

数据科学基础知识 Fundamentals

第一部分:概率论 Probability theory

ISEP第二年
2023-2024

基于 Nathalie Colin 和 Jean-Claude Guillerot 教授的课程

概率论

第四次会议（2023 年 10 月 20 日）：

第 5 章（之二）：A 的变换
实值随机变量

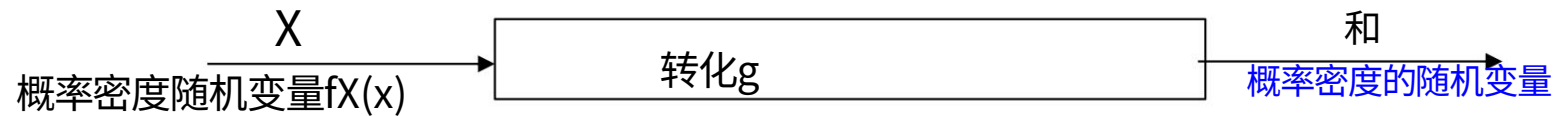
第 6 章：二维随机
变量

第5章 (之二): 实数的变换 随机变量

- 随机变量的确定性函数 · 密度 $f_X(x)$ 的变换 · 示例
 - 特殊情况

问题描述

设 X 为密度 $f_X(x)$ 和 g 确定性实函数的随机变量
变量 X 的。



: 已知 _____

未知

变换: 确定性函数, $Y = g(X)$

变换示例: $Y = X^2$; $Y = aX + B$

可能会发生三种情况:

· X 和 Y 之间一一对应

↔ X

· 非一一对应

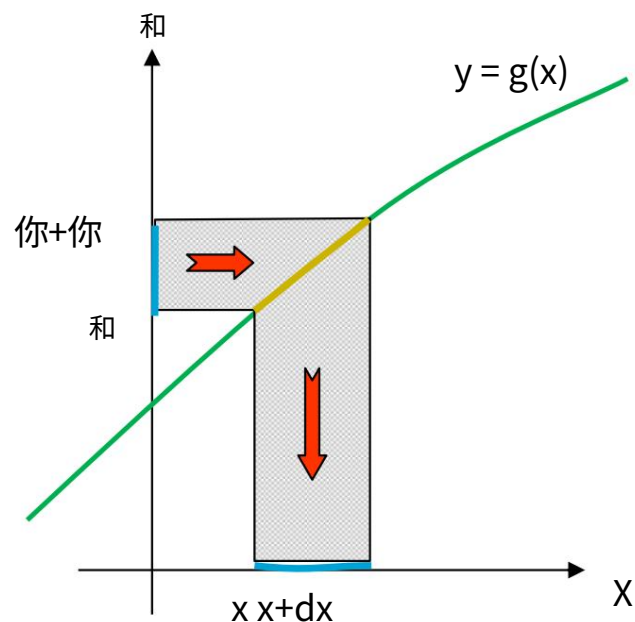
↔ (x_i) 对于 $i = 1, 2, \dots, n$

· 特殊情况

情况1:X和Y一一对应

子情况 1.1:

如果 g 单调递增



事件:

$$\{y \leq Y \leq y + dy\} \Leftrightarrow \{x \leq X \leq x + dx\}$$

创建者:

鉴于以下事实:

y 固定只有一个先行词 x

$dy > 0$ 且 $dx > 0$

$$\begin{aligned} \blacksquare P\{y \leq Y \leq y + dy\} &= P\{x \leq X \leq x + dx\} \\ \blacksquare f_Y(y) dy &= f_X(x) dx \end{aligned}$$

那么, Y 的密度为:

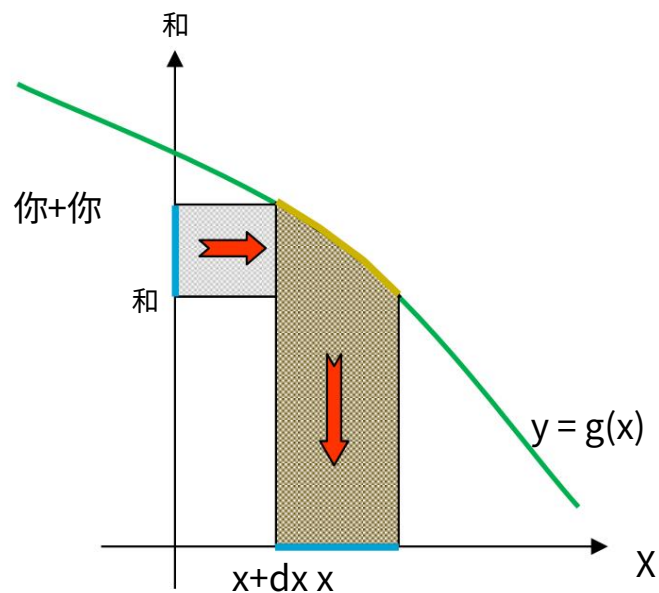
$$f_Y(y) = f_X(x) \frac{dx}{dy} = f_X(g^{-1}(y)) \frac{dx}{dy}$$

其中 g^{-1} 是 g 的反函数

情况1:X和Y一一对应

子情况 1.2:

如果g单调递减



事件:

$$\{y \leq Y \leq y + dy\} \Leftrightarrow \{x + dx \leq X \leq x\}$$

创建者:

