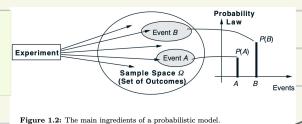


Probabilistic Model 概率模型



1. 是什么：对一个underlying process, i.e. experiment的不确定性的量化描述 (quantitative description)

只考虑一个实验，3次抛硬币被视为一个a simple experiment, 而不是三个实验

2. 包括2个元素：

(1) Sample Space/样本空间 Ω

Def. The set of all possible outcomes (ω) of an experiment.

(2) Probability Law/概率律

• Def. 假设我们已经确定了与实验相关的样本空间。

Intuitively, 这个法则规定了任意结果/outcome或任意事件/Event (a subset of the sample space; a collection of outcomes) 的 "likelihood". More precisely, 概率率为每一个事件 A 赋予一个数值 $P(A)$, i.e. the probability of A .

→ 满足的公理:= 概率公理

$$P(A) = \frac{2}{3} \text{ 可粗略理解为在大量重复实验中, 事件 } A \text{ 会在大约 } \frac{2}{3} \text{ 的实验中出现.}$$

• Discrete Probability Law

If the sample space consists of a finite number of possible outcomes, then the probability law is specified by the probabilities of the events that consist of a single element. In particular, the probability of any event $\{s_1, s_2, \dots, s_n\}$ is the sum of the probabilities of its elements:

$$P(\{s_1, s_2, \dots, s_n\}) = P(\{s_1\}) + P(\{s_2\}) + \dots + P(\{s_n\}).$$

概率律由单个元素组成的事件的概率指定。

Special Case: All single-element events have the same probability: $P(A) = \frac{\# \text{elements of } A}{n}$

3. 与 random variable 的联系。

在许多概率模型中, outcomes 具有数值性质/numerical nature. It's outcomes correspond to stock prices.

在其它实验中, outcome 并非数值, 但它们可能与某些 numerical values of interest 相关联. 比如实验是 Select students from a given population, we may wish to consider their GPA. 处理这些数值时, 可以通过随机变量的概念来为它们分配概率。

→ 符合以下3个特点的试验叫随机试验

{ 可重复: 可在相同条件下重复进行

{ 多样性: 每次试验的可能结果不止一个, 且能事先明确试验的所有结果 has a well-defined set of possible outcomes infinitely repeated

{ 随机性: 一次试验前不确定哪个结果会实现 has more than one possible outcome, and deterministic if it has only one

• 样本空间 Ω 必须满足的条件:

{ mutually exclusive / 互斥: 与该骰子相关的样本空间不能既包含 "1 or 3" 又包含 "1 or 4" 作为可能的 outcome. 否则当掷出1时, 实验 outcome 不唯一

{ collectively exhaustive / 完全穷尽: 无论实验发生什么情况, 我们总能获得包含在样本空间内的 outcome.

{ At the right granularity: 根据研究目的确定适当的细颗粒度. 如对于 "Coin Flip" 随机试验可定义



• Probability Axioms / 柯尔莫果洛夫公理 / Kolmogorov Axioms

Let $A = \text{Pot}(\Omega)$ (Event Space/die Menge aller Ereignisse über Ω /所有可能事件的集合).

Nonnegativity/非负性 $0 \leq P(A) \leq 1 \quad \forall A \in A$

Additivity/可加性

If A and B are two disjoint events: $P(A \cup B) = P(A) + P(B)$

If Ω has infinite number of elements and A_1, A_2, \dots is the sequence of disjoint events:

Normalization/归一化 $P(\Omega) = 1$

$$P(A_1 \cup A_2 \cup \dots) = P(A_1) + P(A_2) + \dots$$

e.g. The experiment of rolling a pair of 4-sided dice.

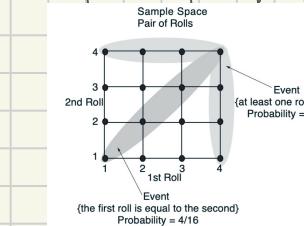


Figure 1.4: Various events in the experiment of rolling a pair of 4-sided dice, and their probabilities, calculated according to the discrete uniform law.

→ Some properties of Probability Laws

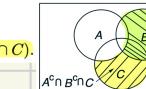
Consider a probability law, and let A , B , and C be events.

(a) If $A \subset B$, then $P(A) \leq P(B)$.

(b) $P(A \cup B) = P(A) + P(B) - P(A \cap B)$.

(c) $P(A \cup B) \leq P(A) + P(B)$.

(d) $P(A \cup B \cup C) = P(A) + P(A^c \cap B) + P(A^c \cap B^c \cap C)$.



Example 1.5. Romeo and Juliet have a date at a given time, and each will arrive at the meeting place with a delay between 0 and 1 hour, with all pairs of delays being equally likely. The first to arrive will wait for 15 minutes and will leave if the other has not yet arrived. What is the probability that they will meet?

Let us use as sample space the square $\Omega = [0, 1] \times [0, 1]$, whose elements are the possible pairs of delays for the two of them. Our interpretation of "equally likely" pairs of delays is to let the probability of a subset of Ω be equal to its area. This probability law satisfies the three probability axioms. The event that Romeo and Juliet will meet is the shaded region in Fig. 1.5, and its probability is calculated to be 7/16.

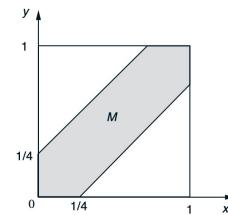


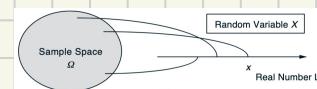
Figure 1.5: The event M that Romeo and Juliet will arrive within 15 minutes of each other (cf. Example 1.5) is

$$M = \{(x, y) \mid |x - y| \leq 1/4, 0 \leq x \leq 1, 0 \leq y \leq 1\},$$

and is shaded in the figure. The area of M is 1 minus the area of the two unshaded triangles, or $1 - (3/4) \cdot (3/4) = 7/16$. Thus, the probability of meeting is 7/16.

Discrete Random Variables

离散随机变量

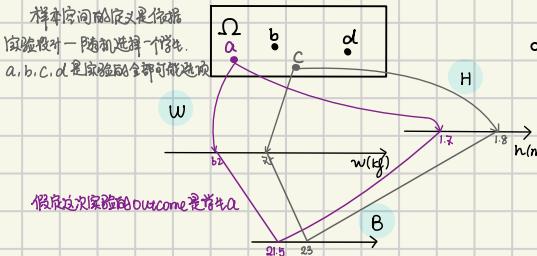


1. Def. $X: \Omega \rightarrow \mathbb{R}$

A random variable X is a real-valued function from the outcome of the experiment to the real numbers. X is called discrete if its range (the set of values that it can take) is finite or at most countably infinite.

• 定义理解:

Our probabilistic experiment is to pick a student at random according to some probability law. 被抽中的学生的概率可能不同, 这些概率值是根据某种概率分布来设置的
 $\Omega = \{a, b, c, d\}$, and then record their weights and heights. 身高体重值不是样本空间的一部分, 只是学生属性的一个描述。是选择某个学生后所记录的一个附带信息。



- 用 W 来表示 "Weight" 这个抽象概念, 一旦我们确定哪位学生被选中, 就能确定 W 具体的取值。

从这个意义上说, W is a function that assigns a numerical value (w) to each possible outcome (a) of the experiment.

- 我们还可以根据已有的随机变量 W 和 H 创建新的随机变量 / A function of a random variable defines another random variable.

例如身体质量指数 Body Mass Index $B = \frac{W}{H^2}$ 也是定义在样本空间上的一个函数 (一旦确定 outcome, BMI 也随之确定)



随机过程的结果是 Ω 中的元素 W , 但是人们往往关注的不是结果本身 $W \in \Omega$, 不是对其产生过程的精确描述 / die exakte Beschreibung des Zustandekommens, 而是函数值 $X(W)$. 从这个意义上说, 随机变量代表了对研究对象 / Untersuchungsgegenstand 的关注。

• e.g.

I toss a coin five times. This is a random experiment and the sample space can be written as

$$\Omega = \{\text{TTTTT}, \text{TTTTH}, \dots, \text{HHHHH}\}.$$

Note that here the sample space Ω has $2^5 = 32$ elements. Suppose that in this experiment, we are interested in the number of heads. We can define a random variable X whose value is the number of observed heads. The value of X will be one of 0, 1, 2, 3, 4 or 5 depending on the outcome of the random experiment.

2. Framework to describe discrete random variable.

(1) PMF/Probability Mass Function / 概率质量函数

• Def.

描述了离散随机变量 X 所有可能取值的概率。

PMF of X , denoted p_X is the "probability law" or "probability distribution" of a discrete random variable X .

