

高斯分布(正态分布)

一、极大似然估计.

Data: $X = (x_1, x_2, \dots, x_N)^T = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_N \end{pmatrix}_{N \times p}$, $x_i \in \mathbb{R}^p$, $x_i \sim N(\mu, \Sigma)$

全参数 $\theta = (\mu, \Sigma)$

则极大似然估计MLE: $\theta_{MLE} = \arg \max_{\theta} P(X | \theta)$
为简化计算, 全 $p=1$, 则 $\theta = (\mu, \sigma^2)$

推导:

$\because x_i$ 独立同分布, 则

$$\begin{aligned} \log P(X | \theta) &= \log \prod_{i=1}^N P(x_i | \theta) = \sum_{i=1}^N \log P(x_i | \theta) \\ &= \sum_{i=1}^N \log \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x_i - \mu)^2}{2\sigma^2}\right) \\ &= \sum_{i=1}^N \left[\log \frac{1}{\sqrt{2\pi}\sigma} + \log \frac{1}{\sigma} - \frac{(x_i - \mu)^2}{2\sigma^2} \right] \end{aligned}$$

$$\begin{aligned} \mu_{MLE} &= \arg \max_{\mu} \log P(X | \theta) \\ &= \arg \max_{\mu} \sum_{i=1}^N \frac{(x_i - \mu)^2}{2\sigma^2} \\ &= \arg \min_{\mu} \sum_{i=1}^N (x_i - \mu)^2 \end{aligned}$$

对 μ 求导:

$$\begin{aligned} \frac{\partial}{\partial \mu} \sum_{i=1}^N (x_i - \mu)^2 &= \sum_{i=1}^N 2(x_i - \mu) \cdot (-1) = 0 \\ \sum_{i=1}^N (x_i - \mu) &= 0 \\ \underbrace{\sum_{i=1}^N x_i}_{N \cdot \mu} - \underbrace{\sum_{i=1}^N \mu}_{N \cdot \mu} &= 0 \end{aligned}$$

$\mu_{MLE} = \frac{1}{N} \sum_{i=1}^N x_i$ 即样本的均值.

iid 即独立同分布.

P1:

$$P(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

P2:

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

$$\sigma^2_{MLE} = \arg \max_{\sigma} p(x|\theta) = \arg \max_{\sigma} \underbrace{\left(-\log \sigma - \frac{1}{2\sigma^2} \sum_{i=1}^N (x_i - \mu)^2 \right)}_L$$

$$\frac{\partial L}{\partial \sigma} = \sum_{i=1}^N \frac{1}{\sigma} - \frac{1}{2} (x_i - \mu)^2 \cdot (-2)\sigma^{-3} = 0$$

$$\sum_{i=1}^N \left[\frac{1}{\sigma} + (x_i - \mu)^2 \sigma^{-3} \right] = 0$$

$$\sum_{i=1}^N \left[-\sigma^2 + (x_i - \mu)^2 \right] = 0$$

$$-\sum_{i=1}^N \sigma^2 + \sum_{i=1}^N (x_i - \mu)^2 = 0$$

$$\sum_{i=1}^N \sigma^2 = \sum_{i=1}^N (x_i - \mu)^2$$

$$\boxed{\sigma^2_{MLE} = \frac{1}{N} \sum_{i=1}^N (x_i - \mu_{MLE})^2}$$

σ^2_{MLE} 是有偏估计， μ_{MLE} 是无偏估计

$$\sigma^2_{MLE} : E[\sigma^2_{MLE}] = \frac{N-1}{N} \sigma^2, \text{真正的方差偏应为 } \hat{\sigma} = \frac{1}{N-1} \sum_{i=1}^N (x_i - \mu_{MLE})^2$$

$$\mu_{MLE} : E[\mu_{MLE}] = \frac{1}{N} \sum_{i=1}^N E[x_i] = \frac{1}{N} \sum_{i=1}^N \mu = \frac{1}{N} \cdot N \cdot \mu = \mu$$

