

# Probability for Computer Scientists

36-218

CMU

Created by Lorenzo Xiao

November 8, 2022

## 1 Sample Space and Probability

### 1.1 Probabilistic Model

#### 1.1.1 Elements of a Probabilistic Model

- Sample Space, the set of all possible outcomes of an experiment.
  - collectively exhaustive, in the sense that no matter what happens in the experiment, we always obtain an outcome that has been included in the sample space .
- The probability law, assigns to a set  $A$  of possible outcome(**event**) a nonnegative  $P(A)$  that encodes the likelihood of the elements of  $A$ .

#### 1.1.2 Sequential Models, Continuous Models, and discrete model

- Events like coin tossing is often described in sequential models such as trees.
- Independent events where the probabilities of single-element is enough to characterize the probability law uses a Discrete model.
  - 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:

$$P(\{s_1, s_2, s_3, \dots, s_n\}) = P(s_1) + P(s_2) + P(s_3) + \dots + P(s_n)$$

- If the sample space consists of  $n$  possible outcomes which are equally likely, then the probability of any event  $A$  is given by:

$$P(A) = \frac{\text{size of } A}{n}$$

- Events that the probabilities of single-element may not be sufficient to characterize the probability law will use a Continuous model

### 1.1.3 Properties of Probability Laws

- $A \subset B \rightarrow P(A) \leq P(B)$
- $P(A \cup B) = P(A) + P(B) - P(A \cap B)$
- $P(A \cup B) \leq P(A) + P(B)$
- $P(A \cup B \cup C) = P(A) + P(A^c \cap B) + P(A^c \cap B^c \cap C)$

## 1.2 Conditional Model

### 1.2.1 Properties of Conditional Probability

- The probability of an event  $A$ , given an event  $B$  with  $P(B) > 0$  is:

$$P(A | B) = \frac{P(A \cap B)}{P(B)}$$

- Conditional probabilities can also be viewed as a probability law on a new universe  $B$ , because all of the conditional probability is concentrated on  $B$ .

### 1.2.2 Using Conditional Probability for Modeling

When constructing probabilistic models for experiments that have a sequential character, it is often natural and convenient to first specify conditional probabilities and then use them to determine unconditional probabilities.

- Multiplication rule :

$$P(\cap_{i=1}^n A_i) = P(A_1) P(A_2 | A_1) P(A_3 | A_1 \cap A_2) \dots P(A_n | \cap_{i=1}^{n-1} A_i)$$

## 1.3 Total Probability Theorem and Baye's Rule

### 1.3.1 Total Probability Theorem

Let  $A_1, \dots, A_n$  be disjoint events that form a partition of the sample space and assume that  $P(A_i) > 0$  for all  $i$ . Then, for any event  $B$ , we have:

$$P(B) = P(A_1 \cap B) + P(A_2 \cap B) + \dots + P(A_n \cap B)$$

### 1.3.2 Bayes' Rule

Let  $A_1, \dots, A_n$  be disjoint events that form a partition of the sample space, and assume that  $P(A_i) > 0$  for all  $i$ :

$$P(A_i | B) = \frac{P(A_i)P(B | A_i)}{P(B)} = \frac{P(A_i)P(B | A_i)}{P(A_1 \cap B) + P(A_2 \cap B) + \dots + P(A_n \cap B)}$$

## 1.4 Independence and Counting

### 1.4.1 Independence

- Two events  $A$  and  $B$  are said to be independent if:

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

- If  $A$  and  $B$  are independent, so are  $A$  and  $B^C$ .
- Independence does not imply conditional independence, and vice versa.

### 1.4.2 The Counting principle

Consider a process that consists of  $r$  stages, suppose that :

- There are  $n_1$  possible results at the first stage
- For every possible result, there are  $n_2$  possible result at the second stage
- The total number of possible results of the  $r$ -stage process is:  $n_1 n_2 n_3 \dots n_r$

### 1.4.3 Counting result

- Permutations of  $n$ :  $n!$
- $k$ -permutation of  $n$ :  $\frac{n!}{(n-k)!}$
- Combination of  $k$  out of  $n$  objects:  $\binom{n}{k} = \frac{n!}{k!(n-k)!}$
- Partition of  $n$  object in  $r$  group:  $\frac{n!}{n_1! n_2! \dots n_r!}$

## 2 Discrete Random Variables

### 2.1 Basic Concepts

#### 2.1.1 Main Concepts Related to Random Variables

- A **random variable** is a real-valued function of the outcome of the experiment.
- A **function of a random variable** defines another random variable.
- We can associate with each random variable certain “averages” of interest, such as the **mean** and the **variance**.
- A random variable can be conditioned on an event or on another random variable.
- There is a notion of **independence** of a random variable from an event or from another random variable.