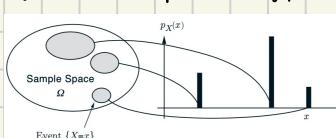
对于 X 的一个可能的数值 x , $p_X(x)$ 给出了事件 $\{X=x\}$ 的概率, 即 X 取到该值 x 的概率。 $p_X(x) = P(\{X=x\}) = P(X=x)$

• PMF Properties

$\sum_x p_X(x) = 1$ 事件 $\{X=x\}$ 互不相容 (必须满足的条件), X 遍历 X 所有可能的数值. 根据 Additivity 和 Normalization 定理可推导得出

• 计算

对于 X 的每个可能取值 x , collect 所有可能使事件 $\{X=x\}$ 发生的 possible outcomes; 将这些结果的概率相加以获得 $p_X(x)$



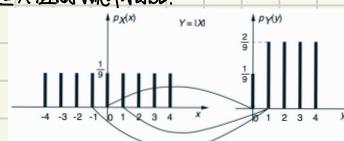
• PMF for Functions of Random Variables $Y = g(X)$

Y 也是一个随机变量, since it provides a numerical value for each possible outcome.

$$P_Y(y) = \sum_{\{x | g(x)=y\}} p_X(x)$$

从已知 $p_X(x)$ 的 PMF 推导出 $p_Y(y)$: 为计算 Y 在某一特定值 y 处的概率 $P_Y(y)$, 寻找到所有使得 $g(x)=y$ 的 x 值, 并将这些 x 值的概率相加。

e.g.



$$\text{Let } Y = ax + b \quad P_Y(y) = P_Y(y = ax + b) = P_X\left(\frac{y-b}{a}\right)$$

(2) Mean and Variance

Expectation / Mean / Expected Value $E[X]$

Motivation

多次旋转一个 fortune wheel，每次旋转出现代表金钱奖励的数字 m_i ($i \in [n]$) 的概率是 p_i 。

What is the amount of money that you "expect" to get "per spin"? Average in large number of independent repetitions of the experiment.

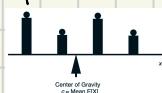
$$\text{total amount received} = m_1 k_1 + m_2 k_2 + \dots + m_n k_n$$

$$\text{total amount received per spin} = \frac{m_1 k_1 + m_2 k_2 + \dots + m_n k_n}{n}$$

$$= m_1 \cdot \frac{k_1}{n} + m_2 \cdot \frac{k_2}{n} + \dots + m_n \cdot \frac{k_n}{n}$$

$$= m_1 \cdot p_1 + m_2 \cdot p_2 + \dots + m_n \cdot p_n \quad (p_i := \text{relative frequency})$$

Interpretation



假设有 n 个 bar，在每个点 x 处放置一个 weight $p_x(x)$ ，且 $p_x(x) > 0$ 。

则在重心 C 处，左侧权重的力矩之和等于右侧权重的力矩之和 $\sum_x (x - c) p_x(x) = 0$, or $c = \sum_x x p_x(x)$

The center of gravity is the point where you would place your finger to balance the diagram so that it doesn't tip in either direction. $E[X]$ 反映了 X 概率分布的中心位置。

Properties

If $X \geq 0$, then $E[X] \geq 0$

$X \geq 0$ i.e. $\forall w: X(w) \geq 0$

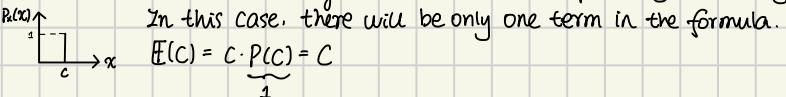
If $a \leq X \leq b$, then $a \leq E[X] \leq b$

$a \leq X \leq b$ i.e. $\forall w: a \leq X(w) \leq b$

Sketch proof: $E[X] = \sum_x x \cdot P_x(x) \geq \sum_x a \cdot P_x(x) = a \cdot \sum_x P_x(x) = a \cdot 1 = a$

If c is a constant, $E[c] = c$

We can think of constant as being a random variable of a very special type.



Variance $\text{Var}(X)$

The Variance provides a measure of dispersion of X around its mean

Standard deviation of X $s_x = \sqrt{\text{Var}(X)}$

计算

$$E[X] = \sum_x x \cdot P_x(x)$$

X 可能取值的概率加权平均

• Expected Value Rule for Functions of Random Variables $E[g(X)] = \sum_x g(x) \cdot P_x(x)$

$$\text{方法 I. Averaging over } Y \quad E[Y] = \sum_y y \cdot P_Y(y)$$

range over the values of Y one at a time. $E[Y] = 3 \cdot (0.1 + 0.2) + 4 \cdot (0.3 + 0.4)$

$$\text{方法 II. Averaging over } X \quad E[Y] = E[g(X)] = \sum_x g(x) \cdot P_x(x)$$

range over the values of X one at a time, and take into account their individual contributions.

$$E[Y] = 3 \cdot 0.1 + 3 \cdot 0.2 + 4 \cdot 0.3 + 5 \cdot 0.3$$

10% of the time, X is going to be 2, and when that happens, Y takes on the value of 3.

Proof: $\sum_x g(x) \cdot P_x(x)$ fix a particular value of y , and add over all those x 's that correspond to that particular y .

$$= \sum_{\substack{y: x \text{ gives } y}} g(x) \cdot P_x(x)$$

$$= \sum_{\substack{y: x \text{ gives } y}} y \cdot P_x(x)$$

$$= \sum_y y \sum_{\substack{x: x \text{ gives } y}} P_x(x)$$

将所有赋予同一 y 值的 x 的概率都相加

$$= \sum_y y \cdot P(y)$$

$$= E[Y]$$

$$\text{应用: } E[X^2] = \sum_x x^2 \cdot P_x(x)$$

• Linearity of Expectation $E[aX+b] = a \cdot E[X] + b$

In general, $E[g(X)] \neq g(E[X])$ 当 g 是线性函数时, 等式成立.

计算

$$\text{var}(X) = E[(X - \mu)^2] = \frac{g(x) = (x - \mu)^2}{E[g(x)] = \sum_x g(x) \cdot P_x(x)} \sum_x (x - \mu)^2 \cdot P_x(x)$$

$$\text{var}(X) = E[(X - \mu)^2] = E[X^2 - 2\mu X + \mu^2] = E[X^2] - E[2\mu X] + E[\mu^2] = E[X^2] - 2\mu E[X] + \mu^2 = E[X^2] - (E[X])^2$$

• Linearity of Variance $\text{Var}(aX+b) = a^2 \text{Var}(X)$

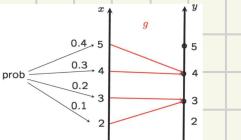
$$E[Y] = E[X+b] = \mu + b$$

$$\text{var}(Y) = E[(Y - E[Y])^2] = E[(X+b - (\mu+b))^2] = E[(X-\mu)^2] = \text{var}(X)$$

$$\text{Let } Y = aX$$

$$E[Y] = E[aX] = a\mu$$

$$\text{var}(Y) = E[(Y - E[Y])^2] = E[(ax - a\mu)^2] = E[a^2(X - \mu)^2] = a^2 \cdot E[(X - \mu)^2] = a^2 \cdot \text{var}(X)$$

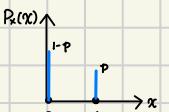


伯努利分布 / 0-1 分布 / Bernoulli Distribution

Parameter $p \in [0, 1]$

- Def.

$$X = \begin{cases} 1 & \text{w.p. } p \\ 0 & \text{w.p. } 1-p \end{cases}$$



- 期望 $E[X] = 1 \cdot p + 0 \cdot (1-p) = p$

If X is the indicator of an event A , $X = IA$
 $E[IA] = P(A)$

- 方差 $\text{Var}(X) = \sum_x (x - E[X])^2 p_X(x)$

$$= (1-p)^2 \cdot p + (0-p)^2 \cdot (1-p)$$

$$= p \cdot p^2$$

$$= p(1-p)$$

$$\text{Var}(X) = E[X^2] - (E[X])^2 = E[X] - (E[X])^2 \quad (1^2 = 1, 0^2 = 0)$$



$$= p \cdot p^2$$

$$= p(1-p)$$

Interpretation: 方差是对随机变量的 uncertainty 的度量，即对 randomness 的度量。只有 coin 的随机性最大，即正反两面出现的概率都是 $\frac{1}{2}$ 时 $p=1-p=\frac{1}{2}$ ，coin 才公平。

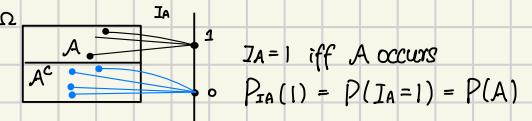
Bernoulli random variable 的方差的验证证了这种直觉；

另一方面，在极端情况下，若 $p=0$ ，即 coin 总是朝上（或朝下），我们没有任何随机性，相应地，方差为 0。

- Application

- Models a trial that results in 2 alternative outcomes: success / failure, Heads / Tail, etc.

- Connect events and random variable, i.e., the indicator random variable IA of an event A



$\Rightarrow IA$ is a Bernoulli random variable, with parameter p equals to the probability of the event of interest.

IA allows us to translate a manipulation of events to a manipulation of random variables.

