

# Gaussian Distribution

---

1. MLE

2. PDF view

3. margin Distri

4. Conditional Distr

---

---

---

---



# Gaussian Distribution MLE

$$\text{Data: } X = (x_1, x_2, \dots, x_N)^T = \begin{pmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_N^T \end{pmatrix}_{N \times p}$$

$x_i \in \mathbb{R}^p \quad x_i \stackrel{\text{iid}}{\sim} N(\mu, \Sigma), \theta = (\mu, \Sigma)$

$$\text{MLE} = \hat{\theta}_{\text{MLE}} = \arg \max_{\theta} P(x|\theta)$$

For  $p=1, \theta = (\mu, \sigma^2)$

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

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

$$\mu_{\text{MLE}} = \arg \max_{\mu} \log P(x|\theta)$$

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

$$\frac{d}{d\mu} \sum_{i=1}^N (x_i-\mu)^2 = \sum_{i=1}^N 2 \cdot (x_i-\mu) \cdot (-1) = 0$$

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

$$\underbrace{\mu_{\text{MLE}} = \frac{1}{N} \sum_{i=1}^N x_i}_{\rightarrow \text{无偏}}$$

$$\begin{aligned} E[\mu_{\text{MLE}}] &= \frac{1}{N} \sum_{i=1}^N E[x_i] \\ &= \frac{1}{N} \cdot N \cdot \mu = \mu \end{aligned}$$

$$\hat{\theta}_{MLE} = \arg \max_{\theta} \log P(X|\theta)$$

$$= \arg \max_{\theta} \sum_{i=1}^N \left[ -\log \theta - \frac{(x_i - \mu)^2}{2\theta^2} \right]$$

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

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

$$\begin{aligned} N\theta^2 &= \sum_{i=1}^N (x_i - \mu)^2 \\ \text{so } \hat{\theta}_{MLE} &= \frac{1}{N} \sum_{i=1}^N (x_i - \mu_{MLE})^2 \end{aligned}$$

有偏估计

$$\text{Because } E[\hat{\theta}_{MLE}] = \frac{N-1}{N} \theta^2$$

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

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

$$\begin{aligned} E\left[\frac{1}{N} \sum_{i=1}^N x_i^2 - \mu^2\right] &= E\left[\frac{1}{N} \sum_{i=1}^N (x_i^2 - \mu^2)\right] = \frac{1}{N} \sum_{i=1}^N E(x_i^2 - \mu^2) \\ &= \frac{1}{N} \sum_{i=1}^N [E(x_i^2) - \mu^2] = \frac{1}{N} \sum_{i=1}^N \underbrace{\text{Var}(x_i)}_{\theta^2} = \theta^2 \end{aligned}$$

$$\begin{aligned}
 & E[(\bar{X}_{MLE}^2 - \bar{u}^2)] \\
 &= E(\bar{X}_{MLE}^2) - E(\bar{u}^2) \\
 &= \bar{E}(\bar{X}_{MLE}^2) - \bar{u}^2 \\
 &= \underbrace{\bar{E}(\bar{X}_{MLE}^2) - E(\bar{X}_{MLE}^2)}_{Var(\bar{X}_{MLE})} = \frac{1}{N^2} \sum_{i=1}^N Var[\sum_{i=1}^N x_i] \\
 &= \frac{1}{N^2} \sum_{i=1}^N \frac{1}{N} \sum_{i=1}^N Var(x_i) = \frac{1}{N^2} \cdot \sum_{i=1}^N 6^2 \\
 &= \frac{1}{N^2} \cdot N \cdot 6^2 = \frac{6^2}{N} \\
 &= \frac{1}{N} 6^2
 \end{aligned}$$

$$So E[\bar{6}_{MLE}^2] = \frac{N-1}{N} 6^2 \neq 6^2$$

$\bar{6}_{MLE}$  是  $6^2$  的一個有偏估計。

$$\hat{6} = \frac{1}{N-1} \sum_{i=1}^N (x_i - \bar{X}_{MLE})^2 \rightarrow \text{无偏估计.}$$

