

The Normal Distribution

Introduction

Mass-produced items should conform to a specification. Usually, a mean is aimed for but due to random errors in the production process a tolerance is set on deviations from the mean. For example if we produce piston rings which have a target mean internal diameter of 45 mm then realistically we expect the diameter to deviate only slightly from this value. The deviations from the mean value are often modelled very well by the **normal distribution**. Suppose we decide that diameters in the range 44.95 mm to 45.05 mm are acceptable, then what proportion of the output is satisfactory? In this Section we shall see how to use the normal distribution to answer questions like this.

1. The normal distribution

The normal distribution is the most widely used model for the distribution of a random variable. There is a very good reason for this. Practical experiments involve measurements and measurements involve errors. However you go about measuring a quantity, inaccuracies of all sorts can make themselves felt. For example, if you are measuring a length using a device as crude as a ruler, you may find errors arising due to:

- the calibration of the ruler itself;
- parallax errors due to the relative positions of the object being measured, the ruler and your eye;
- rounding errors;
- 'guesstimation' errors if a measurement is between two marked lengths on the ruler.
- mistakes.

If you use a meter with a digital readout, you will avoid some of the above errors but others, often present in the design of the electronics controlling the meter, will be present. Errors are unavoidable and are usually the sum of several factors. The behaviour of variables which are the sum of several other variables is described by a very important and powerful result called the Central Limit Theorem which we will study later in this Workbook. For now we will quote the result so that the importance of the normal distribution will be appreciated.

The central limit theorem

Let X be the sum of n independent random variables $X_i, i = 1, 2, \dots, n$ each having a distribution with mean μ_i and variance σ_i^2 ($\sigma_i^2 < \infty$), respectively, then the distribution of X has expectation and variance given by the expressions

$$E(X) = \sum_{i=1}^n \mu_i \quad \text{and} \quad V(X) = \sum_{i=1}^n \sigma_i^2$$

and becomes **normal** as $n \rightarrow \infty$.

Essentially we are saying that a quantity which represents the combined effect of a number of variables will be approximately normal no matter what the original distributions are provided that $\sigma^2 < \infty$. This statement is true for the vast majority of distributions you are likely to meet in practice. This is why the normal distribution is crucially important to engineers. A quotation attributed to Prof. G. Lippmann, (1845-1921, winner of the Nobel prize for Physics in 1908) 'Everybody believes on the law of errors, experimenters because they think it is a mathematical theorem and mathematicians because they think it is an experimental fact.'

You may think that anything you measure follows an approximate normal distribution. Unfortunately this is not the case. While the heights of human beings follow a normal distribution, weights do not. Heights are the result of the interaction of many factors (outside one's control) while weights principally depend on lifestyle (including how much and what you eat and drink!) In practice, it is found that weight is skewed to the right but that the square root of human weights is approximately normal.

The probability density function of a normal distribution with mean μ and variance σ^2 is given by the formula

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-(x-\mu)^2/2\sigma^2}$$

This curve is always bell-shaped with the centre of the bell located at the value of μ . See Figure 1. The height of the bell is controlled by the value of σ . As with all normal distribution curves it is symmetrical about the centre and decays as $x \rightarrow \pm\infty$. As with any probability density function the area under the curve is equal to 1.

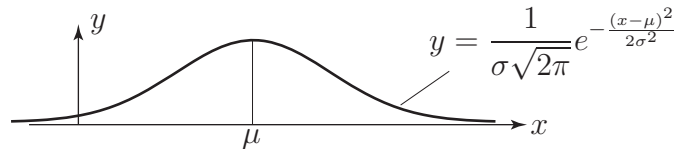


Figure 1

A normal distribution is completely defined by specifying its mean (the value of μ) and its variance (the value of σ^2 .) The normal distribution with mean μ and variance σ^2 is written $N(\mu, \sigma^2)$. Hence the distribution $N(20, 25)$ has a mean of 20 and a standard deviation of 5; remember that the second parameter is the variance which is the square of the standard deviation.