$$\begin{aligned} \sigma_{MLE}^2 &= \frac{1}{N} \sum_{i=1}^N (x_i - \mu_{MLE})^2 = \frac{1}{N} \sum_{i=1}^N (x_i^2 - 2x_i\mu_{MLE} + \mu_{MLE}^2) = \frac{1}{N} \sum_{i=1}^N x_i^2 - \underbrace{\frac{1}{N} \sum_{i=1}^N 2x_i\mu_{MLE}}_{2\cdot\mu_{MLE}\cdot\mu_{MLE}} + \underbrace{\frac{1}{N} \sum_{i=1}^N \mu_{MLE}^2}_{\mu_{MLE}^2} \\ &= \frac{1}{N} \sum_{i=1}^N x_i^2 - \mu_{MLE}^2 \end{aligned}$$

$$\begin{aligned} E[\sigma^2_{MLE}] &= E\left[\frac{1}{N} \sum_{i=1}^N x_i^2 - \mu_{MLE}^2\right] = E\left[\frac{1}{N} \sum_{i=1}^N x_i^2 - (\mu_{MLE}^2 - \mu^2)\right] \\ &= E\left(\frac{1}{N} \sum_{i=1}^N x_i^2 - \mu^2\right) - E(\mu_{MLE}^2 - \mu^2) \\ &= E\left[\frac{1}{N} \sum_{i=1}^N (x_i^2 - \mu^2)\right] - [E(\mu_{MLE}^2) - E(\mu^2)] \\ &= \frac{1}{N} \sum_{i=1}^N E(x_i^2 - \mu^2) - [E(\mu_{MLE}^2) - \mu^2] \\ &= \frac{1}{N} \sum_{i=1}^N (E(x_i^2) - \bar{x}^2) - [E(\mu_{MLE}^2) - \bar{\mu}^2(\mu_{MLE})] \\ &= \frac{1}{N} \sum_{i=1}^N \sigma^2 - D(\mu_{MLE}) \end{aligned}$$

$$D(\mu_{MLE}) = D\left(\frac{1}{N} \sum_{i=1}^N x_i\right) = \frac{1}{N^2} \sum_{i=1}^N D(x_i) = \frac{1}{N^2} \sum_{i=1}^N \sigma^2 = \frac{1}{N} \sigma^2$$

$$\therefore E[\sigma^2_{MLE}] = \sigma^2 - \frac{1}{N} \sigma^2 = \frac{N-1}{N} \sigma^2$$

高维情形: $x \in \mathbb{R}^p$, $x = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_p \end{pmatrix}$, $\mu = \begin{pmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_p \end{pmatrix}$, $\Sigma = \begin{pmatrix} \sigma_{11} & \sigma_{12} & \cdots & \sigma_{1p} \\ \sigma_{21} & \sigma_{22} & \cdots & \sigma_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{p1} & \sigma_{p2} & \cdots & \sigma_{pp} \end{pmatrix}$, $p \times p$

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

$(x-\mu)^T \Sigma^{-1} (x-\mu)$, 马氏距离, (x 和 μ 之间的)

Σ 是协方差矩阵. 如果 Σ 为 I (单位阵), 则马氏距离=欧氏距离

$$\Sigma = U \Lambda U^T, \quad UU^T = U^T U = I, \quad \Lambda = \text{diag}(\lambda_i) \rightarrow \text{对角阵}$$

$$= (u_1, u_2, \dots, u_p) \begin{pmatrix} \lambda_1 & & \\ & \lambda_2 & \\ & & \lambda_p \end{pmatrix} \begin{pmatrix} u_1^T \\ u_2^T \\ \vdots \\ u_p^T \end{pmatrix}$$

$$= (u_1, u_2, \dots, u_p) \begin{pmatrix} u_1^T \\ u_2^T \\ \vdots \\ u_p^T \end{pmatrix}$$

$$= \sum_{i=1}^p \lambda_i u_i u_i^T \quad \because U \text{正交}, \text{则 } U^T = U^{-1}$$

$$\text{Q1) } \Sigma^{-1} = (U \Lambda U^T)^{-1} = (U^T)^{-1} \Lambda^{-1} U^{-1} = U \Lambda^{-1} U^T, \quad \Lambda^{-1} = \text{diag}\left(\frac{1}{\lambda_i}\right), i=1, \dots, p.$$

