

# Primer on Probability Theory

## 一、Probability Density Functions 概率密度函数

### 1.1 Definitions

#### 1.1.1 PDF Definition

我们缩写概率密度函数为PDF

我们设一个随机变量 $x$ ，令 $p(x)$ 为 $x$ 在 $[a,b]$ 上的PDF。

那么，我们用 $Pr$ 代表概率，则：

$$Pr(c \leq x \leq d) = \int_c^d p(x)dx$$

#### 1.1.2 条件概率密度函数

我们还可以类似得到条件概率，令 $p(x|y)$ 是 $x$ 的概率密度函数( $x \in [a,b]$ )，而条件 $y \in [r,s]$ ，我们有：

$$\text{任取 } y, \quad \int_a^b p(x|y)dx = 1$$

(我们要明确一个思想，条件概率情况下，我们对随机变量进行积分，那么在随机变量的所有取值范围之内积分，其值必为1，不需要管条件是什么)

#### 1.1.3 联合概率密度函数

我们还可以类似得到联合概率密度函数(joint probability densities)， $N$ 维情况下：我们有

$$\int_a^b p(x)dx = \int_{a_N}^{b_N} \dots \int_{a_2}^{b_2} \int_{a_1}^{b_1} p(x_1, x_2, \dots, x_N)dx_1 dx_2 \dots dx_N = 1$$

其中， $a=(a_1,a_2,\dots,a_N), b=(b_1,b_2,\dots,b_N)$

## 1.2 贝叶斯公式

首先，我们易得联合概率密度：

$$p(x, y) = p(x|y)p(y) = p(y|x)p(x)$$

我们重新排列上述公式即可得到贝叶斯公式：

$$p(x|y) = \frac{p(y|x)p(x)}{p(y)}$$

若我们有先验概率密度 $p(x)$ 以及传感器模型 $p(y|x)$ ，我们可以用其来推断给定测量值的状态的后验 $p(x|y)$ 。

我们对上式的分母进行扩展：

$$p(x|y) = \frac{p(y|x)p(x)}{\int p(y|x)p(x)dx}$$

我们进行分母的推导：如下

$$p(\mathbf{y}) = p(\mathbf{y}) \underbrace{\int p(\mathbf{x}|\mathbf{y}) d\mathbf{x}}_1 = \int p(\mathbf{x}|\mathbf{y})p(\mathbf{y}) d\mathbf{x}$$

$$= \int p(\mathbf{x}, \mathbf{y}) d\mathbf{x} = \int p(\mathbf{y}|\mathbf{x})p(\mathbf{x}) d\mathbf{x},$$

## 1.3 Moments 矩

*a.k.a also known as 亦称为*

当处理许多动态分布(a.k.a 密度函数)，我们通常关注较少的特性，称作moments of mass(e.g. mass, center of mass质心, inertia matrix惯性矩阵).

第零个概率矩总是1，一阶矩就是期望 $\mu$ ：

$$\mu = E[x] = \int xp(x)dx$$

$E[\cdot]$ 为期望算子，我们还有：

$$E[F(x)] = \int F(x)p(x)dx$$

但是我们必须这样解释上式：

$$E[F(x)] = [E[f_{ij}(x)]] = [\int f_{ij}p(x)dx]$$

我们如此定义二阶矩：（是我们熟知的协方差矩阵） $\Sigma$ ：

$$\Sigma = E[(x - \mu)(x - \mu)^T]$$

三阶矩和四阶矩分别称作偏度(skewness)和峰度(kurtosis)。

## 1.4 Sample Mean and Covariance

假设我们有一个随机变量 $x$ ，以及它的概率密度函数 $p(x)$ ，我们可以从密度函数进行抽样，记作：

$$x_{meas} = p(x)$$

样本，也可以称作随机变量的一次**实现(realization)**，我们可以直观地将其理解为一个测量值。若我们抽取的样本中有 $N$ 个值，来估计随机变量 $x$ 的**均值(mean)**和**方差(covariance)**，如下：

$$\mu_{meas} = \frac{1}{N} \sum_{i=1}^N x_{i,meas}$$

$$\Sigma_{meas} = \frac{1}{N-1} \sum_{i=1}^N (x_{i,meas} - \mu_{meas})(x_{i,meas} - \mu_{meas})^T$$

