

## 9.5: The Normal Distribution

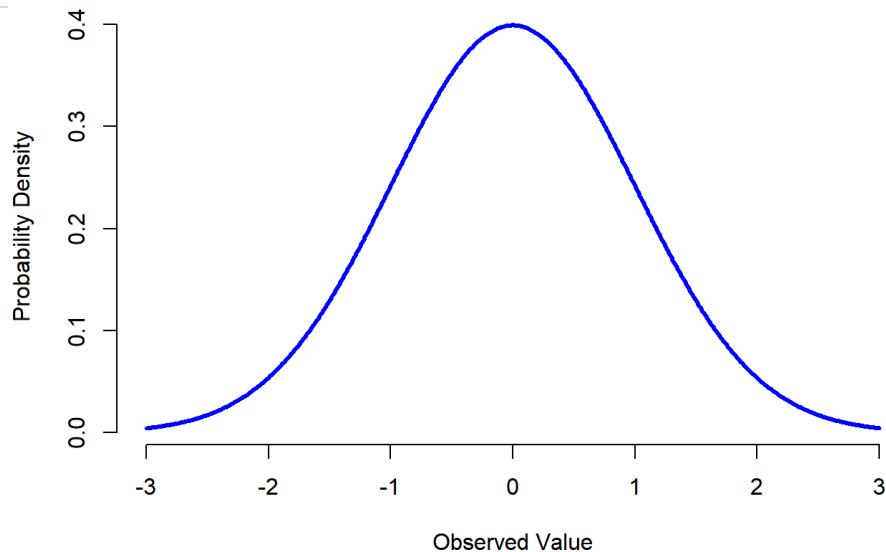


Figure 9.6: {The normal distribution with mean  $\mu=0$  and standard deviation  $\sigma=1$ . The x-axis corresponds to the value of some variable, and the y-axis tells us something about how likely we are to observe that value. However, notice that the y-axis is labelled “Probability Density” and not “Probability”. There is a subtle and somewhat frustrating characteristic of continuous distributions that makes the y axis behave a bit oddly: the height of the curve here isn’t actually the probability of observing a particular x value. On the other hand, it is true that the heights of the curve tells you which x values are more likely (the higher ones!).

While the binomial distribution is conceptually the simplest distribution to understand, it’s not the most important one. That particular honour goes to the **normal distribution**, which is also referred to as “the bell curve” or a “Gaussian distribution”. A normal distribution is described using two parameters, the mean of the distribution  $\mu$  and the standard deviation of the distribution  $\sigma$ . The notation that we sometimes use to say that a variable X is normally distributed is as follows:

$$X \sim \text{Normal}(\mu, \sigma)$$

Of course, that’s just notation. It doesn’t tell us anything interesting about the normal distribution itself. As was the case with the binomial distribution, I have included the formula for the normal distribution in this book, because I think it’s important enough that everyone who learns statistics should at least look at it, but since this is an introductory text I don’t want to focus on it, so I’ve tucked it away in Table 9.2. Similarly, the R functions for the normal distribution are `dnorm()`, `pnorm()`, `qnorm()` and `rnorm()`. However, they behave in pretty much exactly the same way as the corresponding functions for the binomial distribution, so there’s not a lot that you need to know. The only thing that I should point out is that the argument names for the parameters are `mean` and `sd`. In pretty much every other respect, there’s nothing else to add.