MLE 給高斯分布帶來偏差:



从概率密度角度：(多维高斯)

$$x \sim N(\mu, \Sigma), x \in \mathbb{R}^P, \mu \in \mathbb{R}^P$$

PDF的值是与此相关的

$$\text{此时 } x \text{ 的 PDF: } \frac{1}{(2\pi)^{\frac{P}{2}} |\Sigma|^{\frac{1}{2}}} \exp \left( -\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu) \right)$$

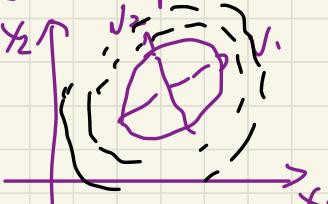
$$x = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_P \end{pmatrix}, \mu = \begin{pmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_P \end{pmatrix}, \Sigma = \begin{pmatrix} \sigma_{11} & \sigma_{12} & \cdots & \sigma_{1P} \\ \sigma_{21} & \sigma_{22} & \cdots & \vdots \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{P1} & \sigma_{P2} & \cdots & \sigma_{PP} \end{pmatrix}$$

Σ: 正定的  
半正定(一般)

$$(x-\mu)^T \Sigma^{-1} (x-\mu) = D$$

D为 x 与 μ 的马氏距离  
当三为工时, D 为欧式距离

$$\Sigma = U \Lambda U^T, U U^T = U^T U = I, \Lambda = \text{diag}(\lambda_1, \dots, \lambda_P)$$

$$\begin{aligned} \therefore \Sigma &= (U_1, U_2, \dots, U_P) \begin{pmatrix} \lambda_1 & & & \\ & \ddots & & \\ & & \lambda_P & \\ & & & U_1^T \\ & & & U_2^T \\ & & & \vdots \\ & & & U_P^T \end{pmatrix} \\ &= (\lambda_1 U_1, \lambda_2 U_2, \dots, \lambda_P U_P) \\ &\approx \sum_{i=1}^P \lambda_i U_i U_i^T \quad \Sigma^{-1} = (U \Lambda U^T)^{-1} = U^{-1} \Lambda^{-1} U^T = \sum_{i=1}^P \frac{1}{\lambda_i} U_i U_i^T \end{aligned}$$

$$\text{故 } \Sigma^{-1} = \sum_{i=1}^P \frac{1}{\lambda_i} U_i U_i^T \quad \therefore (x-\mu)^T \Sigma^{-1} (x-\mu) = (x-\mu)^T \sum_{i=1}^P \frac{1}{\lambda_i} U_i U_i^T (x-\mu)$$

$$\text{原式} = \sum_{i=1}^P \frac{1}{\lambda_i} (x-\mu)^T U_i U_i^T (x-\mu)$$

$$\therefore y = \begin{pmatrix} y_1 \\ \vdots \\ y_P \end{pmatrix}, y_i = (x-\mu)^T U_i \rightarrow \text{意思是 } x \text{ 离中心 } \mu \text{ 的 } U_i \text{ 距离}.$$

$$U \equiv \sum_{i=1}^P y_i \frac{1}{\lambda_i} y_i^T = \sum_{i=1}^P \frac{y_i^2}{\lambda_i} \rightarrow \text{在 } U \text{ 轴上画椭圆}$$

当 x 取不同值时, 圆不同

的值是, 是一个山形的凹面

形成等高线。

形成等高线:

局限性：

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

①  $\sum p \times p \rightsquigarrow \frac{p^2 + p}{2}$

参数个数： $O(p^2)$

高维时，计算太复杂

2  $\rightsquigarrow$  对角阵。 $(\lambda_1 \cdots \lambda_p)$

$\downarrow \lambda_1 = \lambda_2 = \cdots = \lambda_p$

各向同性

factor analysis  $\rightarrow$  2 对角矩阵

P-PCA  $\rightarrow$  2 各向同性

$\hat{\theta} = \arg \max P(x|\theta)$

$= \arg \max \log P(x|\theta)$

$= \arg \max \log [P(x_1|\theta) P(x_2|\theta) \cdots P(x_N|\theta)]$

$= \arg \max \sum_{i=1}^N \log P(x_i|\theta)$

$\hat{\mu} = \arg \max \sum_{i=1}^N \log \left( \frac{1}{N\sqrt{2\pi}} \exp\left(-\frac{(x_i - \mu)^2}{2\sigma^2}\right) \right)$

$\hat{\mu} = \arg \max \sum_{i=1}^N \frac{1}{N} \left[ \log \frac{1}{\sqrt{2\pi}} + \log \frac{1}{\sigma} - \frac{(x_i - \mu)^2}{2\sigma^2} \right]$

$\hat{\mu} = \frac{1}{N} \sum_{i=1}^N x_i$

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

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

② 一阶概率分布不足以描述 data。

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

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

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

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

$$\frac{\sigma^2}{\sigma^{-2}} = N$$

~~只取第一个~~

$$\hat{\mu} = \frac{1}{N} \sum_{i=1}^N x_i$$

无偏估计

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

无偏估计

高維高斯分布求邊緣概率和條件概率

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

$x \in \mathbb{R}^p$ , r.v.

$$x = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_p \end{pmatrix}, \mu = \begin{pmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_p \end{pmatrix}, \Sigma = \begin{pmatrix} \sigma_{11} & \sigma_{12} & \cdots & \sigma_{1p} \\ \sigma_{21} & \ddots & & \\ \vdots & & \ddots & \\ \sigma_{p1} & \sigma_{p2} & \cdots & \sigma_{pp} \end{pmatrix}_{p \times p}$$

已知:

$$x = \begin{pmatrix} x_a \\ x_b \end{pmatrix}_{m+n \times p}, \mu = \begin{pmatrix} \mu_a \\ \mu_b \end{pmatrix}, \Sigma = \begin{pmatrix} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ba} & \Sigma_{bb} \end{pmatrix}_{m+n \times m+n}$$

\*  $P(X_a), P(X_b | X_a)$  西方統計  $\rightarrow PRML$   
 $P(X_a), P(X_a | X_b)$

定理: 已知  $x \sim N(\mu, \Sigma)$ ,  $x \in \mathbb{R}^p$ .

$$Y = AX + B$$

$$(2) Y \sim N(A\mu + B, A\Sigma A^T)$$

$$\text{so: } x_a = \underbrace{(I_m, 0_n)}_A \begin{pmatrix} x_a \\ x_b \end{pmatrix}$$

$$E[x_a] = (I_m, 0_n) \begin{pmatrix} \mu_a \\ \mu_b \end{pmatrix} = \mu_a$$

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

$$\text{So } x_a \sim N(\mu_a, \Sigma_{aa})$$

Now for  $X_{b|X_a}$

we define a new variable  $X_{b,a}$ .

$$X_{b,a} = X_b - \bar{\Sigma}_{ba} \bar{\Sigma}_{aa}^{-1} X_a$$
$$= \underbrace{(-\bar{\Sigma}_{ba} \bar{\Sigma}_{aa}^{-1}, I_n)}_{n \times n} \begin{pmatrix} X_a \\ X_b \end{pmatrix}$$

$$E[X_{b,a}] = A \cdot \begin{pmatrix} u_a \\ u_b \end{pmatrix}$$

$$= (-\bar{\Sigma}_{ba} \bar{\Sigma}_{aa}^{-1}, I_n) \begin{pmatrix} u_a \\ u_b \end{pmatrix}$$
$$= -\bar{\Sigma}_{ba} \bar{\Sigma}_{aa}^{-1} u_a + u_b$$