注意到，协方差中用 $N-1$ 而不是用 $N$ ，这与**贝塞尔校正(Bessel's correction)**有关。

## 1.5 Statistically Independent, Uncorrelated

我们有两个随机变量 $x, y$ ，我们说它们**统计学独立(statistically independent)**，若二者联合概率密度函数满足：

$$p(x, y) = p(x)p(y)$$

我们称变量**不相关(uncorrelated)**，若有：

$$E[xy^T] = E[x]E[y]^T$$

若随机变量统计学独立，那么它们不相关，反之则不一定成立。但出于简化计算的目的，我们经常会直接认为(假设)不相关的随机变量是统计独立的。

## 1.6 Normalized Product

我们翻译为：归一化积

如果 $p_1(x)$ 和 $p_2(x)$ 是随机变量 $x$ 的两个不同的概率密度函数，那么它们的归一化积 $p(x)$ 定义为：

$$p(x) = \eta p_1(x)p_2(x) \\ \eta = \left( \int p_1(x)p_2(x)dx \right)^{-1}$$

其中 $\eta$ 是一个常值的归一化因子，用于确保 $p(x)$ 满足全概率公理。

在贝叶斯理论中，我们可以用归一化积来融合随机变量 $x$ 的多次独立估计：假设 $y_1, y_2$ 为两次独立测量， $x$ 为待估计随机变量，则有

$$p(x|y_1, y_2) = \eta p(x|y_1)p(x|y_2)$$

我们讲上式左侧利用贝叶斯公式展开：

$$p(x|y_1, y_2) = \frac{p(y_1, y_2|x)p(x)}{p(y_1, y_2)}$$

假定 $y_1, y_2$ 独立，那么有：

$$p(y_1, y_2|x) = p(y_1|x)p(y_2|x) = \frac{p(x|y_1)p(y_1)}{p(x)} \frac{p(x|y_2)p(y_2)}{p(x)}$$

代入上上式我们得到：

$$\eta = \frac{p(y_1)p(y_2)}{p(y_1, y_2)p(x)}$$

若令先验 $p(x)$ 为均匀分布(常数)，那么 $\eta$ 也是一个常量。

## 1.7 Shannon and Mutual Information

我们翻译为：香农和互信息

当我们在估计某一随机变量的概率密度函数时，我们也希望知道对某个量(比如均值)有多么不确定。那么，描述不确定性的方法有：**计算事件的负熵(negative entropy)或香浓信息量(Shannon information)**，记作 $H$ ，如下：

$$H(x) = -E[\ln p(x)] = - \int p(x) \ln p(x) dx$$

另一个重要量则是**互信息(mutual information)**， $I(x, y)$ ：

$$I(x, y) = E[\ln(\frac{p(x, y)}{p(x)p(y)})] = \int \int p(x, y) \ln(\frac{p(x, y)}{p(x)p(y)}) dx dy$$

互信息衡量的是已知一个随机变量的信息之后，另一个随机变量不确定性的减少了多少。当 $x$ 和 $y$ 统计学独立时，我们可以得到：

$$\begin{aligned}
 I(\mathbf{x}, \mathbf{y}) &= \iint p(\mathbf{x}) p(\mathbf{y}) \ln \left( \frac{p(\mathbf{x})p(\mathbf{y})}{p(\mathbf{x})p(\mathbf{y})} \right) d\mathbf{x} d\mathbf{y} \\
 &= \iint p(\mathbf{x}) p(\mathbf{y}) \underbrace{\ln(1)}_0 d\mathbf{x} d\mathbf{y} = 0.
 \end{aligned}$$

而当 $x$ 和 $y$ 不独立时，我们有 $I(x,y) \geq 0$ ，且：

$$I(x, y) = H(x) + H(y) - H(x, y)$$

上式关联了香农信息和互信息。

## 1.8 Cramer-Rao Lower Bound and Fisher Information

我们翻译为：克拉美罗下界和费歇尔信息量

假定我们有一个确定的参数 $\theta$ ，其可以影响随机变量 $x$ 的概率密度，我们用条件概率来描述：

$$p(x|\theta)$$

之后，我们抽样：

$$x_{meas} < -p(x|\theta)$$

这里， $x_{meas}$ 有时可以被称作随机变量 $x$ 的一个实现(realization)，我们可以看作为测量值。

接下来，**克拉美罗下界(Cramer-Rao lower bound(CRLB))**指出：参数真实值 $\theta$ 的任意**无偏估计** $\hat{\theta}$ (基于观测值 $x_{meas}$ )的协方差，可以由**费歇尔信息矩阵(Fisher information matrix)  $I(\mathbf{x}|\theta)$** 来定义边界：

$$\text{cov}(\hat{\theta}|x_{meas}) = E[(\hat{\theta} - \theta)(\hat{\theta} - \theta)^T] \geq I^{-1}(x|\theta)$$