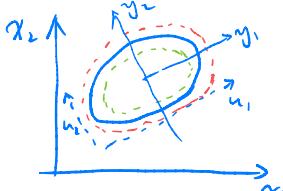
$$= \sum_{i=1}^p u_i \frac{1}{\lambda_i} u_i^T$$

$$\text{Q2) } (x-\mu)^T \Sigma^{-1} (x-\mu) = (x-\mu)^T \sum_{i=1}^p u_i \frac{1}{\lambda_i} u_i^T (x-\mu) = \sum_{i=1}^p (x-\mu)^T u_i \frac{1}{\lambda_i} u_i^T (x-\mu)$$

$$\sum y_i = (x-\mu)^T u_i, \quad y_i = \sum_{i=1}^p y_i \frac{1}{\lambda_i} y_i^T = \boxed{\frac{y_i^2}{\lambda_i}}$$

如果 $p=2$, 那么 $(x-\mu)^T \Sigma^{-1} (x-\mu) = \frac{y_1^2}{\lambda_1} + \frac{y_2^2}{\lambda_2}$, 令该值为 1.

观察 $\frac{y_1^2}{\lambda_1} + \frac{y_2^2}{\lambda_2} = 1$, 其为椭圆曲线.



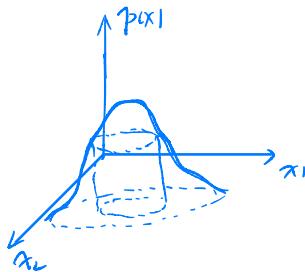
$y_i = (x-\mu)^T u_i$ 表示 x 到中心点 μ , $(x-\mu)$ 这个向量在 u_i 上的投影, 即将 $(x-\mu)$ 在 u_i 上的投影作为 y_i :

我们发现当 $\frac{y_1^2}{\lambda_1} + \frac{y_2^2}{\lambda_2} = 1$ 取不同值时, 设为 r, 对于某个固定的 r, 都能得到不同的椭圆曲线, 是一个等高线.

$$\text{对于 } X \sim N(\mu, \Sigma) = \frac{1}{(2\pi)^{\frac{p}{2}} |\Sigma|^{\frac{1}{2}}} \exp \left(-\frac{1}{2} \frac{(x-\mu)^T \Sigma^{-1} (x-\mu)}{\Delta} \right)$$

它的概率密度函数 $p(x) = \frac{1}{(2\pi)^{\frac{p}{2}} |\Sigma|^{\frac{1}{2}}} \exp(-\frac{1}{2}\Delta)$, 对于 $p(x)$ 来讲, 只有 Δ 是自由量, μ 与 Σ 是模型的参数.

$0 \leq p(x) \leq 1$, 当 $p(x)$ 取某一值时, 总会有一个 $\Delta = r$, 与之对应.



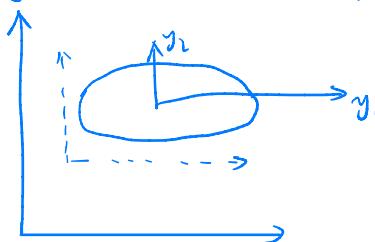
局限性:

① $\Sigma_{p \times p}$: 元素共 $\frac{p(p+1)}{2}$ 个, 又因其对称, 则共 $\frac{p(p+1)}{2}$ 个参数, 为 $O(p^2)$.
因此, 面对实际问题时, 尤其高维的问题, Δ 会很大,
参数会过多.

因此我们尝试对方差矩阵做简化.

$$\Sigma = \begin{pmatrix} \lambda_1 & \\ & \ddots & \lambda_p \end{pmatrix}, \text{ 只在对角线有值, 其余全为 } 0.$$

那么 Σ 就可以不用做特征值分解, 其中的 $y_i = (x - \mu)^T u_i$, 也就变成了计算 $(x - \mu)^T x_i$, 图象变为↓, 其方向与 x_1, x_2 一致.



若 $\lambda_1 = \lambda_2 = \dots = \lambda_p$, 则 A 为圆形, 称为各向同性.

