

## Lecture 1 – Review of Probability

2021 年 2 月 22 日

Lecturer: 张驰豪

Scribe: 杨宽

## 课程信息

课程大纲:

- 马尔可夫链 Markov Chains: discrete / continuous
- 泊松过程 Poisson Process
- 鞅 Martingale
- 布朗运动 Brownian Motion
- AI / 大数据 / 机器学习中的算法应用

参考资料:

- Richard Durrett, *Essentials of Stochastic Processes*
- Sheldon M. Ross, *Introduction to Probability Models*
- <http://www.stat.yale.edu/~pollard/Courses/251.spring2013/Handouts/Chang-notes.pdf>

## 1 Probability and Random Variable

**Definition 1** (Probability space). A *probability space* consists of a 3-ary tuple  $(\Omega, \mathcal{F}, \Pr[\cdot])$ :

- $\Omega$  is a set of “outcomes” (countable or uncountable);
- $\mathcal{F} \subseteq 2^\Omega$  is a  $\sigma$ -algebra (a set of all possible “events”) on which we can define probability, and here we say  $\mathcal{F}$  is a  $\sigma$ -algebra if  $\mathcal{F}$  satisfies
  - $\emptyset \in \mathcal{F}$ ,
  - $A \in \mathcal{F} \implies A^c \in \mathcal{F}$ ,
  - a countable sequence of sets  $A_1, \dots, A_n, \dots \in \mathcal{F} \implies \bigcup_{i=1}^{\infty} A_i \in \mathcal{F}$ ;
- Probability  $\Pr[\cdot] : \mathcal{F} \rightarrow [0, 1]$  is a function s.t.

1.  $\Pr[\emptyset] = 0$ ,
2.  $\Pr[\Omega] = 1$ ,
3. if  $A_1, \dots, A_n, \dots \in \mathcal{F}$  are disjoint, then  $\Pr[\cup_{i=1}^{\infty} A_i] = \sum_{i=1}^{\infty} \Pr[A_i]$ .

**Example 2** (6-face dice).  $\Omega = [6] = \{1, 2, 3, 4, 5, 6\}$ ,  $\mathcal{F} = 2^{[6]}$ ,  $\Pr[i] = 1/6$ .

Generally, in a *discrete space*,  $\Omega$  is countable. Define  $\mathcal{F} = 2^{\Omega}$  and  $\hat{p}: \Omega \rightarrow [0, 1]$  s.t.  $\sum_{\omega \in \Omega} \hat{p}(\omega) = 1$ , then

$$\forall A \in \mathcal{F}, \quad \Pr[A] \triangleq \sum_{\omega \in A} \hat{p}(\omega).$$

**Question.** How to define a probability on  $\mathbb{R}$ ?

Or, what do we mean by drawing a uniform real in  $(0, 1)$ ?

**Example 3** (Uniform real in  $(0, 1)$ ). Define the probability of uniformly drawing real numbers as follows:

- $\Omega = (0, 1)$ ;
- $\mathcal{F}$  is the  $\sigma$ -algebra consisting of all “Borel sets” on  $(0, 1)$ , namely the collection of subsets of  $(0, 1)$  obtained from all open intervals by repeatedly taking *countable* unions and complements;
- $\forall$  interval  $I = (a, b)$ ,  $\Pr[I] = (b - a)$ . (*Lebesgue measure*)

*Remark.*  $\mathcal{F}$  is called the Borel algebra, which is the smallest  $\sigma$ -algebra containing all open intervals. All Borel sets are measurable. The existence of non-Borel set is independent of ZF (Zermelo-Fraenkel set theory).