其中，无偏估计指：

$$E[\hat{\theta} - \theta] = 0$$

费歇尔信息矩阵：

$$\mathbf{I}(\mathbf{x}|\theta) = E \left[ \left( \frac{\partial \ln p(\mathbf{x}|\theta)}{\partial \theta} \right)^T \left( \frac{\partial \ln p(\mathbf{x}|\theta)}{\partial \theta} \right) \right].$$

## 2、Gaussian Probability Density Functions

我们先明确一下协方差矩阵：定义：

设  $X = (X_1, X_2, \dots, X_N)^T$  为  $n$  维随机变量，称矩阵

$$C = (c_{ij})_{n \times n} = \begin{pmatrix} c_{11} & c_{12} & \dots & c_{1n} \\ c_{21} & c_{22} & \dots & c_{2n} \\ \vdots & \vdots & & \vdots \\ \vdots & \vdots & & \vdots \\ c_{n1} & c_{n2} & \dots & c_{nn} \end{pmatrix}$$

为  $n$  维随机变量  $X$  的协方差矩阵 (covariance matrix)，也记为  $D(X)$ ，其中

$$c_{ij} = \text{Cov}(X_i, X_j), i, j = 1, 2, \dots, n$$

为  $X$  的分量  $X_i$  和  $X_j$  的协方差 (设它们都存在)。

那么， $x$  为向量表示，那么我们还可以如此表示协方差：

$$\text{Cov}(x) = E[(x - x_m)(x - x_m)^T]$$

我们举个例子：

例如，二维随机变量  $(X_1, X_2)$  的协方差矩阵为

$$C = \begin{pmatrix} c_{11} & c_{12} \\ c_{21} & c_{22} \end{pmatrix}$$

其中  $c_{11} = E[X_1 - E(X_1)]^2$ ,  $c_{12} = E[X_1 - E(X_1)][X_2 - E(X_2)]$

$c_{21} = E[X_2 - E(X_2)][X_1 - E(X_1)]$ ,  $c_{22} = E[X_2 - E(X_2)]^2$

我们写成矩阵形式：

$$\text{Cov}((x_1, x_2)^T) = E \left[ \begin{bmatrix} x_1 - x_{1m} \\ x_2 - x_{2m} \end{bmatrix} \begin{bmatrix} x_1 - x_{1m} & x_2 - x_{2m} \end{bmatrix} \right]$$

结果与展开式相同。

## 2.1 Definitions

高斯概率密度函数，一维：

$$p(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2} \frac{(x-\mu)^2}{\sigma^2}}$$

高维表示： $x \in \mathbb{R}^N$

$$p(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{\sqrt{(2\pi)^N \det \boldsymbol{\Sigma}}} \exp \left( -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right)$$

其中， $\boldsymbol{\mu} \in \mathbb{R}^N$ ，为均值； $\boldsymbol{\Sigma} \in \mathbb{R}^{(N \times N)}$  是协方差矩阵 (正定对称矩阵)。

因此我们有：

$$\boldsymbol{\mu} = E[\mathbf{x}] = \int_{-\infty}^{\infty} \mathbf{x} \frac{1}{\sqrt{(2\pi)^N \det \boldsymbol{\Sigma}}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right) d\mathbf{x}, \quad (2.37)$$

and

$$\begin{aligned} \boldsymbol{\Sigma} &= E[(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T] \\ &= \int_{-\infty}^{\infty} (\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T \frac{1}{\sqrt{(2\pi)^N \det \boldsymbol{\Sigma}}} \\ &\quad \times \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right) d\mathbf{x}. \end{aligned} \quad (2.38)$$

高维情况下，我们写  $\mathbf{x} \sim N(0, \mathbf{I})$ ，表示  $\mathbf{x} \in \mathbb{R}^N$  服从高维高斯分布，其中， $\mathbf{I}$  为  $N \times N$  阶单位矩阵。

## 2.2 Isserlis' Theorem

多维高斯分布的高阶矩比较难计算，但是它们在后续的学习中也很重要。因此我们利用 Isserlis 理论来计算高斯分布的高阶矩。

假设  $\mathbf{x} = (x_1, x_2, \dots, x_{2M}) \in \mathbb{R}^{(2M)}$ ，通常来讲：

$$E[x_1 x_2 x_3 \dots x_{2M}] = \sum \prod E[x_i x_j]$$

这表明：计算  $2M$  个变量乘积的期望，可以首先计算所有两两不同的变量的乘积的期望，然后把计算出来这些期望做乘积。这样的组合有：

$$\frac{2M!}{2^M M!}$$

种，最后将这些乘积的值求和即可。例如， $M=2$  时

$$E[x_i x_j x_k x_l] = E[x_i x_j] E[x_k x_l] + E[x_i x_k] E[x_j x_l] + E[x_i x_l] E[x_j x_k]$$

