## Chapter 0

# **Review of Probability**

In this chapter, we will review some concepts from probability theory and linear algebra that will be useful in the rest of the course. An excellent resource for review of probability is [1], which also has many examples.

## 0.1 What is probability?

Probability is a branch of mathematics that deals with sets, and functions that assign real values to these sets, in a way that certain axioms are satisfied. Note that this may or may not correspond to our models of the real world. In that sense, probability is similar to geometry, number theory, etc.

#### **Definitions:**

Assuming an experiment with different possible outcomes, we will use the following definitions.<sup>1</sup>

- $\Omega$ : the sample space, the set of all possibilities (*outcomes*)
- $E \subseteq \Omega$ : events
- p: A function from subsets of  $\Omega$  to  $\mathbb{R}_{\geq 0}$ . p(E) is the probability of the event E.

#### Axioms:

- $p(E) \ge 0$  for all  $E \subseteq \Omega$ .
- $p(\Omega) = 1$
- $p(E_1 \cup E_2) = p(E_1) + p(E_2)$  if  $E_1 \cap E_2 = \emptyset$ .

Based on these axioms, many theorems and other results can be proven. For  $A, B \subseteq \Omega$ :

- If  $A \subseteq B$ , then  $p(A) \leq p(B)$ .
- $p(A \cup B) = p(A) + p(B) p(A \cap B).$

<sup>&</sup>lt;sup>1</sup>These definitions and the following axioms are simplified. We cannot always assign probability to all subsets of  $\Omega$ . Also, for the third axiom, for any **countable** sequence of mutually exclusive events  $E_1, E_2, \ldots$ , we require that  $p(\bigcup_{i=1}^{\infty}) = \sum_{i=1}^{\infty} p(E_i)$ .

And then more definitions for basic concepts:

- Two events A and B are *independent*, denoted  $A \perp B$ , if  $p(A \cap B) = p(A)p(B)$ .
- If  $p(B) \neq 0$ , the *conditional probability* of A given B is defined as

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

• Random variables, distributions, expected value, ...

What these theorems and definitions 'mean' depends on what we think probability means.

#### Interpretations of probability

How do we assign probability to events? What does it mean, for example, to say that p(E) = 1/3?

- Frequentist interpretation: Assuming that there is a "random" experiment that can be repeated many times for which E is an event, the relative "frequency" of E occurring is 1/3.
  - Probability of heads for a fair coin: p(H) = 1/2. As odds this is represented as 1:1 (happening : not happening)
  - Probability **distribution** of the number of children ( $\leq 18$ ) of a randomly chosen American household:

	p(0)	p(1)	p(2)	p(3+)
1970	0.442	0.182	0.174	0.203
2008	0.541	0.195	0.169	0.095

- Bayesian interpretation: probability indicates the degree of belief in a way that is consistent with the axioms. This allows us to consider events that are, strictly-speaking, not random.
  - p(Heads) = 1/2 (both Bayesian and frequentist)
  - p(Stock market will hit a certain threshold this year)
  - p(Nuclear war this century)
  - p(A certain person is guilty of a given crime)

Different interpretations lead to different approaches to problems, potentially leading to different real-life decisions.

### 0.2 Sets and their sizes

Finding the probability of an event is easiest when all outcomes are equally likely. In such cases, if we can measure the size of the set of desirable outcomes A, dividing that by the size of the sample space, will yield the probability,

$$p(A) = \frac{|A|}{|\Omega|}$$

where |A| denotes the size of the set A.

**Definition 1.** A set A is **finite** if there is a natural number n such that the number of elements in A is less than n. Otherwise, it is **infinite**. If the elements of A can be counted, i.e., there is a one-to-one function from A to natural numbers, then A is **countable**. Otherwise, it is **uncountable**. A countable set may be finite (e.g.,  $\{1, 5, 6\}$ ) or infinite (e.g., integers, prime numbers, rational numbers).

