

① Derivation of GDA model

Comparison with Discriminative algorithm:

In Disc we have to classify $y=0$ or $y=1$ (binary example) given features x . Logistic Reg tried to find a straight line (Decision boundary) that separates the two levels. While, GDA looks at how $y=1$ and $y=0$ are made up of and based the prediction of a example by comparing if it fits $y=1$ or $y=0$ properties.

Disc Algo: Find $P(y|x)$ directly

GDA: Find $P(x|y)$ and $P(y)$ then use Bayes rule to find $P(y|x)$
 ↙ makes strong assumption
 e.g., $P(x|y=0)$, $P(x|y=1)$
 $P(y)$ → class prior

Use Bayes rule, $P(y|x) = \frac{P(x|y) P(y)}{P(x)}$

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

Model:

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

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

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

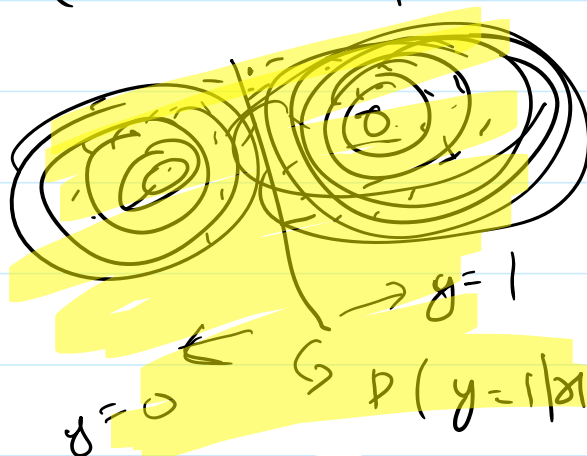
$$x|y=1 \sim \mathcal{N}(\mu_1, \Sigma)$$

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

} → find prob mass/density funct

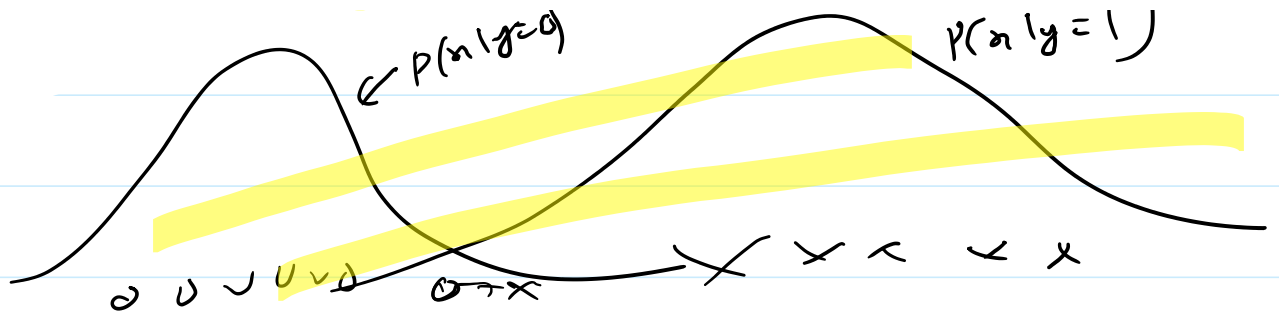
$$\begin{aligned}
 \ell(\phi, \mu_0, \mu_1, \Sigma) &= \log \prod_{i=1}^n P(x^{(i)}, y^{(i)}; \phi, \mu_0, \mu_1, \Sigma) \\
 &= \log \prod_{i=1}^n P(x^{(i)}|y^{(i)}; \mu_0, \mu_1, \Sigma) P(y^{(i)}; \phi)
 \end{aligned}$$

Take derivative of ℓ wrt parameters $(\phi, \mu_0, \mu_1, \Sigma)$ to find the parameter that maximize the likelihood. (Done in Problem Set below)

