

Math 154: Probability Theory, Lecture Notes

KEVIN YANG

CONTENTS

1. Week 1, starting Tue. Jan. 23, 2024	2
1.1. Probability spaces and events	2
1.2. Conditional probability	3
1.3. Independence	4
1.4. Some examples	5
2. Week 2, starting Tue. Jan. 30, 2024	6
2.1. Random variables	6
2.2. Independence of random variables	8
2.3. Expectation	9
2.4. Variance and higher moments	11
2.5. Cauchy-Schwarz and Hölder inequalities	13

1. WEEK 1, STARTING TUE. JAN. 23, 2024

1.1. Probability spaces and events.

Definition 1.1. Take a set Ω . A σ -algebra \mathcal{F} is a collection of subsets of Ω such that

- $\Omega, \emptyset \in \mathcal{F}$.
- If $\{A_n\}_{n=1}^\infty$ is a collection of sets in \mathcal{F} , then $\bigcup_{n=1}^\infty A_n \in \mathcal{F}$ and $\bigcap_{n=1}^\infty A_n \in \mathcal{F}$.

Sets in \mathcal{F} are called *events*. A probability measure \mathbb{P} on (Ω, \mathcal{F}) is a function $\mathbb{P} : \mathcal{F} \rightarrow [0, 1]$ such that

- $\mathbb{P}(\emptyset) = 0$ and $\mathbb{P}(\Omega) = 1$
- If $\{A_n\}_{n=1}^\infty$ is a pairwise disjoint collection of sets in \mathcal{F} , then $\mathbb{P}(\bigcup_{n=1}^\infty A_n) = \sum_{n=1}^\infty \mathbb{P}(A_n)$.
- If $\{E_n\}_{n=1}^\infty$ are in \mathcal{F} and $E_1 \subseteq E_2 \subseteq \dots$, then $\mathbb{P}(E_n) \rightarrow \mathbb{P}(\bigcup_{k=1}^\infty E_k)$.
- If $\{B_n\}_{n=1}^\infty$ are in \mathcal{F} and $B_1 \supseteq B_2 \supseteq \dots$, then $\mathbb{P}(B_n) \rightarrow \mathbb{P}(\bigcap_{n=1}^\infty B_n)$.
- **The previous two bullet points are necessary parts of the definition. They must follow**

The data $(\Omega, \mathcal{F}, \mathbb{P})$ is called a *probability space*.

Example 1.2. A coin is tossed. In this case, $\Omega = \{H, T\}$ (heads or tails). We can take $\mathcal{F} = 2^\Omega$. It contains $\{H, T\}$ (the coin lands heads or tails), $\{H\}$ (the coin lands heads), $\{T\}$ (the coin lands tails), and \emptyset (the coin lands neither heads or tails). We have $\mathbb{P}(H) = 1 - \mathbb{P}(T)$, and $\mathbb{P}(\{H, T\}) = 1$ and $\mathbb{P}(\emptyset) = 0$. If it is a fair coin, then $\mathbb{P}(H), \mathbb{P}(T) = \frac{1}{2}$.

Example 1.3. A six-sided dice is thrown. $\Omega = \{1, 2, 3, 4, 5, 6\}$. We can take $\mathcal{F} = 2^\Omega$. **In general, if Ω is finite, one should always take $\mathcal{F} = 2^\Omega$.** If $X \in \mathcal{F}$ has size 1, then $\mathbb{P}(X) = \frac{1}{6}$. Then, use the additivity property to extend all of \mathbb{P} . (For example, $\mathbb{P}(\{1, 2\}) = \frac{1}{6} + \frac{1}{6} = \frac{1}{3}$.)

Lemma 1.4. (1) $\mathbb{P}(A^C) = 1 - \mathbb{P}(A)$, where $A^C = \Omega \setminus A$.

(2) If $B \supseteq A$, then $\mathbb{P}(B) = \mathbb{P}(A) + \mathbb{P}(B \setminus A) \geq \mathbb{P}(A)$.

(3) If $A_1, \dots, A_n \in \mathcal{F}$, then

$$\begin{aligned} \mathbb{P}(\bigcup_{i=1}^n A_i) &= \sum_{i=1}^n \mathbb{P}(A_i) - \sum_{i < j} \mathbb{P}(A_i \cap A_j) + \sum_{i < j < k} \mathbb{P}(A_i \cap A_j \cap A_k) - \dots \\ &\quad + (-1)^{n+1} \mathbb{P}(A_1 \cap \dots \cap A_n). \end{aligned}$$

For $n = 2$, this reduces to $\mathbb{P}(A \cup B) = \mathbb{P}(A) + \mathbb{P}(B) - \mathbb{P}(A \cap B)$.

(4) If $A_1, \dots, A_n, \dots \in \mathcal{F}$, then $\mathbb{P}(\bigcup_{n=1}^\infty A_n) \leq \sum_{n=1}^\infty \mathbb{P}(A_n)$. **This is the union bound**

Proof. Take the sequence $A_1 = A$ and $A_2 = A^C$ (and $A_n = \emptyset$ for all $n \geq 3$). We have $\mathbb{P}(A) + \mathbb{P}(A^C) = 1$, so point (1) follows. For point (2), write $B = A \cup (B \setminus A)$. Set $A_1 = A$, $A_2 = B \setminus A$, and $A_n = \emptyset$ for $n \geq 3$. Thus $\mathbb{P}(A) + \mathbb{P}(B \setminus A) = \mathbb{P}(A \cup (B \setminus A)) = \mathbb{P}(B)$, so point (2) follows. We will not prove point (3), since it is not really useful, but it's the same general principle as point (2). For point (4), we first define an auxiliary sequence $B_n = A_n \setminus \bigcup_{k=1}^{n-1} A_k$ and $B_1 = A_1$. Then B_n are pairwise disjoint. So $\mathbb{P}(\bigcup_{n=1}^\infty B_n) = \sum_{n=1}^\infty \mathbb{P}(B_n)$. But $\bigcup_{n=1}^\infty B_n = \bigcup_{n=1}^\infty A_n$, and $B_n \subseteq A_n$, so $\mathbb{P}(B_n) \leq \mathbb{P}(A_n)$, and point (4) follows. \square