If A is finite, we define its size (aka, cardinality) as the number of elements. This requires us to be able to count:

- Sum rule: If an action can be performed in m ways and another action can be performed in n, and further if we can choose which action to perform, in total we have m + n options.
- **Product rule:** If an action can be performed in m ways and another action can be performed in n, and further if we must perform both actions, in total we have  $m \times n$  options.
- **Permutations:** The number of ways we can arrange *n* objects is  $n! = 1 \times 2 \times \cdots \times n$ .
- Combinations: The number of ways we can choose k objects from a set of n objects is

$$\binom{n}{k} = \frac{n!}{k!(n-k)!}$$

**Exercise 2.** Prove that  $\binom{n}{x}x = n\binom{n-1}{x-1}$ .

**Exercise 3.** How many 8-bit bytes are there? How many of these have exactly 3 ones? If we pick a random byte, what is the probability that it has exactly 3 ones? What is the probability that it has 6 or more consecutive ones?

**Exercise 4.** How many binary sequences of length n that end with one are there with exactly k ones?

If the sample space has an infinite, even uncountable, number of outcomes, we may still be able to think of the outcomes as equally likely. For example, if we pick a random number between 0 and 1 (doing this is pretty difficult if not impossible), we may assume all outcomes are equally likely. In such cases, the size of the set can be measured via length, area, volume, etc.

**Exercise 5.** A random number in the interval [0, 1] is chosen. What is the probability that it is more than 1/2 but less than 2/3? What is the probability that it is equal to 1/2? What is the probability that it is rational(optional)?

**Exercise 6.** A random point is chosen in a square of unit side. What is the probability that it is inside the circle of diameter one inscribed in the square? What is the probability that it is on the circle?

## 0.3 Random variables and distributions

A random variable is a function that assigns real values to outcomes in  $\Omega$ . In most cases, there is a very natural mapping. For example, let X denote the number showing on a dice. Now X is a random variable, mapping each outcome of the form "the dice shows i" to the real number i. For this reason, the fact that random variables are really functions is often overlooked. Information about the probabilities of different outcomes is given by the **distribution** of the random variable.

A random variable is **discrete** if there are a countable number of possibilities (could be infinite but countable, like natural numbers). They can also be **continuous** (uncountable number of outcomes, defined over the real line or some subset of some Euclidean space).

Examples: a random variable that is 1 if heads shows when a given coin is filliped and it is 0 otherwise (discrete, finite); the arrival time of a plane in seconds from midnight ; the number of people buying a specific product; ...

#### 0.3.1 Discrete distributions

The distribution of a discrete random variable X is given by its **probability mass function** (pmf) denoted by  $p_X(x)$ , where

$$p_X(x) = p(X = x).$$

Clearly,  $p_X(x) \ge 0$  for all x and

$$\sum_{x} p_X(x) = 1.$$

If clear from the context, we drop the X in the subscript.

**Exercise 7.** A red die and a blue die are rolled. Let X denote the number showing on the red die and Y denote the sum of the two dice. Find the pmf of X and the pmf of Y.

**Exercise 8.** Two cards are drawn at random from a standard deck of 52 cards and let Z denote the number of Aces drawn. Find the pmf of Z.

Exercise 9 (Poisson Distribution). The number of times an event occurs in a given interval of times is often assumed to have a Poisson distribution (with good reason). The RV W has Poisson distribution with parameter  $\lambda$  if

$$p_X(k) = \frac{\lambda^k e^{-\lambda}}{k!}, \quad k \ge 0.$$

#### 0.3.2 Continuous distributions

The distribution of a continuous random variable X is given by its **probability distribution function** (pdf)  $p_X(x)$ , also sometimes denoted  $f_X(x)$ . Roughly speaking,

$$p\left(x - \frac{dt}{2} \le X \le x + \frac{dt}{2}\right) = p_X(x)dt.$$

For two real numbers a, b,

$$p(a \le X \le b) = \int_{a}^{b} p_X(x) dx.$$

For any pdf, we have  $p_X(x) \ge 0$  and

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

**Exercise 10** (Exponential distribution). An exponential random variable X with parameter  $\lambda$  has distribution

$$f(x) = \lambda e^{-\lambda x}, \quad x \ge 0.$$

For  $\lambda = 1$ , the probability that X is between 1 and 1.1 is around  $e^{-1} = 0.37 \times 0.1 = 0.037$ .



