

Today we shall discuss a measure of how close a random variable tends to be to its expectation. But first we need to see how to compute the average of a product.

## Independence of Random Variables

Two random variables  $X$  and  $Y$  over the same probability space are independent if, for every possible values  $a, b$ , the events  $X = a$  and  $Y = b$  are independent.

**Theorem 20.1:** *If  $X$  and  $Y$  are independent random variables, then*

$$\mathbf{E}[XY] = \mathbf{E}[X] \cdot \mathbf{E}[Y] \quad (1)$$

**Proof:**

$$\begin{aligned} \mathbf{E}[XY] &= \sum_v v \cdot \Pr[XY = v] = \sum_v v \cdot \sum_{a,b:ab=v} \Pr[X = a \wedge Y = b] \\ &= \sum_{a,b} a \cdot b \cdot \Pr[X = a \wedge Y = b] = \sum_{a,b} a \cdot b \cdot \Pr[X = a] \cdot \Pr[Y = b] \\ &= \left( \sum_a a \cdot \Pr[X = a] \right) \cdot \left( \sum_b b \cdot \Pr[Y = b] \right) \\ &= \mathbf{E}[X] \cdot \mathbf{E}[Y] \end{aligned}$$

□

Notice how we used the assumption that the random variables are independent. If the random variables are not independent, then Equation (1) may fail.

Suppose for example that you are looking for a summer programming job, and that you have determined that there is a 50% chance that you will be offered \$70 an hour and a 50% chance that you will be offered \$50 an hour. Also, there is a 50% chance that you will be offered to work for 30 hours a week and a 50% chance that you will be offered to work for 40 hours a week. If the two random variables  $P$  (for hourly pay) and  $H$  (for weekly hours) were independent, then the expectation of your weekly pay would be  $\mathbf{E}[P \cdot H] = \mathbf{E}[P] \cdot \mathbf{E}[H] = 60 \cdot 35 = 2,100$ . If, however, the situation is that there is a 50% chance of being offered \$70 an hour for 30 hours, and a 50% chance of being offered \$50 an hour for 40 hours, then the average weekly pay is  $\frac{1}{2} \cdot 30 \cdot 70 + \frac{1}{2} \cdot 50 \cdot 40 = 2,050$ .

# Variance

**Question:** At each time step, I flip a fair coin. If it comes up Heads, I walk one step to the right; if it comes up Tails, I walk one step to the left. How far do I expect to have traveled from my starting point after  $n$  steps?

Denoting a right-move by  $+1$  and a left-move by  $-1$ , we can describe the probability space here as the set of all words of length  $n$  over the alphabet  $\{\pm 1\}$ , each having equal probability  $\frac{1}{2^n}$ . Let the r.v.  $X$  denote our position (relative to our starting point 0) after  $n$  moves. Thus

$$X = X_1 + X_2 + \cdots + X_n,$$

$$\text{where } X_i = \begin{cases} +1 & \text{if } i\text{th toss is Heads;} \\ -1 & \text{otherwise.} \end{cases}$$

Now obviously we have  $\mathbf{E}X = 0$ . The easiest rigorous way to see this is to note that  $\mathbf{E}X_i = (\frac{1}{2} \times 1) + (\frac{1}{2} \times (-1)) = 0$ , so by linearity of expectation  $\mathbf{E}X = \sum_{i=1}^n \mathbf{E}X_i = 0$ . Thus after  $n$  steps, my expected position is 0! But of course this is not very informative, and is due to the fact that positive and negative deviations from 0 cancel out.

What the above question is really asking is: What is the expected value of  $|X|$ , our *distance* from 0? Rather than consider the r.v.  $|X|$ , which is a little awkward due to the absolute value operator, we will instead look at the r.v.  $X^2$ . Notice that this also has the effect of making all deviations from 0 positive, so it should also give a good measure of the distance traveled. However, because it is the *squared* distance, we will need to take a square root at the end.

