

Section 2 Math Basics

高斯分布

在机器学习中很重要.

例如 Linear Gaussian Model

$$Y = \underbrace{AX + B}_{\text{线性变换}} + \varepsilon$$

△ 从高斯分布中采样得到

多维高斯分布的定义:

① 若 p 维随机变量 $X = (X_1, X_2, \dots, X_p)^T$ 的密度函数为

$$f(x_1, x_2, \dots, x_p) = \frac{1}{(2\pi)^{\frac{p}{2}} |\Sigma|^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right\}$$

其中 $x = (x_1, x_2, \dots, x_p)^T$. μ 是 p 维向量, Σ 是 p 阶半正定阵.

则 $X \sim N(\mu, \Sigma)$.

但是当 $|\Sigma| = 0$ 时, Σ^{-1} 不存在, 上述定义失效, 则有定义②

② 设 X_1, X_2, \dots, X_p 为 p 个独立的标准正态随机变量.

$$\text{令 } Y = \begin{pmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_m \end{pmatrix} = A_{m \times p} \begin{pmatrix} X_1 \\ X_2 \\ \vdots \\ X_p \end{pmatrix} + \mu_{m \times 1}$$

则称 Y 为 m 维正态随机向量, 记 $Y \sim N_m(\mu, \Sigma)$

其中 $E(Y) = \mu$, $\text{Var}(Y) = AA^T$, 显然, $\Sigma = AA^T$ 分解一般不唯一

这个定义的优点在于, 表明了多元正态分布是由多个独立的标准正态分布线性组合得到的, 所以多元正态分布的一些性质可以由一元正态分布得到.

多元高斯分布的性质:

① 设 $X = (X_1, X_2, \dots, X_p)^T \sim N_p(\mu, \Sigma)$. 若 Σ 为对角阵,

则 X_1, X_2, \dots, X_p 相互独立

② 设 $X \sim N_p(\mu, \Sigma)$. $A \in \mathbb{R}^{s \times p}$, $d \in \mathbb{R}^{s \times 1}$, 则

$$AX + d \sim N_s(A\mu + d, A\Sigma A^T)$$

即正态随机变量的线性变换仍然是正态的.

③ 设 $X \sim N_p(\mu, \Sigma)$. 将 X, μ, Σ 作如下划分:

$$X = \begin{pmatrix} X^{(1)} \\ X^{(2)} \end{pmatrix}_q^{p-q}, \quad \mu = \begin{pmatrix} \mu^{(1)} \\ \mu^{(2)} \end{pmatrix}_q^{p-q}, \quad \Sigma = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix}_q^{p-q}$$

$$\text{则 } X^{(1)} \sim N_q(\mu^{(1)}, \Sigma_{11}), \quad X^{(2)} \sim N_{p-q}(\mu^{(2)}, \Sigma_{22})$$

$$\Sigma_{12} = \text{Cov}(X^{(1)}, X^{(2)})$$

注意:

(1) 多元正态分布的任意元边缘分布仍为正态分布, 反之则不真.

(2) $\Sigma_{12} = \text{Cov}(X^{(1)}, X^{(2)})$, $\Sigma_{12} = 0$ 表示 $X^{(1)}, X^{(2)}$ 不相关, 由定义可知:

此时 $X^{(1)}$ 与 $X^{(2)}$ 独立, 故对于多元正态分布来说不相关与独立是等价的.

参数估计:

$$\text{Data: } X = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_N \end{pmatrix}_{N \times p}.$$

其中 N 表示样本数, p 表示样本维度

$$x_i \in \mathbb{R}^p$$

$$x_i \stackrel{\text{i.i.d.}}{\sim} N_p(\mu, \Sigma)$$

$$\theta = (\mu, \Sigma).$$

极大似然估计 (MLE):

$$\text{假设 } p = 1 \text{ (一维)}. \theta = (\mu, \sigma^2), p(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right).$$

$$\theta = \arg \max_{\theta} P(X|\theta).$$

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

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

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

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

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

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

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

$$\Rightarrow \mu_{MLE} = \frac{1}{N} \sum_{i=1}^N x_i$$

→ 无偏估计.

