# the normal distribution

- 1. The normal distribution
- 2. The normal approximation to the binomial distribution
- 3. Sums and differences of random variables

# Learning outcomes

In a previous Workbook you learned what a continuous random variable was. Here you will examine the most important example of a continuous random variables: the normal distribution. The probabilities of the normal distribution have to be determined numerically. Tables of such probabilities, referring to a simplified normal distribution, the standard normal distribution, which has mean 0 and variance 1 will be used to determine probabilities of the general normal distribution. Finally you will learn how to deal with combinations of random variables which is an important statistical tool applicable to many engineering situations.

# Time allocation

You are expected to spend approximately five hours of independent study on the material presented in this workbook. However, depending upon your ability to concentrate and on your previous experience with certain mathematical topics this time may vary considerably.

# The Normal Distribution





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.

Prerequisites	① be familiar with the basic properties of probability
Before starting this Section you should	② be familiar with continuous random variables
After completing this Section you should be able to	<ul> <li>✓ recognise the shape of the frequency curve for the normal distribution and the standard normal distribution</li> <li>✓ be able to calculate probabilities using</li> </ul>
	<ul><li>the standard normal distribution</li><li>✓ recognise key areas under the frequency</li></ul>

curve

# 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 rule 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, ..., n each having a distribution with mean  $\mu_i$  and variance  $\sigma_i^2$  ( $\sigma_i^2 < \infty$ ), respectively, then the distribution of X has expectation and variance given by the expressions

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

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

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) summarizes the situation:

'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 how much 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  $\mu$  and variance  $\sigma^2$  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  $\mu$ . The height of the bell is controlled by the value of  $\sigma$ . As with all normal distribution curves it is symmetrical about the centre and decays exponentially as  $x \to \pm \infty$ . As with any probability density function the area under the curve is equal to 1. See Figure 1.



Figure 1

A normal distribution is completely defined by specifying its mean (the value of  $\mu$ ) and its variance (the value of  $\sigma^2$ ). The normal distribution with mean  $\mu$  and variance  $\sigma^2$  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.



A normal distribution has mean  $\mu$  and variance  $\sigma^2$ . 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  $\mu$  and  $\sigma^2$  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  $\mu$  and  $\sigma$ .



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 bistribution 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:



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** 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?

#### Solution

Here,  $\mu = 45$  and  $\sigma^2 = 0.000625$  so that  $\sigma = 0.025$ . Hence Z = (X - 45)/0.025 is the required transformation.

**Example** When the random variable X takes values between 44.95 and 45.05, between which values does the random variable Z lie?

#### Solution

When X = 45.05,  $Z = \frac{45.05-45}{0.025} = 2$ When X = 44.95,  $Z = \frac{44.95-45}{0.025} = -2$ Hence Z lies between -2 and 2.



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?

Your solution	
	I bas č.0 n sewted z sil $Z$ sona H
	$1 = \frac{01}{01} = Z$ , $0101 = X$ nəhw
	$\ddot{c}.0 = \frac{\ddot{c}}{01} = Z$ , $\ddot{c}001 = X$ nəhw
	The transformation is $Z = \frac{0001-X}{10}$ .

# 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 this Section. 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 x > 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** What is the probability that Z takes values between 0 and 1.9? (Please refer to the table of normal probabilities on page 15).

#### Solution

The row beginning '1.9' and the column headed '0' is the appropriate choice and its entry is 4713. This is to be read as 0.4713 (we omitted the '0.' in each entry for clarity) The interpretation is that the probability that Z takes values between 0 and 1.9 is 0.4713.

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

#### Solution

This time we want the row beginning 1.9 and the column headed '6'. The entry is 4750 so that the required probability is 0.4750.

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

#### Solution

There is no entry corresponding to 1.965 so we take the average of the values for 1.96 and 1.97. (This linear interpolation is not strictly correct but is acceptable).