离散型均匀分布 / discrete uniform distribution

Parameters: $a, b \quad a \leq b$

- Setting

Experiment: Pick one of $a, a+1, \dots, b$ at random.
all equally likely

Sample Space: $\{a, a+1, \dots, b\}$ $b-a+1$ possible values

Random Variable: $X: X(w)=w$



a $a+1$ \dots $\frac{a+b}{2}$ b
the center of the symmetry

- 期望 $E[X] = a \cdot \frac{1}{b-a+1} + \dots + b \cdot \frac{1}{b-a+1} = \frac{(a+b) \cdot (b-a+1)}{2} \cdot \frac{1}{b-a+1} = \frac{a+b}{2}$

- 方差 $\text{Var}(X) = E[X^2] - (E[X])^2 = \frac{1}{b-a+1} (a^2 + \dots + b^2) - \left(\frac{a+b}{2}\right)^2 = \frac{1}{12} (b-a)(b-a+2)$

- Application

Model of: Complete ignorance, no reason to believe that one value is more likely than the other.

二项分布 / Binomial Distribution

Parameters: positive integer n , $P \in [0, 1]$

- Setting

Experiment: n independent tosses of a coin, $P(\text{head}) = p$

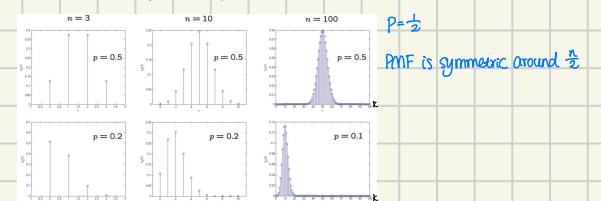
Describe the # of successes in n independent Bernoulli trials.

Sample Space: Set of sequences of H and T, of length n .

Random Variable: X , # Heads observed

- Def.

$$P_X(k) = \binom{n}{k} p^k \cdot (1-p)^{n-k}, \quad k = 0, 1, \dots, n$$



e.g.

$P(X=2) = P(HHT) + P(HTH) + P(THH)$
 $= 3 \cdot p^2 \cdot (1-p)$
 $= \left(\frac{3}{2}\right) p^2 \cdot (1-p)$

- Calculation Preparation

Let X be the indicator variable

$$X_i = \begin{cases} 1 & \text{if } i\text{-th trial is a success} \\ 0 & \text{otherwise} \end{cases} \quad E[X_i] = p$$

successes in n dependent trials $X = X_1 + \dots + X_n$

- 每次试验都是独立的且具有相同的成功概率 p 。

因此每个 X_i 都有相同的分布, X_i is a Bernoulli random variable with parameter p

- X 是 n 次独立试验中成功的总次数, 因此 X 是多个随机变量 X_i 的和。
 X is a binomial random variable with parameters n and p

- 期望 $E[X] = np$

- 方差 $\text{Var}(X) = \text{Var}(X_1) + \dots + \text{Var}(X_n)$

$$= n \cdot \text{Var}(X_1)$$

$$= n \cdot p(1-p)$$

几何分布 / Geometric Distribution

Parameter $p: 0 < p \leq 1$

- Setting

Experiment: infinite many independent tosses of a coin

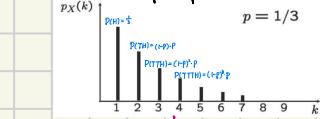
$P(\text{head}) = p$. Describe the # of trials until a success

Sample Space: Set of infinite sequences of H and T.

Random Variable: X , # tosses until the first Heads.

- Def

$$P_X(k) = ((1-p))^{k-1} \cdot p \quad k \in \mathbb{N}_0, \quad p \in (0, 1]$$

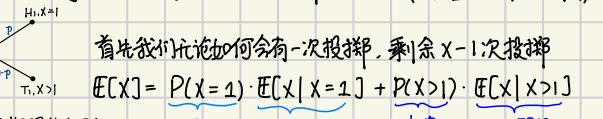


- 期望 $E[X] = \frac{1}{p}$

- 法 I: Formula plug in

$$E[X] = \sum_{k=1}^{\infty} k \cdot P_X(k) = \sum_{k=1}^{\infty} k \cdot (1-p)^{k-1} \cdot p$$

- 法 II: 随机变量 X 的期望值包括两部分： (期望的线性性质)



(全期望值定理) $\Rightarrow E[X] = p + (1-p) \cdot (1+E[X])$
 $\Rightarrow E[X] = \frac{1}{p}$

We've wasted one try, and we are back where we started.

Interpretation: 当 p 较小时，意味着观察到 head 的概率较小也就是说我们要等待较长时间才能第一次看到正面。随机变量 X 的值也会更平均。

- 方差 $\text{Var}(X) = \frac{1-p}{p^2}$

几何分布具有无记忆性，即使已经进行了 k 次失败的试验 ($X>k$)，剩余试验次数 $X-k$ 仍服从与原几何分布相同的分布。

$$E[X^2] = \underbrace{P(X=1) \cdot E[X^2 | X=1]}_1 + \underbrace{P(X>1) \cdot E[X^2 | X>1]}_{1-p} = \underbrace{E[(X+1)^2]}_{= E[X^2 + 2X + 1]} = \underbrace{\frac{2}{p^2} - \frac{1}{p}}$$

$$\text{Var}(X) = E[X^2] - E[X]^2 = \frac{1-p}{p^2}$$

(3) Joint PMF / 联合概率质量函数

Motivation

当我们同时处理多个与同一实验相关的随机变量时，如随机变量X和Y，单独的 P_X 和 P_Y 无法提供X和Y相互作用的信息。

Def.

The joint PMF of X and Y is defined by

$$P_{XY}(x,y) = P(X=x, Y=y)$$

↑
PMF
Subscript: the random variables we are dealing with

表示随机变量X取值x的同时随机变量Y取值y的概率。

Properties

$\sum_x \sum_y P_{XY}(x,y) = 1$ 所有可能的(x,y) pair 的概率总和为1。因为联合PMF覆盖了随机变量X和Y在其整个取值范围内的组合。

$$P_X(x) = P(X=x)$$

$$= \sum_y P(X=x, Y=y)$$

$$= \sum_y P_{XY}(x,y)$$

$$P_Y(y) = \sum_x P_{XY}(x,y)$$

the marginal distribution of X

事件{ $X=x$ }的概率是 disjoint events { $X=x, Y=y$ } 在Y遍历Y的所有值的并集

- The marginal probability is the probability of a single event occurring, independent of other events.

		Joint PMF $P_{XY}(x,y)$ in tabular form			
		0	1/20	1/20	1/20
y	4	0	1/20	1/20	1/20
	3	1/20	2/20	3/20	1/20
2	1/20	2/20	3/20	1/20	7/20
	1	1/20	1/20	1/20	0
		1	2	3	4
		3/20	6/20	8/20	3/20
		Row Sums: Marginal PMF $P_{Y X}$			
		Column Sums: Marginal PMF $P_{X Y}$			

PMF for Functions of Multiple Random Variables $Z = g(X, Y)$

$$P_Z(z) = P(Z=z) = P(g(X, Y)=z) = \sum_{\{(x,y) | g(x,y)=z\}} P_{XY}(x,y)$$

Expected Value Rule for Functions of Multiple Random Variables $E[g(X, Y)] = \sum_x \sum_y g(x, y) \cdot P_{XY}(x, y)$

Interpretation:

函数 $g(X, Y)$ 的输出值依赖于特定的 (x, y) 值对。每一对 (x, y) 发生的概率由它们的联合PMF $P_{XY}(x, y)$ 给出。

每一对 (x, y) 对函数期望值的贡献是 $g(x, y) \cdot P_{XY}(x, y)$

最后累加所有 (x, y) 对。

- $E[aX+bY+c] = aE[X] + bE[Y] + c$ is always true, no matter X, Y, Z are independent or not.

(4) Conditional PMF / 条件概率质量函数

Motivation

条件概率让我们能够“divide-and-conquer”，将复杂的 probability models 分解为更简单的 Submodels，然后逐个分析。

比如给定 an experiment, a corresponding sample space and a probability law. 假设我们知道某个特定事件 B 已经发生，那么我们可以只关注那些在事件 B 发生的前提下可能发生的事情。只关注样本空间中与事件 B 相交的部分，缩小了样本空间。

$P(A|B)$ captures the partial information that event B provides about event A.

Conditional Probability

$$\text{Def. } P(A|B) = \frac{P(A \cap B)}{P(B)} \quad P(B) > 0$$

Conditional Probabilities form a legitimate Probability law

Non-negativity	Additivity	Normalization
\checkmark $P(A_1 \cup A_2 B) = \frac{P((A_1 \cup A_2) \cap B)}{P(B)}$ $= \frac{P((A_1 \cap B) \cup (A_2 \cap B))}{P(B)}$ $= P(A_1 B) + P(A_2 B)$	For any two disjoint events A_1, A_2 : $P(\Omega B) = \frac{P(\Omega \cap B)}{P(B)} = \frac{P(B)}{P(B)} = 1$ 所有概率度量都集中在事件 B 的范围内。 相当于将整个概率空间缩小到 B。	i.e. 条件概率度量在 new universe B 内的概率律

Calculation Rules

I. Multiplication Rule

$$P\left(\bigcap_{i=1}^n A_i\right) = P(A_1) \cdot P(A_2 | A_1) \cdot P(A_3 | A_1 \cap A_2) \cdots P(A_n | \bigcap_{i=1}^{n-1} A_i)$$

Setting: Experiments that have sequential character

乘积规则反映了每个事件在前一个事件已经发生的基础上发生的依赖性。每一步的成功都依赖于前一步的成功。这种类型的问题通常用树状图表示。

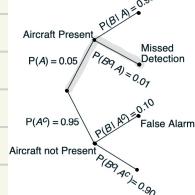
e.g.

