

Generative Learning algorithms

Discriminative vs generative

discriminative: learn $p(y|x)$ directly or mapping directly from x to y .

generative: $p(y|x) = \frac{p(x|y)p(y)}{p(x)} \Rightarrow \text{model } p(x|y)$

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

Gaussian discriminative analysis (GDA)

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

where $\mu \in \mathbb{R}^d$, $\Sigma \in \mathbb{R}^{d \times d}$ (covariance matrix).

假设特征 x 是连续随机变量, 则 GDA 模型为:

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

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

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

\therefore log likelihood is

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

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

$$= \sum_{i=1}^n \log p(x^{(i)} | y^{(i)}; \mu_0, \mu_1, \Sigma) + \sum_{i=1}^n \log p(y^{(i)}; \phi)$$

$$\therefore \frac{\partial l}{\partial \phi} = \frac{\partial}{\partial \phi} \sum_{i=1}^n \log p(y^{(i)}; \phi)$$

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

$$= \frac{\partial}{\partial \phi} \sum_{i=1}^n [y^{(i)} \log \phi + (1-y^{(i)}) \log(1-\phi)]$$

$$= \frac{\partial}{\partial \phi} \sum_{i=1}^n [y^{(i)} \log \phi + \log(1-\phi) - y^{(i)} \log(1-\phi)]$$

$$= \sum_{i=1}^n \left(y^{(i)} \frac{1}{\phi} + \frac{1}{1-\phi} (-1) - y^{(i)} \frac{1}{1-\phi} (-1) \right)$$

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

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

$$\therefore \frac{\partial l}{\partial \mu_0} = \frac{\partial}{\partial \mu_0} \log \exp \left[-\frac{1}{2} (x - \mu_0)^T \Sigma^{-1} (x - \mu_0) \right]$$

$$= -\frac{1}{2} \frac{\partial}{\partial \mu_0} \sum_{i=1}^n (x^{(i)} - \mu_0)^T \Sigma^{-1} (x^{(i)} - \mu_0)$$

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

$$= -\frac{1}{2} ((\Sigma^{-1})^T - \Sigma^{-1}) (x^{(i)} - \mu_0) = 0$$

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

$$\text{同理: } \frac{\partial l}{\partial \mu_1} = 0 \Rightarrow \mu_1 = \frac{\sum_{i=1}^n I\{y^{(i)} = 1\} x^{(i)}}{\sum_{i=1}^n I\{y^{(i)} = 1\}}$$

$$S = \Sigma^{-1}$$

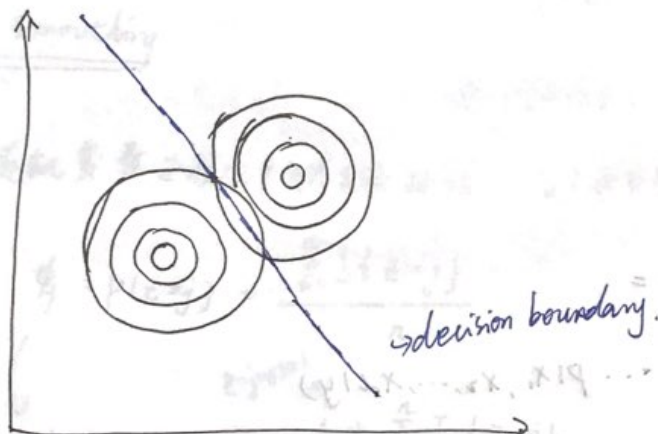
$$\therefore \frac{\partial L}{\partial \Sigma} = \frac{\partial}{\partial \Sigma} - \frac{1}{2} \sum_{i=1}^n [\log |S|^{-1} + (x^{(i)} - \mu_{y^{(i)}})^T S (x^{(i)} - \mu_{y^{(i)}})]$$

$$= \frac{\partial}{\partial \Sigma} - \frac{1}{2} \sum_{i=1}^n [-\log |S| + (x^{(i)} - \mu_{y^{(i)}})^T S (x^{(i)} - \mu_{y^{(i)}})]$$

$$= -\frac{1}{2} \sum_{i=1}^n [-2 + (x^{(i)} - \mu_{y^{(i)}}) (x^{(i)} - \mu_{y^{(i)}})^T]$$

$$= \frac{1}{2} \sum_{i=1}^n 2 - \frac{1}{2} \sum_{i=1}^n (x^{(i)} - \mu_{y^{(i)}}) (x^{(i)} - \mu_{y^{(i)}})^T = 0$$

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



GDA and logistic regression

当 $p(x|y)$ 是变量, 高斯分布时, $p(y|x)$ 满足逻辑回归的公式:

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

⇒ 当 $p(y|x)$ 满足高斯分布时 ~~不是~~ GDA 优于 LR.