#### 0.3.3 Cumulative distribution functions

Cumulative distribution functions (CDF) are defined for both discrete and continuous RVs as  $F_X(x) = p_X(X \le x)$  and can be found via summation or integration:

$$F_X(x) = \sum_{k \le x} p_X(x)$$
$$F_X(x) = \int_{-\infty}^x p_X(t) dt$$

**Example 11.** The CDF of the exponential RV in Example 10 with  $\lambda = 2$  is given by



#### 0.3.4 Expected value

The expected value or the mean  $\mathbb{E}[X]$  of a random variable X with distribution p(x) is given by

$$\begin{split} \mathbb{E}[X] &= \sum_{x} x p(x), \\ \mathbb{E}[X] &= \int_{-\infty}^{\infty} x p(x) dx. \end{split}$$

One way to think about the expected value is as the average of a large number of experiments. For example, if a game pays X each time you play with probability distribution p(x), if you play the game many times, on average you will win  $\mathbb{E}[X]$  per game. That is if you play n times and n is large,

$$\frac{1}{n}(x_1 + x_2 + \dots + x_n) \simeq \mathbb{E}[X].$$

**Exercise 12.** Find the expected value of the discrete and continuous RVs in the examples above. **Exercise 13.** Find  $\mathbb{E}[1]$ .

#### Expectation of functions of random variables

For an RV X and a function f(x) it follows immediately from the definition that

$$\mathbb{E}[f(X)] = \sum_{x} f(x)p(x),$$

$$\mathbb{E}[f(X)] = \int_{-\infty}^{\infty} f(x)p(x)dx.$$
(1)

**Exercise 14.** A random variable X has distribution

$$p_X(-1) = 0.1, \ p_X(0) = 0.2, \ p_X(1) = 0.3, \ p_X(2) = 0.4.$$

Find  $\mathbb{E} X$ . Let  $Y = X^2$ . Find  $\mathbb{E} Y$ , both by finding the distribution of Y and by using (5).

#### Linearity of expectation

For a RV X, functions f(x) and g(x), and real numbers a and b,

$$\mathbb{E}[af(X) + bg(X)] = a \mathbb{E}[f(X)] + b \mathbb{E}[g(X)],$$

which can be proven easily from the definition of expectation. **Example 15.**  $\mathbb{E}[(X-a)^2] = E[X^2 - 2aX + a^2] = \mathbb{E}[X^2] - 2a\mathbb{E}X + a^2.$ 

Consider a collection of random variables  $X_1, X_2, \ldots, X_n$ . By the linearity of expectation

$$\mathbb{E}\left[\sum_{i=1}^{n} X_{i}\right] = \sum_{i=1}^{n} \mathbb{E} X_{i}.$$
(2)

If all variables are identically distributed, then

$$\mathbb{E}\left[\sum_{i=1}^{n} X_{i}\right] = n \mathbb{E} X_{1}.$$
(3)

**Example 16.** In a class of *n* students, what is the expected number of pairs of students who have the same birthday? To find this, for two students *i* and *j*, let  $X_{ij}$  be equal to 1 if they share a birthday and 0 otherwise and let  $X = \sum_{i=1}^{n-1} \sum_{j=i+1}^{n} X_{ij}$ . Now,

$$\mathbb{E} X = \binom{n}{2} \mathbb{E} X_{12} = \binom{n}{2} \Pr(X_{12} = 1) = \binom{n}{2} \frac{1}{365} \simeq \frac{n^2}{730}.$$
 (4)

In particular, if  $n = \sqrt{730} \simeq 27$  students are enough to have on average one pair with the same birthday. With n = 60 students, there should be around 5 such pairs.

#### Variance

Suppose someone offers you a game in which your expected winning is \$100. Will you accept? Which game would you play?

- You always win exactly \$100.
- You win \$0 with probability 1/2 and \$200 with probability 1/2.
- You win \$1200 with probability 1/2 and lose \$1000 with probability 1/2.

All three have the same mean. So what's different between them?

The mean helps us represent a distribution with one value, which describes the average behavior of the RV. But as this example shows, the behavior around the mean is also important. Denoting the mean of X by  $\mu_X$ , the variability around the mean is captured to a degree by the variance  $\operatorname{Var}[X]$ ,

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

The variance gives a sense of how far X is from its mean X, on average. The standard deviation,  $\sigma_X$ , is defined as

$$\sigma_X = \sqrt{\operatorname{Var}[X]},$$

and the variance is usually denoted as  $\sigma_X^2$ . Exercise 17. Prove that

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

Exercise 18. Find the mean and variance of each of the following RVs [1]:

- X + c
- aX
- aX + c
- $\frac{X \mu_X}{\sigma_X}$  (called the standardized version of X)

#### 0.3.5 Common distributions

We denote by  $X \sim \text{Dist}(a, b, ...)$ , where a, b, ..., are the parameters of the distribution.