#### 2.1.2 Concepts Related to Discrete Random Variables

- A **discrete random variable** is a real-valued function of the outcome of the experiment that can take a finite or countably infinite number of values.
- A discrete random variable has an associated **probability mass function (PMF)** which gives the probability of each numerical value that the random variable can take.
- A **function of a discrete random variable** defines another discrete random variable, whose PMF can be obtained from the PMF of the original random variable.

### 2.2 PMF

#### 2.2.1 Calculation of the PMF of a Random Variable X

For each possible value of  $x$  of  $X$ :

1. Collect all the possible outcomes that give rise to the event
2. Add their probability to obtain  $\rho(x)$

#### 2.2.2 Binomial Distribution

A binomial experiment possesses the following properties:

1. consists of a fixed number of identical trials.
2. Each trial can either be fail or success

3. the probability of success of a single trial is maintained trial to trial as  $q = 1 - p$
4. Independent
5. random variable of interest is  $Y$ , the number of success observed during the  $n$  trials.

$$P(y) = \binom{n}{y} p^y q^{n-y}$$

Let  $Y$  be a binomial random variable based on  $n$  trials and success probability  $p$ . Then:

$$\mu = E(Y) = np \quad \text{and} \quad \sigma^2 = V(Y) = npq$$

### 2.2.3 Geometric Distribution

A random variable  $Y$  is said to have a **geometric probability distribution** if and only if

$$p(y) = q^{y-1}p$$

Let  $Y$  be is a random variable with a geometric distribution:

$$\mu = E(Y) = \frac{1}{p} \quad \text{and} \quad \sigma^2 = V(Y) = \frac{1-p}{p^2}$$

### 2.2.4 Poisson probability distribution

A random variable  $Y$  is said to have a Poisson probability distribution if and only if:

$$p(y) = \frac{\lambda^y}{y!} e^{-\lambda}$$

Let  $Y$  be is a random variable with a Poisson probability distribution:

$$\mu = E(Y) = \lambda \quad \text{and} \quad \sigma^2 = V(Y) = \lambda$$

## 2.3 Expectation, Mean, and Variance

### 2.3.1 Expectation

We define the expected value (also called the expectation or the mean) of a random variable  $X$  with PMF  $\rho_X$  by:

$$E[X] = \sum_x x \rho_X(x)$$

### 2.3.2 Expected Value Rule

Let  $X$  be a random variable with PMF  $\rho_X$ , and let  $g(X)$  be a function of  $X$ . Then, the expected value of the random variable  $g(X)$  is given by :

$$E[g(X)] = \sum_x g(x) \rho_X(x)$$

### 2.3.3 Variance

The variance  $\text{var}(X)$  of a random variable  $X$  is defined by:

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

## 2.4 Joint PMF of Multiple Function

Let  $X$  and  $Y$  be random variables associated with the same experiment.

- The joint PMF  $\rho_{x,y}$  of  $X$  and  $Y$  is defined by:

$$\rho_{x,y}(x, y) = P(X = x, Y = y)$$

- The marginal PMFs of  $X$  and  $Y$  can be obtained from the joint PMF, using the formulas:

$$\rho_x(x) = \sum_y \rho_{x,y}(x, y), \quad \rho_y(y) = \sum_x \rho_{x,y}(x, y)$$

- A function  $g(X, Y)$  of  $X$  and  $Y$  defines another random variable, and:

$$E[g(X, Y)] = \sum_x \sum_y g(x, y) \rho_{x,y}(x, y)$$

- The above have natural extensions to the case where more than two random variables are involved.