Lemma 1.5. Let $\{A_n\}_{n=1}^\infty$ be in \mathcal{F} . Then $(\cup_{n=1}^\infty A_n)^C = \cap_{n=1}^\infty A_n^C$ and $(\cap_{n=1}^\infty A_n)^C = \cup_{n=1}^\infty A_n^C$. *One can take $A_n = \emptyset$ or $A_n = \Omega$ for all $n \geq N$ for some N to take finite unions and intersections.*

Proof. Take $x \in (\cup_{n=1}^\infty A_n)^C$. Thus, $x \notin A_n$ for any n . So $x \in A_n^C$ for all n , which means $x \in \cap_{n=1}^\infty A_n^C$. Now, take $x \in \cap_{n=1}^\infty A_n^C$, so $x \notin A_n$ for all n . This means $x \notin \cup_{n=1}^\infty A_n$, thus $x \in (\cup_{n=1}^\infty A_n)^C$. This shows $(\cup_{n=1}^\infty A_n)^C = \cap_{n=1}^\infty A_n^C$. The other claim follows by the same argument. \square

Example 1.6. Let $A, B \in \mathcal{F}$. Suppose $\mathbb{P}(A) = \frac{3}{4}$ and $\mathbb{P}(B) = \frac{1}{3}$. We can bound $\mathbb{P}(A \cap B)$ as follows. First,

$$\mathbb{P}(A \cap B) = \mathbb{P}(A) + \mathbb{P}(B) - \mathbb{P}(A \cup B).$$

We know $\mathbb{P}(A \cup B) \leq 1$, so $\mathbb{P}(A \cap B) \geq \frac{3}{4} + \frac{1}{3} - 1 = \frac{1}{12}$. Also, we know $\mathbb{P}(A \cup B) \geq \mathbb{P}(A)$, so $\mathbb{P}(A \cap B) \leq \frac{3}{4} + \frac{1}{3} - \frac{3}{4} = \frac{1}{3}$.

1.2. Conditional probability.

Definition 1.7. Take $B \in \mathcal{F}$ so that $\mathbb{P}(B) > 0$. The *conditional probability of A given B* is

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

The idea is that one takes Ω , and restricts to a smaller probability space with set B . The σ -algebra is just given by taking \mathcal{F} and intersecting with B (feel free to try to show that this is a σ -algebra). $\mathbb{P}(\cdot|B)$ is the “natural” probability measure on this probability space.

Example 1.8. Two fair dice are thrown. Condition on the first showing 3. What is the probability that the sum of the two rolls is > 6 ? Let A be the event where the sum of the two rolls is > 6 and B is the event where the first roll is a 3. We have

$$\mathbb{P}(A|B) = \frac{\mathbb{P}(A \cap B)}{\mathbb{P}(B)} = \frac{\mathbb{P}(A \cap B)}{\frac{1}{6}}.$$

Note that $A \cap B$ is the event where the second roll is 4, 5, 6, and the first roll is a 3. In particular, there are 3 outcomes out of 36 that are okay, so the probability of $\mathbb{P}(A \cap B) = \frac{3}{36}$. This shows $\mathbb{P}(A|B) = \frac{1}{2}$.

Example 1.9. A coin is flipped twice **independently**. What is the probability that both are heads, given that one is a heads. It is not $\frac{1}{2}$. Indeed, let A be the event of two heads, and B is the event where one is a heads. There are four total outcomes, three of which have at least one heads. So $\mathbb{P}(B) = \frac{3}{4}$. On the other hand, $A \cap B$ is just the event of two heads, so its probability is $\frac{1}{4}$. This shows $\mathbb{P}(A|B) = \frac{\mathbb{P}(A \cap B)}{\mathbb{P}(B)} = \frac{1}{3}$.

Lemma 1.10 (Law of total probability). We say that $B_1, \dots, B_n \in \mathcal{F}$ form a partition of Ω if they are pairwise disjoint, positive probability, and $\cup_{i=1}^n B_i = \Omega$. For any partition B_1, \dots, B_n and any event A , we have

$$\mathbb{P}(A) = \sum_{i=1}^n \mathbb{P}(A|B_i)\mathbb{P}(B_i).$$

In particular, for any events A, B (where $B \neq \Omega, \emptyset$), we have $\mathbb{P}(A) = \mathbb{P}(A|B)\mathbb{P}(B) + \mathbb{P}(A|B^C)\mathbb{P}(B^C)$.

Proof. Since B_1, \dots, B_n is a partition, the collection $A \cap B_1, \dots, A \cap B_n$ are disjoint and $\cup_{k=1}^n A \cap B_k = A$. (To see this, note that clearly $A \cap B_k \subseteq A$, so it suffices to show that $A \subseteq \cup_{k=1}^n A \cap B_k$. Take $x \in A$. Then $x \in \Omega$, and since B_1, \dots, B_n is a partition, we know $x \in B_k$ for some k . Thus $x \in A \cap B_k$, and thus $x \in \cup_{k=1}^n A \cap B_k$.) From the first sentence, we get $\mathbb{P}(A) = \mathbb{P}(\cup_{k=1}^n A \cap B_k) = \sum_{k=1}^n \mathbb{P}(A \cap B_k)$. By the definition of conditional probability, we have $\mathbb{P}(A \cap B_k) = \mathbb{P}(A|B_k)\mathbb{P}(B_k)$. Combining the previous two sentences finishes the proof. \square

Theorem 1.11 (Bayes' formula). *This will be helpful for the homework* For any events A, B of positive probability, we have $\mathbb{P}(A|B) = \frac{\mathbb{P}(B|A)\mathbb{P}(A)}{\mathbb{P}(B)}$.