$$\begin{aligned} \hat{\sigma}_{MLE}^2 &= \underset{\sigma^2}{\operatorname{argmax}} \sum_{i=1}^N \left(-\frac{1}{2} \log \sigma^2 - \frac{(x_i - \mu)^2}{2\sigma^2} \right) \\ &= \underset{\sigma^2}{\operatorname{argmin}} \sum_{i=1}^N \left(\frac{1}{2} \log \sigma^2 + \frac{(x_i - \mu)^2}{2\sigma^2} \right) \end{aligned}$$

$$\frac{\partial}{\partial \sigma^2} \sum_{i=1}^N \left(\frac{1}{2} \log \sigma^2 + \frac{(x_i - \mu)^2}{2\sigma^2} \right) = \sum_{i=1}^N \left(\frac{1}{2} \cdot \frac{1}{\sigma^2} - \frac{(x_i - \mu)^2}{2\sigma^4} \right) = 0.$$

$$\sum_{i=1}^N \left(\frac{1}{\sigma^2} - \frac{(x_i - \mu)^2}{\sigma^4} \right) = 0$$

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

$$\Rightarrow \hat{\sigma}_{MLE}^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2, \text{ 其中 } \mu = \mu_{MLE}$$



有偏估计. $E[\hat{\sigma}_{MLE}^2] = \frac{1}{N} \sum_{i=1}^N E[(x_i - \mu)^2]$

$$= \frac{1}{N} \sum_{i=1}^N D[x_i - \mu] + (E[x_i - \mu])^2$$

$$D[x_i - \mu] = D[x_i - \frac{1}{N} \sum x_j]$$

$$= D\left[\frac{N-1}{N} x_i - \frac{1}{N} \sum x_j\right]$$

$$= \frac{(N-1)^2}{N^2} \sigma^2 + \frac{N-1}{N^2} \sigma^2$$

$$= \frac{N-1}{N} \sigma^2$$

$$DX = EX^2 - (EX)^2$$

↓
 由于 μ 是 \bar{X} ,
 也就是 μ_{MLE} , 而不是 μ
 造成了有偏

$$E[x_i - \mu] = EX_i - E\mu = 0.$$

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

显然, 将 $\hat{\sigma}_{MLE}^2$ 乘上 $\frac{N}{N-1}$, 就能无偏

$$\text{即, } \hat{\sigma}^2 = \frac{1}{N-1} \sum_{i=1}^N (x_i - \mu_{MLE})^2 \text{ 是无偏估计}$$

概率密度函数

如前面定义所示:

$$f(x_1, x_2, \dots, x_p) = \frac{1}{(2\pi)^{\frac{p}{2}} |\Sigma|^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu) \right\}$$

= 二次型.

马氏距离: $D_M(x, y) = \sqrt{(x-y)^T \Sigma^{-1} (x-y)}$

单个数据点的马氏距离:

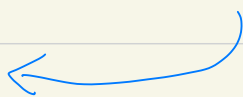
$$D_M(x) = \sqrt{(x-\mu)^T \Sigma^{-1} (x-\mu)}$$

上面两式中, Σ 为协方差阵, μ 为均值

当 $\Sigma = I$ 时, 马氏距离和欧氏距离等价
具体内容参见以下链接

<https://zhuanlan.zhihu.com/p/46626607>

x 和 μ 之间的马氏距离.



基于马氏距离: 可以进行以下推导.

$$\Sigma = U \Lambda U^T, \text{ 其中 } U U^T = U^T U = I, \Lambda = \text{diag}(\lambda_i)_{i=1, \dots, p}.$$

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

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

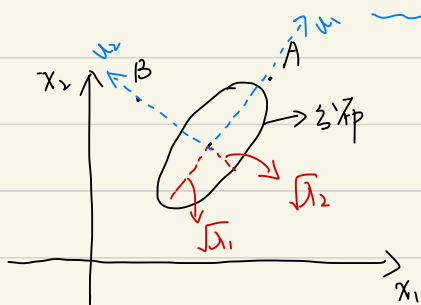
$$\text{所以, } \Sigma^{-1} = (U \Lambda U^T)^{-1} = (U^T)^{-1} \Lambda^{-1} U^{-1} = U \Lambda^{-1} U^T = \sum_{i=1}^p u_i \frac{1}{\lambda_i} u_i^T$$