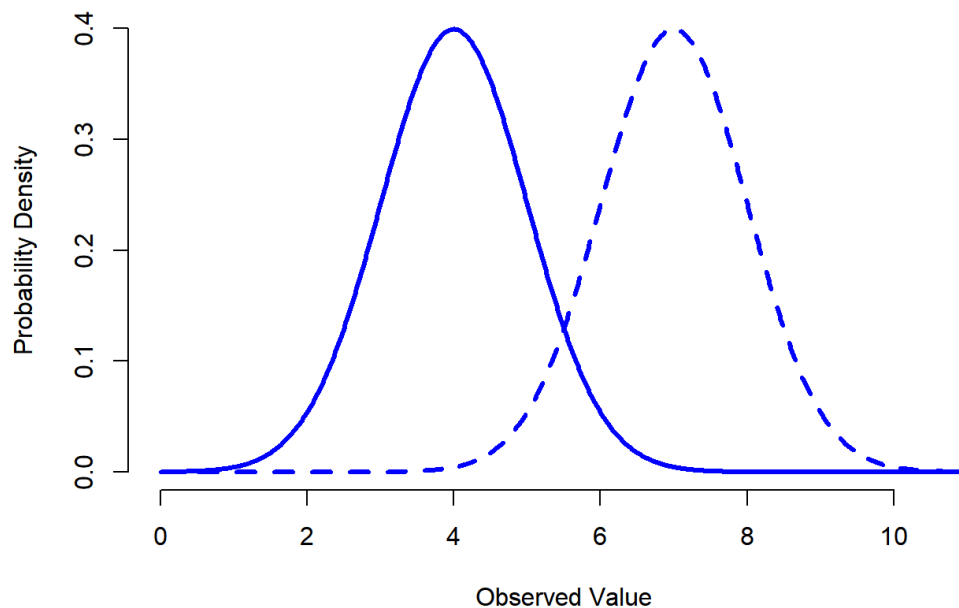


Figure 9.7: An illustration of what happens when you change the mean of a normal distribution. The solid line depicts a normal distribution with a mean of  $\mu=4$ . The dashed line shows a normal distribution with a mean of  $\mu=7$ . In both cases, the standard deviation is  $\sigma=1$ . Not surprisingly, the two distributions have the same shape, but the dashed line is shifted to the right.

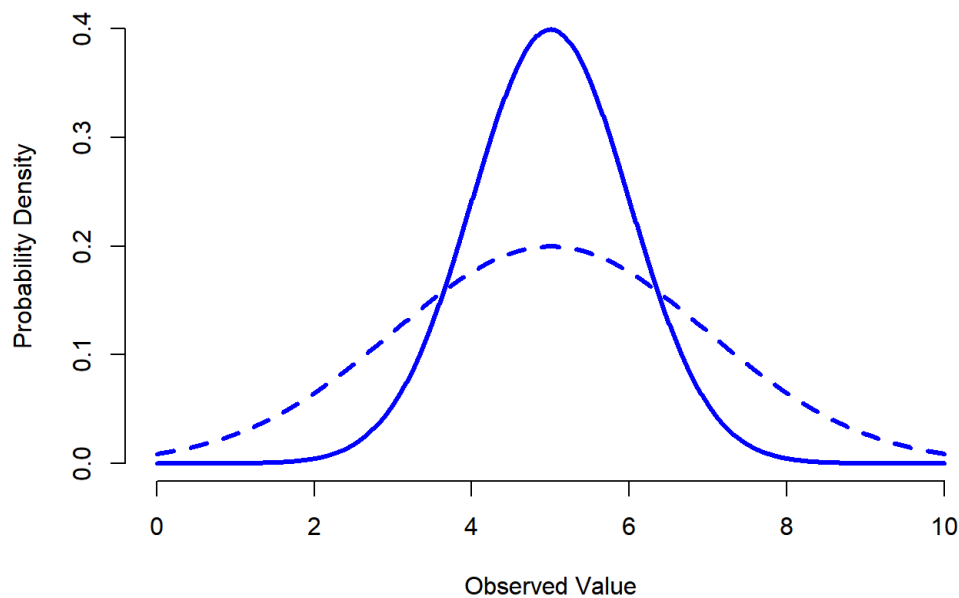


Figure 9.8: An illustration of what happens when you change the standard deviation of a normal distribution. Both distributions plotted in this figure have a mean of  $\mu=5$ , but they have different standard deviations. The solid line plots a distribution with standard deviation  $\sigma=1$ , and the dashed line shows a distribution with standard deviation  $\sigma=2$ . As a consequence, both distributions are “centred” on the same spot, but the dashed line is wider than the solid one.