我们利用这个可以推导出一些有用的结果。

假设我们有  $\mathbf{x} \sim N(0, \boldsymbol{\Sigma}) \in \mathbb{R}^N$ 。我们计算下式：

$$E[\mathbf{x} (\mathbf{x}^T \mathbf{x})^p \mathbf{x}^T]$$

其中  $p$  是一个非负整数。当  $p=0$  时，我们有

$$E[\mathbf{x} \mathbf{x}^T] = \boldsymbol{\Sigma}$$

当  $p=1$  时，我们有

$$E[\mathbf{x} \mathbf{x}_1^T \mathbf{x}_1 \mathbf{x}^T] = E[\mathbf{x} \mathbf{x}^T]$$

## 2.3 联合高斯概率密度函数，分解与推断

我们有联合随机变量服从高斯分布  $(\mathbf{x}, \mathbf{y})$ ，写为，其 PDF 写为：

$$p(\mathbf{x}, \mathbf{y}) = N\left(\begin{bmatrix} \mu_x \\ \mu_y \end{bmatrix}, \begin{bmatrix} \Sigma_{xx} & \Sigma_{xy} \\ \Sigma_{yx} & \Sigma_{yy} \end{bmatrix}\right)$$

我们注意到： $\Sigma_{yx} = \Sigma_{xy}^T$ 。这可以将一个联合概率密度拆分为两个概率密度的乘积(条件概率乘边缘概率)， $p(\mathbf{x}, \mathbf{y}) = p(\mathbf{x}|\mathbf{y})p(\mathbf{y})$ 。特别地，对于高斯分布，我们可以用 **Schur complement** 舒贝尔定理。

$$\begin{bmatrix} \Sigma_{xx} & \Sigma_{xy} \\ \Sigma_{yx} & \Sigma_{yy} \end{bmatrix} = \begin{bmatrix} 1 & \Sigma_{xy} \Sigma_{yy}^{-1} \\ 0 & 1 \end{bmatrix} \begin{bmatrix} \Sigma_{xx} - \Sigma_{xy} \Sigma_{yy}^{-1} \Sigma_{yx} & 0 \\ 0 & \Sigma_{yy} \end{bmatrix} \begin{bmatrix} 1 & 0 \\ \Sigma_{yy}^{-1} \Sigma_{yx} & 1 \end{bmatrix}$$

我们对等式两侧取逆矩阵得到：

$$\begin{bmatrix} \Sigma_{xx} & \Sigma_{xy} \\ \Sigma_{yx} & \Sigma_{yy} \end{bmatrix}^{-1} = \begin{bmatrix} 1 & 0 \\ -\Sigma_{yy}^{-1} \Sigma_{yx} & 1 \end{bmatrix} \times \begin{bmatrix} (\Sigma_{xx} - \Sigma_{xy} \Sigma_{yy}^{-1} \Sigma_{yx})^{-1} & 0 \\ 0 & \Sigma_{yy}^{-1} \end{bmatrix} \begin{bmatrix} 1 & -\Sigma_{xy} \Sigma_{yy}^{-1} \\ 0 & 1 \end{bmatrix}.$$

之后，我们看联合高斯分布种，指数部分(前面可以认为是常数)：

$$\begin{aligned} & \left( \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} - \begin{bmatrix} \mu_x \\ \mu_y \end{bmatrix} \right)^T \begin{bmatrix} \Sigma_{xx} & \Sigma_{xy} \\ \Sigma_{yx} & \Sigma_{yy} \end{bmatrix}^{-1} \left( \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} - \begin{bmatrix} \mu_x \\ \mu_y \end{bmatrix} \right) \\ &= \left( \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} - \begin{bmatrix} \mu_x \\ \mu_y \end{bmatrix} \right)^T \begin{bmatrix} 1 & 0 \\ -\Sigma_{yy}^{-1} \Sigma_{yx} & 1 \end{bmatrix} \begin{bmatrix} (\Sigma_{xx} - \Sigma_{xy} \Sigma_{yy}^{-1} \Sigma_{yx})^{-1} & 0 \\ 0 & \Sigma_{yy}^{-1} \end{bmatrix} \\ & \quad \times \begin{bmatrix} 1 & -\Sigma_{xy} \Sigma_{yy}^{-1} \\ 0 & 1 \end{bmatrix} \left( \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} - \begin{bmatrix} \mu_x \\ \mu_y \end{bmatrix} \right) \\ &= (\mathbf{x} - \mu_x - \Sigma_{xy} \Sigma_{yy}^{-1} (\mathbf{y} - \mu_y))^T (\Sigma_{xx} - \Sigma_{xy} \Sigma_{yy}^{-1} \Sigma_{yx})^{-1} \\ & \quad \times (\mathbf{x} - \mu_x - \Sigma_{xy} \Sigma_{yy}^{-1} (\mathbf{y} - \mu_y)) + (\mathbf{y} - \mu_y)^T \Sigma_{yy}^{-1} (\mathbf{y} - \mu_y), \quad (2.52) \end{aligned}$$

