

$$P(x_1, \dots, x_d | y), \quad x_j \in \{1, 2, \dots, m\}$$

$$\# \text{ parameters} \equiv m^d \times k$$

ex:- Text classificatⁿ. $X = x_1, \dots, x_d$
 $x_j =$ whether doc. contains j^{th} word in vocab or not.

$$y \in \{\text{topics} \in \{\text{sports, science, ...}\} \} \quad (k=500)$$

$$d = 50,000$$

$$\# \text{ parameters} = 2^{50000} \times k$$

Naive Bayes assumption

$$P(x_1, \dots, x_d | y) = \prod_{j=1}^d P(x_j | y)$$

ie. conditioned on the class variable, the attributes are independent of each other. This is quite reasonable and holds approximately in many real applications.

~~$$P(x_1, \dots, x_d) = \prod_{i=1}^d P(x_i)$$~~

Doc classificatⁿ using NB

$$P(x_1, \dots, x_d | y) = \prod_{j=1}^d P(x_j | y)$$

For each subject or class, a prob. of occurrence of each word in the vocab.

$$\theta_{kj} = P(x_j | y = k) \Rightarrow \text{Learned during training}$$

$$\therefore \# \text{ of parameters} = d \times k$$

$$X[0,1,0,0,0,1,0,1] \quad d=8$$

$$P(X|y=k) = (1-\theta_{k1})\theta_{k2}(1-\theta_{k3})\dots(1-\theta_{kd})\theta_{k1}$$

$$\therefore P(y=k|X) \propto P(y=k) \times P(X|y=k)$$

Then normalise it $\prod_{j=1}^d \frac{1}{N} (1-\theta_{kj}) \theta_{kj}^{x_{kj}}$

Gaussian Discriminant Analysis (GDA)

Only works for real valued attr^b.

$$P(x_1, \dots, x_d | y) \sim N(\mu_y, \Sigma_{d \times d})$$

$$(x_1, \dots, x_d \in \mathbb{R}^d) \quad \mu_y \in \mathbb{R}^d$$

Multidimensional gaussian

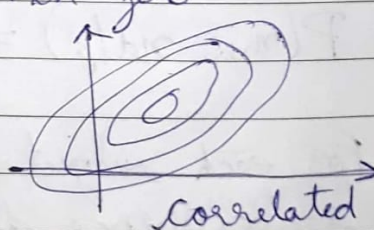
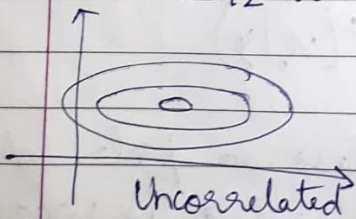
$$[x_1, x_2] \sim N\left([0.5, 1], \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{bmatrix} = \Sigma\right)$$

When x_1 & x_2 are un-correlated

σ_1, σ_2 : std. dev. of attrib. 1 & 2.

When x_1 & x_2 are correlated,

Σ_{12} will be non-zero



Σ_{je} : covariance bⁿ attributes x_j & x_e

Σ_{jj} : variance of attribute x_j

$$x_1, \dots, x_d \sim N(\mu, \Sigma_{d \times d})$$

\downarrow
 $d \times 1$

$$\therefore P(x_1, \dots, x_d) = \frac{1}{\sqrt{2\pi} |\Sigma|^{d/2}} e^{-\frac{(x-\mu)^T \Sigma^{-1} (x-\mu)}{2}}$$

Suppose $d=2$,

$$\begin{aligned} \log P(x_1, x_2) &= -\frac{1}{2} [x_1, x_2] \begin{bmatrix} \sigma_1^2 & \sigma_{12}^2 \\ \sigma_{12}^2 & \sigma_2^2 \end{bmatrix}^{-1} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + \dots \\ &= -\frac{1}{2} [x_1, x_2] \begin{bmatrix} p_{11} & p_{12} \\ p_{21} & p_{22} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \\ &= -\frac{1}{2} [p_{11}x_1^2 + p_{21}x_2x_1 + p_{12}x_1x_2 + p_{22}x_2^2] \end{aligned}$$

$$P(x_1, \dots, x_d | Y=k) \sim N(\mu_y, \Sigma_y)$$

\uparrow \uparrow
 d parameters d^2 parameters
 k classes

Parameter estimation

For each class k ,

$$\hat{\mu}_k \equiv \sum_{\substack{i=1 \\ y_i^o = k}}^N x_i^o$$

$$\sum_{i=1}^N [y_i^o = k] = N_k \leftarrow \begin{array}{l} \text{No. of obs.} \\ \text{in } k^{\text{th}} \text{ class.} \end{array}$$

$$\hat{\Sigma}_k = \frac{1}{N_k} \sum_{i=1}^N (x_i^o - \hat{\mu}_k)(x_i^o - \hat{\mu}_k)^T$$

$N_k - 1$ for unbiased

Alternately,

σ_j^2 = Variance of j^{th} attributes

σ_{je} = co-variance b/w attributes j & e

$$\Sigma_K = \begin{bmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{1j} & \sigma_2^2 \end{bmatrix} \text{ --- Symmetric matrix}$$

Linear Discriminant Analysis (LDA)

Assume Σ_K is the same for all K classes.

$$\Sigma_1 = \Sigma_2 = \dots = \Sigma_K$$

$$P(x_1, \dots, x_d | y = K) = N(\mu_K, \Sigma)$$

8/08

The estimate of the shared covariance is

$$\Sigma = \frac{\sum_{i=1}^N (X^i - \mu_{y_i})(X^i - \mu_{y_i})^T}{N-1}$$

Decision Boundary

$$k=1,2 \quad \mu_1, \Sigma_1, \mu_2, \Sigma_2, \quad P(y=1) \equiv \pi_1 \\ P(y=2) \equiv \pi_2 \\ \pi_1 + \pi_2 = 1$$