鉴于以下事实:

y固定只有一个先行词 x

$dy > 0$ 但 $dx < 0$

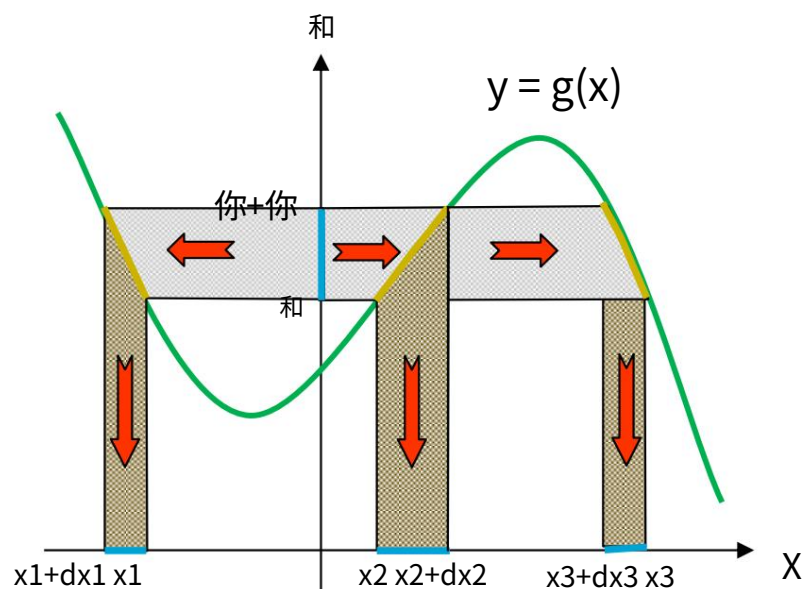
- $P\{y \leq Y \leq y + dy\} = P\{x + dx \leq X \leq x\}$
- $f_Y(y) dy = -f_X(x) dx$

那么,Y的密度为:

$$f_Y(y) = f_X(x) \left| \frac{dx}{dy} \right| = f_X(g^{-1}(y)) \left| \frac{dx}{dy} \right|$$

情况2:X和Y之间的对应关系不是一一对应 (1/2)

X 中可能有多个值 x_i 对应于同一值 y



y 可以有多个先行词

$$y = g(x_i) \quad i = 1, 2, \dots$$

(例如 (x_1, x_2, x_3))

$dy > 0$ 但 $dx_1 < 0$ 、 $dx_2 > 0$ 、 $dx_3 < 0$

事件:

创建者:

$$\{y \leq Y \leq y + dy\} \Leftrightarrow \{x_1 + dx_1 \leq X \leq x_1\} \text{ 或者 } \{x_2 \leq X \leq x_2 + dx_2\} \text{ 或者 } \{x_3 + dx_3 \leq X \leq x_3\}$$

情况2:X和Y之间不是一对一的对应关系 (2/2)

$$\{y \leq Y \leq y + dy\} \Leftrightarrow \{x_1 + dx_1 \leq X \leq x_1\} \text{或者} \{x_2 \leq X \leq x_2 + dx_2\} \text{或者} \{x_3 + dx_3 \leq X \leq x_3\}$$

自事件发生以来

$$\{x_1 + dx_1 \leq X \leq x_1\}, \{x_2 \leq X \leq x_2 + dx_2\}, \{x_3 + dx_3 \leq X \leq x_3\}$$

是不相交的

$$f_Y(y)|dy| = \sum_i f_X(x_i) |dx_i|$$

根据公理3:

$$f_Y(y) = \sum_i f_X(x_i) \left| \frac{dx}{dy} \right|_{x=x_i}$$

其中 x_i 是方程 $x = g$ 的解

$^{-1}(y)$ 当 y 固定时。

示例 (1/2) :

设 X 为密度为 $f_X(x)$ 的随机变量，
变换 $Y = X^2$ 。目的是计算 $f_Y(y)$?

反函数为：

$$x = \pm \sqrt{y} \quad y > 0$$

导数的计算：

$$dy = 2 \cdot x \cdot dx$$

$$\left[g_i^{-1}(y) \right]' = \frac{dx}{dy} = \pm \frac{1}{2\sqrt{y}}$$

密度 $f_Y(y)$ 将为：

$$f_Y(y) = \frac{1}{2\sqrt{y}} \cdot \left[f_X(\sqrt{y}) + f_X(-\sqrt{y}) \right]$$

示例 (2/2) :

应用: $f_X(x)$ 是瑞利密度:

$$f_X(x) = \frac{x}{\sigma^2} \cdot \exp\left(-\frac{x^2}{2\sigma^2}\right) \quad x \geq 0$$

$$f_X(x) = 0 \quad \text{别处}$$

警告: $x \geq 0$ 只有正根属于定义域

$$f_Y(y) = \frac{1}{2\sqrt{y}} \cdot \frac{\sqrt{y}}{\sigma^2} \cdot \exp\left(-\frac{(\sqrt{y})^2}{2\sigma^2}\right) = \frac{1}{2\sigma^2} \cdot \exp\left(-\frac{y}{2\sigma^2}\right)$$