## 2.5 Conditioning and Independence

### 2.5.1 Conditional PMFs

- Conditional PMFs are similar to ordinary PMFs, but pertain to a universe where the conditioning event is known to have occurred.
- The conditional PMF of  $X$  given an event  $A$  with  $P(A) > 0$ , is defined by:

$$\rho_{x|A}(x) = P(X = x | A)$$

and satisfies  $\sum_x \rho_{x|A}(x) = 1$

- If  $A_1, \dots, A_n$  are disjoint events form a partition of the sample space, with  $P(A_i) > 0$ , for all  $i$ , then

$$\rho_{x|A}(x) = \sum_{i=1}^n \rho_{x|A_i}(x)$$

- The conditional PMF of  $X$  given  $Y = y$  is related to the joint PMF by

$$\rho_{X,Y}(x, y) = \rho_Y(Y) \rho_{X|Y}(x | y)$$

- The conditional PMF of  $X$  given  $Y$  can be used to calculate the marginal PMF of  $X$  through the formula

$$\rho_X(x) = \sum_y \rho_Y(Y) \rho_{X|Y}(x | y)$$

- There are natural extensions of the above involving more than two random variables.

### 2.5.2 Conditional Expectation

- The conditional expectation of  $X$  given an event  $A$  with  $P(A) > 0$ , is defined by

$$E[X | A] = \sum_x x \rho_{X|A}(x)$$

- The conditional expectation of  $X$  given a value  $y$  of  $Y$  is defined by

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

- If  $A_1, \dots, A_n$  be disjoint events that form a partition of the sample space, with  $P(A_i) > 0$  for all  $i$ , then

$$E[X] = \sum_{i=1}^n \rho(A_i) E[X | A_i]$$

- We have

$$E[X] = \sum_{i=1}^n \rho_y(y) E[X | Y = y]$$

## 3 General Random Variables

### 3.1 Continuous Random Variables and PDFs

#### 3.1.1 Continuous Function

A random variable  $X$  is called continuous if there is a nonnegative function  $f_X$ , called the probability density function of  $X$ , or PDF for short, such that:

$$P(a \leq X \leq b) = \int_a^b f_X(x) dx$$

can be interpreted as area under the graph of PDF

### 3.1.2 PDF Properties

- $f_X(x) dx \geq 0$  for all x
- $\int_{-\infty}^{\infty} f_X(x) dx = 1$
- if  $\delta$  is very small, then  $P([x, x + \delta]) \approx f_X(x) \cdot \delta$
- For any subset B of the real line:  $P(X \in B) = \int_B f_X(x) dx$

### 3.1.3 Expectation

The **expected value** or **expectation** or **mean** of a continuous random variable  $X$  is defined by

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

### 3.1.4 Properties of Expectations

- The expected value rule for a function  $g(X)$  has the form:

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

- The variance of  $X$  is defined by

$$\text{var}(X) = E[(X - E[X])^2] = \int_{-\infty}^{\infty} (x - E[X])^2 f_X(x) dx$$

- If  $Y = aX + b$ :

$$E[Y] = aE[X] + b \quad \text{var}(Y) = a^2 \text{var}(X)$$

- Exponential\*

## 3.2 Cumulative Distribution Functions

### 3.2.1 Properties of a CDF

The CDF  $F_X$  of a random variable  $X$  is defined by:

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

and has several properties:



- monotonically non-decreasing
- $\lim_{x \rightarrow \infty} F_X(x) = 1, \lim_{x \rightarrow -\infty} F_X(x) = 0$
- If  $X$  is discrete,  $F_X$  will be a piecewise constant function of  $x$ 
  - If  $X$  is discrete and takes integer values, the PMF and the CDF can be obtained from each other by summing or differencing:

$$F_X(k) = \sum_{i=-\infty}^k \rho_X(i)$$

$$\rho_X(k) = P(X \leq k) - P(X \leq k-1) = F_X(k) - F_X(k-1), \forall k \in \mathbb{Z}$$

- If  $X$  is continuous,  $F_X$  will be a piecewise continuous function of  $x$ 
  - If  $X$  is continuous, the PDF and the CDF can be obtained from each other by integration or differentiation:

$$F_X(k) = \int_{-\infty}^k f_X(t) dt \quad f_X(x) = \frac{dF_X}{dx}(x)$$

### 3.3 Normal Random Variables

#### 3.3.1 Normality Preserved by Linear Transformations

If  $X$  is a normal random variable with mean  $\mu$  and variance  $\sigma^2$ , and if  $a \neq 0, b$  are scalars, then the random variable

$$Y = aX + b$$

is also normal with mean and variance:

$$E[Y] = a\mu + b \quad \text{var}(Y) = a^2\sigma^2$$

#### 3.3.2 CDF Calculation for a Normal Random Variable

For a normal random variable with mean  $\mu$  and variance, we use a two-step procedure.