这是两个二次项之和，因此可以进行拆分（指数+ - > ×），因此我们有：

$$\begin{aligned} p(x, y) &= p(x|y)p(y) \\ p(x|y) &= N(\mu_x + \Sigma_{xy} \Sigma_{yy}^{-1} (\mathbf{y} - \mu_y), \Sigma_{xx} - \Sigma_{xy} \Sigma_{yy}^{-1} \Sigma_{yx}) \\ p(y) &= N(\mu_y, \Sigma_{yy}) \end{aligned}$$

我们看到，因子 $p(x|y), p(y)$ 都是高斯概率密度函数，更进一步：如果我们知道 $y$ 的值(比如，其观测值)，那么我们可以解算出 $x$ 的似然值，通过给定 $y$ 值以及利用上式计算 $p(x|y)$ 的值。

这也是高斯估计的重要部分：我们从先验概率 $x \sim N(\mu_x, \Sigma_{xx})$ 入手，基于测量值 $y_{\text{meas}}$ 来缩小(调整、逼近)。上式我们看到，均值改变、方差变小 ( $p(x|y)$ )，在测量值下，出现 $x$ 状态的可能性，通过测量值 $y_{\text{meas}}$ 来修正)。

## 2.4 Statistically Independent , Uncorrelated

我们翻译为：统计学独立、不相关

高斯概率密度函数的情况下，统计学独立和不相关是等价的。如同上一节：我们设 $p(x, y)$ ，若 $x, y$ 独立，那么有：

$$\begin{aligned} p(x, y) &= p(x)p(y) \\ \text{then, } p(x|y) &= p(x) = N(\mu_x, \Sigma_{xx}) \\ \text{then, } \Sigma_{xy} &= 0 \\ \text{then, } \Sigma_{xy} &= E[(x - \mu_x)(y - \mu_y)^T] = E[xy^T] - E[x]E[y]^T \end{aligned}$$

那么，我们就得到了不相关成立的条件：

$$E[xy^T] = E[x]E[y]^T$$

## 2.5 Linear Change of Variables

我们翻译为：高斯分布随机变量的线性变换

假设，我们有服从高斯分布的随机变量 $x$ ：

$$x \in R^N \sim N(\mu_x, \Sigma_{xx})$$

并且，我们还有另一个随机变量 $y \in R^M$ ：

$$y = Gx$$
$$G \in R^{M \times N}$$

其中， $G$ 为常数矩阵。

我们想研究随机变量 $y$ 的统计特性，那么最简单的方法就是计算均值和方差：

$$\mu_y = E[y] = E[Gx] = GE[x] = G\mu_x$$
$$\Sigma_{yy} = E[(y - \mu_y)(y - \mu_y)^T] = GE[(x - \mu_x)(x - \mu_x)^T]G^T = G\Sigma_{xx}G^T$$

因此，我们得到：

$$y \sim N(\mu_y, \Sigma_{yy}) = N(G\mu_x, G\Sigma_{xx}G^T)$$

另外一种方法，我们假设这个映射是**单射**，意思是两个 $x$ 值不可能和同一个 $y$ 值对应；我们通过假定一个更严格的条件来简化单射条件，即 $G$ 是可逆的(因此 $M=N$ )，根据全概率公理：