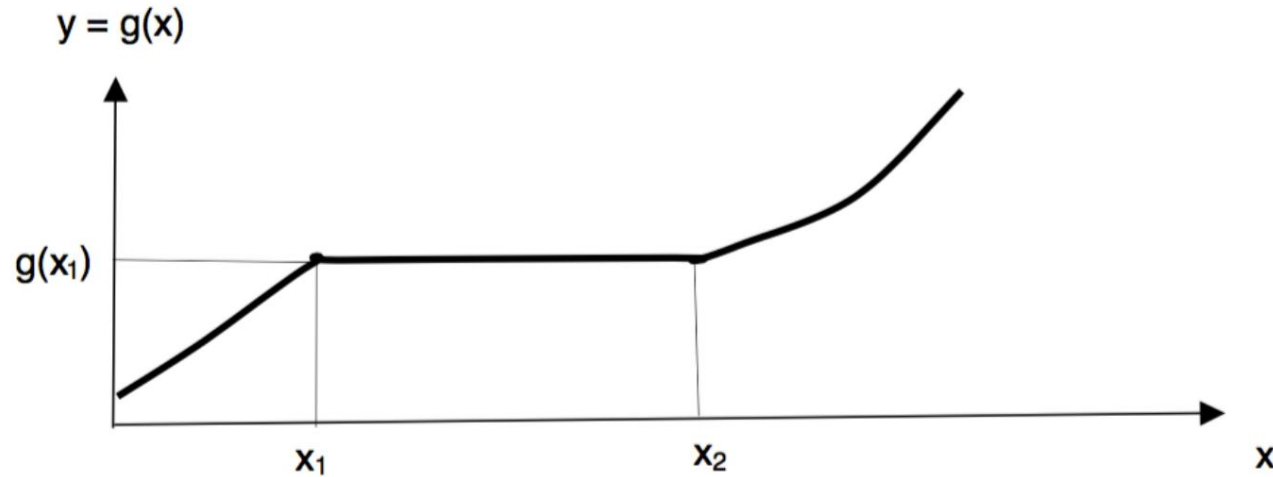
$$y \geq 0$$

它是一个指数分布, 参数为 “#!

—

特殊情况 1: $g(x)$ 在区间内为常数

让我们假设 g 是连续的, 对于除了区间 $[x_1, x_2]$ 之外的所有 x 来说, g 是不递减的



$$g(x_1) = g(x_2)$$

$$P\{Y = g(x_1)\} = P\{x_1 \leq X \leq x_2\} = \int_{x_1}^{x_2} f_X(x) \cdot dx$$

区间 $[x_1, x_2]$ 中 x 的所有值都转换为 $Y = g(x_1)$ 的单个值。

并且 $P(Y = g(x_1)) \geq 0$ (Y 属于长度为 0 的区间的概率)。所以:

Ø 密度, $f_Y(y)$ 是点 $y = g(x_1)$ 处的狄拉克分布

Ø CDF 的不连续性, $F_Y(y)$

$$F_Y\{y = g(x_1^+)\} = F_Y\{y = g(x_1^-)\} + P\{Y = g(x_1)\}$$

具体案例示例 (1/2)

随机变量 X 在区间 $[-1, 1]$ 内服从均匀分布。

考虑转型

$$f_Y(y) = \begin{cases} \frac{1}{2} & \text{如果 } y \geq 0 \\ 0 & \text{如果 } y < 0 \end{cases}$$

$$\text{对于 } x < 0: \quad f_X(x) = 0 \Rightarrow f_Y(y) = 0 \quad \text{对于 } y < 0$$

$$\text{对于 } x \geq 0: \quad f_Y(y) = f_X(x) = \frac{1}{2}$$

因此, Y 的概率密度函数为:

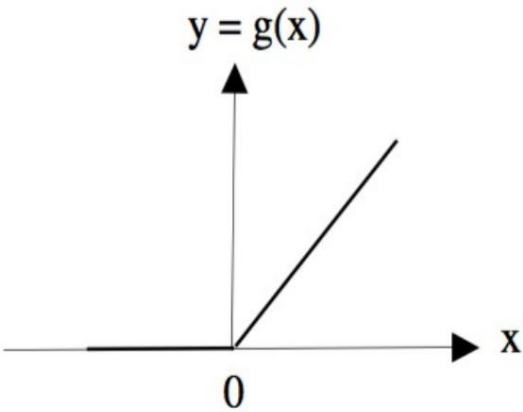
$$f_Y(y) = \begin{cases} \frac{1}{2} & \text{如果 } y \geq 0 \\ 0 & \text{如果 } y < 0 \end{cases}$$

$$f_Y(y) = \frac{1}{2}, \quad 0 \leq y \leq 1$$

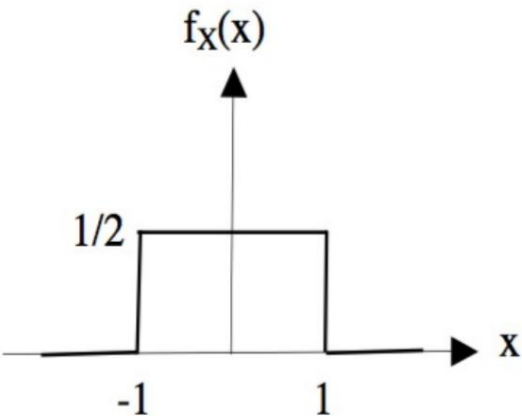
Y 的 CDF 为:

$$F_Y(y) = \begin{cases} 0 & \text{如果 } y < 0 \\ \frac{y}{2} & \text{如果 } 0 \leq y \leq 1 \\ 1 & \text{如果 } y > 1 \end{cases}$$

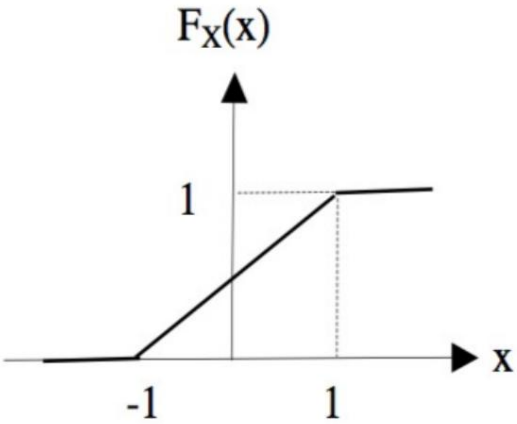
特殊情况示例 (2/2)