Proof. It suffices to combine $\mathbb{P}(A \cap B) = \mathbb{P}(A|B)\mathbb{P}(B)$ and $\mathbb{P}(A \cap B) = \mathbb{P}(B|A)\mathbb{P}(A)$. Indeed, this implies $\mathbb{P}(A|B)\mathbb{P}(B) = \mathbb{P}(B|A)\mathbb{P}(A)$. Now, divide by $\mathbb{P}(B)$ on both sides (which one can do because B has positive probability!). \square

1.3. Independence.

Definition 1.12. We say events A, B are *independent* if $\mathbb{P}(A \cap B) = \mathbb{P}(A)\mathbb{P}(B)$. *Independent and disjoint are totally different notions!* This is the same as $\mathbb{P}(A|B) = \mathbb{P}(A)$.

We say a family of events $\{A_i\}_{i=1}^\infty$ are *jointly independent* if $\mathbb{P}(\cap_{i=1}^n A_i) = \prod_{i=1}^n \mathbb{P}(A_i)$. We say it is *pairwise independent* if A_i, A_j are independent for all $i \neq j$.

Example 1.13. Let $\Omega = \{abc, acb, cab, cba, bca, bac, aaa, bbb, ccc\}$. Each element in Ω occurs with probability $\frac{1}{9}$. Let A_k be the event where the k -th letter (for $k = 1, 2, 3$) is a . We know that A_1, A_2, A_3 are pairwise independent. Indeed, $A_1 \cap A_2$ is the event where the first and second letter are both a . Thus, $A_1 \cap A_2 = \{aaa\}$, so $\mathbb{P}(A_1 \cap A_2) = \frac{1}{9}$. Note that $\mathbb{P}(A_1)\mathbb{P}(A_2) = \frac{1}{3} \cdot \frac{1}{3} = \frac{1}{9}$. Similar arguments apply to A_1, A_3 and A_2, A_3 (try it!).

But, A_1, A_2, A_3 are not jointly independent. Indeed, $A_1 \cap A_2 \cap A_3 = \{aaa\}$, so its probability is $\frac{1}{9}$. But $\mathbb{P}(A_1)\mathbb{P}(A_2)\mathbb{P}(A_3) = \frac{1}{3} \cdot \frac{1}{3} \cdot \frac{1}{3} = \frac{1}{27}$.

Example 1.14. We pick a card uniformly at random from a deck of 52. Each has probability $\frac{1}{52}$. Let A be the event where a king is picked, and B is the event where a spade is picked. Then $\mathbb{P}(A) = \frac{4}{52} = \frac{1}{13}$, and $\mathbb{P}(B) = \frac{1}{4}$. Also, $\mathbb{P}(A \cap B) = \frac{1}{52}$. So A, B are independent.

Lemma 1.15. If A, B are independent, then A^C, B are independent and A^C, B^C are independent.

Proof. We claim $\mathbb{P}(A^C \cap B) + \mathbb{P}(A \cap B) = \mathbb{P}(B)$. (This follows because $A^C \cap B$ and $A \cap B$ are disjoint and union to B .) Since A, B are independent, this implies $\mathbb{P}(A^C \cap B) = \mathbb{P}(B) - \mathbb{P}(A \cap B) = \mathbb{P}(B) - \mathbb{P}(A)\mathbb{P}(B) = (1 - \mathbb{P}(A))\mathbb{P}(B)$. But $\mathbb{P}(A^C) = 1 - \mathbb{P}(A)$, so we get $\mathbb{P}(A^C \cap B) = \mathbb{P}(A^C)\mathbb{P}(B)$, which means A^C, B are independent. To show that A^C, B^C are independent, use the first result (but replace A by B and B by A^C). \square

Example 1.16. Two fair dice are rolled independently. Let A be the event where the sum of the rolls is 7. Let B be the event where the first roll is 1. Then A, B are independent.

Indeed, $\mathbb{P}(A|B) = \frac{1}{6}$ (since a six is needed on the second roll). But $\mathbb{P}(A) = \frac{6}{36}$, since for any value of the first roll, there is exactly one value of the second roll to realize A . If we change 7 to 1, then A, B are no longer independent.

Definition 1.17. Fix an event B with positive probability. We say that A_1, A_2 are *conditionally independent* (given/conditioning on B) if $\mathbb{P}(A_1 \cap A_2|B) = \mathbb{P}(A_1|B)\mathbb{P}(A_2|B)$.

Lemma 1.18. Fix B . Then A_1, A_2 are conditionally independent given B if and only if $\mathbb{P}(A_1|A_2, B) = \mathbb{P}(A_1|B)$.