Key Point 1

A normal distribution has mean μ and variance σ^2 . A random variable X following this distribution is usually denoted by $N(\mu, \sigma^2)$ and we often write

$$X \sim N(\mu, \sigma^2)$$

Clearly, since μ and σ^2 can both vary, there are infinitely many normal distributions and it is impossible to give tabulated information concerning them all.

For example, if we produce piston rings which have a target mean internal diameter of 45 mm then we may realistically expect the actual diameter to deviate from this value. Such deviations are well-modelled by the normal distribution. Suppose we decide that diameters in the range 44.95 mm to 45.05 mm are acceptable, we may then ask the question 'What proportion of our manufactured output is satisfactory?'

Without tabulated data concerning the appropriate normal distribution we cannot easily answer this question (because the integral used to calculate areas under the normal curve is intractable.)

Since tabulated data allow us to apply the distribution to a wide variety of statistical situations, and we cannot tabulate all normal distributions, we tabulate only one - the standard normal distribution - and convert all problems involving the normal distribution into problems involving the standard normal distribution.

2. The standard normal distribution

At this stage we shall, for simplicity, consider what is known as a standard normal distribution which is obtained by choosing particularly simple values for μ and σ .

Key Point 2

The **standard normal distribution** has a mean of zero and a variance of one.

In Figure 2 we show the graph of the standard normal distribution which has probability density function $y = \frac{1}{\sqrt{2\pi}} e^{-x^2/2}$

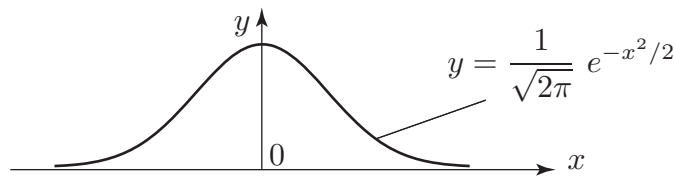


Figure 2: The standard normal distribution curve

The result which makes the standard normal distribution so important is as follows:

Key Point 3

If the behaviour of a continuous random variable X is described by the distribution $N(\mu, \sigma^2)$ then the behaviour of the random variable $Z = \frac{X - \mu}{\sigma}$ is described by the standard normal distribution $N(0, 1)$.

We call Z the **standardised normal variable** and we write

$$Z \sim N(0, 1)$$



Example 1

If the random variable X is described by the distribution $N(45, 0.000625)$ then what is the transformation required to obtain the standardised normal variable?



Example 2

When the random variable $X \sim N(45, 0.000625)$ takes values between 44.95 and 45.05, between which values does the random variable Z lie?



The random variable X follows a normal distribution with mean 1000 and variance 100. When X takes values between 1005 and 1010, between which values does the standardised normal variable Z lie?

3. Probabilities and the standard normal distribution

Since the standard normal distribution is used so frequently a table of values has been produced to help us calculate probabilities - located at the end of the Workbook. It is based upon the following diagram:

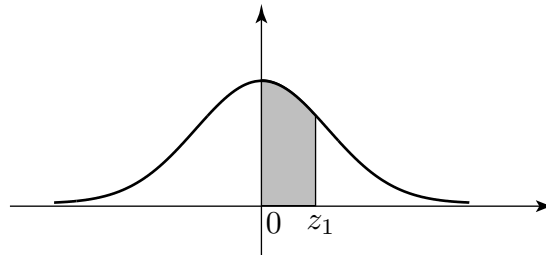


Figure 3

Since the total area under the curve is equal to 1 it follows from the symmetry in the curve that the area under the curve in the region $Z > 0$ is equal to 0.5. In Figure 3 the shaded area is the probability that Z takes values between 0 and z_1 . When we 'look-up' a value in the table we obtain the value of the shaded area.



Example 3

What is the probability that Z takes values between 0 and 1.9? (Refer to the table of normal probabilities at the end of the Workbook.)