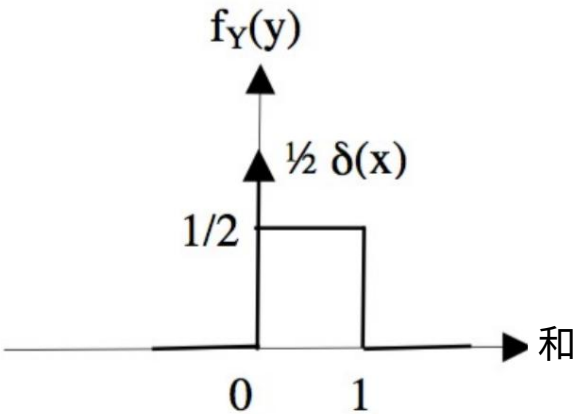
转型



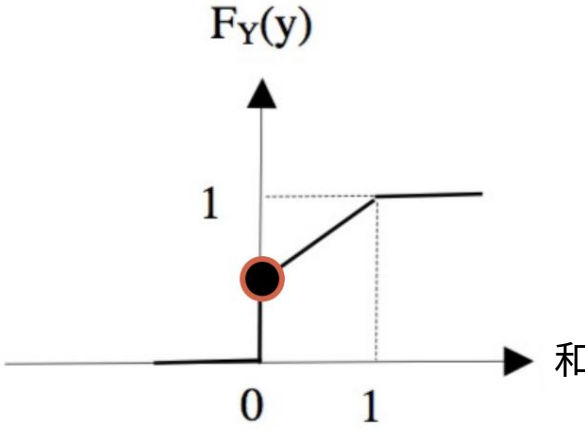
X 的密度



X 的累积分布函数



Y 的密度



Y 的 CDF

第6章 (开头): 二维随机变量

- 二维离散随机变量
- 二维连续随机变量

第6章 (开头): 二维随机变量

- 二维离散随机变量
- 二维连续随机变量

一维随机变量的提醒

sional random variable

一维随机变量 X 完全由以下数据之一定义：

» 其累积分布函数 (CDF)

$F_X(x)$

» 其概率密度函数

$f_X(x)$

» 其特色功能

$\varphi_X(t)$

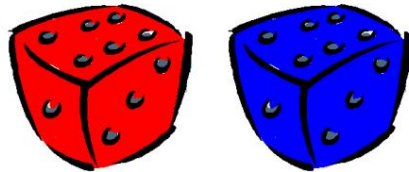
由其矩 m_n , n 描述

通过变换函数创建新的随机变量的可能性。

所有这些概念都将扩展到两个随机的
变量, (X, Y)

示例:实验

掷两个骰子,一红一蓝



	1	2	3	4	5	6
1	(1,1)	(1,2)	(1,3)	(1,4)	(1,5)	(1,6)
2	(2,1)	(2,2)	(2,3)	(2,4)	(2,5)	(2,6)
3	(3,1)	(3,2)	(3,3)	(3,4)	(3,5)	(3,6)
4	(4,1)	(4,2)	(4,3)	(4,4)	(4,5)	(4,6)
5	(5,1)	(5,2)	(5,3)	(5,4)	(5,5)	(5,6)
6	(6,1)	(6,2)	(6,3)	(6,4)	(6,5)	(6,6)

红色骰子($l=1,2,3,4,5,6$),蓝色骰子($k=1,2,3,4,5,6$)

$$P\{\omega_n\} = \frac{1}{36}, \quad n = 1 \text{ à } 36$$

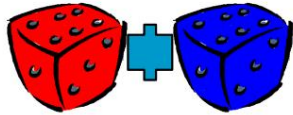
- 36 种可能的结果
:(出现的数字
在每个骰子的上表面)
- 骰子是

对称 => 等概率
(所有结果具有相同的概率)

二维随机变量的定义 (1/2)


为了定义随机变量对 (X,Y) , 我们定义两个映射
从 Ω 到实数集:

X (红色骰子) Y (2 个骰子的总和)






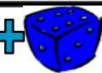
$$\begin{aligned} X(\omega_n) &= X(l, k) = l = x_i \\ Y(\omega_n) &= Y(l, k) = l + k = y_j \end{aligned}$$

该对 (X,Y) 的可能结果:



	1	2	3	4	5	6
1	(1,2)	(1,3)	(1,4)	(1,5)	(1,6)	(1,7)
2	(2,3)	(2,4)	(2,5)	(2,6)	(2,7)	(2,8)
3	(3,4)	(3,5)	(3,6)	(3,7)	(3,8)	(3,9)
4	(4,5)	(4,6)	(4,7)	(4,8)	(4,9)	(4,10)
5	(5,6)	(5,7)	(5,8)	(5,9)	(5,10)	(5,11)
6	(6,7)	(6,8)	(6,9)	(6,10)	(6,11)	(6,12)