$$\text{Var}[X_{b,a}] = [-\bar{\Sigma}_{ba} \bar{\Sigma}_{aa}^{-1}, I_n] \begin{pmatrix} \bar{\Sigma}_{aa} & \bar{\Sigma}_{ab} \\ \bar{\Sigma}_{ba} & \bar{\Sigma}_{bb} \end{pmatrix} \begin{pmatrix} -\bar{\Sigma}_{ba} \bar{\Sigma}_{aa}^{-1} \\ I_n \end{pmatrix}$$
$$= (0 \ \bar{\Sigma}_{bb} - \bar{\Sigma}_{ba} \bar{\Sigma}_{aa}^{-1} \bar{\Sigma}_{ab}) \begin{pmatrix} -\bar{\Sigma}_{ba} \bar{\Sigma}_{aa}^{-1} \\ I_n \end{pmatrix} = \underbrace{\bar{\Sigma}_{bb} - \bar{\Sigma}_{ba} \bar{\Sigma}_{aa}^{-1} \bar{\Sigma}_{ab}}_{I_n}$$

$$\text{So } X_{b,a} \sim N(-\bar{\Sigma}_{ba} \bar{\Sigma}_{aa}^{-1} u_a + u_b)$$

Σaa first Schur  
complementary

because

$$X_{b,a} = X_b - \bar{\Sigma}_{ba} \bar{\Sigma}_{aa}^{-1} X_a$$

$$X_b = X_{b,a} + \bar{\Sigma}_{ba} \bar{\Sigma}_{aa}^{-1} X_a$$

$$X_b | X_a = (X_{b,a} + \bar{\Sigma}_{ba} \bar{\Sigma}_{aa}^{-1} X_a) / \bar{\Sigma}_{aa}$$

Here we need prove  $X_a \perp X_{b-a}$

**Theorem:** If  $X \sim N(\mu, \Sigma)$ , then  $Mx \perp Nx \Leftrightarrow M\Sigma N^T = 0$

**Proof.**

$$X \sim N(\mu, \Sigma)$$

$$\text{then } Mx \sim N(M\mu, M\Sigma M^T)$$

$$Nx \sim N(N\mu, N\Sigma N^T)$$

$$\text{so } \text{Cov}(Mx, Nx)$$

$$= E[(Mx - M\mu)(Nx - N\mu)^T]$$

$$= E[M(X-\mu) - (X-\mu)N^T]$$

$$= M \cdot E[(X-\mu)(X-\mu)^T] - MN^T$$

$$= M \cdot \Sigma - N^T$$

Because  $Mx \perp Nx$  and they are both Gaussian.

$$\therefore \text{Cov}(Mx, Nx) = M\Sigma N^T = 0$$

**Note:** ① 独立一定不相关，相关不一定独立。

② 反之两个变量分布，那么“相关”了“独立”。

$$\begin{aligned}
 \text{Then: } X_{b,a} &= X_b - \bar{z}_{ba} \bar{z}_{aa}^{-1} X_a \\
 &= \underbrace{\left( -\bar{z}_{ba} \bar{z}_{aa}^{-1} \quad I_n \right)}_N \underbrace{\left( \begin{array}{c} X_a \\ X_b \end{array} \right)}_X \\
 X_a &= \underbrace{\left( I_n \quad 0 \right)}_N \underbrace{\left( \begin{array}{c} X_b \\ X_s \end{array} \right)}_X \\
 \text{Cov}(X_{ba}, X_b) &= \left( -\bar{z}_{ba} \bar{z}_{aa}^{-1} \quad 2n \right) \begin{pmatrix} \bar{z}_{aa} \bar{z}_{ss} \\ \bar{z}_{ba} \bar{z}_{ss} \end{pmatrix} \begin{pmatrix} I_n \\ 0 \end{pmatrix} \\
 &= 0
 \end{aligned}$$