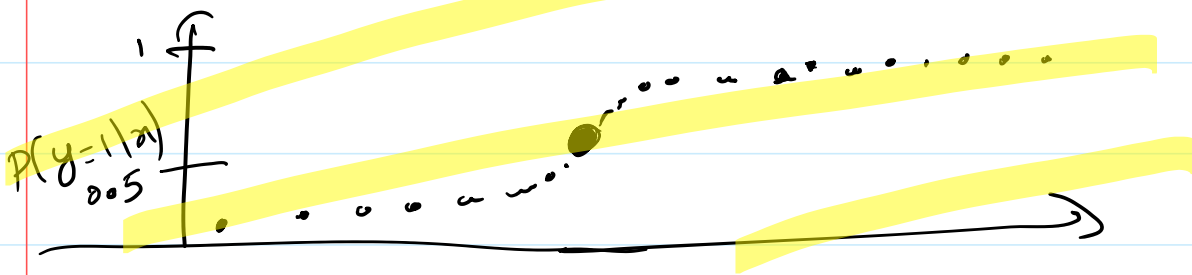


Find gaussian controlled by μ_0, μ_1 to shift with shape shape Σ .

GDA and logistic



$P(y=1) = 0.5$ Since data is half and half



Generative

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

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

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

Discriminative

(logistic)

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

~~\Rightarrow~~

(B) Problem Set 1 (C) (d) (e)

(c) [5 points] Recall that in GDA we model the joint distribution of (x, y) by the following equations:

$$p(y) = \begin{cases} \phi & \text{if } y = 1 \\ 1 - \phi & \text{if } y = 0 \end{cases}$$

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

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

where ϕ , μ_0 , μ_1 , and Σ are the parameters of our model.

Suppose we have already fit ϕ , μ_0 , μ_1 , and Σ , and now want to predict y given a new point x . To show that GDA results in a classifier that has a linear decision boundary, show the posterior distribution can be written as

$$p(y = 1 \mid x; \phi, \mu_0, \mu_1, \Sigma) = \frac{1}{1 + \exp(-(\theta^T x + \theta_0))},$$

where $\theta \in \mathbb{R}^n$ and $\theta_0 \in \mathbb{R}$ are appropriate functions of ϕ , Σ , μ_0 , and μ_1 .

Answer:

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

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

$$1 + \frac{\exp(-\frac{1}{2} (x-w_1)^T \Sigma^{-1} (x-w_1)) (1-\phi)}{\exp(-\frac{1}{2} (x-w_0)^T \Sigma^{-1} (x-w_0)) (\phi)}$$

$$x = e^{\left(-\frac{1}{2}(x-a)^T \Sigma^{-1}(x-a)\right) - \left(-\frac{1}{2}(x-a)^T \Sigma^{-1}(x-a)\right)} e^{\log\left(\frac{1-\phi}{\phi}\right)}$$

$$\alpha = \rho \left(-\frac{1}{2} (x_{cl})^T \Sigma^{-1} (x_{cl}) + \frac{1}{2} (x_{ub})^T \Sigma^{-1} (x_{ub}) \right) e^{\log\left(\frac{1-\phi}{\phi}\right)}$$

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

So we can take

$$\theta = (\mu_1 - \mu_0)^T \Sigma^{-1}, \theta_0 = \frac{1}{2}(\mu_0^T \Sigma^{-1} \mu_0 - \mu_1^T \Sigma^{-1} \mu_1) - \log \frac{1-\phi}{\phi}$$

\Rightarrow 'e'
for coding part *
- - - - -

Using $(A+B)C = (AC) + (BC)$

Also $A^T + B^T = (A+B)^T$

(d) [7 points] For this part of the problem only, you may assume n (the dimension of x) is 1, so that $\Sigma = [\sigma^2]$ is just a real number, and likewise the determinant of Σ is given by $|\Sigma| = \sigma^2$. Given the dataset, we claim that the maximum likelihood estimates of the parameters are given by

$$\begin{aligned}\phi &= \frac{1}{m} \sum_{i=1}^m 1\{y^{(i)} = 1\} \\ \mu_0 &= \frac{\sum_{i=1}^m 1\{y^{(i)} = 0\} x^{(i)}}{\sum_{i=1}^m 1\{y^{(i)} = 0\}} \\ \mu_1 &= \frac{\sum_{i=1}^m 1\{y^{(i)} = 1\} x^{(i)}}{\sum_{i=1}^m 1\{y^{(i)} = 1\}} \\ \Sigma &= \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \mu_{y^{(i)}})(x^{(i)} - \mu_{y^{(i)}})^T\end{aligned}$$

The log-likelihood of the data is

$$\begin{aligned}\ell(\phi, \mu_0, \mu_1, \Sigma) &= \log \prod_{i=1}^m p(x^{(i)}, y^{(i)}; \phi, \mu_0, \mu_1, \Sigma) \\ &= \log \prod_{i=1}^m p(x^{(i)} | y^{(i)}; \mu_0, \mu_1, \Sigma) p(y^{(i)}; \phi).\end{aligned}$$

By maximizing ℓ with respect to the four parameters, prove that the maximum likelihood estimates of ϕ , μ_0 , μ_1 , and Σ are indeed as given in the formulas above. (You may assume that there is at least one positive and one negative example, so that the denominators in the definitions of μ_0 and μ_1 above are non-zero.)