1. Standardize  $X$
2. Read the cdf value from the table

$$P(X \leq x) = P\left(\frac{X - \mu}{\sigma} \leq \frac{x - \mu}{\sigma}\right) = \Phi\left(\frac{x - \mu}{\sigma}\right)$$

### 3.4 Joint PDFs of Multiple Random Variables

- **Joint PDF** is used to calculate probabilities:

$$\iint_{(x,y) \in b} f_{X,Y}(x,y) dx dy$$

- **Marginal PDF** of X and Y can be obtained by the joint PDF as:

$$f_{X,Y}(x) = \iint_{(x,y) \in b} f_{X,Y}(x,y) dy \quad f_{X,Y}(y) = \iint_{(x,y) \in b} f_{X,Y}(x,y) dx$$

- **Joint CDF** is defined by  $f_{X,Y}(x,y) = P(X \leq x, Y \leq y)$ , and determines the joint PDF through:

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

for each  $(x,y)$  where the joint cdf is continuous.

- A function  $g(x,y)$  of X and Y defines a new random variable, and

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

- The above have natural extensions to the case where more than two random variables are involved.

### 3.5 Conditioning

#### 3.5.1 Conditional PDF Given an Event

- The conditional PDF  $f_{X|A}$  of a continuous random variable X, given an event A with  $P(A) > 0$ , satisfies

$$P(X \in B | A) = \int_B f_{X|A}(x) dx$$

- If A is a subset of the real line with  $P(X \in A) > 0$ , then

$$f_{X|\{X \in A\}}(x) = \begin{cases} \frac{f_X(x)}{P(X \in A)} & x \in A \\ 0 & \text{otherwise} \end{cases}$$

- Let  $A_1, A_2, \dots, A_n$  be disjoint events that form a partition of the sample space and assume that  $P(A_i) > 0$ :

$$f_X(x) = \sum_{i=1}^n P(A_i) f_{X|A_i}(x)$$

### 3.5.2 Conditional PDF Given a Random Variable

- The joint, marginal, and conditional PDFs are related to each other by the formulas:

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

### 3.5.3 Conditional Expectations

- Definitions: The conditional expectation of  $X$  given the event  $A$  is defined by:

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

The conditional expectation of  $X$  given that  $Y = y$  is defined by:

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

- The expected value rule: For a function  $g(X)$ , we have

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

The conditional expectation of  $X$  given that  $Y = y$  is defined by:

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

- Total expectation theorem: Let  $A_1, A_2, \dots, A_n$  be disjoint events that form a partition of the sample space and assume that  $P(A_i) > 0, \forall i$ :

$$E[X] = \sum_{i=1}^n P(A_i) E[X|A_i]$$

Similarly,

$$E[X] = \int_{-\infty}^{\infty} E[X|Y=y] f_Y(y) dy$$

### 3.5.4 Independence of Continuous Random Variables

- $X$  and  $Y$  are independent if

$$f_{x,y}(x,y) = f_X(x) f_Y(y), \forall x,y$$

- If  $X$  and  $Y$  are independent, then:

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

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

### 3.6 The Continuous Bayes' Rule

- If  $X$  is a continuous random variable, we have

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

and

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

- If  $N$  is a discrete random variable, we have

$$f_Y(y) P(N = n | Y = y) = \rho_N(n) f_{Y|N}(y|n)$$

resulting in:

$$P(N = n | Y = y) = \frac{\rho_N(n) f_{Y|N}(y|n)}{f_Y(y)} = \frac{\rho_N(n) f_{Y|N}(y|n)}{\sum_i \rho_N(i) f_{Y|N}(y|i)}$$

- There are similar formulas for  $P(A | Y = y)$  and  $f_{Y|A}(y)$ .