Example 1.9. Radar detection. If an aircraft is present in a certain area, a radar correctly registers its presence with probability 0.99. If it is not present, the radar falsely registers an aircraft presence with probability 0.10. We assume that an aircraft is present with probability 0.05. What is the probability of false alarm (a false indication of aircraft presence), and the probability of missed detection (nothing registers, even though an aircraft is present)?

A sequential representation of the sample space is appropriate here, as shown in Fig. 1.8. Let A and B be the events

A = {an aircraft is present},

B = {the radar registers an aircraft presence},

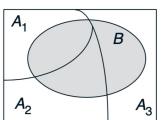


II. Total Probability Theorem / 全概率定理

$$P(B) = P(A_1 \cap B) + \cdots + P(A_n \cap B) \\ = P(A_1) \cdot P(B|A_1) + \cdots + P(A_n) \cdot P(B|A_n)$$

Setting: 事件的生成由多个互斥的途径决定

Let A_1, \dots, A_n be partition of Ω , i.e. each possible outcome is included in one and only one of the events A_1, \dots, A_n

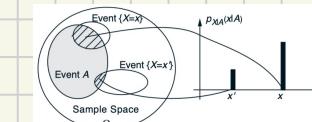


The probability that B occurs is a weighted average of its conditional probability under each scenario/event A_i , where each scenario is weighted according to its unconditional probability.

Conditional PMF

Conditioning a random variable X on an Event A.

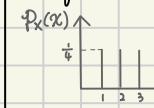
$$P_{X|A}(x) = P(X=x | A) = \frac{P(x=x \cap A)}{P(A)}$$



- Interpretation: 不同的 x 值对应的事件 $\{x=x\} \cap A$ 之间互斥，且这些事件的并集是 A，即 $P(A) = \sum_x P(\{x=x\} \cap A)$

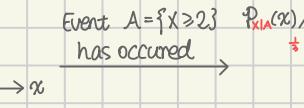
结合上面公式可得 $\sum_x P_{X|A}(x) = 1$

e.g.



$$E[X] = \frac{a+b}{2} = 2.5$$

$$\text{Var}(X) = \frac{(b-a)(b-a+2)}{12} = \frac{5}{4}$$

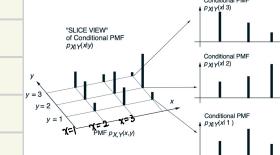


$$E[X|A] = \frac{a+b}{2} = 3$$

$$\text{Var}(X|A) = \frac{1}{3}(4-3)^2 + \frac{1}{3}(3-3)^2 + \frac{1}{3}(2-3)^2 = \frac{2}{3}$$

Conditioning a random variable X on another random variable Y.

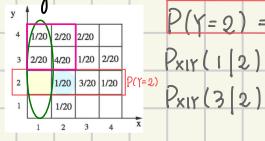
$$P_{X|Y}(x|y) = P(X=x | Y=y) = \frac{P(X=x, Y=y)}{P(Y=y)} = \frac{p_{XY}(x,y)}{p_Y(y)}$$



- Interpretation: 当固定 y 并考虑 $P_{X|Y}(x|y)$ 时，实际上是在考察一个以 y 为变量的函数，i.e. 固定 Y=y 的条件下，X 取每一个可能的值的概率。

因此，我们实际上有一个由所有可能的 y 值定义的 PMF family，每一个不同的 y 值都会定义一个不同的 $P_{X|Y}(x|y)$ 且 $\sum_x P_{X|Y}(x|y) = 1$ 。无论 y 取何值，X 的条件分布都是合法的概率分布。

e.g.



$$P(Y=2) = \frac{1}{20}$$

$$P_{X|Y}(1|2) = \frac{1}{20} \quad P_{X|Y}(2|2) = \frac{1}{20} \quad P_{X|Y}(3|2) = \frac{3}{20} \quad P_{X|Y}(4|2) = \frac{1}{5}$$

$$X, Y \text{ are independent?} - \text{NO. } P_X(1) = \frac{3}{20} \neq P_{X|Y}(1|2) = 0$$

what if we condition on $X \leq 2$ and $Y \geq 3$? — Yes.



Calculation Rules

I. Multiplication Rule

$$P_{X,Y,Z}(x,y,z) = P_X(x) \cdot P_{Y|X}(y|x) \cdot P_{Z|X,Y}(z|x,y)$$

II. Conditional Expectation / 条件期望

条件期望与普通期望的概念在本质上是相同的，计算方式都是加权随机变量的所有可能值。

但条件期望不是在原有的概率框架内进行计算，而是在一个通过条件改变了概率分布的新宇宙中进行。

$$E[X|A] = \sum_x x \cdot P_{X|A}(x|A) \quad E[X] = \sum_x x \cdot P_X(x) \quad E[g(x)] = \sum_x g(x) \cdot P_X(x)$$

$$E[X|Y=y] = \sum_x x \cdot P_{X|Y}(x|y)$$

$$E[g(x)|Y=y] = \sum_x g(x) \cdot P_{X|Y}(x|y)$$

$$E[g(x)|A] = \sum_x g(x) \cdot P_{X|A}(x|A)$$

III. Total Probability Theorem / 全概率定理

$$\text{let } B = \{X=x\} \quad P_X(x) = P(A_1) \cdot P_{X|A_1}(x) + \cdots + P(A_n) \cdot P_{X|A_n}(x)$$

IV. Total Expectation Theorem / 全期望值定理

The total probability theorem is true for all x's : $\sum_x x \cdot P_X(x) = \underbrace{\sum_{A_i} P(A_i) \cdot \sum_x x \cdot P_{X|A_i}(x)}_{E[X|A_i]} + \cdots + \sum_{A_n} P(A_n) \cdot \sum_x x \cdot P_{X|A_n}(x) = E[X|A]$

$$E[X] = \sum_{A_i} P(A_i) \cdot E[X|A_i]$$

$$= \sum_y P_Y(y) \cdot E[X|Y=y]$$

$$= E[E[X|Y]]$$

(5) Independence

① Independence of two Events A, B $P(A \cap B) = P(A) \cdot P(B)$

- $P(A|B)$ captures the partial information that event B provides about event A.

When the occurrence of B provides no information and does not alter the probability that A has occurred, i.e., $P(A|B) = P(A)$, we say that A is independent of B.

- By definition $P(A|B) = \frac{P(A \cap B)}{P(B)}$, in case of independent A and B, $P(A \cap B) = P(A) \cdot P(B)$ 我们采用此公式作为独立性的定义, 因为即使 $P(B) = 0$, 公式依然成立, 而 $P(A|B)$ is undefined.

- independent \neq disjoint

独立事件: 一个事件A的发生不影响另一个事件B的发生概率。独立事件可以联合发生, 其联合概率是各自概率的乘积。

Independence is not easily visualized in terms of the sample space.

不相交事件: 两个事件不能同时发生. $P(A \cap B) = 0$.

如果两个不相交的事件A, B都有正概率, 即 $P(A) > 0, P(B) > 0$, 那么它们绝对不可能是独立的.

② Conditional Independence of two Events A, B given an event C

- $P(A \cap B|C) = P(A|C) \cdot P(B|C)$ (*)

条件独立性表明, 在已知事件C发生的情况下, 事件A和B的发生相互独立。

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

$$\text{推导: } P(A \cap B|C) = \frac{P(A \cap B \cap C)}{P(C)} = \frac{P(C) \cdot P(B|C) \cdot P(A|B \cap C)}{P(C)} = P(B|C) \cdot P(A|B \cap C)$$

(Def. Conditional Probability)
(Multiplication Rule)

$$P(A|C) \cdot P(B|C) = P(B|C) \cdot P(A|B \cap C) \quad (*)$$

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

在事件C发生的情况下, 事件B的发生对事件A的发生概率没有额外的影响。

从信息论的角度来看, All the informations that B provides regarding the occurrence of A are completely encapsulated in C. 某些看似重要的条件实际上不提供额外信息。

- unconditionally independent \Rightarrow conditionally independent

2 independent fair coin tosses

$$H_1 = \{\text{1st toss is a head}\}$$

$$H_2 = \{\text{2nd toss is a head}\}$$

$$D = \{\text{the two tosses have different results}\}$$

- H_1 and H_2 are unconditionally independent

- But NOT conditionally independent:

$$P(H_1 \cap H_2 | D) = 0$$

$$P(H_1 | D) = P(H_2 | D) = \frac{1}{2}$$

$$P(H_1 \cap H_2 | D) \neq P(H_1 | D) \cdot P(H_2 | D)$$

若已知抛出了正一面, 那么 H_1 发生时必有 H_2 不发生

- conditionally independent \Rightarrow unconditionally independent

2 Coins ● $P(\text{Head}) = 0.99$ ● $P(\text{Head}) = 0.01$

Choose one of the two at random

Proceed with 2 independent coin tosses

$$H_i = \{i\text{-th toss is a head}\}$$

$$B = \{\text{the blue coin was selected}\}, P(B) = \frac{1}{2}$$

- H_1 and H_2 are conditionally independent:

选定硬币后, 两次抛掷都是独立的:

$$P(H_1 \cap H_2 | B) = P(H_1 | B) \cdot P(H_2 | B) = 0.99 \cdot 0.99$$

- But NOT unconditionally independent:

若没有给定硬币信息, 只知道 H_1 发生, 我们自然会倾向于认为选择的是蓝色硬币, 这将直接提高对 H_2 发生的预期。

$$P(H_1) = P(B) \cdot P(H_1 | B) + P(B^c) \cdot P(H_1 | B^c) = \frac{1}{2} \cdot 0.99 + \frac{1}{2} \cdot 0.01 = \frac{1}{2}$$

$$P(H_1 \cap H_2) = P(B) \cdot P(H_1 \cap H_2 | B) + P(B^c) \cdot P(H_1 \cap H_2 | B^c) = \frac{1}{2} \cdot 0.99^2 + \frac{1}{2} \cdot 0.01^2 = 0.4999$$

③ Independence of a Random Variable from an Event $P(X=x \text{ and } A) = P(X=x) \cdot P(A) = P_x(x) \cdot P(A)$, for all x

- Interpretation I

$$P_{x|A}(x) = P_x(x), \text{ for all } x \quad \text{无论 } A \text{ 是否发生, } X \text{ 的概率分布保持不变}$$

$$\text{Using the definition of the conditional PMF } P_{x|A}(x) = \frac{P(\{x\} \cap A)}{P(A)}$$

- Interpretation II

$$P(A) = P(A | X=x), \text{ for all } x \quad \text{对事件 } A \text{ 的概率判断不受随机变量 } X \text{ 取值的影响。}$$

④ Independence of Random Variables $P_{x,y}(x,y) = P_x(x) \cdot P_y(y)$, for all x, y

- Interpretation: $P_{x,y}(x,y) = P_x(x)$, for all y with $P_y(y) > 0$, and all x

- If X and Y are independent variables, then:

$$E[XY] = E[X] \cdot E[Y]$$

$$\text{Var}(X+Y) = \text{Var}(X) + \text{Var}(Y)$$

Proof: Expected value rule for functions of multiple random variables

$$E[g(X,Y)] = \sum_x \sum_y g(x,y) \cdot P_{x,y}(x,y) \quad \text{Let } g(x,y) = XY$$

$$E[XY] = \sum_x \sum_y x \cdot y \cdot P_x(x) \cdot P_y(y) \quad \text{By independence}$$

$$= \left(\sum_x x \cdot P_x(x) \right) \left(\sum_y y \cdot P_y(y) \right)$$

$$E[g(X) \cdot h(Y)] = E[g(X)] \cdot E[h(Y)]$$

Proof: Let $E[X] := a$ $E[Y] := b$ $E[X+Y] = a+b$

$$\begin{aligned} \text{Var}(X+Y) &= E[(X+E[X])(Y+E[Y])] \\ &= E[(X-a)+(Y-b)]^2 \\ &= E[(X-a)^2 + 2 \cdot (X-a) \cdot (Y-b) + (Y-b)^2] \\ &= E[(X-a)^2] + 2 \cdot E[(X-a) \cdot (Y-b)] + E[(Y-b)^2] \\ &= \text{Var}(X) + 2 \cdot E[X-a] \cdot E[Y-b] + \text{Var}(Y) \\ &= \text{Var}(X) + \text{Var}(Y) \end{aligned}$$

Covariance: $E[(X - E[X])(Y - E[Y])] = 0$

$$\text{Proof: } = E[(X-E[X]) \cdot E[(Y-E[Y])]] = 0$$

- In general: $\text{Var}(X+Y) \neq \text{Var}(X) + \text{Var}(Y)$

$$\text{e.g. } X = Y \quad \text{Var}(X+Y) = \text{Var}(2X) = 4 \cdot \text{Var}(X)$$

$$X = -Y \quad \text{Var}(X+Y) = \text{Var}(0) = 0$$

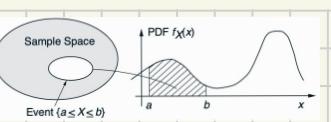
In general $E[X+Y] = E[X] + E[Y]$

⑤ Conditional Independence of two Random Variables X, Y given an event A

$$P_{x,y|A}(x,y) = P_{x|A}(x) \cdot P_{y|A}(y), \text{ for all } x, y$$

Again, equivalent to $\underbrace{P_{x,y|A}(x|y)}_{\{Y=y\} \text{ 和事件 } A \text{ 同时发生的条件下, } X \text{ 取值为 } x \text{ 的概率}} = P_{x|A}(x) \text{ for all } x \text{ and } y \text{ s.t. } P_{y|A}(y) > 0$

Continuous Random Variables 连续随机变量



1. Def.

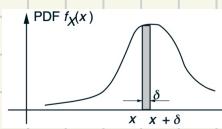
A random variable X is continuous if its possible values comprise either a single interval or a union of disjoint intervals on the number line.

A random variable X is continuous if its probability law can be described by a PDF.

2. Framework to describe continuous random variable

(1) PDF / Probability Density Function / 概率密度函数

• Intuition



We can think of the bars in the PMF as point masses with positive weight, each sitting on top of a specific numerical value. To calculate the probability that the random variable falls within a specific interval, we add up all the masses that sit on top of the numerical values within that interval.

In the case of PDF, there's still a total of one unit of probability mass assigned to the possible values of the random variable. However, this mass is spread across the entire real line, rather than being concentrated at specific points. The distribution of this mass is not uniform: some parts of the real line contain more mass per unit length than others. 我们可以想象有1磅雪落在314数线上。PDF告诉我们有峰值点上积累的雪的高度，然后我们通过计算曲线下面积来得到在某个区间上积累的整体雪量的重量。

• PDF Properties

$$\textcircled{1} \quad P(a \leq X \leq b) = \int_a^b f_X(x) dx \quad [a, b] \subseteq \Omega \quad \text{PDF 用于确定小区间的概率}$$

$$\textcircled{2} \quad f_X(x) \geq 0, \text{ for every } x$$

$$\textcircled{3} \quad \int_{-\infty}^{+\infty} f_X(x) dx = P(-\infty < X < +\infty) = 1 \quad \text{The entire area under the graph of the PDF must be equal to 1.}$$

$$\textcircled{4} \quad P(X=a) = 0 \quad \text{在连续取值范围内选取一点的概率极小.}$$

◦ For an interval $[x, x+\Delta]$ with very small length Δ , we have: 由此可以看出, PDF $f_X(x)$ 的单位是

$$P([x, x+\Delta]) = \int_x^{x+\Delta} f_X(t) dt \approx f_X(x) \cdot \Delta \Rightarrow f_X(x) \approx \frac{P([x, x+\Delta])}{\Delta} \quad \text{"Probability mass per unit length"}$$

连续随机变量 X 取值在某个非常小区间内的概率 \approx PDF在该点的取值 \times 小区间的长度. 随着这个区间长度趋向于0,

$$P([x, x+\Delta]) \rightarrow 0 \Rightarrow P(X=a) = 0$$

PDF在某一点 x 处的值 是连续随机变量 X 在点 x 附近取值概率密度的度量, 而不是概率本身.

◦ 或代入公式理解: $P(a \leq X \leq a) = \int_a^a f_X(x) dx = 0$

$$\textcircled{5} \quad P(a < X < b) = P(a \leq X < b) = P(a < X \leq b) = P(a \leq X \leq b)$$

◦ PDFs don't have to be continuous functions.

- 一旦我们得到了PDF, 可以不再关注其样本空间的构成. PDF为我们提供了计算和推导连续随机变量行为的所有必要信息
- 之后我们默认 $\Omega = \mathbb{R}$. 若 $\Omega \neq \mathbb{R}$, 我们可以通过将不属于 Ω 的 X 值在 \mathbb{R} 上的概率密度定义为0来将模型扩展到整个 \mathbb{R} .

(2) CDF / Cumulative Density Function / (累积)分布函数

• Motivation

一些性质在离散随机变量和连续随机变量之间是共通的. CDF提供了一个统一的框架来表示随机变量的分布情况: $\{X \leq x\}$ is always an event and therefore has a well-defined probability.

CDF是一个单变量函数. 它给出随机变量取≤某个具体值 x 的概率.