Example 4

What is the probability that Z takes values between 0 and 1.96?



Example 5

What is the probability that Z takes values between 0 and 1.965?



What are the probabilities that Z takes values between

- (a) 0 and 2 (b) 0 and 2.3 (c) 0 and 2.33 (d) 0 and 2.333?

Note from Table 1 that as Z increases from 0 the entries increase, rapidly at first and then more slowly, toward 5000 i.e. a probability of 0.5. This is consistent with the shape of the curve.

After $Z = 3$ the increase is quite slow so that we tabulate entries for values of Z rising by increments of 0.1 instead of 0.01 as in the rest of Table 1.

4. Calculating other probabilities

In this Section we see how to calculate probabilities represented by areas other than those of the type shown in Figure 3.

Case 1

Figure 4 illustrates what we do if both Z values are positive. By using the properties of the standard normal distribution we can organise matters so that any required area is always of 'standard form'.

Figure 4



Example 6

Find the probability that Z takes values between 1 and 2.

Case 2

The following diagram illustrates the procedure to be followed when finding probabilities of the form $P(Z > z_1)$.

Figure 5



Example 7

What is the probability that $Z > 2$?

Case 3

Here we consider the procedure to be followed when calculating probabilities of the form $P(Z < z_1)$. Here the shaded area is the sum of the left-hand half of the total area and a 'standard' area.

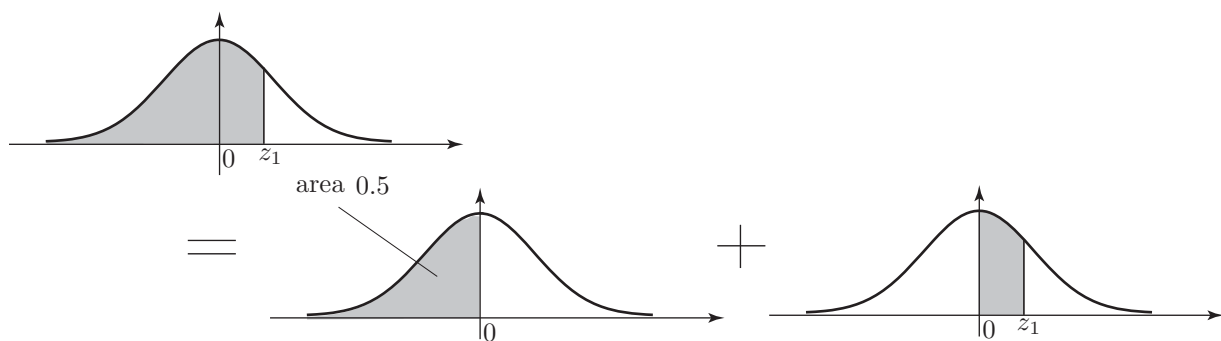


Figure 6



Example 8

What is the probability that $Z < 2$?

Case 4

Here we consider what needs to be done when calculating probabilities of the form $P(-z_1 < Z < 0)$ where z_1 is positive. This time we make use of the symmetry in the standard normal distribution curve.

Figure 7



Example 9

What is the probability that $-2 < Z < 0$?

Case 5

Finally we consider probabilities of the form $P(-z_2 < Z < z_1)$. Here we use the sum property and the symmetry property.