### 3.7 Covariance and Correlation

#### 3.7.1 Covariance

The covariance of two random variable  $X$  and  $Y$ , denoted by  $\text{cov}(X, Y)$ , is defined by:

$$\text{cov}(X, Y) = E[(X - E[X])(Y - E[Y])] = E[XY] - E[X]E[Y]$$

When  $\text{cov}(X, Y) = 0$ , two functions are **uncorrelated**. If  $X$  and  $Y$  are independent, they are uncorrelated. We also have:

$$\text{var}(X + Y) = \text{var}(X) + \text{var}(Y) + 2\text{cov}(X, Y)$$

#### 3.7.2 Correlation coefficient

The correlation coefficient  $\rho(X, Y)$  of two random variables  $X$  and  $Y$  that have nonzero variances is defined as:

$$\rho(X, Y) = \frac{\text{cov}(X, Y)}{\sqrt{\text{var}(X) \text{var}(Y)}}, \quad -1 \leq \rho(X, Y) \leq 1$$

## 4 Limit Theorem

### 4.1 Markov and Chebyshev Inequalities

#### 4.1.1 Markov Inequality

$$P(X \geq a) \leq \frac{E[X]}{a}, \quad \forall a > 0$$

#### 4.1.2 Chebyshev Inequality

$$P(|X - \mu| \geq c) \leq \frac{\sigma^2}{c^2}, \forall c > 0$$

### 4.2 Central Limit Theorem

#### 4.2.1 Central Limit Theorem

Let  $X_1, X_2, \dots, X_n$  be a sequence of independent identically distributed random variables with common mean  $\mu$  and variance  $\sigma^2$ :

$$Z_n = \frac{X_1, X_2, \dots, X_n - n\mu}{\sigma\sqrt{n}}$$

Then, the CDF of  $Z_n$  converges to standard normal CDF:

$$\Phi(z) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^z e^{-\frac{x^2}{2}} dx$$

in the sense that:

$$\lim_{n \rightarrow \infty} P(Z_n \leq z) = \Phi(z), \forall z$$

#### 4.2.2 Normal Approximation Based on the Central Limit Theorem

Let  $S_n = X_1 + \dots + X_n$ , where the  $X_i$  are independent identically distributed random variables with mean  $\mu$  and variance  $\sigma^2$ . If  $n$  is large, the probability  $P(S_n \leq c)$  can be approximated by treating  $S_n$  as normal through:

1. calculate the mean  $n\mu$  and the variance  $n\sigma^2$  of  $S_n$ .
2. calculate the normalized z-value  $z = \frac{(c - n\mu)}{\sigma\sqrt{n}}$
3. Use approximation

$$P(S_n \leq c) \approx \Phi(z)$$

#### 4.2.3 De Moivre-Laplace Approximation to the Binomial

If  $S_n$  is a binomial random variable with parameters  $n$  and  $p$ ,  $n$  is large and  $k, l$  are nonnegative integers, then:

$$P(k \leq S_n \leq l) = \Phi\left(\frac{l + \frac{1}{2} - np}{\sqrt{np(1-p)}}\right) - \Phi\left(\frac{k - \frac{1}{2} - np}{\sqrt{np(1-p)}}\right)$$

### 4.3 Laws of Large Numbers

#### 4.3.1 Weak Laws of Large Numbers

Let  $X_1, X_2, \dots, X_n$  be independent identically distributed random variables with common mean  $\mu$ . For every  $\epsilon > 0$ :

$$P(|M_n - \mu| \geq \epsilon) \rightarrow 0, \text{ as } n \rightarrow \infty$$

#### 4.3.2 Strong Law of Large Numbers

Let  $X_1, X_2, \dots, X_n$  be a sequence of independent identically distributed random variables with common mean  $\mu$ . Then, the sequence of sample means  $M_n = \frac{(X_1 + \dots + X_n)}{n}$  converges to  $\mu$ :

$$P\left(\lim_{n \rightarrow \infty} \frac{(X_1 + \dots + X_n)}{n} = \mu\right) = 1$$

#### 4.3.3 Convergence with Probability

Let  $X_1, X_2, \dots, X_n$  be a sequence of random variables. Let  $c$  be a real number. We say that  $Y_n$  converges to  $c$  with a probability 1 if:

$$P\left(\lim_{n \rightarrow \infty} Y_n = c\right) = 1$$