Under what X is,

$$P(y=1|X) > P(y=2|X) ?$$

$$\therefore \log P(y=1|x) > \log P(y=2|x)$$

$$\therefore \log P(y=1) P(x|y=1) > \log P(y=2) P(x|y=2)$$

$$\therefore \log \pi_1 + \log \frac{1}{\sqrt{2\pi} |\Sigma_1|^{d/2}} e^{-\frac{1}{2}(x-M_1)^T \Sigma_1^{-1} (x-M_1)}$$

$$\log \pi_1 - \log |\Sigma_1|^{d/2} - \frac{1}{2}(x-M_1)^T \Sigma_1^{-1} (x-M_1) > \log \pi_2 - \log |\Sigma_2|^{d/2}$$

$$\rightarrow \text{const. \& indep of } x. \quad -\frac{1}{2}(x-M_2)^T \Sigma_2^{-1} (x-M_2)$$

$$\therefore (x-M_2)^T \Sigma_2^{-1} (x-M_2) - (x-M_1)^T \Sigma_1^{-1} (x-M_1) > \text{const.}$$

$$\begin{aligned} & -x \Sigma_1^{-1} x + 2M_1^T \Sigma_1^{-1} x - M_1^T \Sigma_1^{-1} M_1 + x \Sigma_2^{-1} x \\ & - 2M_2^T \Sigma_2^{-1} x + M_2^T \Sigma_2^{-1} M_2 > \text{const.} \end{aligned}$$

$$\therefore -x^T [\Sigma_1^{-1} - \Sigma_2^{-1}] x + 2(M_1^T \Sigma_1^{-1} - M_2^T \Sigma_2^{-1}) x > \text{const.}$$

When $\Sigma_1 = \Sigma_2$; the quadratic term vanishes.
The 2nd term,

$$2(M_1^T - M_2^T) \Sigma^{-1} x = a_1 x_1 + a_2 x_2 + \dots + a_d x_d$$

↑
linear in x .

Decision boundary

$$a_1 x_1 + a_2 x_2 + \dots + a_d x_d > \text{const}^n$$

When $\Sigma_1 \neq \Sigma_2$

$$X^T (\Sigma_1^{-1} - \Sigma_2^{-1}) X = \sum_{l=1, j=1}^{d, d} (p_{xj1} - p_{xj2}) x_l x_j$$

quadratic decision boundary

+ linear terms

Boundary for naive Bayes

$$\sum_{j=1}^d \left(\frac{1}{\sigma_{j1}^2} - \frac{1}{\sigma_{j2}^2} \right) x_j^2 + \text{linear terms}$$

The max. likelihood training method

Identify the form of the distribⁿ that characterizes

Let θ = the parameters of distribⁿ.

The max. likelihood estimates of θ are the ones which maximize the prob^l of the training data.

Typically assume that training samples are iid.

$D \equiv \{$

generative $P(x)$

log. li

\equiv

During LL

$\hat{\theta}^{ML}$

Apply

$P(x)$

$\theta =$

$|a|$

LL

$$D \equiv \{(x^1, y_1), (x^2, y_2), \dots, (x^N, y_N)\}$$

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

log. likelihood of D given $P(x, y|\theta)$

$$\equiv LL(\theta, D) = \sum_{i=1}^N \log P(x^i, y_i|\theta)$$

During training, find θ for which $LL(\theta, D)$ is maximized.

$$\hat{\theta}^{ML} = \arg \max_{\theta} LL(\theta, D)$$

Apply ML training for LDA

$$P(y=1) = \pi_1, P(y=2) = \pi_2, \mu_1, \mu_2, \Sigma$$

$$\theta \equiv \{ \pi_1, \pi_2, \underbrace{\mu_1}_{d \times 1}, \underbrace{\mu_2}_{d \times 1}, \underbrace{\Sigma}_{d \times d} \}$$

$$|\theta| = d^2 + 2d + 2$$

$$LL(\theta, D) = \sum_{i=1}^N \log P(x^i, y_i|\theta)$$

$$= \sum_{i=1}^N \log P(y_i) P(x^i|y_i, \theta)$$

$$= \sum_{i=1, y_i=1}^N \log \pi_1 + \log \frac{1}{\sqrt{2\pi} |\Sigma|^{d/2}} e^{-\frac{1}{2} (x^i - \mu_1)^T \Sigma^{-1} (x^i - \mu_1)}$$

$$+ \sum_{i=1, y_i=2}^N \log \pi_2 + \log \frac{1}{\sqrt{2\pi} |\Sigma|^{d/2}} e^{-\frac{1}{2} (x^i - \mu_2)^T \Sigma^{-1} (x^i - \mu_2)}$$

$$\begin{aligned}
 & \therefore N_1 \log \pi_1 + \sum_{\substack{i=1 \\ y_i=1}}^N \left(-\frac{1}{2} \right) (X_i^p - M_1)^T \Sigma^{-1} (X_i^p - M_1) \\
 & + \sum_{\substack{i=1 \\ y_i=2}}^N \left(-\frac{1}{2} \right) (X_i^p - M_2)^T \Sigma^{-1} (X_i^p - M_2) \\
 & + N_2 \log \pi_2 - \frac{Nd \log |\Sigma|}{2} + \text{const}^n
 \end{aligned}$$

$$\therefore \max_{\pi_1, \pi_2, M_1, M_2, \Sigma} \textcircled{LL}(\pi_1, \pi_2, M_1, M_2, \Sigma) \text{ s.t. } \pi_1 + \pi_2 = 1$$

$\rightarrow \nabla$. with each parameters $\Delta = 0$
 if LL is concave