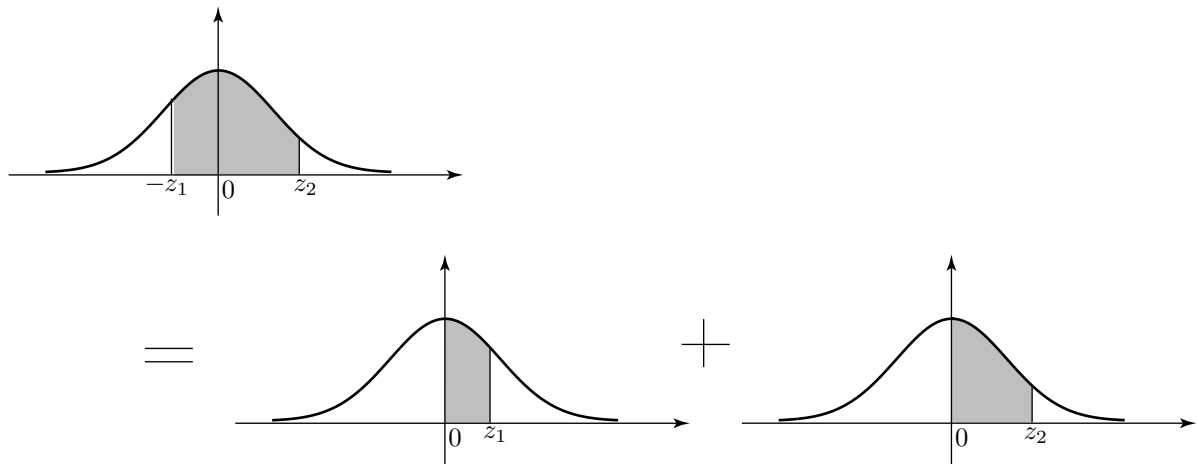


Figure 8



Example 10

What is the probability that $-1 < Z < 2$?



Find the following probabilities.

- | | |
|--------------------------|-------------------------|
| (a) $P(0 < Z < 1.5)$ | (b) $P(Z > 1.8)$ |
| (c) $P(1.5 < Z < 1.8)$ | (d) $P(Z < 1.8)$ |
| (e) $P(-1.5 < Z < 0)$ | (f) $P(Z < -1.5)$ |
| (g) $P(-1.8 < Z < -1.5)$ | (h) $P(-1.5 < Z < 1.8)$ |

(A simple sketch of the standard normal curve will help.)

5. The cumulative distribution function

We know that the normal probability density function $f(x)$ is given by the formula

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-(x-\mu)^2/2\sigma^2}$$

and so the cumulative distribution function $F(x)$ is given by the formula

$$F(x) = \frac{1}{\sigma\sqrt{2\pi}} \int_{-\infty}^x e^{-(u-\mu)^2/2\sigma^2} du$$

In the case of the cumulative distribution for the standard normal curve, we use the special notation $\Phi(z)$ and, substituting 0 and 1 for μ and σ^2 , we obtain

$$\Phi(z) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^z e^{-u^2/2} du$$

The shape of the curve is essentially 'S' -shaped as shown in Figure 9. Note that the curve runs from $-\infty$ to $+\infty$. As you can see, the curve approaches the value 1 asymptotically.

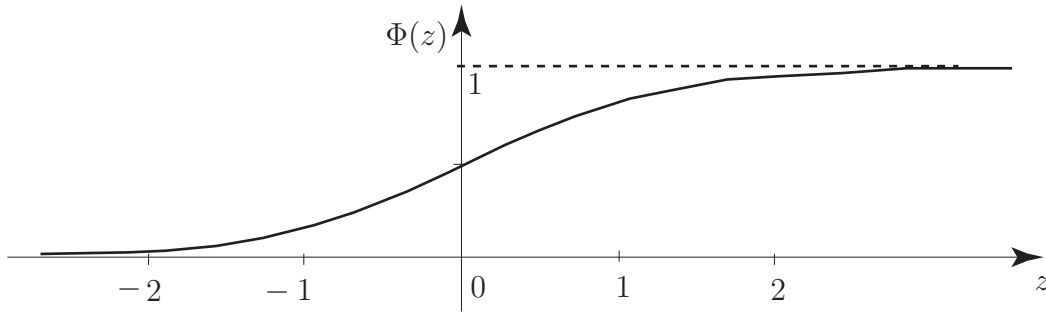


Figure 9

Comparing the integrals

$$F(x) = \frac{1}{\sigma\sqrt{2\pi}} \int_{-\infty}^x e^{-(u-\mu)^2/2\sigma^2} du \quad \text{and} \quad \Phi(z) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^z e^{-v^2/2} dv$$