$$\int_{-\infty}^{\infty} p(x) dx = 1$$

一个小区域内的 $x$ 映射到 $y$ 上，变为：

$$dy = |\det G| dx$$

代入：

$$\begin{aligned} 1 &= \int_{-\infty}^{\infty} p(\mathbf{x}) d\mathbf{x} \\ &= \int_{-\infty}^{\infty} \frac{1}{\sqrt{(2\pi)^N \det \Sigma_{xx}}} \exp \left( -\frac{1}{2} (\mathbf{x} - \mu_x)^T \Sigma_{xx}^{-1} (\mathbf{x} - \mu_x) \right) d\mathbf{x} \\ &= \int_{-\infty}^{\infty} \frac{1}{\sqrt{(2\pi)^N \det \Sigma_{xx}}} \\ &\quad \times \exp \left( -\frac{1}{2} (\mathbf{G}^{-1}\mathbf{y} - \mu_x)^T \Sigma_{xx}^{-1} (\mathbf{G}^{-1}\mathbf{y} - \mu_x) \right) |\det \mathbf{G}|^{-1} d\mathbf{y} \\ &= \int_{-\infty}^{\infty} \frac{1}{\sqrt{(2\pi)^N \det \mathbf{G} \det \Sigma_{xx} \det \mathbf{G}^T}} \\ &\quad \times \exp \left( -\frac{1}{2} (\mathbf{y} - \mathbf{G}\mu_x)^T \mathbf{G}^{-T} \Sigma_{xx}^{-1} \mathbf{G}^{-1} (\mathbf{y} - \mathbf{G}\mu_x) \right) d\mathbf{y} \\ &= \int_{-\infty}^{\infty} \frac{1}{\sqrt{(2\pi)^N \det(\mathbf{G}\Sigma_{xx}\mathbf{G}^T)}} \\ &\quad \times \exp \left( -\frac{1}{2} (\mathbf{y} - \mathbf{G}\mu_x)^T (\mathbf{G}\Sigma_{xx}\mathbf{G}^T)^{-1} (\mathbf{y} - \mathbf{G}\mu_x) \right) d\mathbf{y}, \end{aligned}$$



值得注意的是：若 $M < N$ ，那么线性映射就不是单射了，我们无法通过定积分变量代换的方法，求得 $y$ 的分布。

但是，若 $\text{rank}(G)=M$ ，我们还是可以求。

## 2.6 Normalized Product of Gaussians

翻译为：高斯分布的归一化积

高斯概率密度函数一个很有用的性质： $K$ 个高斯概率密度函数的归一化积还是高斯概率密度函数。

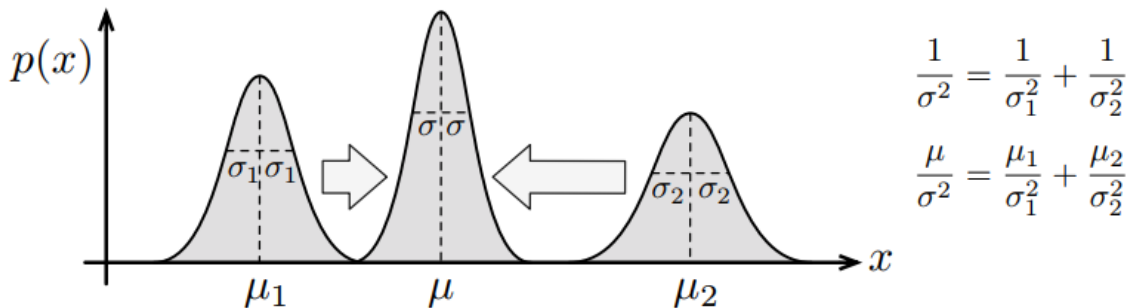
$$\begin{aligned} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right) \\ \equiv \eta \prod_{k=1}^K \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1}(\mathbf{x} - \boldsymbol{\mu}_k)\right), \quad (2.68) \end{aligned}$$

其中：

$$\boldsymbol{\Sigma}^{-1} = \sum_{k=1}^K \boldsymbol{\Sigma}_k^{-1}, \quad (2.69a)$$

$$\boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} = \sum_{k=1}^K \boldsymbol{\Sigma}_k^{-1} \boldsymbol{\mu}_k, \quad (2.69b)$$

$\eta$ 是归一化常数，保证其满足概率公理。高斯概率密度函数在融合过程中会起到作用，如下图所示：



## 2.7 Sherman-Morrison-Woodbury Identity

Identity 译作 等式

Sherman-Morrison-Woodbury 恒等式有时合称为矩阵求逆引理。这个等式是从一个恒等式衍生出来的四个不同的等式。

对于可逆矩阵，我们可以将它分解为一个下三角——对角——上三角（LDU）形式，或上三角——对角——下三角（UDL）形式，如下所示：

$$\begin{aligned}
& \begin{bmatrix} \mathbf{A}^{-1} & -\mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix} \\
&= \begin{bmatrix} \mathbf{1} & \mathbf{0} \\ \mathbf{CA} & \mathbf{1} \end{bmatrix} \begin{bmatrix} \mathbf{A}^{-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{D} + \mathbf{CAB} \end{bmatrix} \begin{bmatrix} \mathbf{1} & -\mathbf{AB} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \quad (\text{LDU}) \\
&= \begin{bmatrix} \mathbf{1} & -\mathbf{BD}^{-1} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \begin{bmatrix} \mathbf{A}^{-1} + \mathbf{BD}^{-1}\mathbf{C} & \mathbf{0} \\ \mathbf{0} & \mathbf{D} \end{bmatrix} \begin{bmatrix} \mathbf{1} & \mathbf{0} \\ \mathbf{D}^{-1}\mathbf{C} & \mathbf{1} \end{bmatrix}. \quad (\text{UDL})
\end{aligned} \tag{2.72}$$