Let's calculate  $\mathbf{E}X^2$ :

$$\begin{aligned} \mathbf{E}X^2 &= \mathbf{E}(X_1 + X_2 + \cdots + X_n)^2 \\ &= \mathbf{E}\sum_{i=1}^n X_i^2 + \sum_{i \neq j} \mathbf{E}X_i X_j \\ &= \sum_{i=1}^n \mathbf{E}X_i^2 + \sum_{i \neq j} \mathbf{E}X_i X_j \end{aligned}$$

In the last line here, we used linearity of expectation. To proceed, we need to compute  $\mathbf{E}X_i^2$  and  $\mathbf{E}X_i X_j$  (for  $i \neq j$ ). Let's consider first  $X_i^2$ . Since  $X_i$  can take on only values  $\pm 1$ , clearly  $X_i^2 = 1$  always, so  $\mathbf{E}X_i^2 = 1$ . What about  $\mathbf{E}X_i X_j$ ? Since  $X_i$  and  $X_j$  are *independent*, it is the case that  $\mathbf{E}X_i X_j = \mathbf{E}X_i \mathbf{E}X_j = 0$ .

Plugging these values into the above equation gives

$$\mathbf{E}X^2 = (n \times 1) + 0 = n.$$

So we see that our expected squared distance from 0 is  $n$ . One interpretation of this is that we might expect to be a distance of about  $\sqrt{n}$  away from 0 after  $n$  steps. However, we have to be careful here: we **cannot** simply argue that  $\mathbf{E}|X| = \sqrt{\mathbf{E}X^2} = \sqrt{n}$ . (Why not?) We will see shortly (see Chebyshev's Inequality below) how to make precise deductions about  $|X|$  from knowledge of  $\mathbf{E}X^2$ .

For the moment, however, let's agree to view  $\mathbf{E}X^2$  as an intuitive measure of "spread" of the r.v.  $X$ . In fact, for a more general r.v. with expectation  $\mathbf{E}X = \mu$ , what we are really interested in is  $\mathbf{E}(X - \mu)^2$ , the expected squared distance *from the mean*. In our random walk example, we had  $\mu = 0$ , so  $\mathbf{E}(X - \mu)^2$  just reduces to  $\mathbf{E}X^2$ .

**Definition 20.1 (variance):** For a r.v.  $X$  with expectation  $\mathbf{E}X = \mu$ , the variance of  $X$  is defined to be

$$\text{Var}(X) = \mathbf{E}(X - \mu)^2.$$

The square root  $\sqrt{\text{Var}(X)}$  is called the standard deviation of  $X$ .

The point of the standard deviation is merely to “undo” the squaring in the variance. Thus the standard deviation is “on the same scale as” the r.v. itself. Since the variance and standard deviation differ just by a square, it really doesn’t matter which one we choose to work with as we can always compute one from the other immediately. We shall usually use the variance. For the random walk example above, we have that  $\text{Var}(X) = n$ , and the standard deviation of  $X$  is  $\sqrt{n}$ .

The following easy observation gives us a slightly different way to compute the variance that is simpler in many cases.

**Theorem 20.2:** For a r.v.  $X$  with expectation  $\mathbf{E}X = \mu$ , we have  $\text{Var}(X) = \mathbf{E}X^2 - \mu^2$ .

**Proof:** From the definition of variance, we have

$$\text{Var}(X) = \mathbf{E}(X - \mu)^2 = \mathbf{E}X^2 - 2\mu X + \mu^2 = \mathbf{E}X^2 - 2\mu\mathbf{E}X + \mu^2 = \mathbf{E}X^2 - \mu^2.$$

In the third step here, we used linearity of expectation.  $\square$

Let’s see some examples of variance calculations.

1. **Fair die.** Let  $X$  be the score on the roll of a single fair die. Recall from an earlier lecture that  $\mathbf{E}X = \frac{7}{2}$ . So we just need to compute  $\mathbf{E}X^2$ , which is a routine calculation:

$$\mathbf{E}X^2 = \frac{1}{6} (1^2 + 2^2 + 3^2 + 4^2 + 5^2 + 6^2) = \frac{91}{6}.$$

Thus from Theorem 22.1

$$\text{Var}(X) = \mathbf{E}X^2 - (\mathbf{E}X)^2 = \frac{91}{6} - \frac{49}{4} = \frac{35}{12}.$$