② 高斯分布本身无法表达模型, 不能确切表达模型.
可用混合模型解决 (多个高斯)

P5 边缘概率和条件概率

$$\text{设 } X = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_p \end{pmatrix} \quad \mu = \begin{pmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_p \end{pmatrix} \quad \Sigma = \begin{pmatrix} \sigma_{11} & \cdots & \sigma_{1p} \\ \sigma_{21} & \cdots & \vdots \\ \vdots & \ddots & \vdots \\ \sigma_{p1} & \cdots & \sigma_{pp} \end{pmatrix}$$

$$\text{已知: } X = \begin{pmatrix} x_a \\ x_b \end{pmatrix} \xrightarrow{\text{m}} \begin{pmatrix} x_a \\ x_b \end{pmatrix} \xrightarrow{\text{n}} \begin{pmatrix} m_a \\ m_b \end{pmatrix} \quad \mu = \begin{pmatrix} \mu_a \\ \mu_b \end{pmatrix} \quad \Sigma = \begin{pmatrix} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ba} & \Sigma_{bb} \end{pmatrix}$$

$$P(X_a), P(X_b | X_a)$$

~~$P(X_b), P(X_a | X_b)$~~

定理: $X \sim N(\mu, \Sigma), y = AX + B, A \in \mathbb{R}^{n \times p}, y \sim N(A\mu + B, A\Sigma A^T)$

$$P(X_a): \quad X_a = \underbrace{(I_m, 0)}_A \underbrace{\begin{pmatrix} x_a \\ x_b \end{pmatrix}}_X, \quad E[X_a] = (I_m, 0) \begin{pmatrix} \mu_a \\ \mu_b \end{pmatrix} = \mu_a$$

$$\text{Var}[X_a] = (I_m, 0) \begin{pmatrix} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ba} & \Sigma_{bb} \end{pmatrix} \begin{pmatrix} I_m \\ 0 \end{pmatrix} = (\Sigma_{aa}, \Sigma_{ab}) \begin{pmatrix} I_m \\ 0 \end{pmatrix} = \Sigma_{aa}$$

$$\text{从而 } X_a \sim N(\mu_a, \Sigma_{aa})$$

$$P(X_b | X_a): \quad \hat{x}_b = X_b - \Sigma_{ba} \Sigma_{aa}^{-1} X_a$$

$$\left\{ \begin{array}{l} \mu_{b|a} = \mu_b - \Sigma_{ba} \Sigma_{aa}^{-1} \mu_a \\ \Sigma_{b|a} = \Sigma_{bb} - \Sigma_{ba} \Sigma_{aa}^{-1} \Sigma_{ab} \end{array} \right.$$

$$X_{ba} = (-\Sigma_{ba} \Sigma_{aa}^{-1} \quad I_n) \begin{pmatrix} x_a \\ x_b \end{pmatrix}$$

$$\text{则 } E[X_{b|a}] = (-\Sigma_{ba} \Sigma_{aa}^{-1} \quad I_n) \begin{pmatrix} \mu_a \\ \mu_b \end{pmatrix} = \mu_b - \Sigma_{ba} \Sigma_{aa}^{-1} \mu_a = \mu_{b|a}.$$

$$\begin{aligned} \text{Var}[X_{b.a}] &= (-\Sigma_{ba}\Sigma_{aa}^{-1} \quad I_n) \begin{pmatrix} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ba} & \Sigma_{bb} \end{pmatrix} \begin{pmatrix} -\Sigma_{ac}\Sigma_{ba}^T \\ I_n \end{pmatrix} \\ &= (0 \quad \Sigma_{bb} - \Sigma_{ba}\Sigma_{aa}^{-1}\Sigma_{ab}) \begin{pmatrix} -\Sigma_{ac}\Sigma_{ba}^T \\ I \end{pmatrix} = \Sigma_{bb} - \Sigma_{ba}\Sigma_{aa}^{-1}\Sigma_{ab} = \Sigma_{bb.a} \end{aligned}$$

$\therefore X_{b.a} \sim N(\mu_{b.a}, \Sigma_{bb.a})$

$$X_b = X_{b.a} + \underbrace{\Sigma_{ba}\Sigma_{aa}^{-1}X_a}_B$$

由上式可知 $X_{b.a} \perp \text{独立于 } X_a$, 且 $X_b | X_a = X_{b.a}$

$$\therefore E[X_b | X_a] = \mu_{b.a} + \Sigma_{ba}\Sigma_{aa}^{-1}X_a$$

$$\text{Var}[X_b | X_a] = \text{Var}[X_{b.a}] = \Sigma_{bb.a}$$

$$\therefore X_b | X_a \sim N(\mu_{b.a} + \Sigma_{ba}\Sigma_{aa}^{-1}X_a, \Sigma_{bb.a})$$

P6.

已知