$$\hookrightarrow X_{s,a} \perp X_b$$

$$\begin{aligned}
 \hookrightarrow X_s | X_a &= X_{s-a} | X_a + \bar{z}_{sa} \bar{z}_{aa}^{-1} X_a | X_a \\
 &= \underbrace{X_{s-a}}_A + \underbrace{\bar{z}_{sa} \bar{z}_{aa}^{-1} X_a}_B
 \end{aligned}$$

$$\hookrightarrow E[X_b | X_a] = \bar{u}_{b,a} + \bar{z}_{ba} \bar{z}_{aa}^{-1} X_a.$$

$$\text{Var}(X_b | X_a) = \bar{s}_{bb,a}$$

$$\hookrightarrow X_s | X_a \sim N(\bar{u}_{s,a} + \bar{z}_{sa} \bar{z}_{aa}^{-1} X_a, \bar{s}_{ss,a})$$

已知:  $P(x) = N(x | \mu, \Lambda^{-1})$   
 $P(y|x) = N(y | Ax + b, L^{-1})$

precision matrix  
 $\Rightarrow L^{-1}$   
 covariance matrix

求  $P(y) \cdot P(x|y)$

(1)      (2)

解: Assume: 
$$\begin{cases} y = Ax + b + \xi, \quad x, y, \xi \sim r.v., \quad \xi \perp x \\ \xi \sim N(0, L^{-1}) \end{cases}$$

Then (1)  $E[y] = E[Ax + b + \xi]$

$$= E[Ax] + b$$

$$= A\mu + b$$

$\text{Var}[y] = \text{Var}[Ax + b + \xi]$

$$= \text{Var}[Ax + \xi]$$

$$= A\Lambda^{-1}A^T + L^{-1}$$

$\checkmark$  from the theorem above.

so  $y \sim N(A\mu + b, L^{-1} + A\Lambda^{-1}A^T)$

(2) Let  $z = \begin{pmatrix} x \\ y \end{pmatrix} \sim N\left(\begin{bmatrix} \mu \\ A\mu + b \end{bmatrix}, \begin{bmatrix} \Lambda^{-1} & ? \\ ? & L^{-1} + A\Lambda^{-1}A^T \end{bmatrix}\right)$

$\Delta = \text{cov}(x, y)$

$$= E[(x - \bar{x})(y - \bar{y})^T]$$

$$= E[(x - \mu)(y - A\mu - b)^T]$$

$$= E[(x - \mu)(Ax + b + \xi - A\mu - b)^T]$$

$$\begin{aligned}
&= \bar{E}[(x-u)(Ax-Au + \varepsilon)^T] \\
&= \bar{E}[(x-u)(Ax-Au)^T] + \bar{E}[(x-u)\varepsilon^T] \\
&= \bar{E}[(x-u)(Ax-Au)^T] + \underbrace{\bar{E}[x-u]\varepsilon^T}_0 \\
&= \bar{E}[(x-u)(Ax-Aw)^T] \\
&= \bar{E}[(x-u)(x-w)^TA^T] \\
&= \underbrace{\bar{E}[(x-u)(x-w)^T]}_{\text{var}(x)} A^T \\
&= A^{-1}A^T
\end{aligned}$$

Then based on

$$P(x_a|x_b) = N(\mu_{a|b} + \bar{\Sigma}_{ab}\bar{\Sigma}_{bb}^{-1}x_b, \bar{\Sigma}_{a|b})$$

$$\begin{aligned}
\bar{E}[x|y] &= \mu + A^{-1}A^T(L^{-1} + A\Lambda^{-1}A^T)^{-1}(y - A\mu - b) \\
\text{var}[x|y] &= \Lambda^{-1} - \Lambda^{-1}A^T(L^{-1} + A\Lambda^{-1}A^T)^{-1}A\Lambda^{-1}
\end{aligned}$$