Loosely speaking, the CDF $F_X(x)$ "accumulates" probability "up to" the value x . $F_X(x) = P(X \leq x)$

• Def. and Properties

Discrete Random Variable X	Continuous Random Variable X
$F_X(x) = \sum_{k \leq x} P_X(k)$ <p>$F_X(x)$ has a piecewise constant and staircase-like form.</p> $P_X(k) = P(X \leq k) - P(X \leq k-1) = F_X(k) - F_X(k-1), \quad \forall k \in \mathbb{Z}$ <p>每当离散随机变量 X 的PMF在某个值 k 上有正值 $p_X(k)$,</p> <p>$F_X(k)$ 在 k 处会有一个跳跃 其大小等于 $p_X(k)$</p>	$F_X(x) = \int_{-\infty}^x f_X(t) dt$ <ul style="list-style-type: none"> ◦ $F_X(x)$ has a continuously varying form. ◦ The 2-step procedure to find a general PDF $f_X(y) = g(X)$ of a continuous random variable: <ol style="list-style-type: none"> Find the CDF of Y: $F_Y(y) = P(Y \leq y)$ Differentiate: $f_Y(y) = \frac{dF_Y(y)}{dy}$ <p>$F_X(x)$ is monotonically non-decreasing : $x \leq y \Rightarrow F_X(x) \leq F_X(y)$</p> <p>$x \rightarrow \infty, F_X(x) \rightarrow 1$ $x \rightarrow -\infty, F_X(x) \rightarrow 0$</p>

• PDF for Functions of Random Variables

$y = ax + b$ $f_Y(y) = \frac{1}{ a } f_X(\frac{y-b}{a})$ st. PDF integrates to 1	$y = g(x)$, g is strictly monotonic and differentiable $f_Y(y) = f_X(g(y)) \left \frac{dg(y)}{dy} \right $
$(a > 0)$ $F_Y(y) = P(Y \leq y)$ $= P(ax+b \leq y)$ $= P(X \leq \frac{y-b}{a})$ $= F_X(\frac{y-b}{a})$ $f_Y(y) = f_X(\frac{y-b}{a}) \cdot \frac{1}{a}$	$(a < 0)$ $F_Y(y) = P(Y \leq y)$ $= P(ax+b \leq y)$ $= P(X > \frac{y-b}{a})$ $= 1 - P(X \leq \frac{y-b}{a})$ $= 1 - F_X(\frac{y-b}{a})$ $f_Y(y) = -f_X(\frac{y-b}{a}) \cdot \frac{1}{a}$

An important fact is that a monotonic function can be "inverted" in the sense that there is some function h , called the inverse of g , such that for all $x \in I$, we have $y = g(x)$ if and only if $x = h(y)$.

g monotonically increasing g monotonically decreasing

$F_Y(y) = P(g(x) \leq y)$ $= P(X \leq h(y))$ $= F_X(h(y))$ $f_Y(y) = f_X(h(y)) \cdot \frac{dh(y)}{dy}$	$F_Y(y) = P(g(x) \leq y)$ $= P(X \geq h(y))$ $= 1 - F_X(h(y))$ $f_Y(y) = -f_X(h(y)) \cdot \frac{dh(y)}{dy}$
---	--

(3) Joint PDF / 联合概率质量函数

• Def

Two continuous random variables associated with a common experiment are jointly continuous (联合连续) if they can be described by a joint PDF $f_{x,y}$ that satisfies:

$$\begin{cases} f_{x,y}(x,y) \geq 0 \\ P((x,y) \in B) = \iint_{(x,y) \in B} f_{x,y}(x,y) dx dy \end{cases} \quad \text{for every subset } B \text{ of the two-dimensional plane.}$$

• Interpretation

我们不直接定义联合PDF为某个具体的概率值，而是通过定义联合PDF的功能来间接定义其概念，i.e. 联合PDF被用以计算概率，本身不表示概率。

具体来说，联合PDF描述两个连续随机变量 X, Y 落在任意给定的二维区域 B 内的概率。

$f_{x,y}(x,y)$ 在点 (x,y) 的值还可以想像为点 (x,y) 上方的高度，这个高度代表了随机变量 (X,Y) 落在此点的概率密度。

若在点 (x,y) 附近有一个非常小的矩形区域 $\Delta x \times \Delta y$ ，那么 X 和 Y 同时落在这个小矩形内的概率约为

$$f_{x,y}(x,y) \cdot \Delta x \cdot \Delta y$$

将这种几何理解扩展到二维区域 B 时，计算 $P((X,Y) \in B)$ 的过程本质上是计算由 $f_{x,y}(x,y)$ 形成的曲面与区域 B 所围成的体积。

这个体积代表了在区域 B 上对 $f_{x,y}(x,y)$ 和 $dx \cdot dy$ 的乘积进行整体积分。

• Properties

$$\textcircled{1} \quad P(a \leq X \leq b, c \leq Y \leq d) = \int_c^d \int_a^b f_{x,y}(x,y) dx dy$$

B is a rectangular of the form $B = [a, b] \times [c, d]$

$$\textcircled{2} \quad \int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} f_{x,y}(x,y) dx dy = 1$$

$$\textcircled{3} \quad \text{area}(B) = 0 \Rightarrow P((X,Y) \in B) = 0$$

e.g. B 是一条曲线

如果矩形的边长都是某个很小的数值 δ :

$$P(a \leq X \leq a+\delta, c \leq Y \leq c+\delta) = \int_c^{c+\delta} \int_a^{a+\delta} f_{x,y}(x,y) dx dy \approx f_{x,y}(a,c) \cdot \delta^2$$

$f_{x,y}(a,c) \approx \frac{P(a \leq X \leq a+\delta, c \leq Y \leq c+\delta)}{\delta^2}$ 由此可以看出, $f_{x,y}(a,c)$ 的单位是 "Probability per unit area in the vicinity of (a,c) "

④ X, Y are continuous $\Rightarrow X$ and Y are jointly continuous.

e.g. $X=Y$

可以确定实验的结果将会落在 $X=Y$ 这条线上, i.e. 所有的概率都集中在这条线上，而这条线的面积为0。

我们在面积为0的集合上得到3正确的概率，看似与性质③相矛盾，但这实际上意味着 X 和 Y 并不是联合连续的。

Joint continuity requires more than individual continuity of the random variables involved; it requires a genuine two-dimensional spread of probability.

• From the joint to the marginals

$$f_x(x) = \int f_{x,y}(x,y) dy$$

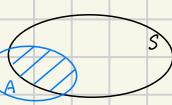
$$f_y(y) = \int f_{x,y}(x,y) dx$$

e.g.

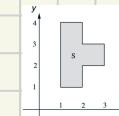
$$(i) \quad f_{x,y}(x,y) = \begin{cases} \frac{1}{\text{area of } S} & \text{if } (x,y) \in S, \\ 0 & \text{otherwise.} \end{cases}$$

For any set $A \subset S$, the probability that the experimental value of (X,Y) lies in A is

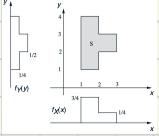
$$P((X,Y) \in A) = \iint_{(x,y) \in A} f_{x,y}(x,y) dx dy = \frac{1}{\text{area of } S} \iint_{(x,y) \in A \cap S} dx dy = \frac{\text{area of } A \cap S}{\text{area of } S}.$$



(ii) We are told that the joint PDF of the random variables X and Y is a constant C on the set S shown in Fig. and is zero outside. Find the value of C and the marginal PDFs of X and Y .



The area of the set S is equal to 4 and, therefore, $f_{x,y}(x,y) = C = 1/4$, for $(x,y) \in S$. To find the marginal PDF $f_x(x)$ for some particular x , we integrate (with respect to y) the joint PDF over the vertical line corresponding to that x . The resulting PDF is shown in the figure. We can compute f_y similarly.



• Expected Value Rule for Functions of Multiple Random Variables $E[g(x,y)] = \int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} g(x,y) \cdot f_{x,y}(x,y) dx dy$

• $E[aX+bY+c] = aE[X] + bE[Y] + c$

(4) Joint CDF / 联合分布函数

- Def

If X and Y are two continuous random variables associated with a common experiment, we define their joint CDF by:

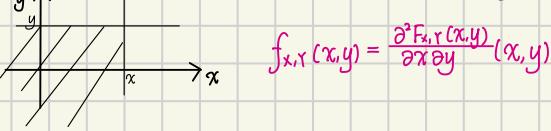
$$F_{X,Y}(x,y) = P(X \leq x, Y \leq y) = \int_{-\infty}^y \int_{-\infty}^x f_{X,Y}(s,t) ds dt$$

PDF \rightarrow CDF

- CDF \rightarrow PDF

$$F_{X,Y}(x,y) \text{ 对 } y \text{ 偏导: } \frac{\partial}{\partial y} F_{X,Y}(x,y) = \int_{-\infty}^x f_{X,Y}(s,y) ds$$

$$\text{对第一步的结果再对 } x \text{ 偏导: } \frac{\partial}{\partial x} (\int_{-\infty}^x f_{X,Y}(s,y) ds) = f_{X,Y}(x,y)$$



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

$X \leq x$ 且 $Y=y$ 条件下的 marginal PDF

X 和 Y 在点 (x,y) 处的 joint PDF

(5) Mean and Variance

$$\mathbb{E}[X] = \int_{-\infty}^{+\infty} x f_X(x) dx$$

离散型随机变量 X : $\mathbb{E}[X] = \sum_x x \cdot p_x(x)$ (An integral is just a limiting form of a sum)