Instead of focusing on the maths, let’s try to get a sense for what it means for a variable to be normally distributed. To that end, have a look at Figure 9.6, which plots a normal distribution with mean  $\mu=0$  and standard deviation  $\sigma=1$ . You can see where the name “bell curve” comes from: it looks a bit like a bell. Notice that, unlike the plots that I drew to illustrate the binomial distribution, the picture of the normal distribution in Figure 9.6 shows a smooth curve instead of “histogram-like” bars. This isn’t an arbitrary choice: the normal distribution is continuous, whereas the binomial is discrete. For instance, in the die rolling example from the last section, it was possible to get 3 skulls or 4 skulls, but impossible to get 3.9 skulls. The figures that I drew in the previous section reflected this fact: in Figure 9.3, for instance, there’s a bar located at  $X=3$  and another one at  $X=4$ , but there’s nothing in between. Continuous quantities don’t have this constraint. For instance, suppose we’re talking about the weather. The temperature on a pleasant Spring day could be 23 degrees, 24 degrees, 23.9 degrees, or anything in between since temperature is a continuous variable, and so a normal distribution might be quite appropriate for describing Spring temperatures.<sup>145</sup>

With this in mind, let's see if we can't get an intuition for how the normal distribution works. Firstly, let's have a look at what happens when we play around with the parameters of the distribution. To that end, Figure 9.7 plots normal distributions that have different means, but have the same standard deviation. As you might expect, all of these distributions have the same "width". The only difference between them is that they've been shifted to the left or to the right. In every other respect they're identical. In contrast, if we increase the standard deviation while keeping the mean constant, the peak of the distribution stays in the same place, but the distribution gets wider, as you can see in Figure 9.8. Notice, though, that when we widen the distribution, the height of the peak shrinks. This has to happen: in the same way that the heights of the bars that we used to draw a discrete binomial distribution have to *sum* to 1, the total *area under the curve* for the normal distribution must equal 1. Before moving on, I want to point out one important characteristic of the normal distribution. Irrespective of what the actual mean and standard deviation are, 68.3% of the area falls within 1 standard deviation of the mean. Similarly, 95.4% of the distribution falls within 2 standard deviations of the mean, and 99.7% of the distribution is within 3 standard deviations. This idea is illustrated in Figure ??.

Shaded Area = 68.3%

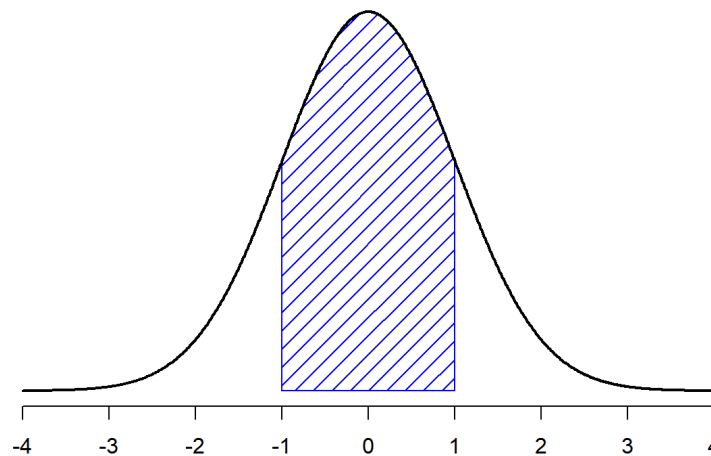


Figure 9.9: The area under the curve tells you the probability that an observation falls within a particular range. The solid lines plot normal distributions with mean  $\mu=0$  and standard deviation  $\sigma=1$ . The shaded areas illustrate "areas under the curve" for two important cases. Here we can see that there is a 68.3% chance that an observation will fall within one standard deviation of the mean

