

IID.

$$(x_1, x_2, \dots, x_n | \theta) = f_0(x_1 | \theta) f_0(x_2 | \theta) \dots f_0(x_n | \theta)$$

多元高斯密度函数:

均值: 0, 方差: 1 的一元正态分布.

$$f(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}}$$

均值 μ , 方差 σ^2 时. 标准化: $z = \frac{x - \mu}{\sigma}$. 标准化后方差为 1. 标准化的意义在于将数据点 x 到均值 μ 的距离转化为数据点 x 到均值的距离等于多少个总体的标准差 σ , 消除了数据分布差异和量纲对概率计算的影响:

$$u) \quad f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{z^2}{2}} = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

高斯分布概率密度计算核心在于计算数据点到中心的距离. 并且除以标准差将这个绝对距离转化为相对距离. 然后通过距离平方的指数衰减计算概率密度.

先从各维度不相关的多元正态分布入手. 数据点, 通过 d 维的列向量描述 $x = [x_1, x_2, \dots, x_d]^T$. 各维度均值方差分别为 $\mu_1, \mu_2, \dots, \mu_d, \sigma_1, \sigma_2, \dots, \sigma_d$ 描述:

$$f(x) = \frac{1}{(\sqrt{2\pi})^d \sigma_1 \sigma_2 \dots \sigma_d} e^{-\frac{1}{2} [(\frac{x_1 - \mu_1}{\sigma_1})^2 + (\frac{x_2 - \mu_2}{\sigma_2})^2 + \dots + (\frac{x_d - \mu_d}{\sigma_d})^2]}$$

前面多出的项是为了让概率之和为 1, 该方程也可写为:

$$f(x) = \frac{1}{(\sqrt{2\pi})\sigma_1} e^{-\frac{1}{2}(\frac{x_1 - \mu_1}{\sigma_1})^2} \frac{1}{(\sqrt{2\pi})\sigma_2} e^{-\frac{1}{2}(\frac{x_2 - \mu_2}{\sigma_2})^2} \dots \frac{1}{(\sqrt{2\pi})\sigma_d} e^{-\frac{1}{2}(\frac{x_d - \mu_d}{\sigma_d})^2}$$

各维度之间不相关的多元正态分布概率密度其实就是各个维度正态分布概率密度函数的乘积. 因为各变量之间互不相关, 联合概率密度等于各自概

率密度的乘积

$$d(x, u) = \left(\frac{x_1 - u_1}{\delta_1}\right)^2 + \left(\frac{x_2 - u_2}{\delta_2}\right)^2 + \dots + \left(\frac{x_d - u_d}{\delta_d}\right)^2$$

$$= [x_1 - u_1, x_2 - u_2, \dots, x_d - u_d] \begin{bmatrix} \frac{1}{\delta_1^2} & 0 & \dots & 0 \\ 0 & \frac{1}{\delta_2^2} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \frac{1}{\delta_d^2} \end{bmatrix} \cdot$$

$$\begin{bmatrix} x_1 - u_1 \\ x_2 - u_2 \\ \vdots \\ x_d - u_d \end{bmatrix} = (\bar{x} - \bar{u})^T \Sigma^{-1} (\bar{x} - \bar{u})$$

$$\therefore f(x) = \frac{1}{(\sqrt{2\pi})^d (\delta_1 \delta_2 \dots \delta_d)} \exp\left(-\frac{1}{2} d(x, u)\right)$$

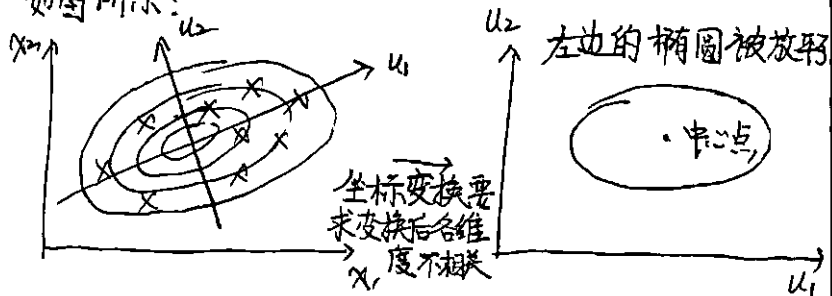
$$= \frac{1}{(2\pi)^{\frac{d}{2}} |\Sigma|^{\frac{1}{2}}} e^{-\frac{1}{2} (\bar{x} - \bar{u})^T \Sigma^{-1} (\bar{x} - \bar{u})}$$