$$\begin{aligned} \text{则 } (x-\mu)^T \Sigma^{-1} (x-\mu) &= (x-\mu)^T \left(\sum_{i=1}^p u_i \frac{1}{\lambda_i} u_i^T \right) (x-\mu) \\ &= \sum_{i=1}^p (x-\mu)^T u_i \frac{1}{\lambda_i} u_i^T (x-\mu) \end{aligned}$$

$$\text{令 } y_i = (x-\mu)^T u_i, \quad Y = (x-\mu)^T U$$

$$\text{则 } (x-\mu)^T \Sigma^{-1} (x-\mu) = \sum_{i=1}^p \frac{y_i^2}{\lambda_i}$$

$$(x-\mu)^T \Sigma^{-1} (x-\mu) = (x-\mu)^T U \Lambda^{-1} U^T (x-\mu) = Y^T \Lambda^{-1} Y$$



将 $(x-\mu)$ 映射到 u_1, u_2, \dots, u_p

如图，如果要判断 A 和 B 属于分布的概率，显然，A 比 B 更符合该分布，但是如果只考虑 $x-\mu$ ，那么 A 和 B 的概率一样。

马氏距离就是将 x_1, x_2 变换到 u_1, u_2 。

然后再沿 u_1, u_2 按特征值缩放。

最后再去评价距离。

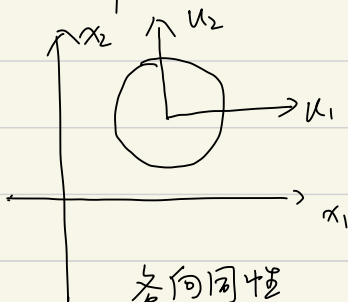
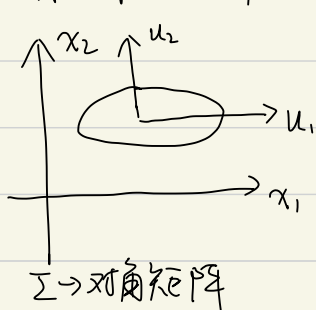
这个 λ_1 和 λ_2 与 σ_1^2 和 σ_2^2 有关系

所以，多维高斯分布的概率密度相当于将协方差考虑在内，使得各维独立，再根据特征值缩放，来衡量概率密度大小。

局限性:

① 参数过多.

Σ 如果是 p 维, Σ 参数个数: $\frac{p^2-p}{2} + p = \frac{p(p+1)}{2}$, 是 $O(p^2)$ 级别的.
- 一般来说, 即使设 Σ 为对角矩阵, 甚至是各向同性的.



例子: factor analysis $\rightarrow \Sigma$ 对角矩阵

p -PCA \rightarrow 各向同性.

② 一般来说, 一个高斯分布无法表示模型

例子: GMM \rightarrow 多个高斯分布

边缘分布与条件概率分布

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

求: $p(x_a)$, $p(x_b|x_a)$

$p(x_b)$, $p(x_a|x_b)$

PRML 中用配方法求解.

性质②证明: $Y = AX + b$.

$$E[Y] = E[AX + b] = AE[X] + b = A\mu + b$$

$$\text{Var}[Y] = \text{Var}[AX + b] = \text{Var}[AX] = A\text{Var}[X]A^T = A\Sigma A^T$$

性质③证明:

由性质④. 得: $x_a = (I_m, 0_n) \begin{pmatrix} x_a \\ x_b \end{pmatrix} + 0$

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

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

则 $x_a \sim N(\mu_a, \Sigma_{aa})$

x_a 的边缘分布求解完成.