shows that

$$v = \frac{u - \mu}{\sigma} \quad \text{and so} \quad dv = \frac{du}{\sigma}$$

and $F(x)$ may be written as

$$\begin{aligned} F(x) &= \frac{1}{\sigma\sqrt{2\pi}} \int_{-\infty}^{(x-\mu)/\sigma} e^{-v^2/2} \sigma dv \\ &= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{(x-\mu)/\sigma} e^{-v^2/2} dv = \Phi\left(\frac{x-\mu}{\sigma}\right) \end{aligned}$$

We already know, from the basic definition of a cumulative distribution function, that

$$P(a < X < b) = F(b) - F(a)$$

so that we may write the probability statement above in terms of $\Phi(z)$ as

$$P(a < X < b) = F(b) - F(a) = \Phi\left(\frac{b-\mu}{\sigma}\right) - \Phi\left(\frac{a-\mu}{\sigma}\right).$$

The value of $\Phi(z)$ is measured from $z = -\infty$ to any ordinate $z = z_1$ and represents the probability $P(Z < z_1)$.

The values of $\Phi(z)$ start as shown below:

z	0.00	0.01	0.02	0.03	0.04	0.05	0.06	0.07	0.08	0.09
0.0	.5000	5040	5080	5120	5160	5199	5239	5279	5319	5359
0.1	.5398	5438	5478	5517	5577	5596	5636	5675	5714	5753
0.2	.5793	5832	5871	5909	5948	5987	6026	6064	6103	6141

You should compare the values given here with the values given for the normal probability integral (Table 1 at the end of the Workbook). Simply adding 0.5 to the values in the latter table gives the values of $\Phi(z)$. You should also note that the diagrams shown at the top of each set of tabulated values tells you whether you are looking at the values of $\Phi(z)$ or the values of the normal probability integral.

Exercises

- If a random variable X has a standard normal distribution find the probability that it assumes a value:
 - less than 2.00
 - greater than 2.58
 - between 0 and 1.00
 - between -1.65 and -0.84
- If X has a standard normal distribution find k in each of the following cases:
 - $P(X < k) = 0.4$
 - $P(X < k) = 0.95$
 - $P(0 < X < k) = 0.1$

6. Applications of the normal distribution

We have, in the previous subsection, noted that the probability density function of a normal distribution X is

$$y = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

This curve is always ‘bell-shaped’ with the centre of the bell located at the value of μ . The height of the bell is controlled by the value of σ . See Figure 10.

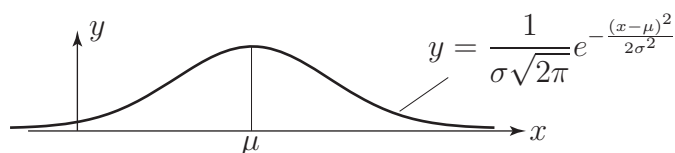


Figure 10

We now show, by example, how probabilities relating to a general normal distribution X are determined. We will see that being able to calculate the probabilities of a standard normal distribution Z is crucial in this respect.



Example 11

Given that the variate X follows the normal distribution $X \sim N(151, 15^2)$, calculate:

- (a) $P(120 \leq X \leq 155)$; (b) $P(X \geq 185)$

We note that, as for any continuous random variable, we can only calculate the probability that

- X lies between two given values;
- X is greater than a given value;
- X is less than a given value.

rather than for individual values.



A worn, poorly set-up machine is observed to produce components whose length X follows a normal distribution with mean 20 cm and variance 2.56 cm. Calculate:

- (a) the probability that a component is at least 24 cm long;
- (b) the probability that the length of a component lies between 19 and 21 cm.