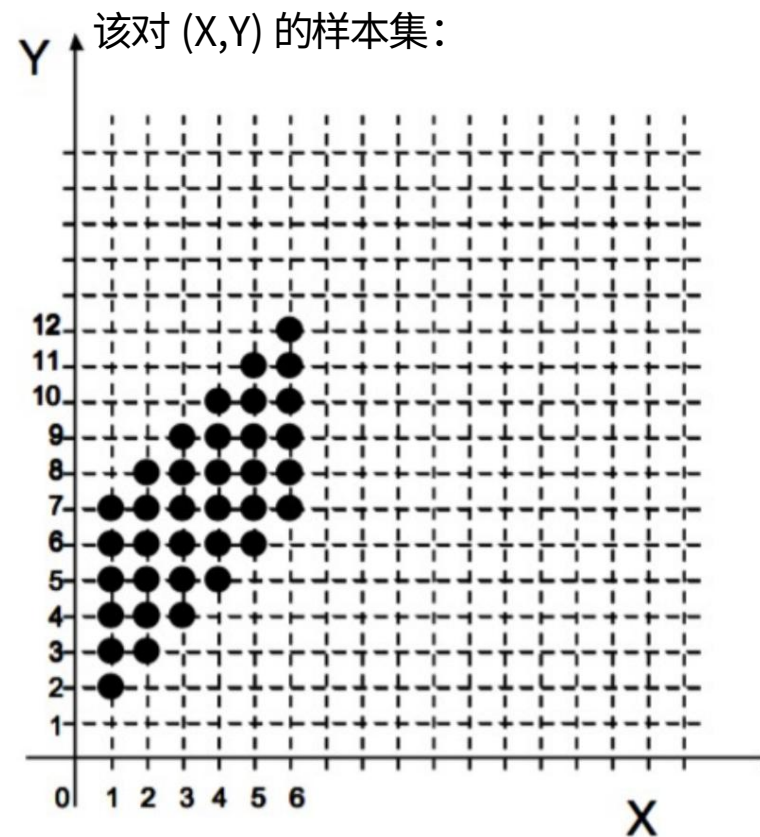
 ,  + )

二维随机变量的定义 (2/2)

X 的可能结果: $x_i = 1, 2, \dots, 6$ (对于 $i=1$ 到 6)

Y 的可能结果: $y_j = 2, 3, \dots, 12$ (对于 $j=1$ 到 11)

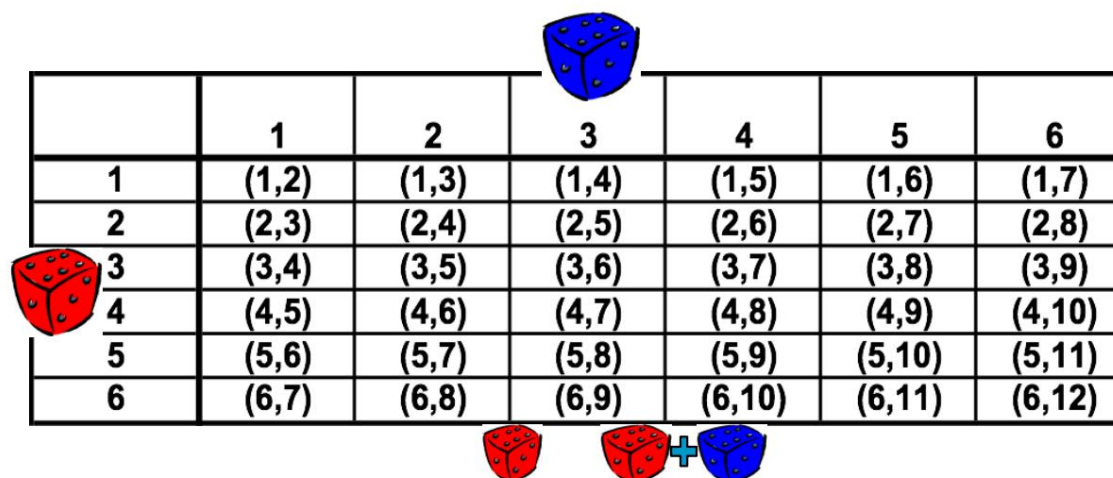
这对夫妇 (X, Y) 不能
取任意值
组合 (x_i, y_j) 。例如, 值 (4,3)
或 (2,10) 永远不会出现!



(X,Y) 对的概率密度: $P(X=x_i \text{ 且 } Y=y_j) = P_{ij}$

这两个映射在两者之间建立了一对一的关系 (双射)

结果 (ω_n) 和有序对 (x_i, y_j) 。



	1	2	3	4	5	6
1	(1,2)	(1,3)	(1,4)	(1,5)	(1,6)	(1,7)
2	(2,3)	(2,4)	(2,5)	(2,6)	(2,7)	(2,8)
3	(3,4)	(3,5)	(3,6)	(3,7)	(3,8)	(3,9)
4	(4,5)	(4,6)	(4,7)	(4,8)	(4,9)	(4,10)
5	(5,6)	(5,7)	(5,8)	(5,9)	(5,10)	(5,11)
6	(6,7)	(6,8)	(6,9)	(6,10)	(6,11)	(6,12)

夫妻的每个值有 36 个结果:

(x_i, y_j)

$$p_{i,j} = P\{\omega_n\} = \frac{1}{36}, \quad \text{对于 } n=1 \text{ 到 } 36$$

列联表

定义 (X, Y) 对概率分布的表格

Y	x1=1	x2=2	x3=3	x4=4	x5=5	x6=6	P Yj
y1=2	1/36	0	0	0	0	0	1/36
y2=3	1/36	1/36	0	0	0	0	2/36
y3=4	1/36	1/36	1/36	0	0	0	3/36
y4=5	1/36	1/36	1/36	1/36	0	0	4/36
y5=6	1/36	1/36	1/36	1/36	1/36	0	5/36
y6=7	1/36	1/36	1/36	1/36	1/36	1/36	6/36
y7=8	0	1/36	1/36	1/36	1/36	1/36	5/36
y8=9	0	0	1/36	1/36	1/36	1/36	4/36
y9=10	0	0	0	1/36	1/36	1/36	3/36
y10=11	0	0	0	0	1/36	1/36	2/36
y11=12	0	0	0	0	0	1/36	1/36
P Xi	6/36	6/36	6/36	6/36	6/36	6/36	1

联合律 (X,Y)

$$P_{ij} \geq 0$$

$$\sum_i \sum_j P_{ij} = 1$$

Y 的边际律

$$P(Y = y_j) = \sum_i P_{ij}$$

X 的边际法则

$$P(X = x_i) = \sum_j P_{ij}$$

表中所有条目的总和
必须始终为 1!