Proof. Suppose conditional independence of A_1, A_2 . Then

$$\begin{aligned}\mathbb{P}(A_1|A_2, B) &= \frac{\mathbb{P}(A_1 \cap A_2 \cap B)}{\mathbb{P}(A_2 \cap B)} \\ &= \frac{\mathbb{P}(A_1 \cap A_2|B)\mathbb{P}(B)}{\mathbb{P}(A_2|B)\mathbb{P}(B)} \\ &= \frac{\mathbb{P}(A_1|B)\mathbb{P}(A_2|B)\mathbb{P}(B)}{\mathbb{P}(A_2|B)\mathbb{P}(B)} \\ &= \mathbb{P}(A_1|B).\end{aligned}$$

Now suppose that $\mathbb{P}(A_1|A_2, B) = \mathbb{P}(A_1|B)$. Then

$$\begin{aligned}\mathbb{P}(A_1 \cap A_2|B) &= \frac{\mathbb{P}(A_1 \cap A_2 \cap B)}{\mathbb{P}(B)} \\ &= \frac{\mathbb{P}(A_1|A_2, B)\mathbb{P}(A_2 \cap B)}{\mathbb{P}(B)} \\ &= \frac{\mathbb{P}(A_1|B)\mathbb{P}(A_2|B)\mathbb{P}(B)}{\mathbb{P}(B)} = \mathbb{P}(A_1|B)\mathbb{P}(A_2|B).\end{aligned}$$

This finishes the proof. \square

Example 1.19. Suppose I have two coins. One is fair, and the other one has probability of heads equal to $\frac{1}{3}$. I choose one of the two coins uniformly at random, and I toss it twice (independently). Let X be the value of the first flip and Y be the value of the second flip. Then X and Y are conditionally independent given that I choose the fair coin. (Same is true if I condition on choosing the non-fair coin.)

1.4. Some examples.

- (1) (Symmetric random walk, “gambler’s ruin”) Let’s play a game. We flip a coin repeatedly. If it lands heads, I get one dollar. If it lands tails, I lose a dollar. (Suppose this is a fair coin for now.) I want to save N dollars, at which point I stop the game, so that I can retire happily. But if I end up with zero dollars at any point, we stop the game, since I can’t play anymore.

Suppose I start with $0 < k < N$ dollars. What is the probability that I win?

- Let $p_k = \mathbb{P}_k(A)$ be the event that I win if we start at k dollars. By the law of total probability, if B is the event that we toss a heads, then

$$\mathbb{P}_k(A) = \mathbb{P}_k(A|B)\mathbb{P}(B) + \mathbb{P}_k(A|B^C)\mathbb{P}(B^C).$$

We have $\mathbb{P}_k(A|B) = p_{k+1}$ and $\mathbb{P}_k(A|B^C) = p_{k-1}$ and $\mathbb{P}(B), \mathbb{P}(B^C) = \frac{1}{2}$. So $p_k = \frac{1}{2}(p_{k+1} + p_{k-1})$. But also $p_0 = 0$ and $p_N = 1$. We will talk later in this class about how to solve this equation efficiently, but one can check that $p_k = 1 - \frac{k}{N}$ solves this equation.

- (2) (Testimonies) We are in court over whether or not Kevin stole the piece of chalk. We have two witnesses Alf and Bob. Alf tells the truth with probability α and Bob commits perjury with probability β . There is no collusion between these two (as in whether Kevin did it or not, their testimonies are independent). Let A be the event where Alf says Kevin stole it, and B be the event where Bob says Kevin stole it. Let T be the event where Kevin stole it. What is probability that Kevin stole it given that Alf and Bob said so, in terms of $\tau = \mathbb{P}(T)$?

- We need to compute $\mathbb{P}(T|A \cap B)$. By Bayes' rule, we have

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

We have $\mathbb{P}(A \cap B|T) = \mathbb{P}(A|T)\mathbb{P}(B|T) = \alpha\beta$, so the numerator is $\alpha\beta\tau$. For the bottom, by the law of total probability, we have

$$\begin{aligned}\mathbb{P}(A \cap B) &= \mathbb{P}(A \cap B|T)\mathbb{P}(T) + \mathbb{P}(A \cap B|T^C)\mathbb{P}(T^C) \\ &= \alpha\beta\tau + (1 - \alpha)(1 - \beta)(1 - \tau).\end{aligned}$$

$$\text{So, } \mathbb{P}(T|A \cap B) = \frac{\alpha\beta\tau}{\alpha\beta\tau + (1 - \alpha)(1 - \beta)(1 - \tau)}.$$

- (3) (Simpson's paradox)

2. WEEK 2, STARTING TUE. JAN. 30, 2024

2.1. Random variables.

Definition 2.1. A *random variable* is a function $X : \Omega \rightarrow \mathbb{R}$ such that for any $x \in \mathbb{R}$, the event $\{X \leq x\}$ is in \mathcal{F} . The function $F(x) := \mathbb{P}(X \leq x)$ is the *distribution function* associated to X .

We say X is *discrete* if it only takes values in a countable set $\{x_1, \dots, x_n, \dots\}$ of \mathbb{R} . We say X is *continuous* if its distribution function can be represented as

$$F(x) = \mathbb{P}(X \leq x) = \int_{-\infty}^x f(u)du,$$

where $f : \mathbb{R} \rightarrow [0, \infty)$ is called the *probability density function* (it needs to be integrable, i.e. $\int_{\mathbb{R}} f(u)du < \infty$).

It is a fact that if X, Y are random variables and $a, b \in \mathbb{R}$, then $aX + bY$ is a random variable!

Lemma 2.2. A distribution function F satisfies

- (1) If $x \leq y$, then $F(x) \leq F(y)$ (even if $x < y$, we can still have $F(x) = F(y)$!)
- (2) $F(x) \rightarrow 0$ as $x \rightarrow -\infty$ and $F(x) \rightarrow 1$ as $x \rightarrow \infty$.
- (3) $F(x + h) \rightarrow F(x)$ as $h \rightarrow 0$ from above.

Proof. (1) If $x \leq y$, then $\{X \leq x\} \subseteq \{X \leq y\}$.