#### **Discrete** distributions

- $X \sim \text{Ber}(p)$ :  $p(X = 1) = p, p(X = 0) = 1 p, \quad \mathbb{E}[X] = p, \quad \text{Var}[X] = p(1 p).$
- $X \sim \operatorname{Bin}(n,p)$ :  $p(x) = \binom{n}{x} p^x (1-p)^{n-x}, \ 0 \le x \le n, \quad \mathbb{E}[X] = np, \quad \operatorname{Var}[X] = np(1-p).$
- $X \sim \text{Geo}(p)$ :  $p(x) = (1-p)^{x-1}p, \ x \ge 1, \quad \mathbb{E}[X] = 1/p, \quad \text{Var}[X] = (1/p)^2 (1/p).$
- $X \sim \text{NegBin}(k,p)$ :  $p(x) = \binom{x-1}{k-1}(1-p)^{x-k}p^k, \ x \ge k, \quad \mathbb{E}[X] = k/p, \quad \text{Var}[X] = k[(1/p)^2 (1/p)].$
- $X \sim \operatorname{Poi}(\lambda)$ :  $p(x) = \frac{\lambda^{x} e^{-\lambda}}{x!}, x \ge 0, \quad \mathbb{E}[X] = \lambda, \quad \operatorname{Var}[X] = \lambda.$
- $X \sim \text{Uni}[a, b]$ :  $p(x) = \frac{1}{b-a+1}, x \in \mathbb{Z}, a \le x \le b, \mathbb{E}[X] = \frac{a+b}{2}, \text{Var}[X] = \frac{(b-a+1)^2-1}{12}.$

**Exercise 19.** Prove that the mean of Bin(n, p) is as given using Exercise 2.

#### Continuous distributions

- $X \sim \text{Uni}(a, b)$ :  $p(x) = \frac{1}{b-a}, x \in (a, b), \quad \mathbb{E}[X] = \frac{a+b}{2}, \quad \text{Var}[X] = \frac{(b-a)^2}{12}.$ •  $X \sim \mathcal{N}(\mu, \sigma^2)$ :  $p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp(-\frac{(x-\mu)^2}{2\sigma^2}), x \in \mathbb{R}, \quad \mathbb{E}[X] = \mu, \quad \text{Var}[X] = \sigma^2.$
- $X \sim \operatorname{Exp}(\lambda)$ :  $p(x) = \lambda e^{-\lambda x}, \ x \ge 0, \quad \mathbb{E}[X] = 1/\lambda, \quad \operatorname{Var}[X] = 1/\lambda^2.$

Sometimes, we drop the normalization constant, that is, the constant by which we divide to ensure that the distribution integrates to 1. This could be because the constant is not important (e.g., in Bayesian inference) or because it is hard to determine. In such cases, we use  $\propto$  to show proportionality rather than equality. We should be careful which of the entities appearing is the *variable*. For example, viewed as a function of x, we have  $f(x) = \frac{\lambda^x e^{-\lambda}}{x!} \propto \frac{\lambda^x}{x!}$  and as a function of  $\lambda$ , we have  $g(\lambda) = \frac{\lambda^x e^{-\lambda}}{x!} \propto \lambda^x e^{-\lambda}$ .

•  $X \sim \text{Beta}(\alpha, \beta)$ :  $p(x) \propto x^{\alpha-1}(1-x)^{\beta-1}, \ 0 \le x \le 1, \quad \mathbb{E}[X] = \frac{\alpha}{\alpha+\beta}, \quad \text{Var}[X] = \frac{\alpha\beta}{(\alpha+\beta)^2(\alpha+\beta+1)}.$ •  $X \sim \text{Gamma}(\alpha, \beta)$ :  $p(x) \propto x^{\alpha-1}e^{-\beta x}, \ x > 0, \quad \mathbb{E}[X] = \frac{\alpha}{\beta}, \quad \text{Var}[X] = \frac{\alpha}{\beta^2}.$ 

**Example 20.** For the distributions given in this section, try changing what the variable is and what the parameters are and check whether another distribution from the list can be obtained with appropriate normalization. For example, Bin(n, p) viewed as a distribution in p is equivalent to Beta(x + 1, n - x + 1).