2. **Biased coin.** Let  $X$  be the number of Heads in  $n$  tosses of a biased coin with Heads probability  $p$  (i.e.,  $X$  has the binomial distribution with parameters  $n, p$ ). We already know that  $\mathbf{E}X = np$ . Writing as usual  $X = X_1 + X_2 + \dots + X_n$ , where  $X_i = \begin{cases} 1 & \text{if } i\text{th toss is Head;} \\ 0 & \text{otherwise} \end{cases}$  we have

$$\begin{aligned} \mathbf{E}X^2 &= \mathbf{E}(X_1 + X_2 + \dots + X_n)^2 \\ &= \sum_{i=1}^n \mathbf{E}X_i^2 + \sum_{i \neq j} \mathbf{E}X_i X_j \\ &= (n \times p) + (n(n-1) \times p^2) \\ &= n^2 p^2 + np(1-p). \end{aligned}$$

In the third line here, we have used the facts that  $\mathbf{E}X_i^2 = p$ , and that  $\mathbf{E}X_i X_j = \mathbf{E}X_i \mathbf{E}X_j = p^2$  (since  $X_i, X_j$  are independent). Note that there are  $n(n-1)$  pairs  $i, j$  with  $i \neq j$ .

Finally, we get that  $\text{Var}(X) = \mathbf{E}X^2 - (\mathbf{E}X)^2 = np(1-p)$ .<sup>1</sup> As an example, for a fair coin the expected number of Heads in  $n$  tosses is  $\frac{n}{2}$ , and the standard deviation is  $\frac{\sqrt{n}}{2}$ .

---

<sup>1</sup>Notice that in fact  $\text{Var}(X) = \sum_i \text{Var}(X_i)$ , and the same was true in the random walk example. This is in fact no coincidence, and it depends on the fact that the  $X_i$  are mutually independent. See Problem 1 on Homework 11. So *in the independent case* we can compute variances of sums very easily. Note that this isn’t the case, however, in example 4.

3. **Number of fixed points.** Let  $X$  be the number of fixed points in a random permutation of  $n$  items (i.e., the number of students in a class of size  $n$  who receive their own homework after shuffling). We saw in an earlier lecture that  $\mathbf{E}X = 1$  (regardless of  $n$ ). To compute  $\mathbf{E}X^2$ , write  $X = X_1 + X_2 + \dots + X_n$ ,

$$\text{where } X_i = \begin{cases} 1 & \text{if } i \text{ is a fixed point;} \\ 0 & \text{otherwise} \end{cases}$$

Then as usual we have

$$\mathbf{E}X^2 = \sum_{i=1}^n \mathbf{E}X_i^2 + \sum_{i \neq j} \mathbf{E}X_i X_j. \quad (2)$$

Since  $X_i$  is an indicator r.v., we have that  $\mathbf{E}X_i^2 = \Pr[X_i = 1] = \frac{1}{n}$ . In this case, however, we have to be a bit more careful about  $\mathbf{E}X_i X_j$ : note that we *cannot* claim as before that this is equal to  $\mathbf{E}X_i \mathbf{E}X_j$ , because  $X_i$  and  $X_j$  are not independent (why not?). But since both  $X_i$  and  $X_j$  are indicators, we can compute  $\mathbf{E}X_i X_j$  directly as follows:

$$\mathbf{E}X_i X_j = \Pr[X_i = 1 \wedge X_j = 1] = \Pr[\text{both } i \text{ and } j \text{ are fixed points}] = \frac{1}{n(n-1)}.$$

[Check that you understand the last step here.] Plugging this into equation (2) we get

$$\mathbf{E}X^2 = (n \times \frac{1}{n}) + (n(n-1) \times \frac{1}{n(n-1)}) = 1 + 1 = 2.$$

Thus  $\text{Var}(X) = \mathbf{E}X^2 - (\mathbf{E}X)^2 = 2 - 1 = 1$ . I.e., the variance and the mean are both equal to 1. Like the mean, the variance is also independent of  $n$ . Intuitively at least, this means that it is unlikely that there will be more than a small number of fixed points even when the number of items,  $n$ , is very large.