Answer:

$$\begin{aligned}\ell &= \log \prod_{i=1}^m 1\{y^{(i)}=0\} p(x^{(i)}|y^{(i)}=0) p(y^{(i)}=0) + 1\{y^{(i)}=1\} p(x^{(i)}|y^{(i)}=1) p(y^{(i)}=1) \\ &= \log \prod_{i=1}^m \left[1\{y^{(i)}=0\} \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x^{(i)}-\mu_0)^2}{2\sigma^2}\right) (1-\phi) + \dots \right] \\ &= \sum_{i=1}^m \log \left[1\{y^{(i)}=0\} \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x^{(i)}-\mu_0)^2}{2\sigma^2}\right) (1-\phi) + \dots \right] \\ &= \sum_{i=1}^m \log \left(1\{y^{(i)}=0\} \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x^{(i)}-\mu_0)^2}{2\sigma^2}\right) + 1\{y^{(i)}=1\} \dots \right) + \log \left(\phi^{n_1} (1-\phi)^{n_0} \right)\end{aligned}$$

$$\frac{\partial \ell}{\partial \phi} = \sum_{i=1}^m \frac{\partial}{\partial \phi} \left[y^{(i)} \log(\phi) + (1-y^{(i)}) \log(1-\phi) \right] = 0$$

$$\sum_{i=1}^m \frac{y^{(i)}}{\phi} + \frac{1-y^{(i)}}{1-\phi} = 0$$

$$\sum_{i=1}^m y^{(i)} - \phi y^{(i)} + \phi - \phi y^{(i)} = 0$$

$$\sum_{i=1}^m y^{(i)} + \phi \sum_{i=1}^m (1 - 2y^{(i)}) = 0 \Rightarrow \phi = \frac{\sum_{i=1}^m y^{(i)}}{\sum_{i=1}^m (2y^{(i)} - 1)}$$

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

$$\Rightarrow \sum_{i=1}^m 1 \left[\text{For } y^{(i)} \right]$$

$$l = \frac{\partial}{\partial \mu_0} \sum_{i=1}^m \log \left(1\{y^{(i)}=0\} \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x^{(i)} - \mu_0)^2}{2\sigma^2}\right) \right)$$

$$\frac{d}{dx} e^{-(x-2)^2} = \frac{1}{e^{(x-2)^2}} \frac{d}{dx} e^{-(x-2)^2}$$

$$= \sum_{i=1}^m 1\{y^{(i)}=0\} \frac{\partial}{\partial \mu_0} \log \left(e^{-\frac{(x^{(i)} - \mu_0)^2}{2\sigma^2}} \right) \quad \frac{d}{dx} \log(e^x) = x \quad \left[\text{log and exp inverse function} \right]$$

$$= \sum_{i=1}^m 1\{y^{(i)}=0\} \frac{\partial}{\partial \mu_0} \left(-\frac{(x^{(i)} - \mu_0)^2}{2\sigma^2} \right) = -1\{y^{(i)}=0\} \sum_{i=1}^m \frac{x^{(i)} - \mu_0}{\sigma^2} = 0$$

$$\text{So } -\sum_{i=1}^m 1\{y^{(i)}=0\} x^{(i)} + \sum_{i=1}^m 1\{y^{(i)}=0\} \mu_0 = 0$$

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

Similarly for σ^2

Note $y = \sigma^2$
 $\sqrt{y} = \sigma$

Bad approximation of GDA on data set \Rightarrow Since the dataset is not Gaussian
 Fix (take $\log(x^{(i)}/s)$)

② Implementation Notes

Algorithm

- ① Compute $\phi, \mu_0, \mu_1, \Sigma$ using above formula
- ② Compute θ using '*'

③ Find $P(y=1|x)$, if $P(y=1|x) \geq 0.5$
predict that label is '1' else predict
that label is 0.