## 0.4 Joint probability distributions

Joint probability distributions allow us to encode information about relationships between quantities, from independence to strongly correlated.

For random variables X and Y, the CDF and the pmf/pdf give their joint distribution, depending on their type,

$$F_{X,Y}(x,y) = \Pr(X \le x, Y \le y),$$

$$p_{X,Y}(x,y) = \Pr(X = x, Y = y),$$

$$p_{X,Y}(x,y)dxdy \simeq \Pr\left(x - \frac{dx}{2} \le X \le x + \frac{dx}{2}, y - \frac{dy}{2} \le Y \le y + \frac{dy}{2}\right),$$
CDF for continuous and discrete pmf for discrete pmf for continuous and discrete pmf for discrete pmf for continuous pdf for continuous p

We can find the distribution for each random variable (in this context these are called the **marginals**) by integration/summation,

$$p_X(x) = \sum_y p(x,y), \qquad \qquad p_X(x) = \int_{-\infty}^{\infty} p_{X,Y}(x,y) dy.$$

#### 0.4.1 Expectation, correlation, and covariance

Given two or more RVs, we may be interested in finding the expected value of a function of these RVs, e.g.,  $\mathbb{E}[XY]$ . In such case, similar to (5), we have

$$\mathbb{E}[f(X,Y)] = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(x,y)p(x,y)dxdy,$$
(5)

and similarly for discrete variables.

The correlation between X and Y is  $\mathbb{E}[XY] = \int \int xyp(x,y)dxdy$ . The covariance Cov(X,Y) and the correlation coefficient  $\rho_{X,Y}$  are defined as

$$\operatorname{Cov}(X,Y) = \mathbb{E}[(X - \mu_X)(Y - \mu_Y)]$$
$$\rho_{X,Y} = \frac{\operatorname{Cov}(X,Y)}{\sigma_X \sigma_Y}.$$

It can be shown that  $-1 \le \rho_{X,Y} \le 1$ . If  $\rho = 0$ , then the random variables are **uncorrelated**. **Exercise 21.** Show that  $\operatorname{Cov}(X, Y) = \mathbb{E}[XY] - \mathbb{E}X\mathbb{E}Y$ .

**Example 22.** The bivariate jointly Gaussian distribution for X, Y with means  $\mu_X$  and  $\mu_Y$ , variances  $\sigma_X$  and  $\sigma_Y$ , and correlation coefficient  $\rho$  is given as

$$p(x,y) = \frac{1}{2\pi\sigma_x\sigma_y\sqrt{1-\rho^2}}e^{-\frac{1}{2(1-\rho^2)}\left[\frac{(x-\mu_X)^2}{\sigma_X^2} + \frac{(y-\mu_Y)^2}{\sigma_Y^2} - \frac{2\rho(x-\mu_X)(y-\mu_Y)}{\sigma_X\sigma_Y}\right]}.$$

Examples of this pdf are given in Figure 1.

**Exercise 23.** For random variables X, Y, Z and constants a, b, c, d, e, prove that

- $\operatorname{Var}(X) = \operatorname{Cov}(X, X)$
- $\operatorname{Cov}(X+Y,Z) = \operatorname{Cov}(X,Z) + \operatorname{Cov}(Y,Z)$
- $\operatorname{Cov}(aX, Y) = a \operatorname{Cov}(X, Y)$
- $\operatorname{Cov}(X, b) = 0$
- $\operatorname{Cov}(aX + bY + c, dZ + e) = ad \operatorname{Cov}(X, Z) + bd \operatorname{Cov}(Y, Z)$

**Exercise 24.** Find the expected value and variance of X and Y from Exercise 7. Find Cov(X, Y).

#### 0.4.2 Independence

Recall that two events A and B are independent iff  $p(A \cap B) = p(A)p(B)$ . Two random variables X and Y are independent, if  $\{X \in S_1\}$  and  $Y \in S_2$  are independent for all sets  $S_1$  and  $S_2$ . This implies that

$$p(x,y) = p(x)p(y).$$
(6)



Figure 1: Bivariate Normal pdfs with  $\mu_X = \mu_Y = 0$ ,  $\sigma_X = \sigma_Y = 1$ , with  $\rho = 0$  (uncorrelated),  $\rho = .5$  (positively correlated), and  $\rho = -.5$  (negatively correlated), respectively.