## Chebyshev's Inequality

We have seen that, intuitively, the variance (or, more correctly the standard deviation) is a measure of "spread", or deviation from the mean. Our next goal is to make this intuition quantitatively precise. What we can show is the following:

**Theorem 20.3: [Chebyshev's Inequality]** For a random variable  $X$  with expectation  $\mathbf{E}X = \mu$ , and for any  $\alpha > 0$ ,

$$\Pr[|X - \mu| \geq \alpha] \leq \frac{\text{Var}(X)}{\alpha^2}.$$

Before proving Chebyshev's inequality, let's pause to consider what it says. It tells us that the probability of any given deviation,  $\alpha$ , from the mean, either above it or below it (note the absolute value sign), is at most  $\frac{\text{Var}(X)}{\alpha^2}$ . As expected, this deviation probability will be small if the variance is small. An immediate corollary of Chebyshev's inequality is the following:

**Corollary 20.4:** For a random variable  $X$  with expectation  $\mathbf{E}X = \mu$ , and standard deviation  $\sigma = \sqrt{\text{Var}(X)}$ ,

$$\Pr[|X - \mu| \geq \beta\sigma] \leq \frac{1}{\beta^2}.$$

**Proof:** Plug  $\alpha = \beta\sigma$  into Chebyshev's inequality.  $\square$

So, for example, we see that the probability of deviating from the mean by more than (say) two standard deviations on either side is at most  $\frac{1}{4}$ . In this sense, the standard deviation is a good working definition of the “width” or “spread” of a distribution.

We should now go back and prove Chebyshev’s inequality. The proof will make use of the following simpler bound, which applies only to *non-negative* random variables (i.e., r.v.’s which take only values  $\geq 0$ ).

**Theorem 20.5: [Markov’s Inequality]** For a non-negative random variable  $X$  with expectation  $\mathbf{E}X = \mu$ , and any  $\alpha > 0$ ,

$$\Pr[X \geq \alpha] \leq \frac{\mathbf{E}X}{\alpha}.$$

**Proof:** From the definition of expectation, we have

$$\begin{aligned} \mathbf{E}X &= \sum_a a \times \Pr[X = a] \\ &\geq \sum_{a \geq \alpha} a \times \Pr[X = a] \\ &\geq \alpha \sum_{a \geq \alpha} \Pr[X = a] \\ &= \alpha \Pr[X \geq \alpha]. \end{aligned}$$

The crucial step here is the second line, where we have used the fact that  $X$  takes on only non-negative values. (Why is this step not valid otherwise?)  $\square$

Now we can prove Chebyshev’s inequality quite easily.

**Proof of Theorem 22.2:** Define the r.v.  $Y = (X - \mu)^2$ . Note that  $\mathbf{E}Y = \mathbf{E}(X - \mu)^2 = \text{Var}(X)$ . Also, notice that the probability we are interested in,  $\Pr[|X - \mu| \geq \alpha]$ , is exactly the same as  $\Pr[Y \geq \alpha^2]$ . (Why?) Moreover,  $Y$  is obviously non-negative, so we can apply Markov’s inequality to it to get

$$\Pr[Y \geq \alpha^2] \leq \frac{\mathbf{E}Y}{\alpha^2} = \frac{\text{Var}(X)}{\alpha^2}.$$

This completes the proof.  $\square$

Let’s apply Chebyshev’s inequality to answer our question about the random walk at the beginning of the lecture. Recall that  $X$  is our position after  $n$  steps, and that  $\mathbf{E}X = 0$ ,  $\text{Var}(X) = n$ . Corollary 22.3 says that, for any  $\beta > 0$ ,  $\Pr[|X| \geq \beta\sqrt{n}] \leq \frac{1}{\beta^2}$ . Thus for example, if we take  $n = 10^6$  steps, the probability that we end up more than 10000 steps away from our starting point is at most  $\frac{1}{100}$ .

Here are a few more examples of applications of Chebyshev’s inequality (you should check the algebra in them):

1. **Coin tosses.** Let  $X$  be the number of Heads in  $n$  tosses of a fair coin. The probability that  $X$  deviates from  $\mu = \frac{n}{2}$  by more than  $\sqrt{n}$  is at most  $\frac{1}{4}$ . The probability that it deviates by more than  $5\sqrt{n}$  is at most  $\frac{1}{100}$ .
2. **Fixed points.** Let  $X$  be the number of fixed points in a random permutation of  $n$  items; recall that  $\mathbf{E}X = \text{Var}(X) = 1$ . Thus the probability that more than (say) 10 students get their own homeworks after shuffling is at most  $\frac{1}{100}$ , however large  $n$  is.