Σ 是协方差矩阵, 里面第 i 行第 j 列元素表示第 i 个变量与第 j 个变量的协方差, 由于假设各个维度不相关 \therefore 协方差只有在对角线位置有值. $\delta_1 \delta_2 \dots \delta_d = |\Sigma|^{\frac{1}{2}}$

维度之间互相关的多元高斯分布

化归思想: 找到倾斜的椭圆分布的长轴方向 u_1 和短轴方向 u_2 . 计算数据点在这两个轴上的坐标 (变换到这两个方向后, 新维度之间显然不相关). 以长轴为 x 轴, 短轴为 y 轴建立新的坐标系, 经过变换后相当于将倾斜的椭圆放平, 此时数据的各个维度之间不相关, 就可以用前面各维度不相关的高斯分布解了。

如图所示:



如果变换的方向 u_1 和 u_2 用列向量 $\bar{u}_1 = \begin{bmatrix} u_1^1 \\ u_1^2 \end{bmatrix}$, $\bar{u}_2 = \begin{bmatrix} u_2^1 \\ u_2^2 \end{bmatrix}$ 那么数据的投影长度可以用点积来计算为:

$\bar{u}_1^T \cdot \bar{X} = [u_1^1, u_1^2] \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$ 投影长度代表了数据点在 u_1 方向上的坐标. 这个过程可以用矩阵变化表示为:

$$\bar{Y} = \begin{bmatrix} \bar{u}_1^T \cdot \bar{X} \\ \bar{u}_2^T \cdot \bar{X} \end{bmatrix} = \begin{bmatrix} \bar{u}_1^T \\ \bar{u}_2^T \end{bmatrix} \bar{X} = U^T \bar{X} \quad \text{其中 } X = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

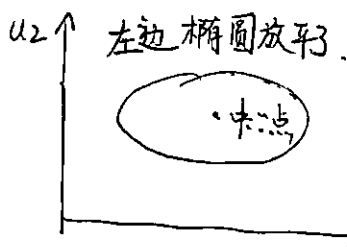
代表源空间的数据坐标. $U = [\bar{u}_1, \bar{u}_2]$ 每一列代表一个投影方向. \bar{u}_1, \bar{u}_2 代表变换方向的单位向量, $\bar{u}_1^T \cdot \bar{X}$ 代表将数据点 \bar{X} 投影到 u_1 方向上的长度.

$\therefore \bar{u}_1, \bar{u}_2$ 都是单位向量, 且相互垂直, $\therefore U$ 是一个正交矩阵 $U^T = U^{-1}$

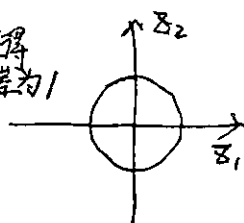
现在数据的各个维度已经去相关, 即可用多元正态分布计算. 首先需将数据标准化(消除量纲)

$$Z = \begin{bmatrix} \frac{Y_1 - \bar{u}_{Y_1}}{\sigma_{Y_1}} \\ \frac{Y_2 - \bar{u}_{Y_2}}{\sigma_{Y_2}} \end{bmatrix} = \begin{bmatrix} \frac{1}{\sigma_{Y_1}}, & 0 \\ 0, & \frac{1}{\sigma_{Y_2}} \end{bmatrix} \begin{bmatrix} Y_1 - \bar{u}_{Y_1} \\ Y_2 - \bar{u}_{Y_2} \end{bmatrix}$$

$$\begin{aligned} &= D(\bar{Y} - \bar{u}_Y) = D(U^T \bar{X} - U^T \bar{u}_X) \\ &= D U^T (\bar{X} - \bar{u}_X) \end{aligned}$$



标准化使得
均值为0, 方差为1



相对距离的平方,

$$d^2(z, u) = z_1^2 + z_2^2 = [z_1 \ z_2] \begin{bmatrix} z_1 \\ z_2 \end{bmatrix} = \bar{z}^T \bar{z}$$

$$d^2(x, u) = \bar{z}^T \bar{z} = (DU^T(\bar{X} - \bar{u}_x))^T (DU^T(\bar{X} - \bar{u}_x)) \\ = (\bar{X} - \bar{u}_x)^T U D^T D U^T (\bar{X} - \bar{u}_x)$$

$$\Sigma_Y^T = D^T D = \begin{bmatrix} \frac{1}{\sigma_{Y_1}^2} & 0 \\ 0 & \frac{1}{\sigma_{Y_2}^2} \end{bmatrix}$$