For two independent random variables, we have

$$\mathbb{E}[XY] = \mathbb{E}[X] \mathbb{E}[Y] \tag{7}$$

and  $\operatorname{Cov}(X, Y) = 0.$ 

**Exercise 25.** Prove (7) using (6).

**Exercise 26.** For two independent RVs X and Y, find Var[X + Y] and  $\mathbb{E}[(X - Y)^2 + 3XY + 5]$  in terms of means and variances of X and Y.

A collection  $X_1, \ldots, X_n$  of random variables that are independent from each other but have the same distribution are called **independent and identically distributed (iid)**. We have

$$p(x_1, \dots, x_n) = \prod_{i=1}^n p(x_i).$$
 (8)

**Exercise 27.** For iid RVs  $X_1, \ldots, X_n$ , let  $E[X_i] = \mu$  and  $Var[X_i] = \sigma^2$ . Show that

$$\mathbb{E}[\bar{X}] = \mu, \qquad \qquad \operatorname{Var}[\bar{X}] = \frac{\sigma^2}{n}. \tag{9}$$

#### 0.4.3 Conditional probability and conditional distributions

For two discrete variables X and Y, the conditional probability distribution of Y given X is given by

$$p_{Y|X}(y|x) = \Pr(Y = y|X = x) = \frac{\Pr(Y = y, X = x)}{\Pr(X = x)} = \frac{p_{X,Y}(x,y)}{p_X(x)}$$

For continuous RVs, we also have  $p_{Y|X}(y|x) = \frac{p_{X,Y}(x,y)}{p_X(x)}$ . In this case, however, we interpret the conditional density as

$$p_{Y|X}(y|x) \simeq \frac{\Pr(y - \epsilon/2 \le Y \le y + \epsilon/2 | x - \epsilon/2 \le X \le x + \epsilon/2)}{\epsilon}$$

This essentially says to find  $p_{Y|X}(y|x)$ , we first assume that X is in a narrow strip around x and then find the density for Y given this assumption.

**Law of total probability.** Let  $A_1, A_2, \ldots, A_n$  be a partition of the sample space. That is  $\bigcup_{i=1}^n A_i = \Omega$  and for all  $i \neq j$ , we have  $A_i \cap A_j = \emptyset$ . For an event  $B_i$ , we have

$$p(B) = \sum_{i=1}^{n} p(B \cap A_i) = \sum_{i=1}^{n} p(B|A_i)p(A_i).$$

In particular, if X can take on  $\{1, 2, ..., n\}$ , then for another RV Y,

$$p_Y(y) = \sum_{x=1}^n p_{Y|X}(y|x)p_X(x).$$

**Chain rule of probability.** For events  $A_1, \ldots, A_n$ , we have

$$p(A_1 \cap A_2 \cap \dots \cap A_n) = p(A_1)p(A_2|A_1)p(A_3|A_1, A_2) \cdots p(A_n|A_1, \dots, A_{n-1}),$$

which can be easily proven by induction. A similar rule holds for random variables  $X_1, \ldots, X_n$ :

$$p(x_1,\ldots,x_n) = p(x_1)p(x_2|x_1)p(x_3|x_1,x_2)\ldots p(x_n|x_1,\ldots,x_{n-1}).$$

Conditional expectations are defined based on conditional distributions, e.g.,

$$\mathbb{E}[X|Y=y] = \sum_{x} x p_{X|Y}(x|y).$$

Exercise 28. Suppose the joint pmf is given as

$p_{X,Y}(x,y)$	x = 0	x = 1
y = 0	0.25	0
y = 1	0.5	0.25

Find p(y|x), p(x|y),  $\mathbb{E}[Y|X=0]$ ,  $\mathbb{E}[Y|X=1]$ ,  $\mathbb{E}[X|Y=0]$ ,  $\mathbb{E}[X|Y=1]$ .