- (2) Let $A_n := \{X \leq -a_n\}$, where $a_n \rightarrow \infty$ is strictly increasing. Then $F(a_n) = \mathbb{P}(A_n)$. But $A_n \supseteq A_m$ for all $m \geq n$. So $F(a_n) = \mathbb{P}(A_n) \rightarrow \mathbb{P}(\cap_{m=1}^{\infty} A_m) = \mathbb{P}(\emptyset) = 0$.
Let $B_n := \{X \geq b_n\}$, where $b_n \rightarrow \infty$ is strictly increasing. Note that $B_n \subseteq B_m$ if $m \geq n$. Then $F(b_n) = \mathbb{P}(B_n) = \mathbb{P}(\cup_{m=1}^{\infty} B_m) = \mathbb{P}(\Omega) = 1$.
- (3) Let $A_n = \{X \leq x + h_n\}$, where $h_n \rightarrow 0$ is strictly decreasing. Then $\cap_{n=1}^{\infty} A_n = \{X \leq x\}$, and $A_n \supseteq A_m$ if $m \geq n$. So $F(x + h_n) = \mathbb{P}(A_n) \rightarrow \mathbb{P}(\cap_{n=1}^{\infty} A_n) = \mathbb{P}(X \leq x) = F(x)$. □

Definition 2.3. Suppose X is a discrete random variable. Its *probability mass function* (or *pmf*) is the function $f : \mathbb{R} \rightarrow [0, 1]$ such that $f(x) = \mathbb{P}(X = x)$. **This is generally much easier to compute than the distribution function!**

Example 2.4 (Bernoulli distribution). Any random variable which is valued in $\{0, 1\}$. For example, the outcome of flipping a coin is Bernoulli, if we interpret heads as 1 and tails as 0. If the probability of heads is p , then its pmf is $p(1) = p$ and $p(0) = 1 - p$ and $p(x) = 0$ for $x \neq 0, 1$. The distribution function is $F(x) = 0$ for all $x < 0$, and $F(x) = 1 - p$ for all $x \in [0, 1)$, and $F(x) = 1$ for all $x \geq 1$.

For shorthand, we write $X \sim \text{Bern}(p)$.

Example 2.5 (Binomial distribution). Let X_1, \dots, X_n be *independent* Bernoulli random variables. Set $Y = X_1 + \dots + X_n$. This is a *binomial* random variable. It is discrete, since it takes values in $\{0, 1, \dots, n\}$. Its probability mass function satisfies $p(x) = 0$ if $x \notin \{0, 1, \dots, n\}$. For any $k \in \{0, 1, \dots, n\}$, $p(k)$ is the probability of flipping exactly k heads. There are $\binom{n}{k} = \frac{n!}{k!(n-k)!}$ ways to choose k out of n flips to be heads. The probability of flipping this particular sequence of heads and tails is $p^k(1 - p)^{n-k}$. So $p(k) = \binom{n}{k} p^k (1 - p)^{n-k}$.

For shorthand, we write $X \sim \text{Bin}(n, p)$.

Example 2.6 (Poisson distribution). X takes values in the set $\{0, 1, 2, \dots\}$. Its pmf is defined to be

$$\mathbb{P}(X = k) = p_{\lambda}(k) = \frac{\lambda^k}{k!} e^{-\lambda}.$$

Here, $\lambda > 0$ is a fixed parameter. Note that

$$\sum_{k=0}^{\infty} p_{\lambda}(k) = \sum_{k=0}^{\infty} \frac{\lambda^k}{k!} e^{-\lambda} = e^{\lambda} e^{-\lambda} = 1,$$

so $p_{\lambda}(\cdot)$ is indeed a probability mass function. For shorthand, we write $X \sim \text{Pois}(\lambda)$.

Example 2.7 (Geometric distribution). Flip a coin repeatedly with probability of heads being p . Let X be the first time that the coin turns up heads. This takes values in $\{1, \dots\}$. Its pmf is $\mathbb{P}(X = k) = p(k) = (1 - p)^{k-1} p$. This is called the geometric distribution, since

$$\sum_{k=1}^{\infty} p(k) = p \sum_{k=0}^{\infty} (1 - p)^k = p \frac{1}{1 - (1 - p)} = 1$$

is a geometric series.

Definition 2.8. A *random vector* of dimension n is a vector $\mathbf{X} = (X_1, \dots, X_n)$ such that $X_1, \dots, X_n : \Omega \rightarrow \mathbb{R}$ are random variables. If X_1, \dots, X_n are discrete random variables, then the pmf of \mathbf{X} is defined to be the function

$$p(x_1, \dots, x_n) = \mathbb{P}(X_1 = x_1, \dots, X_n = x_n).$$

2.2. Independence of random variables.

Definition 2.9. A collection of random variables $X_1, \dots, X_n : \Omega \rightarrow \mathbb{R}$ (i.e. on the same probability space) are *jointly independent* if for open or closed subsets $A_1, \dots, A_n \subseteq \mathbb{R}$, we have

$$\mathbb{P}(\cap_{i=1}^n \{X_i \in A_i\}) = \prod_{i=1}^n \mathbb{P}(X_i \in A_i).$$

We say they are pairwise independent if X_i, X_j are independent for all $i \neq j$.