我们对两侧取逆并且取等号，得到如下四个恒等式：

$$(\mathbf{A}^{-1} + \mathbf{BD}^{-1}\mathbf{C})^{-1} \equiv \mathbf{A} - \mathbf{AB}(\mathbf{D} + \mathbf{CAB})^{-1}\mathbf{CA}, \tag{2.75a}$$

$$(\mathbf{D} + \mathbf{CAB})^{-1} \equiv \mathbf{D}^{-1} - \mathbf{D}^{-1}\mathbf{C}(\mathbf{A}^{-1} + \mathbf{BD}^{-1}\mathbf{C})^{-1}\mathbf{BD}^{-1}, \tag{2.75b}$$

$$\mathbf{AB}(\mathbf{D} + \mathbf{CAB})^{-1} \equiv (\mathbf{A}^{-1} + \mathbf{BD}^{-1}\mathbf{C})^{-1}\mathbf{BD}^{-1}, \tag{2.75c}$$

$$(\mathbf{D} + \mathbf{CAB})^{-1}\mathbf{CA} \equiv \mathbf{D}^{-1}\mathbf{C}(\mathbf{A}^{-1} + \mathbf{BD}^{-1}\mathbf{C})^{-1}. \tag{2.75d}$$

这四个恒等式我们可能在后续经常用到。

## 2.8 Passing a Gaussian through a Nonlinearity

译作：高斯分布随机变量的非线性变换

高斯分布经过一个随机非线性变换之后的情况：

$$p(y) = \int_{-\infty}^{\infty} p(y|x)p(x)dx$$

其中：

$$\begin{aligned}
p(y|x) &= N(g(x), R) \\
p(x) &= N(\mu_x, \Sigma_{xx})
\end{aligned}$$