Example 12

Piston rings are mass-produced. The target internal diameter is 45 mm but records show that the diameters are normally distributed with mean 45 mm and standard deviation 0.05 mm. An acceptable diameter is one within the range 44.95 mm to 45.05 mm. What proportion of the output is unacceptable?



Example 13

If the standard deviation is halved by improved production practices what is now the proportion of unacceptable items?



The resistance of a strain gauge is normally distributed with a mean of 100 ohms and a standard deviation of 0.2 ohms. To meet the specification, the resistance must be within the range 100 ± 0.5 ohms.

(a) What percentage of gauges are unacceptable?

(b) To what value must the standard deviation be reduced if the proportion of unacceptable gauges is to be no more than 0.2%?

First sketch the standard normal curve marking on it the lower and upper values z_1 and z_2 and appropriate areas:



7. Probability intervals - standard normal distribution

We use probability models to make predictions in situations where there is not sufficient data available to make a definite statement. Any statement based on these models carries with it a **risk** of being proved incorrect by events. Notice that the normal probability curve extends to infinity in both directions. **Theoretically** any value of the normal random variable is possible, although, of course, values far from the mean position (zero) are very unlikely.

Consider the diagram in Figure 13:

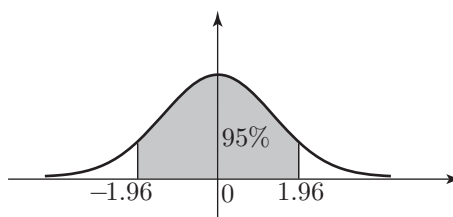


Figure 13

The shaded area is 95% of the total area. If we look at the entry in Table 1 (at the end of the Workbook) corresponding to $Z = 1.96$ we see the value 4750. This means that the probability of Z taking a value between 0 and 1.96 is 0.475. By symmetry, the probability that Z takes a value between -1.96 and 0 is also 0.475. Combining these results we see that

$$P(-1.96 < Z < 1.96) = 0.95 \text{ or } 95\%$$

We say that the 95% probability interval for Z (about its mean of 0) is $(-1.96, 1.96)$. It follows that there is a 5% chance that Z lies outside this interval.



Find the 99% probability interval for Z about its mean, i.e. the value of z_1 in the diagram:



Find the value of Z

- (a) which is exceeded on 5% of occasions
- (b) which is exceeded on 99% of occasions.

8. Probability intervals - general normal distribution

We saw in subsection 3 that 95% of the area under the standard normal curve lay between $z_1 = -1.96$ and $z_2 = 1.96$. Using the formula $Z = \frac{X - \mu}{\sigma}$ in the re-arrangement $X = \mu + Z\sigma$. We can see that 95% of the area under the general normal curve lies between $x_1 = \mu - 1.96\sigma$ and $x_2 = \mu + 1.96\sigma$.

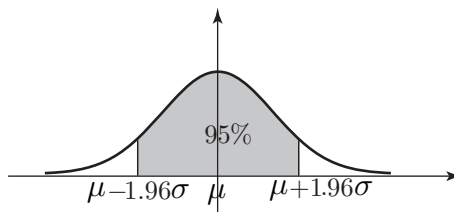


Figure 14



Example 14

Suppose that the internal diameters of mass-produced pipes are normally distributed with mean 50 mm and standard deviation 2 mm. What are the 95% probability limits on the internal diameter of a single pipe?



What is the 99% probability interval for the lifetime of a bulb when the lifetimes of such bulbs are normally distributed with a mean of 2000 hours and standard deviation of 40 hours?

The Normal Approximation to the Binomial Distribution

Introduction

We have already seen that the Poisson distribution can be used to approximate the binomial distribution for large values of n and small values of p provided that the correct conditions exist. The approximation is only of practical use if just a few terms of the Poisson distribution need be calculated. In cases where many - sometimes several hundred - terms need to be calculated the arithmetic involved becomes very tedious indeed and we turn to the normal distribution for help. It is possible, of course, to use high-speed computers to do the arithmetic but the normal approximation to the binomial distribution negates the necessity of this in a fairly elegant way. In the problem situations which follow this introduction the normal distribution is used to avoid very tedious arithmetic while at the same time giving a very good approximate solution.