Lemma 2.10. Let X_1, \dots, X_n be independent discrete random variables with pmfs p_1, \dots, p_n . Then X_1, \dots, X_n are jointly independent if and only if for any $x_1, \dots, x_n \in \mathbb{R}$, we have

$$\mathbb{P}(X_1 = x_1, \dots, X_n = x_n) = \prod_{i=1}^n p_i(x_i).$$

Proof. If X_1, \dots, X_n are jointly independent, just take the formula for joint independence above and set $A_i = \{x_i\}$ for all i . For the other direction, we have

$$\begin{aligned} \mathbb{P}(\cap_{i=1}^n \{X_i \in A_i\}) &= \sum_{x_1 \in A_1, \dots, x_n \in A_n} \mathbb{P}(X_1 = x_1, \dots, X_n = x_n) \\ &= \sum_{x_1 \in A_1} \dots \sum_{x_n \in A_n} p_1(x_1) \dots p_n(x_n) \\ &= \sum_{x_1 \in A_1} p_1(x_1) \dots \sum_{x_n \in A_n} p_n(x_n) \\ &= \mathbb{P}(X_1 \in A_1) \dots \mathbb{P}(X_n \in A_n). \end{aligned}$$

□

Example 2.11. A coin flips heads with probability p and tails with probability $1 - p$. Let X be the number of heads and Y be the number of tails. These are *not* independent. (As for the details why, $\mathbb{P}(\{X = 1\} \cap \{Y = 1\}) = 0$ but $\mathbb{P}(X = 1)\mathbb{P}(Y = 1) = p(1 - p)$.)

Suppose that N is a Poisson random variable of parameter λ (it is independent of the coin). Then X and Y are independent! Indeed,

$$\begin{aligned} \mathbb{P}(X = x, Y = y) &= \mathbb{P}(X = x, Y = y | N = x + y) \mathbb{P}(N = x + y) \\ &= \binom{x+y}{x} p^x (1-p)^y \frac{\lambda^{x+y}}{(x+y)!} e^{-\lambda} \\ &= \frac{(\lambda p)^x (\lambda(1-p))^y}{x!y!} e^{-\lambda}. \end{aligned}$$

Since this has the form of $f(x)f(y)$, this means independence. To see this exactly,

$$\begin{aligned}\mathbb{P}(X = x) &= \sum_{y=0}^{\infty} \mathbb{P}(X = x, Y = y) \\ &= \frac{(\lambda p)^x}{x!} e^{-\lambda p} \sum_{y=0}^{\infty} \frac{(\lambda(1-p))^y}{y!} e^{-\lambda(1-p)} \\ &= \frac{(\lambda p)^x}{x!} e^{-\lambda p}.\end{aligned}$$

In particular, the number of heads and the number of tails are Poisson random variables of parameters λp and $\lambda(1-p)$, and they are independent of each other!

Lemma 2.12 (Convolution formula). *Suppose X, Y are independent discrete random variables that take values in \mathbb{Z} . Let p_X and p_Y be their pmfs. Then $Z = X + Y$ takes values in \mathbb{Z} , and its pmf is*

$$p_Z(z) = \sum_{k \in \mathbb{Z}} p_X(z - k) p_Y(k).$$

Proof. For any $z \in \mathbb{Z}$, the event $\{Z = z\}$ is equal to $\cup_{k \in \mathbb{Z}} \{X = z - k\} \cap \{Y = k\}$. These events in the union are disjoint, since X, Y cannot obtain two values simultaneously. So, by independence, we have

$$\begin{aligned}\mathbb{P}(Z = z) &= \mathbb{P}(\cup_{k \in \mathbb{Z}} \{X = z - k\} \cap \{Y = k\}) = \sum_{k \in \mathbb{Z}} \mathbb{P}(X = z - k, Y = k) \\ &= \sum_{k \in \mathbb{Z}} \mathbb{P}(X = z - k) \mathbb{P}(Y = k).\end{aligned}$$

□

Example 2.13. Take X_1, X_2 independent Bernoullis of parameter p (so $\mathbb{P}(X_1 = 1), \mathbb{P}(X_2 = 1) = p$). Let $Z = X_1 + X_2$. By the convolution formula and the fact that X_1, X_2 cannot attain values other than 0 and 1, we have

$$\begin{aligned}\mathbb{P}(Z = 0) &= \sum_{k \in \mathbb{Z}} \mathbb{P}(X = -k) \mathbb{P}(Y = k) = \mathbb{P}(X = 0) \mathbb{P}(Y = 0) = (1 - p)^2, \\ \mathbb{P}(Z = 1) &= \sum_{k \in \mathbb{Z}} \mathbb{P}(X = 1 - k) \mathbb{P}(Y = k) \\ &= \mathbb{P}(X = 1) \mathbb{P}(Y = 0) + \mathbb{P}(X = 0) \mathbb{P}(Y = 1) = 2p(1 - p), \\ \mathbb{P}(Z = 2) &= \sum_{k \in \mathbb{Z}} \mathbb{P}(X = 2 - k) \mathbb{P}(Y = k) = \mathbb{P}(X = 1) \mathbb{P}(Y = 1) = p^2.\end{aligned}$$

In particular, $Z \sim \text{Bin}(2, p)$!

Lemma 2.14. *If X, Y are independent, then so are $f(X)$ and $g(Y)$ (for any functions f, g).*

2.3. Expectation.

Definition 2.15. Let X be a discrete random variable with pmf p . Its *expectation* is $\mathbb{E}(X) = \sum_{x:p(x)>0} xp(x)$.

Lemma 2.16. (1) If $X \geq 0$ with probability 1, then $\mathbb{E}(X) \geq 0$.

(2) If $a, b \in \mathbb{R}$, then $\mathbb{E}(aX + bY) = a\mathbb{E}(X) + b\mathbb{E}(Y)$ (linearity of expectation; note that X, Y do not have to be independent!).

(3) If $X = c$ with probability 1, then $\mathbb{E}(X) = c$.

Proof. (1) We have $\mathbb{E}(X) = \sum_{x:p(x)>0} xp(x)$. Since $p(x) > 0$ only if $x \geq 0$ by assumption, we know $xp(x) \geq 0$, so $\mathbb{E}(X) \geq 0$.