有条件分配

给定X等于固定值,例如X = x3 = 3,计算当X取值 3时Y的条件概率。 , 我们可以

Y	x1=1	x2=2	x3=3	x4=4	x5=5	x6=6	P Yj
y1=2	1/36	0	0	0	0	0	1/36
y2=3	1/36	1/36	0	0	0	0	2/36
y3=4	1/36	1/36	1/36	0	0	0	3/36
y4=5	1/36	1/36	1/36	1/36	0	0	4/36
y5=6	1/36	1/36	1/36	1/36	1/36	0	5/36
y6=7	1/36	1/36	1/36	1/36	1/36	1/36	6/36
y7=8	0	1/36	1/36	1/36	1/36	1/36	5/36
y8=9	0	0	1/36	1/36	1/36	1/36	4/36
y9=10	0	0	0	1/36	1/36	1/36	3/36
y10=11	0	0	0	0	1/36	1/36	2/36
y11=12	0	0	0	0	0	1/36	1/36
P Xi	6/36	6/36	6/36	6/36	6/36	6/36	1

的条件概率
X 给定 Y 固定

的条件概率
给定 X Y

$$P(Y = y_j | X = x_i) = \frac{P_{ij}}{P_{x_i}} = \dots \frac{1}{6}$$

$$P(X = x_i | Y = y_j) = \frac{P_{ij}}{P_{y_j}} = \dots \frac{1}{3}$$

总结几个离散随机变量 (X,Y)

该对 (X,Y) 具有以下结果:

$$\begin{aligned} X &= x_1, x_2, \dots, x_n, \dots \\ Y &= y_1, y_2, \dots, y_m, \dots \end{aligned}$$

联合概率函数

$$P(X=x_i, Y=y_j) = p_{ij}, \quad i=1,2,\dots, n; j=1,2,\dots, m$$

X的边际分布

$$P(X=x_i) = \sum_{j=1}^m p_{ij} = p_{i\cdot}$$

Y的边际分布

$$P(Y=y_j) = \sum_{i=1}^n p_{ij} = p_{\cdot j}$$

给定 Y 的 X 的条件分布

$$P(X=x_i | Y=y_j) = \frac{p_{ij}}{p_{\cdot j}} = \frac{p_{ij}}{\sum_{i=1}^n p_{ij}}$$

给定 X 的 Y 的条件分布

$$P(Y=y_j | X=x_i) = \frac{p_{ij}}{p_{i\cdot}} = \frac{p_{ij}}{\sum_{j=1}^m p_{ij}}$$

独立

如果对于任何 i, j , 两个变量 X 和 Y 是独立的:

$$(X_i = x_i \cap Y_j = y_j) = P(X_i = x_i) P(Y_j = y_j) \text{ 或者 } P(X_i = x_i | Y_j = y_j) = P(X_i = x_i) \quad \forall x_i, y_j,$$

如果它们是独立的, 则条件分布为

与边际分布相同。

$$(X_i = x_i | Y_j = y_j) = P(X_i = x_i), \quad \forall y_j$$

$$(Y_j = y_j | X_i = x_i) = P(Y_j = y_j), \quad \forall x_i$$

第6章 (开头): 二维随机变量

- 二维离散随机变量
- 二维连续随机变量

累积分布函数 (CDF)