1. The normal approximation to the binomial distribution

A typical problem

An engineering professional body estimates that 75% of the students taking undergraduate engineering courses are in favour of studying of statistics as part of their studies. If this estimate is correct, what is the probability that more than 780 undergraduate engineers out of a random sample of 1000 will be in favour of studying statistics?

Discussion

The problem involves a binomial distribution with a large value of n and so very tedious arithmetic may be expected. This can be avoided by using the normal distribution to approximate the binomial distribution underpinning the problem.

If X represents the number of engineering students in favour of studying statistics, then

$$X \sim B(1000, 0.75)$$

Essentially we are asked to find the probability that X is greater than 780, that is $P(X > 780)$.

The calculation is represented by the following statement

$$P(X > 780) = P(X = 781) + P(X = 782) + P(X = 783) + \cdots + P(X = 1000)$$

In order to complete this calculation we have to find all 220 terms on the right-hand side of the expression. To get some idea of just how big a task this is when the binomial distribution is used, imagine applying the formula

$$P(X = r) = \frac{n(n-1)(n-2) \cdots (n-r+1)p^r(1-p)^{n-r}}{r(r-1)(r-2) \cdots 3.2.1}$$

220 times! You would have to take $n = 1000$, $p = 0.75$ and vary r from 781 to 1000. Clearly, the task is enormous.

Fortunately, we can approximate the answer very closely by using the normal distribution with the *same mean and standard deviation* as $X \sim B(1000, 0.75)$. Applying the usual formulae for μ and σ we obtain the values $\mu = 750$ and $\sigma = 13.7$ from the binomial distribution.

We now have two distributions, $X \sim B(1000, 0.75)$ and (say) $Y \sim N(750, 13.7^2)$. Remember that the second parameter represents the variance. By doing the appropriate calculations, (this is extremely tedious even for one term!) it can be shown that

$$P(X = 781) \approx P(780.5 \leq Y \leq 781.5)$$

This statement means that the probability that $X = 781$ calculated from the binomial distribution $X \sim B(1000, 0.75)$ can be very closely approximated by the area under the normal curve $Y \sim N(750, 13.7^2)$ between 780.5 and 781.5. This relationship is then applied to all 220 terms involved in the calculation.

The result is summarised below:

$$\begin{aligned}
 P(X = 781) &\approx P(780.5 \leq Y \leq 781.5) \\
 P(X = 782) &\approx P(781.5 \leq Y \leq 782.5) \\
 &\vdots \\
 P(X = 999) &\approx P(998.5 \leq Y \leq 999.5) \\
 P(X = 1000) &\approx P(999.5 \leq Y \leq 1000.5)
 \end{aligned}$$

By adding these probabilities together we get

$$\begin{aligned}
 P(X > 780) &= P(X = 781) + P(X = 782) + \cdots + P(X = 1000) \\
 &\approx P(780.5 \leq Y \leq 1000.5)
 \end{aligned}$$

To complete the calculation we need only to find the area under the curve $Y \sim N(750, 13.7^2)$ between the values 780.5 and 1000.5. This is far easier than completing the 220 calculations suggested by the use of the binomial distribution.

Finding the area under the curve $Y \sim N(750, 13.7^2)$ between the values 780.5 and 1000.5 is easily done by following the procedure used previously. The calculation, using the tables on page 15 and working to three decimal places, is

$$\begin{aligned}
 P(X > 780) &\approx P\left(\frac{780.5 - 750}{13.7} \leq Z \leq \frac{1000.5 - 750}{13.7}\right) \\
 &= P(2.23 \leq Z \leq 18.28) \\
 &= P(Z \geq 2.23) \\
 &= 0.013
 \end{aligned}$$

Notes:

1. Since values as high as 18.28 effectively tell us to find the area to the right of 2.33 (the area to the right of 18.28 is so close to zero as to make no difference) we have

$$P(Z \geq 2.23) = 0.0129 \approx 0.013$$

2. The solution given *assumes* that the original binomial distribution can be approximated by a normal distribution. This is not always the case and you must always check that the following conditions are satisfied before you apply a normal approximation. The conditions are:

- $np > 5$
- $n(1 - p) > 5$

You can see that these conditions are satisfied here.