(2) We have

$$\begin{aligned}
\mathbb{E}(aX + bY) &= \sum_z z \mathbb{P}(aX + bY = z) \\
&= \sum_z z \sum_w \mathbb{P}(aX + bY = z | Y = w) \mathbb{P}(Y = w) \\
&= \sum_z z \sum_w \mathbb{P}(aX + bw = z) \mathbb{P}(Y = w) \\
&= \sum_z z \sum_w \sum_s \mathbb{P}(aX + bw = z | X = s) \mathbb{P}(X = s) \mathbb{P}(Y = w) \\
&= \sum_{w,s} (as + bw) \mathbb{P}(X = s) \mathbb{P}(Y = w) \\
&= \sum_s \left(\sum_w (as + bw) \mathbb{P}(Y = w) \right) \mathbb{P}(X = s) \\
&= \sum_s (as + b\mathbb{E}(Y)) \mathbb{P}(X = s) \\
&= a\mathbb{E}(X) + b\mathbb{E}(Y).
\end{aligned}$$

(3) By definition, we have $\mathbb{E}(X) = \sum_{x:p(x)>0} xp(x)$. Only $x = c$ has $p(x) > 0$, so $\mathbb{E}(X) = cp(c) = c$ since $p(c) = 1$. □

Example 2.17. If $X \sim \text{Bern}(p)$, then $\mathbb{E}(X) = p$. If $X \sim \text{Bin}(n, p)$, then $X = Y_1 + \dots + Y_n$ where $Y_i \sim \text{Bern}(p)$, so $\mathbb{E}(X) = np$. If $X \sim \text{Pois}(\lambda)$, then

$$\begin{aligned}
\mathbb{E}(X) &= \sum_{k=0}^{\infty} \frac{k\lambda^k}{k!} e^{-\lambda} \\
&= \sum_{k=1}^{\infty} \frac{\lambda^k}{(k-1)!} e^{-\lambda} \\
&= \lambda \sum_{k=0}^{\infty} \frac{\lambda^k}{k!} e^{-\lambda} = \lambda.
\end{aligned}$$

Now, suppose X has pmf $p(k) = Ak^{-2}$ for $k \geq 1$ (where A is a “normalization constant”, so that $\sum_{k \geq 1} p(k) = 1$). Then $\mathbb{E}X = \sum_{k=1}^{\infty} Ak^{-1} = \infty$.

2.4. Variance and higher moments.

Definition 2.18. Given a random variable X , its *variance* is $\text{Var}(X) = \mathbb{E}(X - \mathbb{E}(X))^2$. Its k -th *moment* (for any $k \geq 0$) is $\mathbb{E}X^k$. We will often take k to be an integer.

Given any random variables X, Y , the *covariance* between X and Y is $\text{Cov}(X, Y) := \mathbb{E}[(X - \mathbb{E}(X))(Y - \mathbb{E}(Y))]$. In particular, we have $\text{Cov}(X, X) = \text{Var}(X)$. We say X, Y are *uncorrelated* if $\mathbb{E}(XY) = \mathbb{E}(X)\mathbb{E}(Y)$.

Lemma 2.19. (1) For any random variables X and Y , we have

$$\text{Var}(X) = \mathbb{E}(X^2) - (\mathbb{E}(X))^2, \quad \text{Cov}(X, Y) = \mathbb{E}(XY) - \mathbb{E}(X)\mathbb{E}(Y).$$

In particular, if X, Y are uncorrelated, then $\text{Cov}(X, Y) = 0$.

- (2) If X, Y are independent, then X, Y are uncorrelated.
(3) If X_1, \dots, X_n and Y_1, \dots, Y_n are random variables, and a_1, \dots, a_n and b_1, \dots, b_n are real numbers, then

$$\text{Cov}\left(\sum_{i=1}^n a_i X_i, \sum_{j=1}^n b_j Y_j\right) = \sum_{i,j=1}^n a_i b_j \text{Cov}(X_i, Y_j).$$

This is often called *bilinearity of the covariance*.

- (4) For any $a \in \mathbb{R}$, we have $\text{Var}(aX) = a^2 \text{Var}(X)$. (In words, the variance is “quadratic”.)
(5) There exists a constant c such that $X = c$ with probability 1 if and only if $\text{Var}(X) = 0$.

Proof. (1) By definition, we have $\text{Cov}(X, Y) = \mathbb{E}[XY - \mathbb{E}(X)Y - X\mathbb{E}(Y) + \mathbb{E}(X)\mathbb{E}(Y)] = \mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y]$, since $\mathbb{E}[\cdot]$ is always a constant (we also use linearity of expectation here). The formula for variance follows by taking $Y = X$.

- (2) If X, Y are independent, then

$$\begin{aligned} \mathbb{E}[XY] &= \sum_z z \mathbb{P}(XY = z) \\ &= \sum_z z \sum_w \mathbb{P}(XY = z | Y = w) \mathbb{P}(Y = w) \\ &= \sum_z z \sum_w \mathbb{P}(wX = z | Y = w) \mathbb{P}(Y = w) \\ &= \sum_z z \sum_w \mathbb{P}(Y = w) \mathbb{P}(wX = z) \\ &= \sum_z z \sum_w \mathbb{P}(Y = w) \sum_s \mathbb{P}(wX = z | X = s) \mathbb{P}(X = s) \\ &= \sum_{w,s} ws \mathbb{P}(Y = w) \mathbb{P}(X = s) \\ &= \sum_w w \mathbb{P}(Y = w) \sum_s \mathbb{P}(X = s) = \mathbb{E}(Y)\mathbb{E}(X). \end{aligned}$$