Shaded Area = 95.4%

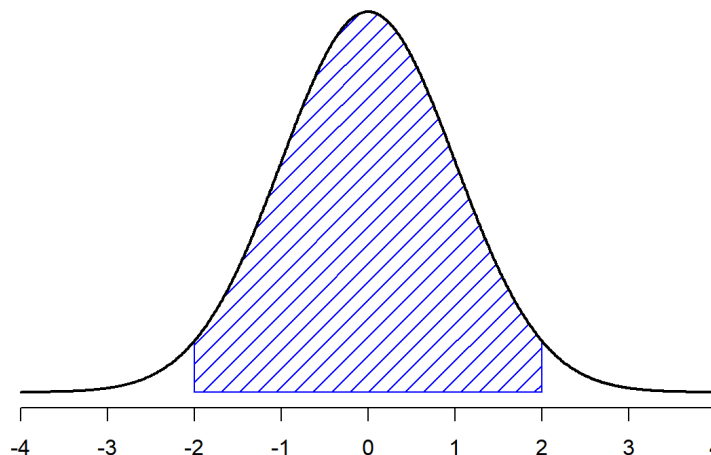


Figure 9.10: The area under the curve tells you the probability that an observation falls within a particular range. The solid lines plot normal distributions with mean  $\mu=0$  and standard deviation  $\sigma=1$ . The shaded areas illustrate "areas under the curve" for two important cases. Here we see that there is a 95.4% chance that an observation will fall within two standard deviations of the mean.

Shaded Area = 15.9%

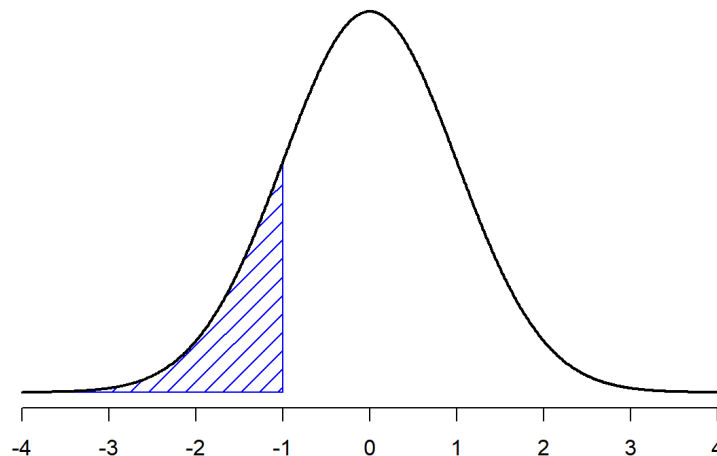


Figure 9.11: Two more examples of the “area under the curve idea”. There is a 15.9% chance that an observation is one standard deviation below the mean or smaller

Shaded Area = 34.1%

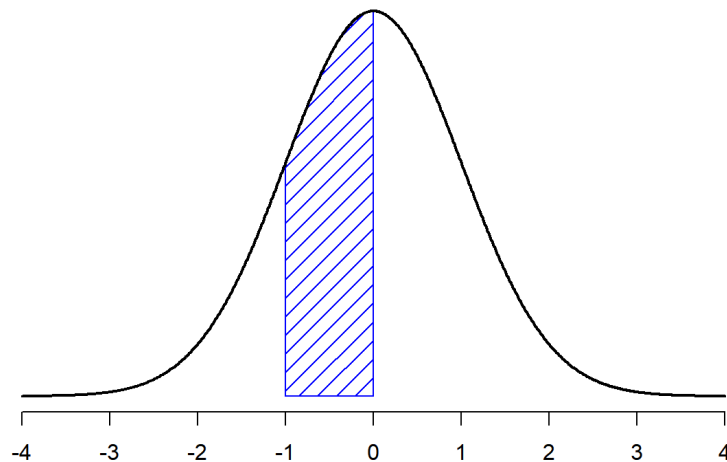


Figure 9.12: There is a 34.1% chance that the observation is greater than one standard deviation below the mean but still below the mean. Notice that if you add these two numbers together you get  $15.9 + 34.1 = 50$ . For normally distributed data, there is a 50% chance that an observation falls below the mean. And of course that also implies that there is a 50% chance that it falls above the mean.

### 9.5.1 Probability density

There’s something I’ve been trying to hide throughout my discussion of the normal distribution, something that some introductory textbooks omit completely. They might be right to do so: this “thing” that I’m hiding is weird and counterintuitive even by the admittedly distorted standards that apply in statistics. Fortunately, it’s not something that you need to understand at a deep level in order to do basic statistics: rather, it’s something that starts to become important later on when you move beyond the basics. So, if it doesn’t make complete sense, don’t worry: try to make sure that you follow the gist of it.