Can be interpreted as: Average in large number of independent repetitions of the experiment.
& "Center of gravity" of the probability

$$\begin{aligned} \text{Var}(X) &= \mathbb{E}[(X - \mu)^2] = \int_{-\infty}^{+\infty} (x - \mu)^2 f_X(x) dx \\ &= \mathbb{E}[X^2] - (\mathbb{E}[X])^2 \end{aligned}$$

- Linearity of Variance $\text{Var}(aX + b) = a^2 \text{Var}(X)$

(6) Conditional PDF / 条件概率密度函数

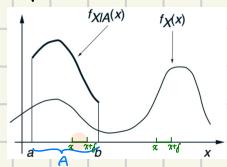
- Conditioning a random variable X on an Event A

$$P(X \in B | A) = \int_B f_{X|A}(x) dx, P(A) > 0, \text{ for any subset } B \text{ of the real line}$$

- Intuition

	PMF	PDF
Ordinary Probability Definition Calculation	$P_x(x) = P(X=x)$ $P(X \in B) = \sum_{x \in B} P_x(x)$	$f_x(x) \cdot f \approx P(x \leq X \leq x+f)$ $P(X \in B) = \int_B f_x(x) dx$
Conditional Probability Definition Calculation	$P_{X A}(x) = P(X=x A)$ $P(X \in B A) = \sum_{x \in B} P_{X A}(x)$	$f_{X A}(x) \cdot f \approx P(x \leq X \leq x+f A)$ $P(X \in B A) = \int_B f_{X A}(x) dx$

- Special Case: Conditional PDF of X , given that $X \in A$ A is a subset of the real line.



$$\begin{aligned} P(x \leq X \leq x+f | X \in A) &\approx f_{X|A}(x) \cdot f \quad (*) \\ &= \frac{P(x \leq X \leq x+f, X \in A)}{P(X \in A)} \\ &= \frac{P(x \leq X \leq x+f)}{P(X \in A)} \approx \frac{f_x(x) \cdot f}{P(X \in A)} \quad (**) \end{aligned}$$

$f_{X|A}(x)$ has exactly the same shape as the unconditional one, except that it's scaled by the constant factor $\frac{1}{P(X \in A)}$

- Calculation Rules

$$\text{I. Total Probability Theorem / 全概率定理} \quad f_X(x) = \sum_{i=1}^n P(A_i) \cdot f_{X|A_i}(x)$$

$$\text{验证: } P(X \leq x) = P(A_1) \cdot P(X \leq x | A_1) + \dots + P(A_n) \cdot P(X \leq x | A_n)$$

$$\Rightarrow F_X(x) = P(A_1) \cdot F_{X|A_1}(x) + \dots + P(A_n) \cdot F_{X|A_n}(x) \quad \text{Take the derivatives of both sides w.r.t. } x$$

$$\Rightarrow f_X(x) = P(A_1) \cdot f_{X|A_1}(x) + \dots + P(A_n) \cdot f_{X|A_n}(x)$$

$$\text{II. Total Expectation Theorem / 全期望值定理} \quad E[X] = \sum_{i=1}^n P(A_i) \cdot E[X | A_i]$$

$$\text{验证: Multiply both sides of } f_X(x) = \sum_{i=1}^n P(A_i) \cdot f_{X|A_i}(x) \text{ by } x \text{ and then integrate from } -\infty \text{ to } +\infty.$$

Application: mixed distributions

Conditioning a random variable X on another $f_{X|Y}(x|y) = \frac{f_{X,Y}(x,y)}{f_Y(y)}$ if $f_Y(y) > 0$

$$\text{推导: } P(X \in A | Y=y) = \int_A f_{X|Y}(x|y) dx$$

$$f_{X|A}(x) \cdot f \approx P(x \leq X \leq x+f | A) \quad \text{where } P(A) > 0 \quad (\text{Conditional Probability Definition})$$

在事件 A 发生的条件下, 小区间 $[x, x+f]$ 的条件概率

我们不使用 $A = \{Y=y\}$, 因 PDF 用于确定小区间内的概率, $P(Y=y) = 0$

所以定义 $A = \{Y \approx y\}$, i.e., $y \leq Y \leq y+\epsilon$

$$\begin{aligned} P(x \leq X \leq x+f | y \leq Y \leq y+\epsilon) &= \frac{P(x \leq X \leq x+f, y \leq Y \leq y+\epsilon)}{P(y \leq Y \leq y+\epsilon)} \\ &= \frac{P(x \leq X \leq x+f)}{P(y \leq Y \leq y+\epsilon)} \end{aligned}$$

$$(f_{X|Y}(x|y) = \frac{f_{X,Y}(x,y)}{f_Y(y)})$$

$f_{X|Y}(x|y) \cdot f$ 表示 Y 落在小区间 $[y, y+\epsilon]$ 时, X 落在小区间 $[x, x+f]$ 的概率.

由于 $f_{X|Y}(x|y) \cdot f$ 不依赖于 ϵ , 可以考虑当 ϵ 趋近于 0 的极限情况, 将到 $P(x \leq X \leq x+f | Y=y) \approx f_{X|Y}(x|y) \cdot f$

$$\text{More generally: } P(X \in A | Y=y) = \int_A f_{X|Y}(x|y) dx$$

我们之前无法定义布给定 0 概率事件 $\{Y=y\}$ 下的条件概率, 但现在可以通过公式 $P(X \in A | Y=y) = \int_A f_{X|Y}(x|y) dx$ 处理这种情况.

In addition, it allows us to view the conditional PDF $f_{X|Y}(x|y)$ (as a function of x) as a description of the probability law of X , given that the event $\{Y=y\}$ has occurred.

- Property

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

当我们把 y 视为 a fixed number, 并将 $f_{X|Y}(x|y)$ 视为单变量 x 的函数: 条件 PDF $f_{X|Y}(x|y)$ 和联合 PDF $f_{X,Y}(x,y)$ 具有相同的形式, 因为 normalizing factor $f_Y(y)$ 不依赖于 x .

- Calculation Rules

$$\text{I. Multiplication Rule} \quad f_{X,Y}(x,y) = f_Y(y) \cdot f_{X|Y}(x|y) = f_X(x) \cdot f_{Y|X}(y|x)$$

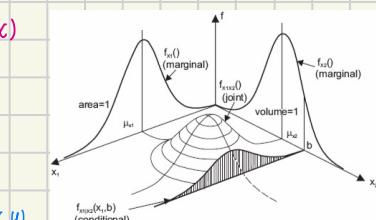
$$\text{II. Conditional Expectation / 条件期望}$$

$$E[X | A] = \int_{-\infty}^{+\infty} x \cdot f_{X|A}(x) dx$$

$$E[X | Y=y] = \int_{-\infty}^{+\infty} x \cdot f_{X|Y}(x|y) dx$$

$$E[g(x) | Y=y] = \int_{-\infty}^{+\infty} g(x) \cdot f_{X|Y}(x|y) dx$$

$$E[g(x) | A] = \int_{-\infty}^{+\infty} g(x) \cdot f_{X|A}(x) dx$$



$$\text{III. Total Probability Theorem / 全概率定理} \quad f_X(x) = \int_{-\infty}^{+\infty} f_Y(y) \cdot f_{X|Y}(x|y) dy$$

$$\text{IV. Total Expectation Theorem / 全期望值定理} \quad E[X] = \int_{-\infty}^{+\infty} f_Y(y) \cdot E[X | Y=y] dy$$

$$\text{验证: } \int_{-\infty}^{+\infty} f_Y(y) \cdot E[X | Y=y] dy \quad \text{代入 } E[X | Y=y] = \int_{-\infty}^{+\infty} x \cdot f_{X|Y}(x|y) dx$$

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

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

$$= \int_{-\infty}^{+\infty} y \int_{-\infty}^{+\infty} f_Y(y) \cdot f_{X|Y}(x|y) dx dy \quad (\text{全概率定理})$$

$$= \int_{-\infty}^{+\infty} x f_X(x) dx$$

$$= E[X]$$

Continuous Uniform Random Variable

Parameters: a, b $a \leq b$

$$f_x(x) = \begin{cases} \frac{1}{b-a}, & \text{if } a \leq x \leq b \\ 0, & \text{otherwise} \end{cases}$$



$$F_x(x) = \begin{cases} 0, & \text{if } x < a \\ \frac{x-a}{b-a}, & \text{if } a \leq x \leq b \\ 1, & \text{if } x > b \end{cases}$$

- $P(X \in I) = \int_{[a,b] \cap I} \frac{1}{b-a} dx = \frac{1}{b-a} \int_{[a,b] \cap I} dx$
 $= \frac{\text{length of } [a,b] \cap I}{\text{length of } I}$

- 期望 $E[X] = \int_{-\infty}^{+\infty} x f_x(x) dx = \int_a^b x \cdot \frac{1}{b-a} dx = \frac{a+b}{2}$

- 方差 $\text{var}(X) = E[X^2] - (E[X])^2 = \frac{(b-a)^2}{12}$

$$E[X^2] = \int_{-\infty}^{+\infty} x^2 f_x(x) dx = \int_a^b x^2 \cdot \frac{1}{b-a} dx = \left(\frac{1}{b-a}\right) \left(\frac{b^3}{3} - \frac{a^3}{3}\right)$$

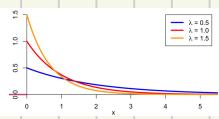
- 标准差 $\sigma = \frac{b-a}{\sqrt{12}}$

$\frac{b-a}{\sqrt{12}}$ is proportional to the width of the uniform distribution. This aligns with the intuition that the standard deviation captures the "width" or dispersion of a particular distribution.