定义：

一维变量提醒：

$$F_X(x) = P(X \leq x)$$

对于二维随机变量 (X, Y) ，CDF 是
事件相交的概率： $\{X \leq x\}$ 和 $\{Y \leq y\}$

$$F_{X,Y}(x, y) = (\{ \leq \cap \{ \leq \})$$

符号：

F :联合累积分布函数

X, Y :随机变量

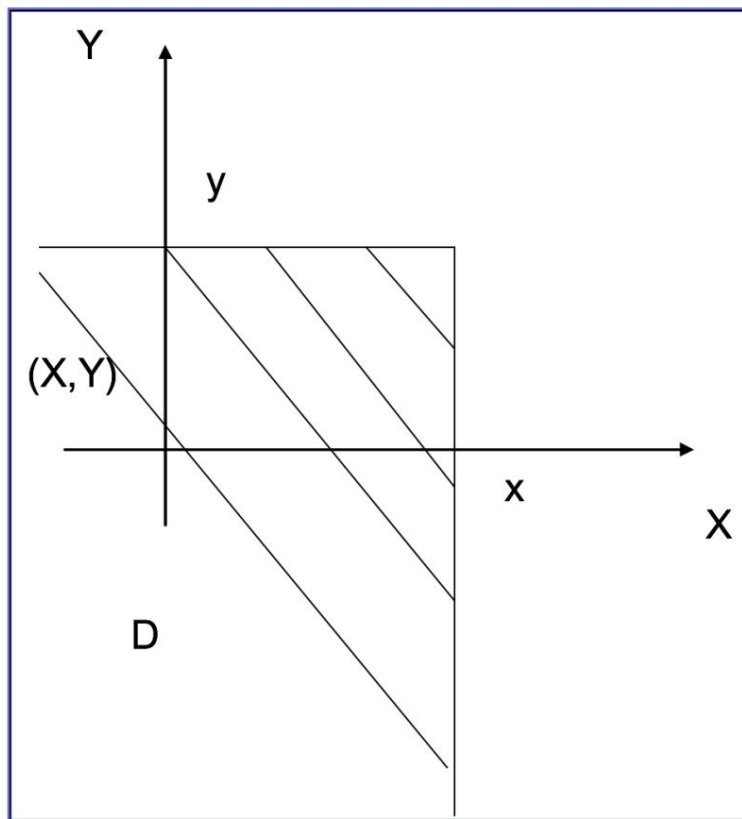
x, y :实际阈值

P :概率

图形表示（非常有用！）：

笛卡尔坐标系中的表示

$F_{XY}(x, y) : (X, Y) \in D$ 的概率



边际累积分布函数

X 的边际累积分布函数

它只是X的累积分布函数

$$F_X(x) = P(\{ \omega : X(\omega) \leq x \}) = P(\{ \omega : X(\omega) \leq x \}, \forall \omega \in \Omega)$$

只需考虑二维累积分布函数并设置：

$y = +\infty$ (或域 $D_{X,Y}$ 中 $y = y_{\max}$)

$$F_X(x) = P(\{ \omega : X(\omega) \leq x, Y(\omega) \leq +\infty \}) = F_{X,Y}(x, +\infty)$$

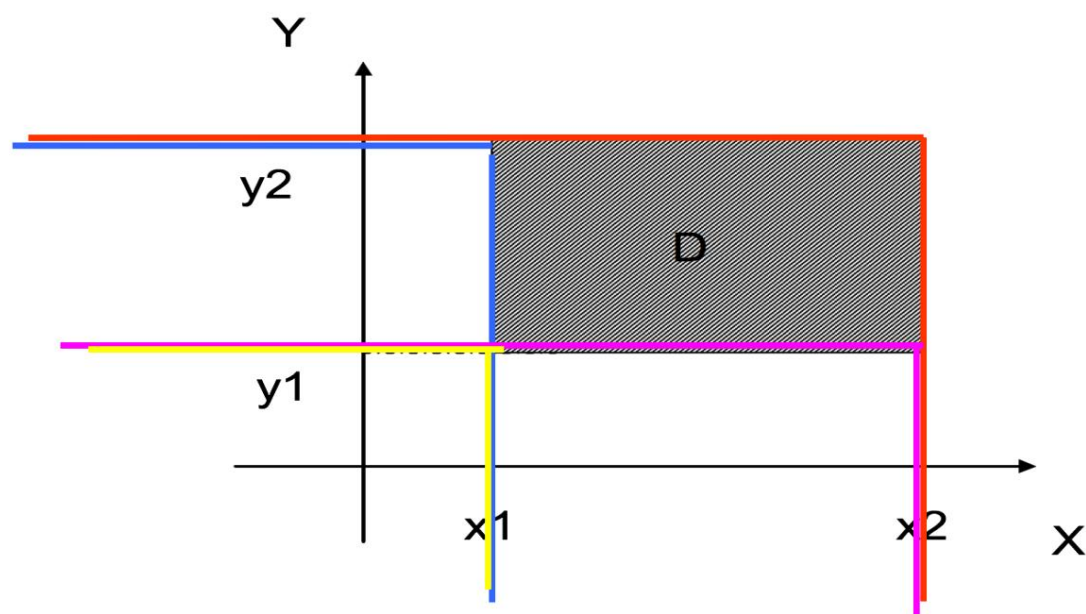
Y 的边际累积分布函数

$$F_Y(y) = F_{X,Y}(+\infty, y)$$

矩形的概率计算（二维区间）

我们的目标是根据 $F_{X,Y}(x,y)$ 计算 (X,Y) 对属于矩形 D 的概率：

$$\{ \varepsilon \in [a, b] \text{ 和 } \varepsilon \in [c, d] \}$$



$$\begin{aligned}
 & \left(\{ \varepsilon \in [a, b] \text{ 和 } \varepsilon \in [c, d] \} \right) \\
 &= F_{X,Y}(b, d) - F_{X,Y}(a, d) - F_{X,Y}(b, c) + F_{X,Y}(a, c)
 \end{aligned}$$

联合概率密度函数

一维情况：

$$f_X(x) = \frac{dF_X(x)}{dx}$$

二维情况：

$$f_{X,Y}(x,y) = \frac{\partial}{\partial y} \left(\frac{\partial F_{X,Y}(x,y)}{\partial x} \right) = \frac{\partial^2 F_{X,Y}(x,y)}{\partial y \partial x}$$

