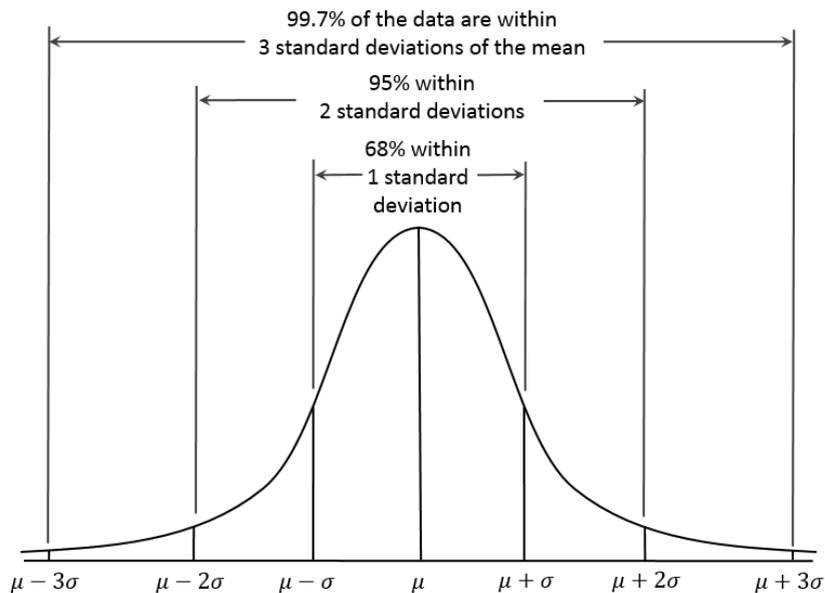


## Normal distribution and Z-scores

### 1. NORMAL DISTRIBUTION

Each random variable has a probability distribution associated to it. One of the most common distributions is the **normal distribution**. The normal distribution is a probability density function represented by a bell shaped curve centred at the population mean. The total area under the curve is 1. If a random variable,  $X$ , is normally distributed we write  $X \approx N(\mu, \sigma)$ , where  $\mu$  is the population mean and  $\sigma$  is standard deviation.



- The area under the curve is 1. This is meant to reflect the probability of obtain any one value in my data at random.
- The curve is symmetric, meaning 50% of the population lie below the mean and 50% lie above.
- The area enclosed within one standard deviation of the mean represents 68% of the population.
- The area enclosed within two standard deviations of the mean represents 95% of the population.
- The area enclosed within three standard deviations of the mean represents 99.7% of the population.

### 2. STANDARDISATION AND Z-SCORES

Simply put, standardisation allows us to realise the data in a more convenient way. This is done by repositioning the data relative the to the standard normal distribution, which is a normal distribution with mean 0 and standard deviation 1, i.e.,  $X \approx N(0, 1)$ . We can standardise any piece of data,  $x$ , using the **Z-score** formula:

$$z = \frac{x - \mu}{\sigma},$$

where  $\mu$  is the population mean and  $\sigma$  is the standard deviation.

Why do we standardise the data? Because it makes things easier! Once we have standardised the data, we can use tables, known as z-tables to simply read off the probability we need. The number relevant number in the z-table (at the end of the document) gives the area to the left of the z-score under the standard normal curve.

**Example 2.1. Question:** Suppose that the average daily calorie intake of British men (in kcal) is normal distributed, with mean 3,100 and standard deviation 125. Calculate the probability that the average daily calorie intake of a randomly-chosen male will be less than 3,190 kcal.

**Solution:**

**Step 1** Calculate the z-score of 3,200:

$$z = \frac{3,190 - 3,100}{125} = 0.72$$

**Step 2** Look up 0.72 in the z-table:

$$0.72 = 0.70 + 0.02.$$

So first we go to the 0.7 row. Then we go across that row until we get to the 0.02 column. That number is our probability. So the probability a randomly selected man's calorie intake is less than 3,190 is 0.7642. Or more mathematically:

$$P(x \leq 3,190) = 0.7642.$$

Alternatively if we wanted to know the probability of a randomly selected man's calorie intake is more than 3,190, we'd go through the same procedure except we would need to add one final step. For the probability of more than 3,190 is 1 minus the probability of less than 3,190! That is

$$P(x \geq 3,190) = 1 - 0.7642 = 0.2358.$$

**Exercise 2.2.** The average score for a class sitting a Physics test is 54, with standard deviation of 7. The average score for the same class in a Chemistry exam is 72, with standard deviation 10. Each set of scores is distributed normally. Which is more likely: that a student in the class

- (1) scores over 77 on the Chemistry test; or
- (2) scores over 58 on the Physics test?