Exponential Random Variable

Parameter: $\lambda > 0$

$$f_x(x) = \begin{cases} \lambda e^{-\lambda x}, & \text{if } x \geq 0 \\ 0, & \text{otherwise} \end{cases}$$



$$F_x(x) = \begin{cases} 1 - e^{-\lambda x}, & \text{if } x \geq 0 \\ 0, & \text{otherwise} \end{cases}$$

$$\begin{aligned} P(X > a) &= \int_a^{\infty} \lambda e^{-\lambda x} dx \\ &= \lambda \int_a^{\infty} e^{-\lambda x} dx \\ &= \lambda \cdot \left(-\frac{1}{\lambda}\right) e^{-\lambda x} \Big|_a^{\infty} \\ &= -e^{-\lambda a} - (-e^{-\lambda a}) \\ &= e^{-\lambda a} \end{aligned}$$

(recall: $\int e^{\alpha x} dx = \frac{1}{\alpha} e^{\alpha x}$)
let $a = -\lambda$

Observation: Let $a=0$, we obtain the integral of the PDF over the entire range of X .

In this case, $P(X > 0) = 1$. Thus, we have verified that the integral of PDF is 1, as it should be.

- 期望 $E[X] = \int_0^{\infty} x \cdot \lambda e^{-\lambda x} dx = \frac{1}{\lambda}$

- 方差 $\text{var}(X) = E[X^2] - (E[X])^2 = \frac{1}{\lambda^2}$

$$E[X^2] = \int_0^{\infty} x^2 \cdot \lambda e^{-\lambda x} dx = \frac{2}{\lambda^2}$$

Observation:

当 λ 较小时, PDF 下降较慢, 这意味着 X 取较大值的概率增加, 因此 X 的均值会更高. 较大值的分布跨越范围更广, 所以方差变大.

Application

Model the amount of time until a piece of equipment breaks down, until a light bulb burns out, or until an accident occurs.

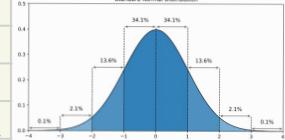
Standard Normal Random Variable

Parameters: $N(0, 1)$

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

recall: $\int_{-\infty}^{+\infty} e^{-x^2/2} dx = \sqrt{2\pi}$

$F_x(x)$: Standard normal table



- 期望 $E[X] = 0$

X^2 is symmetric around 0
 \Rightarrow PDF $f_x(x)$ is symmetric around 0.

- 方差 $\text{var}(X) = 1$

$$\begin{aligned} E[X^2] &= \int_{-\infty}^{\infty} x^2 f_x(x) dx \\ &= (2\pi)^{-1/2} \int_{-\infty}^{\infty} x^2 \exp(-\frac{1}{2}x^2) dx \\ &= (2\pi)^{-1/2} \left\{ \int_{-\infty}^0 x \left(x \exp(-\frac{1}{2}x^2)\right) dx + \int_0^{\infty} x \left(x \exp(-\frac{1}{2}x^2)\right) dx \right\} \\ &= (2\pi)^{-1/2} \left\{ \left[-x \exp(-\frac{1}{2}x^2)\right]_0^0 + \int_0^{\infty} \exp(-\frac{1}{2}x^2) dx + \left[-x \exp(-\frac{1}{2}x^2)\right]_0^{\infty} \right. \\ &\quad \left. + \int_0^{\infty} \exp(-\frac{1}{2}x^2) dx \right\} \quad (\text{integrating by parts}) \\ &= (2\pi)^{-1/2} \left\{ (0-0) + (0-0) + \int_0^{\infty} \exp(-\frac{1}{2}x^2) dx + \int_0^{\infty} \exp(-\frac{1}{2}x^2) dx \right\} \\ &= (2\pi)^{-1/2} \int_{-\infty}^{\infty} \exp(-\frac{1}{2}x^2) dx \\ &= \int_{-\infty}^{\infty} f_x(x) dx = 1 \quad (\text{the integral of a pdf over its support is equal to 1}) \end{aligned}$$

$$E[X]^2 = 0^2 = 0$$

$$\text{Var}[X] = E[X^2] - E[X]^2 = 1 - 0 = 1$$

Application

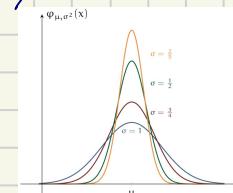
• It models well the additive effect of many independent factors in a variety of contexts. It's important in the theory of probability, e.g. Central Limit Theorem.

Normal Random Variable

Parameters: $N(\mu, \sigma^2)$

$\mu \in \mathbb{R}, \sigma > 0$

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



$F_x(x)$: Standard normal table

$$\begin{aligned} P(X \leq x) &= P\left(\frac{X-\mu}{\sigma} \leq \frac{x-\mu}{\sigma}\right) \\ &= P(Y \leq \frac{x-\mu}{\sigma}) \\ &= \Phi\left(\frac{x-\mu}{\sigma}\right) \end{aligned}$$

Y 表示 X 与均值 μ 之间的偏差. 除以标准差后, Y 以标准差为单位来量化这个偏差. 若 $Y=3$, 则意味着 X 与均值之间有3个标准差的距离.

Y 是一个标准正态随机变量:

$$E[Y] = 0 \quad (\because X \sim \mu)$$

$$\text{Var}(Y) = \frac{1}{\sigma^2} \cdot \text{Var}(X) = 1$$

- 期望 $E[X] = \mu$

$e^{-\frac{(x-\mu)^2}{2\sigma^2}}$ is symmetric around μ

- 方差 $\text{var}(X) = \sigma^2$

$$\text{var}(X) = \frac{1}{\sqrt{2\pi}\sigma} \int_{-\infty}^{\infty} (x-\mu)^2 e^{-(x-\mu)^2/2\sigma^2} dx.$$

Using the change of variables $y = (x-\mu)/\sigma$ and integration by parts, we have

$$\begin{aligned} \text{var}(X) &= \frac{\sigma^2}{\sqrt{2\pi}\sigma} \int_{-\infty}^{\infty} y^2 e^{-y^2/2} dy \\ &= \frac{\sigma^2}{\sqrt{2\pi}\sigma} \left[y \exp(-y^2/2) \right]_{-\infty}^{\infty} + \frac{\sigma^2}{\sqrt{2\pi}\sigma} \int_{-\infty}^{\infty} \exp(-y^2/2) dy \\ &= \frac{\sigma^2}{\sqrt{2\pi}\sigma} \int_{-\infty}^{\infty} e^{-y^2/2} dy \\ &= \sigma^2 \end{aligned}$$

The last equality above is obtained by using the fact

$$\frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{-y^2/2} dy = 1,$$

which is just the normalization property of the normal PDF for the case where $\mu = 0$ and $\sigma = 1$.

- Normality is preserved by linear transformations

Let $X \sim N(\mu, \sigma^2)$ $Y = ax + b$

$$E[Y] = a \cdot E[X] + b = a\mu + b$$

$$\text{Var}(Y) = a^2 \cdot \text{Var}(X) = a^2 \sigma^2$$

$$\text{或利用之前得出的结论 } f_Y(y) = \text{tai} f_X(\frac{y-b}{a}) = \text{tai} \left(\frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(y-b-a\mu)^2}{2\sigma^2}} \right) = \frac{1}{\sqrt{2\pi}\sigma} \left(\frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(y-b-a\mu)^2}{2\sigma^2}} \right)^2 = \frac{1}{\sqrt{2\pi}\sigma} \cdot \frac{1}{\sqrt{2\pi}\sigma} \cdot e^{-\frac{(y-b-a\mu)^2}{2\sigma^2}} = \frac{1}{\sqrt{2\pi}\sigma} \cdot \frac{1}{\sqrt{2\pi}\sigma} \cdot e^{-\frac{(y-b-a\mu)^2}{2\sigma^2}} = \frac{1}{\sqrt{2\pi}\sigma} \cdot \frac{1}{\sqrt{2\pi}\sigma} \cdot e^{-\frac{(y-b-a\mu)^2}{2\sigma^2}}$$

$$Y \sim N(a\mu + b, a^2\sigma^2)$$

若 $a=0$, 则 $Y=b$, 即 Y 退化为一个常数, 但我们仍可使用正态分布的表示方法: $Y \sim (b, 0)$

Application

- In Signal processing and communication engineering, model noise and unpredictable distortions of signals