符号：

f : 概率密度函数

X, Y : 随机变量

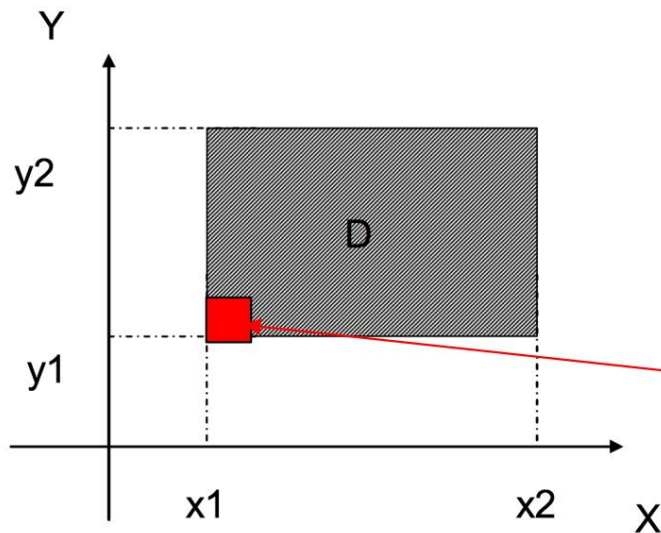
x, y : 实际阈值

$\frac{\partial^2}{\partial y \partial x}$: 相对于 x (阈值) 和 y (阈值) 的导数

概率密度的无穷小解释

让我们计算域 D (矩形) 上的积分:

$$\int_{x_1}^{x_2} \int_{y_1}^{y_2} f_{X,Y}(x,y) dx dy = F_{X,Y}(x_2, y_2) - F_{X,Y}(x_1, y_2) - F_{X,Y}(x_2, y_1) + F_{X,Y}(x_1, y_1)$$



$$(\{ \text{ } < \leq \text{ } \text{ \& , } \text{ } < \leq \text{ } \text{ \& } \})$$

在无穷小域 $dS = dx dy$ 中密度
可以认为是常数

$$x_1 = x, x_2 = x + dx$$

$$y_1 = y, y_2 = y + dy$$

$$(\{ \text{ } < \leq \text{ } + \text{ } , \text{ } < \leq \text{ } + \text{ } \}) = \text{ } , \text{ } (\text{ })$$

解释: (X,Y) 对在点 (x,y) 的邻域属于 dS 的概率与该点的联合密度函数的值成正比。

结论: 对于连续随机变量, $\{X = x, Y = y\}$ 的概率为零
因为 dS 被视为零。

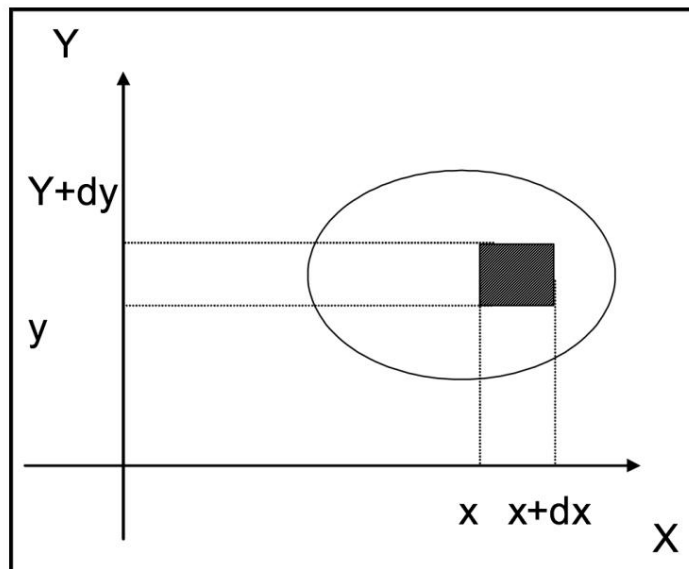
使用密度:

密度允许计算几个随机的概率

变量 (X, Y) 属于任意域 D

$$P(\{(X, Y) \in D\}) = \int_D f_{X,Y}(x, y) dx dy$$

只需将 D 分解为不相交的元素并应用概率公理 3 即可。



积分表示由联合概率密度函数定义的
表面之间的体积

和域 D 。

与累积分布函数的关系：

$$\int_{-\infty}^{x_2} \int_{-\infty}^{y_2} f_{X,Y}(x,y) dx dy = F_{X,Y}(x_2, y_2)$$

通过取 $x_2 = +\infty$, $y_2 = +\infty$, 我们得到：

$$\int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} f_{X,Y}(x,y) dx dy = F_{X,Y}(+\infty, +\infty) = 1$$

特性：_____

- 表面下的体积为 1。
- 密度是非负的。

$$f_{X,Y}(x,y) \geq 0$$

X 和 Y 的边际密度:

X 的边际密度

$$f_X(x) = \int_{-\infty}^{+\infty} f_{X,Y}(x,y) dy$$

Y 的边际密度

$$f_Y(y) = \int_{-\infty}^{+\infty} f_{X,Y}(x,y) dx$$

提示:计算另一个变量在定义域上的积分,使其消失

条件分布

给定 X 的 Y 的条件分布

$$f_Y(y|X=x) = \frac{f_{X,Y}(x,y)}{\int_{-\infty}^{+\infty} f_{X,Y}(x,y) dy} = \frac{f_{X,Y}(x,y)}{f_X(x)}$$

给定 Y 的 X 的条件分布

$$f_X(x|Y=y) = \frac{f_{X,Y}(x,y)}{\int_{-\infty}^{+\infty} f_{X,Y}(x,y) dx} = \frac{f_{X,Y}(x,y)}{f_Y(y)}$$

如果 X 和 Y 独立会发生什么？

提醒:如果事件A和B是独立的:

$$P(A \cap B) = P(A)P(B)$$

如果 $A = \{X \leq x\}$ 且 $B = \{Y \leq y\}$, 则 $P(X \leq x, Y \leq y) = P(X \leq x)P(Y \leq y)$:

$$F_{X,Y}(x,y) = F_X(x)F_Y(y)$$

耦合的联合累积分布函数等于边际累积分布函数的乘积。

密度怎么样？

$$f_{X,Y}(x,y) = f_X(x)f_Y(y)$$

联合密度也是边缘密度的乘积。