求条件概率密度

$$\text{设 } x_{b|a} = x_b - \Sigma_{ba} \Sigma_{aa}^{-1} x_a$$

则由性质② 得:

$$x_{b|a} = (-\Sigma_{ba} \Sigma_{aa}^{-1}, I_n) \begin{pmatrix} x_a \\ x_b \end{pmatrix}$$

$$\mu_{b|a} = E[x_{b|a}] = \mu_b - \Sigma_{ba} \Sigma_{aa}^{-1} \mu_a$$

$$\Sigma_{bb|a} = \text{Var}[x_{b|a}] = (-\Sigma_{ba} \Sigma_{aa}^{-1}, I_n) \begin{pmatrix} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ba} & \Sigma_{bb} \end{pmatrix} \begin{pmatrix} -(\Sigma_{aa}^{-1})^T \Sigma_{ba}^T \\ I_n \end{pmatrix}$$

$$= (0, \Sigma_{bb} - \Sigma_{ba} \Sigma_{aa}^{-1} \Sigma_{ab}) \begin{pmatrix} -(\Sigma_{aa}^{-1})^T \Sigma_{ba}^T \\ I_n \end{pmatrix}$$

$$= \Sigma_{bb} - \Sigma_{ba} \Sigma_{aa}^{-1} \Sigma_{ab} \quad \Sigma_{aa} \text{ is schur complement}$$

有关 schur complement 在相关内参见以下链接.

<https://blog.csdn.net/sheagu/article/details/115771184>

$$\text{则 } x_{b|a} \sim \mathcal{N}(\mu_{b|a}, \Sigma_{bb|a})$$

$$\therefore x_b = x_{b|a} + \Sigma_{ba} \Sigma_{aa}^{-1} x_a \quad \text{这里还是套用 } Y = AX + b.$$

由性质① 得:

$$E[x_b | x_a] = \mu_{b|a} + \Sigma_{ba} \Sigma_{aa}^{-1} x_a$$

$$\text{Var}[x_b | x_a] = \text{Var}[x_{b|a}] = \Sigma_{bb|a}.$$

$$\text{则 } x_b | x_a \sim \mathcal{N}(\mu_{b|a} + \Sigma_{ba} \Sigma_{aa}^{-1} x_a, \Sigma_{bb|a})$$

边缘概率密度求的完毕. x_b 和 $x_a | x_b$ 的求解式也一样

上述过程中，有一步缺少证明：

$$\text{即 } x_b = x_{b \cdot a} + \Sigma_{ba} \Sigma_{aa}^{-1} x_a$$

$$\text{应该写为 } x_b | x_a = x_{b \cdot a} | x_a + \Sigma_{ba} \Sigma_{aa}^{-1} x_a | x_a$$

需要证明 $x_{b \cdot a}$ 与 x_a 独立。

定理：若 $x \sim N(\mu, \Sigma)$ ，则 $MX \perp NX \Leftrightarrow M\Sigma N^T = 0$ 。

$$\begin{aligned} \text{证明： } \text{Cov}(MX, NX) &= E[(MX - E[MX])(NX - E[NX])^T] \\ &= E[(MX - M\mu)(NX - N\mu)^T] \\ &= E[M(x - \mu)(x - \mu)^T N^T] \\ &= M E[(x - \mu)(x - \mu)^T] N^T \\ &= M \Sigma N^T \end{aligned}$$

$$\text{由于 } MX \perp NX \Leftrightarrow \text{Cov}(MX, NX) = 0$$

$$\text{则 } MX \perp NX \Leftrightarrow M\Sigma N^T = 0.$$

$$\text{代入上式： } x_{b \cdot a} = \underbrace{(-\Sigma_{ba} \Sigma_{aa}^{-1})}_M \underbrace{\begin{pmatrix} x_a \\ x_b \end{pmatrix}}_x$$

$$x_a = \underbrace{(\mathbb{I}_m, 0)}_N \underbrace{\begin{pmatrix} x_a \\ x_b \end{pmatrix}}_x.$$

$$\begin{aligned} \Rightarrow M\Sigma N^T &= (-\Sigma_{ba} \Sigma_{aa}^{-1}, \mathbb{I}_n) \begin{pmatrix} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ba} & \Sigma_{bb} \end{pmatrix} \begin{pmatrix} \mathbb{I}_m \\ 0 \end{pmatrix} \\ &= (0, \Sigma_{bb} - \Sigma_{ba} \Sigma_{aa}^{-1} \Sigma_{ab}) \begin{pmatrix} \mathbb{I}_m \\ 0 \end{pmatrix} = 0. \end{aligned}$$

2. $x_{b,a}$ 与 x_a 独立, 即 $x_{b,a} | x_a = x_{b,a}$.