**Exercise 29.** A point is chosen uniformly at random in a triangle with vertices on (0,0), (1,0), (1,1). Let X and Y determine the x and y coordinates of the chosen point. Find p(x|y), p(y|x),  $\mathbb{E}[X|Y=y]$ ,  $\mathbb{E}[Y|X=x]$ .

**Law of iterated expectations.** Consider a function g(x). Instead of a deterministic value for x, we can consider a random value. An example of this was given in Exercise 14 with  $g(x) = x^2$ .

Now let  $g(x) = \mathbb{E}[Y|X = x]$ . This is, of course, a well-defined function. So we can consider  $g(X) = \mathbb{E}[Y|X]$ . Now that we have a random variable, we can compute its expectation, i.e.,  $\mathbb{E}[\mathbb{E}[Y|X]]$ .

**Exercise 30.** A die is rolled, showing X. A coin is then flipped X times resulting in Y heads. Find  $\mathbb{E}[Y]$ ,  $\mathbb{E}[Y|X = x]$ , the pmf of  $\mathbb{E}[Y|X]$ , and  $\mathbb{E}[\mathbb{E}[Y|X]]$ .

It can be shown that

$$\mathbb{E}[\mathbb{E}[Y|X]] = \mathbb{E}[Y], \qquad \qquad \mathbb{E}[\mathbb{E}[Y|X,Z]|Z] = \mathbb{E}[Y|Z]. \tag{10}$$

#### 0.4.4 Bayes' rule

In Exercise 30, the conditional distribution p(y|x) is readily available as

$$p(y|x) = \binom{x}{y} 2^{-x}$$

But what if we are interested in p(x|y)? Since  $p(x|y) = \frac{p(x,y)}{p(y)}$  and p(x,y) = p(y|x)p(x), we have

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

which is called the **Bayes rule**.

**Example 31.** In Exercise 30, we can use the Bayes rule to find p(x|y),

$$p(x|y) = \frac{\binom{x}{y}2^{-x}(1/6)}{\sum_{x'=y}^{6}\binom{x'}{y}2^{-x'}(1/6)} = \frac{\binom{x}{y}2^{-x}}{\sum_{x'=y}^{6}\binom{x'}{y}2^{-x'}}$$

We may ask for example, what are the likeliest value for X if Y = 2. Below,  $p_{X|Y}(x|2)$ , i.e., the conditional distribution of X when Y = 2 is given. We can see that the likeliest values for X are 3, 4.



Bayes' rule is used in *evidential reasoning*, examples of which we will see in the next chapter. In this setting, the goal is to find the probabilities of different causes based on the evidence.

*Bayesian inference* takes its name from Bayes rule. In this setting it is often the case that we know the distribution of data given the parameters. But what we actually have is data and need to find the distribution of the parameters. The Bayes rule allows us to find this conditional distribution, a topic which will discuss in detail later.

## 0.5 Inequalities and limits

#### 0.5.1 Inequalities

#### Markov inequality

Suppose the average length of a blue whale is 22 m and we do not know anything other than the mean of distribution of the lengths of blue whales. Can we say anything about the probability that the length of a randomly chosen blue whale is  $\geq 30$ m? For example, is it possible that this probability is 0.8 or larger? No, since in that case, the average would be  $\geq 0.8 \times 30 = 24$ m. So only knowing the mean enables us to say something about the extremes of the probability distribution.

This observation is formalized via the **Markov inequality**. For a *non-negative* random variable X, we have

$$\Pr(X \ge a) \le \frac{\mathbb{E}\,X}{a}.$$

Exercise 32. Prove the Markov inequality.

A special case of this occurs when X counts something, i.e., it only takes non-negative integer values. Then,

$$\Pr(X \ge 1) = \Pr(X > 0) \le \mathbb{E} X, \qquad \qquad \Pr(X = 0) \ge 1 - \mathbb{E} X.$$

In particular, if the mean  $\mathbb{E} X$  is small, then there is a large probability that X = 0. **Exercise 33** (optional). Provide a bound on the probability that in a random binary sequence of length n, there exists a run (consecutive occurrences) of 1s of length at least  $2 \log_2 n$ ? (The result will tell you that this is unlikely for large n.)