A particular production process used to manufacture ferrite magnets used to operate reed switches in electronic meters is known to give 10% defective magnets on average. If 200 magnets are randomly selected, what is the probability that the number of defective magnets is between 24 and 30?



Example 15

Overbooking of passengers on intercontinental flights is a common practice among airlines. Aircraft which are capable of carrying 300 passengers are booked to carry 320 passengers. If on average 10% of passengers who have a booking fail to turn up for their flights, what is the probability that at least one passenger who has a booking will end up without a seat on a particular flight?

Table 1: The Standard Normal Probability

$Z = \frac{x-\mu}{\sigma}$	0	1	2	3	4	5	6	7	8	9
0	0000	0040	0080	0120	0160	0199	0239	0279	0319	0359
.1	0398	0438	0478	0517	0577	0596	0636	0675	0714	0753
.2	0793	0832	0871	0909	0948	0987	1026	1064	1103	1141
.3	1179	1217	1255	1293	1331	1368	1406	1443	1480	1517
.4	1555	1591	1628	1664	1700	1736	1772	1808	1844	1879
.5	1915	1950	1985	2019	2054	2088	2123	2157	2190	2224
.6	2257	2291	2324	2357	2389	2422	2454	2486	2517	2549
.7	2580	2611	2642	2673	2703	2734	2764	2794	2822	2852
.8	2881	2910	2939	2967	2995	3023	3051	3078	3106	3133
.9	3159	3186	3212	3238	3264	3289	3315	3340	3365	3389
1.0	3413	3438	3461	3485	3508	3531	3554	3577	3599	3621
1.1	3643	3665	3686	3708	3729	3749	3770	3790	3810	3830
1.2	3849	3869	3888	3907	3925	3944	3962	3980	3997	4015
1.3	4032	4049	4066	4082	4099	4115	4131	4147	4162	4177
1.4	4192	4207	4222	4236	4251	4265	4279	4292	4306	4319
1.5	4332	4345	4357	4370	4382	4394	4406	4418	4429	4441
1.6	4452	4463	4474	4484	4495	4505	4515	4525	4535	4545
1.7	4554	4564	4573	4582	4591	4599	4608	4616	4625	4633
1.8	4641	4649	4656	4664	4671	4678	4686	4693	4699	4706
1.9	4713	4719	4726	4732	4738	4744	4750	4756	4761	4767
2.0	4772	4778	4783	4788	4793	4798	4803	4808	4812	4817
2.1	4821	4826	4830	4834	4838	4842	4846	4850	4854	4857
2.2	4861	4865	4868	4871	4875	4878	4881	4884	4887	4890
2.3	4893	4896	4898	4901	4904	4906	4909	4911	4913	4916
2.4	4918	4920	4922	4925	4927	4929	4931	4932	4934	4936
2.5	4938	4940	4941	4943	4946	4947	4948	4949	4951	4952
2.6	4953	4955	4956	4957	4959	4960	4961	4962	4963	4964
2.7	4965	4966	4967	4968	4969	4970	4971	4972	4973	4974
2.8	4974	4975	4976	4977	4977	4978	4979	4979	4980	4981
2.9	4981	4982	4982	4983	4984	4984	4985	4985	4986	4986
	3.0	3.1	3.2	3.3	3.4	3.5	3.6	3.7	3.8	3.9
	4987	4990	4993	4995	4997	4998	4998	4999	4999	4999