Σ_Y 是去相关后数据的协方差矩阵, 因为是对角阵, 它的逆等于对角元素取倒数

\therefore 变换后数据各维度不相关, 也就是变换后的协方差矩阵是对角阵, 即:

$$\Sigma_Y = \begin{bmatrix} \sigma_{Y_1}^2 & 0 \\ 0 & \sigma_{Y_2}^2 \end{bmatrix}$$

$$\Sigma_Y = E[(\bar{Y} - \bar{u}_Y)(\bar{Y} - \bar{u}_Y)^T] \\ = E[U^T(\bar{X} - \bar{u}_X)(\bar{X} - \bar{u}_X)^T U] \\ = U^T E[(\bar{X} - \bar{u}_X)(\bar{X} - \bar{u}_X)^T] U \\ = U^T \Sigma_X U$$

即 $\Sigma_X = U \Sigma_Y U^T$, U 是正交矩阵, $U^{-1} = U^T$

$$\therefore U D^T D U^T = U \Sigma_Y^{-1} U^T$$

$$\text{And } (U \Sigma Y^{-1} U^T)(U \Sigma Y U^T) = I$$

$$\therefore (U \Sigma Y^{-1} U^T) \Sigma X = I$$

$$\therefore U D^T D U^T = \Sigma X^{-1}$$

$$\therefore d^2(x, u) = (\bar{x} - \bar{u}_x) \Sigma_x^{-1} (\bar{x} - \bar{u}_x)$$

我们在计算过程中得到最终零均值, 方差为1的Z, 相当于对原坐标X做了一次变换;

$$Z = D U^T (\bar{X} - \bar{u}_x)$$

因此, 概率密度函数在源空间做全空间积分的时候需要做换元变换, 整体减小了 $|DU^T|$,

$$\therefore U D^T D U^T = \Sigma_x^{-1}, \therefore |DU^T| = \sqrt{|\Sigma_x^{-1}|} = |\Sigma_x|^{-\frac{1}{2}}$$

为保证概率密度函数全空间积分为1, 需要乘上 $|\Sigma_x|^{-\frac{1}{2}}$, 还需要除以 $(\sqrt{2\pi})^d$. 这一项是在计算 e^x 的积分时引入的, 每个维度都会有, 所以是 d 次方

$$f(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} = \frac{1}{(2\pi)^{\frac{d}{2}} |\Sigma_x|^{\frac{1}{2}}} e^{-\frac{(\bar{x} - \bar{u}_x)^T \Sigma_x^{-1} (\bar{x} - \bar{u}_x)}{2}}$$

主要思想, 通过线性变换, 将数据的各维度去相关, 再将去相关后的数据标准化. 在推导概率分布的过程中, 可以消去这个变换, 只需求源空间的协方差矩阵就可以了.

假设平面上有一点A, 这个点客观存在, 一旦A指定, 它的概率大小P就已经确定了. 现在我们添加一个坐标系, 添加坐标系使得 $P(A)$ 可以被量化. $P(A) = f(u_1, u_2)$ $P(A) = f(u_1, u_2)$ 使用其他坐标系量化 $P(A) = f(v_1, v_2)$ $P(A) = f(v_1, v_2)$. 不管使用哪个

坐标系, A点, 概率始终不变. $\therefore f(u, u_y) = f(v, v_y)$

Example

$$\Sigma = \begin{bmatrix} 1 & 0.8 \\ 0.8 & 1 \end{bmatrix}$$

$$U = \begin{bmatrix} \cos\theta & -\sin\theta \\ \sin\theta & \cos\theta \end{bmatrix} \quad U^T = \begin{bmatrix} \cos\theta & \sin\theta \\ -\sin\theta & \cos\theta \end{bmatrix}$$

$$(\Sigma)_{\text{new}} = U^T \Sigma U = \begin{bmatrix} \cos\theta & \sin\theta \\ -\sin\theta & \cos\theta \end{bmatrix} \begin{bmatrix} 1 & 0.8 \\ 0.8 & 1 \end{bmatrix} -$$

$$\begin{bmatrix} \cos\theta & -\sin\theta \\ \sin\theta & \cos\theta \end{bmatrix}$$

$$= \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{bmatrix}$$

$$\Rightarrow \theta = \frac{\pi}{4}$$

$$\therefore (\Sigma)_{\text{new}} = \begin{bmatrix} 1.8 & 0 \\ 0 & 0.2 \end{bmatrix}$$

新的坐标是原坐标系经过 $\theta = \frac{\pi}{4}$ 旋转而来。
在新坐标系下, 输入元素变得不相关, x_1 方向方差为 1.8, 分布比较宽, x_2 方向的方差为 0.2, 分布比较窄。