**Definition 4** (Random variable). A *random variable* is a function or mapping from the probability space to a field. Given  $(\Omega, \mathcal{F}, \Pr[\cdot])$ , a real-valued random variable is a function of  $\Omega$ :

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

**Definition 5** (Distribution (discrete)). For a countable  $\Omega$  and a random variable  $X$ , the distribution of  $X$  is given by

$$\forall a \in \text{Range}(X), \quad \mu(a) = \Pr[X = a] \triangleq \Pr[\{\omega : X(\omega) = a\}] = \Pr[X^{-1}(a)].$$

**Example 6** (Binomial distribution). Toss a *biased* coin (Head with probability  $p$  and Tail with probability  $1 - p$ )  $n$  times. Let  $X$  be the number of Heads, then the distribution of  $X$  is called the *binomial distribution*:

$$\Omega = \{0, 1\}^n, \quad X \sim \text{Binom}(n, p) \iff \Pr[X = k] = \binom{n}{k} \cdot p^k \cdot (1 - p)^{n-k}.$$

**Example 7** (Geometric distribution). Toss a biased coin. Let  $X$  be the number of trials until the first Head, then the distribution of  $X$  is called the *geometric distribution*:

$$\Omega = \{0, 1\}^*, \quad X \sim \text{Geometric}(n, p) \iff \Pr[X = k] = (1 - p)^{k-1} p.$$

**Definition 8** (Distribution (continuous)). For an uncountable  $\Omega$  and a random dx variable  $X$ , if there exists a nonnegative function  $f(x)$  s.t.

$$\Pr[a \leq X \leq b] = \int_a^b f(x) dx,$$

then  $f(x)$  is called the *probability density function* of  $X$ .

The function

$$F(x) = \Pr[X \leq x] = \int_{-\infty}^x f(t) dt$$

is called the *cumulative distribution function* of  $X$ .

**Example 9** (Uniform distribution on  $(a, b)$ ). The probability density function is given by

$$f(x) = \begin{cases} \frac{1}{b-a}, & \text{if } a < x < b; \\ 0, & \text{otherwise.} \end{cases}$$

**Example 10** (Exponential distribution). The probability density function of the exponential distribution with  $\lambda > 0$  is defined as

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

**Example 11** (Gaussian / Standard normal distribution). The probability density function is defined as

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

## 2 Expectation

**Definition 12** (Expectation). Given a probability space  $(\Omega, \mathcal{F}, \Pr[\cdot])$  and a random variable  $X$ , the *expectation* of  $X$  is defined as:

$$\mathbb{E}[X] = \sum_a a \cdot \Pr[X = a] \quad (\text{discrete})$$

$$\mathbb{E}[X] = \int_{-\infty}^{\infty} t \cdot f(t) dt \quad (\text{continuous})$$

**Example 13** (Uniform distribution on  $(a, b)$ ).

$$\mathbb{E}[X] = \int_a^b t \cdot \frac{1}{b-a} dt = \frac{1}{b-a} \cdot \frac{t^2}{2} \Big|_a^b = \frac{b+a}{2}$$

**Example 14** (Exponential distribution).

$$\mathbb{E}[X] = \int_0^\infty t \cdot \lambda e^{-\lambda t} dt = - \int_0^\infty t de^{-\lambda t} = -te^{-\lambda t} \Big|_0^\infty + \int_0^\infty e^{-\lambda t} dt = -\frac{1}{\lambda} \cdot e^{-\lambda t} \Big|_0^\infty = \frac{1}{\lambda}$$

**Definition 15** (Variance). Given a probability space  $(\Omega, \mathcal{F}, \Pr[\cdot])$  and a random variable  $X$ , the *variance* of  $X$  is defined as:

$$\text{Var}[X] = \mathbb{E}[(X - \mathbb{E}[X])^2] = \mathbb{E}[X^2] - (\mathbb{E}[X])^2.$$

**Definition 16** (Independence). Given a probability space  $(\Omega, \mathcal{F}, \Pr[\cdot])$  and two random variables  $X$  and  $Y$ ,  $X$  and  $Y$  are *independent* ( $X \perp Y$ ) if

$$\forall A, B \subseteq \mathbb{R}, \quad \Pr[X \in A \wedge Y \in B] = \Pr[X \in A] \cdot \Pr[Y \in B].$$

**Proposition 1** (Linearity of expectation). If  $X_1, X_2, \dots, X_n$  are  $n$  random variables (not necessarily independent), then

$$\mathbb{E}\left[\sum_{i=1}^n X_i\right] = \sum_{i=1}^n \mathbb{E}[X_i].$$

**Proposition 2.** If  $X_1, X_2, \dots, X_n$  are  $n$  “mutually independent” random variables, then

$$\begin{aligned} \mathbb{E}\left[\prod_{i=1}^n X_i\right] &= \prod_{i=1}^n \mathbb{E}[X_i], \\ \text{Var}\left[\sum_{i=1}^n X_i\right] &= \sum_{i=1}^n \text{Var}[X_i]. \end{aligned}$$

### 3 Conditional Probability and Conditional Expectation

In this section, assume the set of outcomes  $\Omega$  is finite. Then case for infinite  $\Omega$  is more subtle.