#### Chebyshev inequality

If in addition to the mean, we also have the variance, we can use the Chebyshev bound. For a random variable X with mean  $\mu$  and variance  $\sigma^2$ ,

$$\Pr\left(\left|\frac{X-\mu}{\sigma}\right| \ge a\right) \le \frac{1}{a^2}.$$

Exercise 34. Prove the Chebyshev bound using the Markov bound.

**Example 35.** The Chebyshev bound tells us that being k standard deviations away from the mean has probability at most  $1/k^2$ .

k	2	3	4	5	6	7	8	9	10
Probability of deviating	25%	11.1%	6.25~%	4%	2.78%	2.04%	1.56%	1.23%	1%
more than $k \times$ std is $\leq$									

In particular, being 10 standard deviations away from the mean has probability at most 1%.

#### 0.5.2 Limits

Limits in probability provide a way to understand what happens when the number of experiments grows or many random effects accumulate. Limit theorems are beneficial given that we often deal with large volumes of data. The following limit theorems will be helpful to us later in the course.

#### Law of large numbers

Let  $X_1, \ldots, X_n$  be random variables with mean  $\mu$  and variance  $\leq \sigma^2$  and suppose that for each i and j,  $X_i$  and  $X_j$  are uncorrelated (in particular, it is sufficient for them to be independent). Also, let  $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$ . Then, for any  $\epsilon > 0$ ,

$$\Pr(|\bar{X}_n - \mu| \ge \epsilon) \le \frac{\sigma^2}{n\epsilon^2}.$$
(11)

As *n* becomes large the right side becomes smaller and smaller. So for large *n* the probability of  $\bar{X}_n$  being too far from the mean is very small. This is referred to as the **Law of Large Numbers** (LLN). In other words, if we take *n* independent samples from a random variable *X*, then the average of those samples will be close to the mean  $\mathbb{E} X$ ,

$$\frac{1}{n}(x_1 + x_2 + \ldots + x_n) \simeq \mathbb{E}[X],$$

which is what we used to motivate expected value.

**Exercise 36.** Use the Chebyshev inequality to prove LLN when random variables are independent and all have the same variance  $\sigma^2$ .

**Example 37.** Suppose  $X_i \sim \text{Poi}(2)$ ,  $1 \leq i \leq 500$ , and let  $\bar{X}_n$  be the average of the first  $n X_i$ s. Figure 2 shows the plot for  $\bar{X}_n$  for a realization of  $X_i$ s obtained via computer simulation. It is observed that for large values of n,  $\bar{X}_n$  is close to 2, the mean of the Poisson distribution.

#### Central limit theorem

Let  $X_1, X_2, \ldots$  be iid random variables with mean  $\mu$  and variance  $\sigma^2$  and let  $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$ . As  $n \to \infty$ . The **Central Limit Theorem (CLT)** states that

distribution of 
$$\sqrt{n}(\bar{X}_n - \mu) \to \mathcal{N}(0, \sigma^2).$$
 (12)



Figure 2:  $\bar{X}_n$  based on  $X_i \sim \text{Poi}(2)$  as a function of n.

That is, the distribution of  $\sqrt{n}(\bar{X}_n - \mu)$  approaches the distribution of a normal random variable with mean 0 and variance  $\sigma^2$ .

Loosely speaking, the CLT also means  $S_n = \sum_{i=1}^n X_i$  has distribution  $\mathcal{N}(n\mu, n\sigma^2)$ . **Example 38.** Let  $X_i \sim \text{Uni}(0, 1), 1 \leq i \leq n = 10$ . We produce 50,000 samples of  $\bar{X}_n$  (and  $S_n$ ), and plot the normalized histograms for  $\sqrt{n}(\bar{X}_n - \mu)$  and the pdf of  $\mathcal{N}(0, \sigma^2)$  and the normalized histogram for  $S_n$  and the pdf of  $\mathcal{N}(n\mu, n\sigma^2)$  in Figure 3.



Figure 3: The normalized histograms for  $\sqrt{n}(\bar{X}_n - \mu)$  and the pdf of  $\mathcal{N}(0, \sigma^2)$  (on the left) and the normalized histogram for  $S_n$  and the pdf of  $\mathcal{N}(n\mu, n\sigma^2)$  (on the right) for uniform  $X_i$  with  $\mu = 1/2$  and  $\sigma^2 = 1/12$  and with n = 10.