对于 $\Sigma_{ba} \Sigma_{aa}^{-1} x_a | x_a$, 理解如下:

求期望时, $E[\Sigma_{ba} \Sigma_{aa}^{-1} x_a | x_a]$

$$= \Sigma_{ba} \Sigma_{aa}^{-1} E[x_a | x_a]$$

$$= \Sigma_{ba} \Sigma_{aa}^{-1} \sum p(x_a | x_a) \cdot x_a$$

$$= \Sigma_{ba} \Sigma_{aa}^{-1} x_a.$$

求方差时, $\text{Var}[\Sigma_{ba} \Sigma_{aa}^{-1} x_a | x_a]$

$$= \Sigma_{ba} \Sigma_{aa}^{-1} (\Sigma_{ba} \Sigma_{aa}^{-1})^T \text{Var}[x_a | x_a]$$

$$= \Sigma_{ba} \Sigma_{aa}^{-1} (\Sigma_{ba} \Sigma_{aa}^{-1})^T E[(x_a - x_a | x_a)(x_a - x_a | x_a)^T]$$

$$= 0$$

上述为个人理解, 可能有误.

已知 $x, y|x$, 求 $y, x|y$.

已知 $p(x) = N(x|\mu, \Lambda^{-1})$ \rightarrow precision matrix \rightarrow (covariance matrix) $^{-1}$
精度矩阵就是协方差矩阵的逆

$$p(y|x) = N(y|Ax+b, L^{-1})$$

求 $p(y), p(x|y)$

求解 $p(y)$

由已知: $y = Ax + b + \varepsilon$, 其中 $\varepsilon \sim N(0, L^{-1})$ 且 $\varepsilon \perp A$ ε 和 A 独立
 y, x, ε 都是 random variable, A 和 b 看作系数

则|由性质②:

$$E[y] = E[Ax + b + \varepsilon] = A\mu + b.$$

$$\text{Var}[y] = \text{Var}[Ax + b + \varepsilon] = A\Lambda^{-1}A^T + L^{-1}$$

则| $y \sim N(A\mu + b, A\Lambda^{-1}A^T + L^{-1})$, $p(y)$ 求解完毕.

求解 $p(x|y)$

$$\text{构造 } z = \begin{pmatrix} x \\ y \end{pmatrix}, \text{ 则 } \mu_z = \begin{pmatrix} \mu \\ A\mu + b \end{pmatrix}, \Sigma_z = \begin{pmatrix} \Lambda^{-1} & \Delta \\ \Delta^T & L^{-1} + A\Lambda^{-1}A^T \end{pmatrix}$$

$$\begin{aligned} \Delta &= \text{Cov}(x, y) = E[(x - E[x])(y - E[y])^T] \\ &= E[(x - \mu)(Ax + b + \varepsilon - A\mu - b)^T] \\ &= E[(x - \mu)(Ax - A\mu + \varepsilon)^T] \\ &= E[(x - \mu)(Ax - A\mu)^T + (x - \mu)\varepsilon^T] \\ &= E[(x - \mu)(Ax - A\mu)^T] + E[(x - \mu)\varepsilon^T] \end{aligned}$$

其中, 由于 $x \perp \varepsilon$. 则 $E[(x-\mu)\varepsilon^T] = E[\underbrace{(x-\mu)}_0]E[\varepsilon^T] = 0$

$$\begin{aligned}\Delta &= E[(x-\mu)(Ax-A\mu)^T] \\&= E[(x-\mu)(x-\mu)^T A^T] \\&= E[(x-\mu)(x-\mu)^T] A^T \\&= \text{Var}[x] A^T \\&= \Lambda^{-1} A^T\end{aligned}$$

$$\therefore z \sim N\left(\begin{bmatrix} \mu \\ A\mu+b \end{bmatrix}, \begin{bmatrix} \Lambda^{-1} & \Lambda^{-1}A^T \\ A\Lambda^{-1} & L^{-1}+A\Lambda^{-1}A^T \end{bmatrix}\right)$$

沿用上一节求解条件概率率之法, 即可求得 $p(x|y)$.
构造 $x_{b,a}$, 求条件概率