**Definition 17** (Conditional probability). Given a probability space  $(\Omega, \mathcal{F}, \Pr[\cdot])$ , let  $A$  and  $B$  are two events, then the probability of  $A$  conditioned on  $B$  is

$$\Pr[A|B] \triangleq \frac{\Pr[A \cap B]}{\Pr[B]}.$$

Similarly, we would like to define conditional expectation. However, to formally deal with conditional expectation, we should introduce measurable function first.

**Definition 18** (Measurable functions). Let  $\mathcal{F}$  be a  $\sigma$ -algebra and  $X$  be a function. Then  $X$  is  $\mathcal{F}$ -measurable if

$$\forall a \in \text{Range}(X), \quad X^{-1}(a) \in \mathcal{F}.$$

Denote by  $\sigma(X)$  the minimum  $\sigma$ -algebra  $\hat{\mathcal{F}}$  such that  $X$  is  $\hat{\mathcal{F}}$ -measurable.

Now we can define conditional expectation. Note that we only define it in discrete cases here.

**Definition 19** (Conditional expectation (discrete)). Suppose  $X, Y: \Omega \rightarrow \mathbb{R}$  are two random variables, and  $A \subseteq \Omega$  is a event. Then the expectation of  $X$  conditioned on  $A$  is

$$\mathbb{E}[X | A] = \sum_x x \cdot \Pr[X = x | A].$$

Specifically, let  $A = Y^{-1}(a)$  be the event that  $Y = a$ ,

$$\mathbb{E}[X | Y = a] = \sum_b b \cdot \Pr[X = b | Y = a].$$

Given  $Y$ ,  $f_Y = \mathbb{E}[X | Y]$  is a function (a random variable) on  $\Omega$  which satisfies

$$\forall \omega \in \Omega, \quad f_Y(\omega) = \mathbb{E}[X | Y = Y(\omega)].$$

**Proposition 3.** *Conditional expectation has the following two important propositions:*

1.  $\mathbb{E}[X | Y]$  is  $\sigma(Y)$ -measurable;
2.  $\mathbb{E}[\mathbb{E}[X | Y]] = \mathbb{E}[f_Y] = \mathbb{E}[X]$ .

*Proof.* Item 1 is trivial. We only prove item 2 here.

$$\begin{aligned} \mathbb{E}[\mathbb{E}[X | Y]] &= \sum_y \mathbb{E}[X | Y = y] \cdot \Pr[Y = y] \\ &= \sum_x \sum_y x \cdot \Pr[X = x | Y = y] \cdot \Pr[Y = y] \\ &= \sum_x x \cdot \sum_y \Pr[X = x | Y = y] \cdot \Pr[Y = y] \\ &= \sum_x x \cdot \sum_y \Pr[X = x \wedge Y = y] \\ &= \sum_x x \cdot \Pr[X = x] \\ &= \mathbb{E}[x]. \end{aligned}$$

□

**Example 20.** Consider the probability space  $(\Omega, \mathcal{F}, \Pr[\cdot])$  where  $\Omega$  is the set of all Chinese people and  $\Pr[\cdot]$  is the uniform probability.

Let  $X$  and  $Y$  be two random variables that  $X$  is the height of a person and  $Y$  is the gender of a person. Then  $\mathbb{E}[X]$  is the average height of Chinese people.

Let  $f_Y = \mathbb{E}[X | Y] : \Omega \rightarrow \mathbb{R}$  be the random variable that  $f_Y(\omega)$  is the average height of people with the same gender as  $\omega$ . Then  $\mathbb{E}[\mathbb{E}[X | Y]] = \mathbb{E}[X]$ .