Naive Bayes

特征是独立的.

Naive Bayes Assumption.

$$\begin{aligned} p(x_1, \dots, x_d | y) &= \\ &= p(x_1 | y) p(x_2 | y) \dots p(x_d | y) \\ &= p(x_1 | y) p(x_2 | y) \dots p(x_d | y) = \prod_{j=1}^d p(x_j | y) \end{aligned}$$



$$L(\phi_y, \phi_{j|y=0}, \phi_{j|y=1}) = \prod_{i=1}^n p(x^{(i)} | y^{(i)}) p(y^{(i)})$$

where $\phi_y = p(y=1)$, $\phi_{j|y=0} = p(x_j=1 | y=0)$, $\phi_{j|y=1} = p(x_j=1 | y=1)$

$$\Rightarrow \phi_{j|y=1} = \frac{\sum_{i=1}^n \mathbb{I}\{x_j^{(i)}=1 \wedge y^{(i)}=1\}}{\sum_{i=1}^n \mathbb{I}\{y^{(i)}=1\}}, \quad \phi_{j|y=0} = \frac{\sum_{i=1}^n \mathbb{I}\{x_j^{(i)}=1 \wedge y^{(i)}=0\}}{\sum_{i=1}^n \mathbb{I}\{y^{(i)}=0\}}$$

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

预测:

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

$$\propto p(x|y=1) p(y=1)$$

$$= \prod_{j=1}^d p(x_j|y=1) \cdot p(y=1)$$

假设输入的特征 x_j 只取 k 种不同的值, 则 $p(x_j|y)$ 应该建模为多项式分布.

且若 x_j 为连续值, 需要离散化操作.

Laplace smoothing

假设随机变量 z 服从多项式分布, 从 $1, \dots, k$ 中取值. 给定 n 个训练样本

$$\phi = p(z=j) = \frac{\sum_{i=1}^n I\{z=j\}}{n}$$

smoothing.

$$\phi' = p(z=j)' = \frac{1 + \sum_{i=1}^n I\{z=j\}}{k + n}$$

\downarrow
每个维度的取值

Text classification with Naïve

假设 d_i 表示第 i 个文档, c 表示类别且有 K 种取值.

1) Bernoulli model

文档使用 one-hot 方式表示,

$$p(d_i | c) = p(d_i | c) p(c) = \prod_{j=1}^V p(d_{ij} | c) p(c)$$

$$p(c) = \frac{n_k}{n} \rightarrow \text{属于类别 } k \text{ 的文档数} \\ \rightarrow \text{总文档数.}$$

$$p(d_{ij} | c) = \frac{n_k(w_j)}{n_k} \rightarrow \text{属于类别 } k \text{ 且包含单词 } w_j \text{ 的文档数.}$$

laplace smoothing

$$\hat{p}(d_{ij} | c) = \frac{n_k(w_j) + 1}{n_k + 2}$$

2) Multinomial model

文档使用 multi-hot (term-frequency) 方式表示.

$$p(d_i | c) = p(d_i | c) p(c) = \prod_{j=1}^V p(d_{ij} | c) p(c)$$

$$p(c) = \frac{n_k}{n}$$

$$p(d_{ij} | c) = \frac{n w_k(w_j)}{n w_k} \rightarrow \text{属于类别 } k \text{ 的文档中 } w_j \text{ 出现的次数之和} \\ \rightarrow \text{属于类别 } k \text{ 的文档中单词出现的总次数之和}$$

Laplace smoothing

$$\hat{p}(d_{ij} | c) = \frac{n w_k(w_j) + 1}{n w_k + V}$$