这里：g(·)表示一个非线性映射：g:→y，被N(0,R)高斯分布干扰。我们经常需要这类随机非线性变换来对传感器建模。

### 2.8.1 Scalar Deterministic Case via Change of Variables

译作：标量情况下的非线性映射

一个简单的情况：x是一个标量，g(·)是一个非线性函数，且：

$$x \sim N(0, \sigma^2)$$

对于高斯概率密度函数：

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

特殊例子：

$$y = e^x$$

反函数：

$$x = \ln(y)$$

在很小区间我们有：

$$dy = e^x dx$$

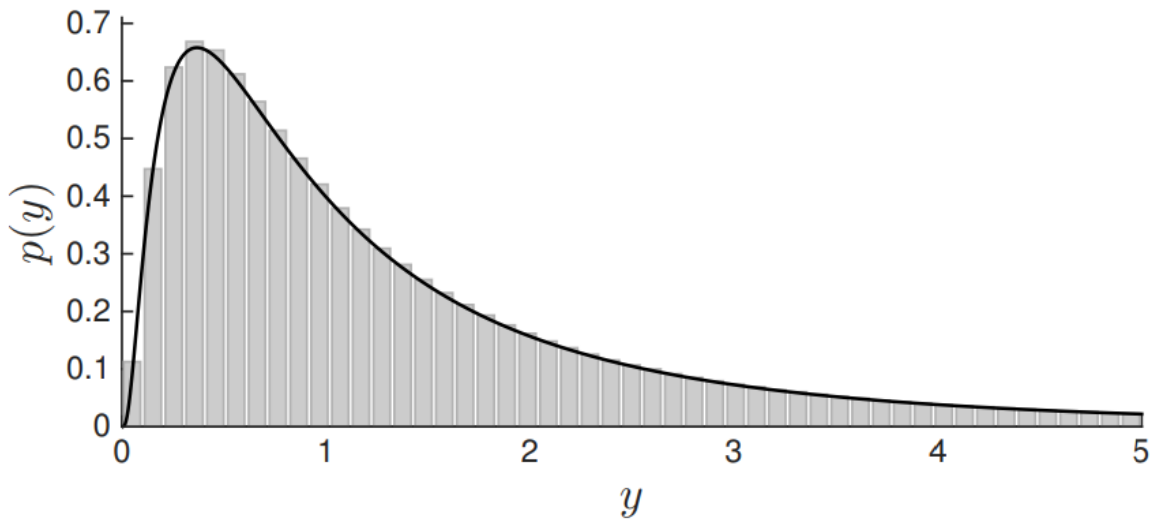
or

$$dx = \frac{1}{y} dy$$

根据上述推导：

$$1 = \int_{-\infty}^{\infty} p(x) dx = \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2} \frac{x^2}{\sigma^2}} dx$$
$$= \int_0^{\infty} \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2} \frac{(\ln(y))^2}{\sigma^2}} \frac{1}{y} dx = \int_0^{\infty} p(y) dy$$

非线性变换后的概率密度函数如下图所示：



## 2.8.2 General Case via Linearization

我们引入**线性化**：

$$g(x) \approx \mu_y + G(x - \mu_x)$$
$$G = \left. \frac{\partial g(x)}{\partial x} \right|_{x=\mu_x}$$
$$\mu_y = g(\mu_x)$$

其中，G为g(·)的雅可比矩阵。

## 2.9 Shannon Information of a Gaussian

译作：高斯分布的香农信息、

在高斯概率密度情况下，我们有如下Shannon information：

$$\begin{aligned}
H(\mathbf{x}) &= - \int_{-\infty}^{\infty} p(\mathbf{x}) \ln p(\mathbf{x}) d\mathbf{x} \\
&= - \int_{-\infty}^{\infty} p(\mathbf{x}) \left( -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) - \ln \sqrt{(2\pi)^N \det \boldsymbol{\Sigma}} \right) d\mathbf{x} \\
&= \frac{1}{2} \ln ((2\pi)^N \det \boldsymbol{\Sigma}) + \int_{-\infty}^{\infty} \frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) p(\mathbf{x}) d\mathbf{x} \\
&= \frac{1}{2} \ln ((2\pi)^N \det \boldsymbol{\Sigma}) + \frac{1}{2} E [(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})], \quad (2.91)
\end{aligned}$$

其中，我们用期望算子来表示第二项。实际上，第二项就是平方**马氏距离(Mahalanobis distance)**的期望值，与欧式距离差一个协方差权重。我们由：

$$x^T A x = \text{tr}(A x x^T)$$

可得：

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

## 2.10 Mutual Information of a Joint Gaussian PDF

译作：联合高斯概率密度函数的互信息

## 2.11 Cramer-Rao Lower Bound Applied to Gaussian PDFs

译作：高斯概率密度函数的克拉美罗下界

# 3、Gaussian Processes

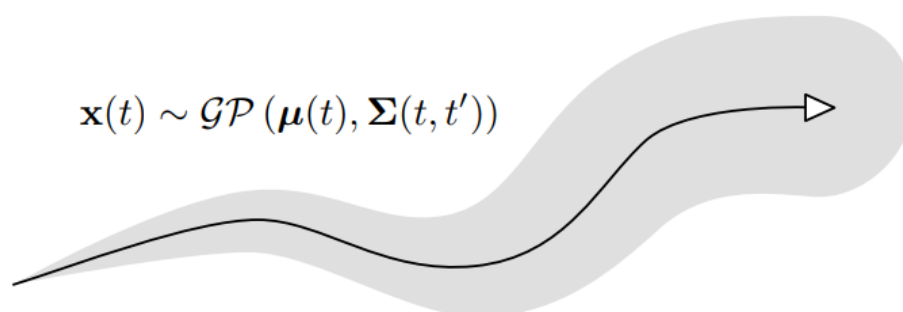
译作：高斯过程

我们将满足高斯分布的变量  $\mathbf{x} \in \mathbb{R}^N$  记为：

$$\mathbf{x} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

我们会大量使用这类随机变量表达离散时间的状态量。接下来我们要着手讨论时间  $t$  上的连续的状态量。

首先，引入**高斯过程(Gaussian processes , GPs)**。下图描述了高斯过程表示的轨迹：



其中，其均值函数为黑色的实线，协方差函数为阴影区域。

我们认为整个轨迹是一个单独的随机变量，其属于一个函数集合。一个函数越接近均值函数，轨迹就越相似。协方差函数通过描述两个时刻 $t, t'$ 的随机变量的相关性来刻画轨迹的平滑程度。我们把这个随机变量函数记为：

$$x(t) \sim GP(\mu(t), \Sigma(t, t'))$$

这表明了连续时间轨迹是一个高斯过程。实际上高斯过程不仅限于表达对于时间是一维的情况，但不需要考虑那么多。

## 4、习题

---

①