Throughout my discussion of the normal distribution, there’s been one or two things that don’t quite make sense. Perhaps you noticed that the y-axis in these figures is labelled “Probability Density” rather than density. Maybe you noticed that I used  $p(X)$  instead of  $P(X)$  when giving the formula for the normal distribution. Maybe you’re wondering why R uses the “d” prefix for functions like `dnorm()`. And maybe, just maybe, you’ve been playing around with the `dnorm()` function, and you accidentally typed in a command like this:

```
dnorm( x = 1, mean = 1, sd = 0.1 )
```

And if you’ve done the last part, you’re probably very confused. I’ve asked R to calculate the probability that  $x = 1$ , for a normally distributed variable with `mean = 1` and standard deviation `sd = 0.1`; and it tells me that the probability is 3.99. But, as we discussed earlier, probabilities *can’t* be larger than 1. So either I’ve made a mistake, or that’s not a probability.

As it turns out, the second answer is correct. What we’ve calculated here isn’t actually a probability: it’s something else. To understand what that something is, you have to spend a little time thinking about what it really *means* to say that  $X$  is a continuous variable. Let’s say we’re talking about the temperature outside. The thermometer tells me it’s 23 degrees, but I know that’s not really true. It’s not *exactly* 23 degrees. Maybe it’s 23.1 degrees, I think to myself. But I know that that’s not really true either, because it might actually be 23.09 degrees. But, I know that... well, you get the idea. The tricky thing with genuinely continuous quantities is that you never really know exactly what they are.

Now think about what this implies when we talk about probabilities. Suppose that tomorrow’s maximum temperature is sampled from a normal distribution with mean 23 and standard deviation 1. What’s the probability that the temperature will be *exactly* 23 degrees? The answer is “zero”, or possibly, “a number so close to zero that it might as well be zero”. Why is this? It’s like trying to throw a dart at an infinitely small dart board: no matter how good your aim, you’ll never hit it. In real life you’ll never get a value of exactly 23. It’ll always be something like 23.1 or 22.99998 or something. In other words, it’s completely meaningless to talk about the probability that the temperature is exactly 23 degrees. However, in everyday language, if I told you that it was 23 degrees outside and it turned out to be 22.9998 degrees, you probably wouldn’t call me a liar. Because in everyday language, “23 degrees” usually means something like “somewhere between 22.5 and 23.5 degrees”. And while it doesn’t feel very meaningful to ask about the probability that the temperature is exactly 23 degrees, it does seem sensible to ask about the probability that the temperature lies between 22.5 and 23.5, or between 20 and 30, or any other range of temperatures.

The point of this discussion is to make clear that, when we’re talking about continuous distributions, it’s not meaningful to talk about the probability of a specific value. However, what we *can* talk about is the probability that the value lies within a particular range of values. To find out the probability associated with a particular range, what you need to do is calculate the “area under the curve”. We’ve seen this concept already: in Figures 9.9 and (fig:sdnorm1b), the shaded areas shown depict genuine probabilities (e.g., in Figure 9.9 it shows the probability of observing a value that falls within 1 standard deviation of the mean).

Okay, so that explains part of the story. I’ve explained a little bit about how continuous probability distributions should be interpreted (i.e., area under the curve is the key thing), but I haven’t actually explained what the `dnorm()` function actually calculates. Equivalently, what does the formula for  $p(x)$  that I described earlier actually mean? Obviously,  $p(x)$  doesn’t describe a probability, but what is it? The name for this quantity  $p(x)$  is a **probability density**, and in terms of the plots we’ve been drawing, it corresponds to the *height* of the curve. The densities themselves aren’t meaningful in and of themselves: but they’re “rigged” to ensure that the *area* under the curve is always interpretable as genuine probabilities. To be honest, that’s about as much as you really need to know for now.<sup>146</sup>

---

This page titled [9.5: The Normal Distribution](#) is shared under a [CC BY-SA 4.0](#) license and was authored, remixed, and/or curated by [Danielle Navarro](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.