(3) We have

$$\begin{aligned}\mathbb{E} \left[\sum_{i=1}^n a_i X_i \sum_{j=1}^n b_j Y_j \right] &= \mathbb{E} \left[\sum_{i,j=1}^n a_i b_j X_i Y_j \right] \\ &= \sum_{i,j=1}^n a_i b_j \mathbb{E}[X_i Y_j]\end{aligned}$$

and

$$\begin{aligned}\mathbb{E} \left[\sum_{i=1}^n a_i X_i \right] \mathbb{E} \left[\sum_{j=1}^n b_j Y_j \right] &= \left\{ \sum_{i=1}^n a_i \mathbb{E}[X_i] \right\} \left\{ \sum_{j=1}^n b_j \mathbb{E}[Y_j] \right\} \\ &= \sum_{i,j=1}^n a_i b_j \mathbb{E}[X_i] \mathbb{E}[Y_j].\end{aligned}$$

Plug this into $\text{Cov}(\sum_{i=1}^n a_i X_i, \sum_{j=1}^n b_j Y_j) = \mathbb{E} \left[\sum_{i=1}^n a_i X_i \sum_{j=1}^n b_j Y_j \right] - \mathbb{E} \left[\sum_{i=1}^n a_i X_i \right] \mathbb{E} \left[\sum_{j=1}^n b_j Y_j \right]$ to get the formula.

- (4) Use part (3) with $n = 1$ and $a_1, b_1 = a$ and $X_1, Y_1 = X$.
(5) If $X = c$ with probability 1, then $\mathbb{E}(X) = c$ and $X - \mathbb{E}(X) = 0$ with probability 1. So $\mathbb{E}[(X - \mathbb{E}(X))^2] = \mathbb{E}[(c - c)^2] = 0$. If $\text{Var}(X) = \mathbb{E}[(X - \mathbb{E}(X))^2] = 0$, then $X = \mathbb{E}(X)$ with probability 1. Indeed, if $X = d$ for $d \neq \mathbb{E}(X)$ with positive probability p , since $(X - \mathbb{E}(X))^2 \geq 0$ with probability 1, we would get $\mathbb{E}[(X - \mathbb{E}(X))^2] \geq p(d - \mathbb{E}(X))^2 > 0$, a contradiction.

□

Example 2.20. Let $X \sim \text{Bern}(p)$. We saw before that $\mathbb{E}X = p$. Now, note that $X^2 = X$, since $X \in \{0, 1\}$, so that $\mathbb{E}X^2 = \mathbb{E}X = p$ as well. Thus, its variance is $\mathbb{E}X^2 - (\mathbb{E}X)^2 = p - p^2$. Now, assume that $X \sim \text{Pois}(\lambda)$. We saw that $\mathbb{E}X = \lambda$. We compute

$$\begin{aligned}\mathbb{E}X^2 &= \sum_{k=0}^{\infty} \frac{k^2 \lambda^k}{k!} e^{-\lambda} = e^{-\lambda} \sum_{k=1}^{\infty} \frac{k^2 \lambda^k}{k!} \\ &= e^{-\lambda} \sum_{k=1}^{\infty} \frac{k \lambda^k}{(k-1)!} \\ &= e^{-\lambda} \lambda \sum_{k=0}^{\infty} \frac{(k+1) \lambda^k}{k!} \\ &= e^{-\lambda} \lambda \frac{d}{d\lambda} \left(\sum_{k=0}^{\infty} \frac{\lambda^{k+1}}{k!} \right) \\ &= e^{-\lambda} \lambda \frac{d}{d\lambda} (\lambda e^{\lambda}) \\ &= \lambda^2 + \lambda.\end{aligned}$$

Hence, the variance of $X \sim \text{Pois}(\lambda)$ is $\text{Var}(X) = \mathbb{E}X^2 - (\mathbb{E}X)^2 = \lambda^2 + \lambda - \lambda^2 = \lambda$. Notice how this does not scale quadratically in λ !

2.5. Cauchy-Schwarz and Hölder inequalities.

Lemma 2.21. *Suppose X, Y are two random variables. Then for any $a > 0$, we have $|\mathbb{E}(XY)| \leq \frac{a^2\mathbb{E}(X^2)}{2} + \frac{\mathbb{E}(Y^2)}{2a^2}$. We also have $|\mathbb{E}(XY)| \leq (\mathbb{E}(X^2))^{1/2}(\mathbb{E}(Y^2))^{1/2}$.*

Proof. For the first inequality, we first note $(aX - \frac{1}{a}Y)^2 = a^2X^2 + \frac{Y^2}{a^2} - 2XY \geq 0$ (it is non-negative because it is the square of something). Thus, $XY \leq \frac{a^2X^2}{2} + \frac{Y^2}{2a^2}$. Now, take expectations to get $\mathbb{E}(XY) \leq \frac{a^2\mathbb{E}(X^2)}{2} + \frac{\mathbb{E}(Y^2)}{2a^2}$. In the case where $\mathbb{E}(XY) \geq 0$, this is the first claim. If $\mathbb{E}(XY) < 0$, use the claim after replacing X by $-X$. To prove the second claim, use the first claim for $a = \sqrt{2} \frac{\sqrt{\mathbb{E}(Y^2)}}{\sqrt{\mathbb{E}(X^2)}}$. \square

Lemma 2.22. *Suppose $p \in [1, \infty) \cup \{\infty\}$ and suppose $\frac{1}{p} + \frac{1}{q} = 1$. Then $|\mathbb{E}(XY)| \leq (\mathbb{E}|X|^p)^{1/p}(\mathbb{E}|Y|^q)^{1/q}$. (Note that if $p = q = 2$, this recovers Cauchy-Schwarz.)*

Proof. It suffices to instead use $XY \leq \frac{a^p|X|^p}{p} + \frac{|Y|^q}{a^q q}$ for any $a > 0$, take expectation, and choose a appropriately. \square