=1; \phi \mu_1, \Sigma) \phi + p(x|y=0; \phi \mu_0, \Sigma)(1-\phi)} = \frac{p(x|y=1, P)}{p(x|y=1, P) \phi + p(x|y=0, P)(1-\phi)} \\ &\stackrel{!}{=} \frac{1}{1 + \left(\frac{1-\phi}{\phi}\right) \left(\frac{p(x|y=0, P)}{p(x|y=1, P)} \right)} \xleftarrow{r} \end{aligned}$$

$$r = \frac{\exp(-\frac{1}{2}(x-\mu_0)^T \Sigma^{-1} (x-\mu_0))}{\exp(-\frac{1}{2}(x-\mu_1)^T \Sigma^{-1} (x-\mu_1))} = \exp \left[-\frac{1}{2} (x-\mu_0)^T \Sigma^{-1} (x-\mu_0) + \frac{1}{2} (x-\mu_1)^T \Sigma^{-1} (x-\mu_1) \right]$$

$$= \exp \left[\frac{1}{2} \left(-(x^T \Sigma^{-1} - \mu_0^T \Sigma^{-1})(x-\mu_0) + (x^T \Sigma^{-1} - \mu_1^T \Sigma^{-1})(x-\mu_1) \right) \right]$$

$$= \exp \left(\frac{1}{2} \left(-x^T \cancel{\Sigma^{-1}} x + x^T \Sigma^{-1} \mu_0 + \mu_0^T \Sigma^{-1} x - \mu_0^T \Sigma^{-1} \mu_0 + \cancel{x^T \Sigma^{-1} x} - \cancel{x^T \Sigma^{-1} \mu_1} \right. \right. \\ \left. \left. - \mu_1^T \Sigma^{-1} x + \mu_1^T \Sigma^{-1} \mu_1 \right) \right)$$

$$= \exp \left(\frac{1}{2} \left((\mu_0^T - \mu_1^T) \Sigma^{-1} x + x^T \Sigma^{-1} (\mu_0 - \mu_1) - \mu_0^T \Sigma^{-1} \mu_0 + \mu_1^T \Sigma^{-1} \mu_1 \right) \right)$$

$$= \exp \left(\frac{1}{2} \left(2(\mu_0^T - \mu_1^T) \Sigma^{-1} x - \mu_0^T \Sigma^{-1} \mu_0 + \mu_1^T \Sigma^{-1} \mu_1 \right) \right)$$

$$p(y=1|x, \phi, \mu_0, \mu_1, \Sigma) = \frac{1}{1 + \exp(\log(\frac{1-\phi}{\phi}))} \exp \left(\frac{(\mu_0^T - \mu_1^T) \Sigma^{-1} x - \frac{1}{2} \mu_0^T \Sigma^{-1} \mu_0 + \frac{1}{2} \mu_1^T \Sigma^{-1} \mu_1}{\log(\frac{1-\phi}{\phi})} \right)$$

$$= \frac{1}{1 + \exp \left(\frac{(\mu_0^T - \mu_1^T) \Sigma^{-1} x - \frac{1}{2} \mu_0^T \Sigma^{-1} \mu_0 + \frac{1}{2} \mu_1^T \Sigma^{-1} \mu_1 + \log(\frac{1-\phi}{\phi})}{\log(\frac{1-\phi}{\phi})} \right)}$$

$$= \frac{1}{1 + \exp \left(- \left((\mu_1 - \mu_0)^T \Sigma^{-1} x + \frac{1}{2} \mu_0^T \Sigma^{-1} \mu_0 - \frac{1}{2} \mu_1^T \Sigma^{-1} \mu_1 + \log(\frac{\phi}{1-\phi}) \right) \right)}$$

$$\text{let } \underline{\theta}^T = (\mu_0 - \mu_1 \Sigma^{-1}), \quad \underline{\theta}_0 = \frac{1}{2}(\mu_0^T \Sigma^{-1} \mu_0 - \mu_1^T \Sigma^{-1} \mu_1) + \log\left(\frac{\phi}{1-\phi}\right)$$

then

$$p(y=1 | x, \phi, \mu_0, \mu_1, \Sigma) = \frac{1}{1 + \exp(-(\underline{\theta}^T x + \underline{\theta}_0))}$$