The two values are 4750 and 4756 with an average of 4753. Hence the required probability is 0.4753.



Your solution

What are the probabilities that Z takes values between (i) 0 and 2 (ii) 0 and 2.3 (iii) 0 and 2.33 (iv) 0 and 2.333?

(i) The entry is 4893; the probability is 0.4772.
(ii) The entry is 4893; the probability is 0.4903.
(iii) The entry is 4901; the probability is 0.4904.
(iv) The entry for 2.33 is 4904, that for 2.34 is 4904.
Linear interpolation gives a value of 4901 + 0.3(4904 - 4901) i.e. about 4902; the probability is 0.4902.

Note from the table 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 0.1 instead of 0.01 as in the rest of the table.

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

#### Solution

Using the table  $P(Z = z_2)$  i.e. P(Z = 2) is 0.4772  $P(Z = z_1)$  i.e. P(Z = 1) is 0.3413. Hence P(1 < Z < 2) = 0.4772 - 0.3413 = 0.1359(Remember that with a continuous distribution, P(Z = 1) is meaningless so that  $P(1 \le Z \le 2)$ is interpreted as P(1 < Z < 2).

#### Case 2

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





**Example** What is the probability that Z > 2?

#### Solution

P(0 < Z < 2) = 0.4772 (from the table) Hence the probability is 0.5 - 0.4772 = 0.0228.

#### 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** What is the probability that Z < 2?

Solution P(Z > 2) = 0.5 + 0.4772 = 0.9772.

#### 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.



by symmetry this shaded area is equal in value to the one above.

Figure 7

**Example** What is the probability that -2 < Z < 0?

#### Solution

The area is equal to that corresponding to P(0 < Z < 2) = 0.4772.

#### 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.





**Example** What is the probability that -1 < Z < 2?

Solution

$$P(-1 < Z < 0) = P(0 < Z < 1) = 0.3413$$
$$P(0 < Z < 2) = 0.4772$$

Hence the required probability, P(-1 < Z < 2) is 0.8185.

Other cases can be handed by a combination of the ideas already used.



Find the following probabilities.

(i)	P(0 < Z < 1.5)	(ii)	P(Z > 1.8)
(iii)	P(1.5 < Z < 1.8)	(iv)	P(Z < 1.8)
(v)	P(-1.5 < Z < 0)	(vi)	P(Z < -1.5)
(vii)	P(-1.8 < Z < -1.5)	(viii)	P(-1.5 < Z < 1.8)

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

Aon. solution Aon.

## 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  $\mu$  and  $\sigma^2$ , 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 the diagram below. Note that the

HELM (VERSION 1: April 9, 2004): Workbook Level 1 39.1: The Normal Distribution

curve runs from  $-\infty$  to  $+\infty$ . As you can see, the curve approaches the value 1 asymptotically.



Comparing the integrals

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

shows that

$$\nu = \frac{u - \mu}{\sigma} \quad \text{and so} \quad d\nu = \frac{du}{\sigma}$$

and F(x) may be written as

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

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

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

Some values of  $\Phi(z)$  are given in the table below. You should compare the values given here with the values given for the normal probability integral on page 15. 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.

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

$Z = \frac{x-\mu}{\sigma}$	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
<u> </u>	3.0	3.1	3.2	3.3	34	3.5	3.6	37	3.8	39
	4987	4990	4993	4995	4997	4998	4998	4999	4999	4999

The Standard Normal Probability Integral

#### **Exercises**

- 1. If a random variable X has a standard normal distribution find the probability that it assumes a value:
  - (a) less than 2.00 (b) greater than 2.58 (c) between 0 and 1.00
  - (d) between -1.65 and -0.84
- 2. If X has a standard normal distribution find k in each of the following cases:

(a) P(X < k) = 0.4 (b) P(X < k) = 0.95 (c) P(0 < X < k) = 0.1

Answers 1. Standard normal probabilities found using tables. 2. ditto

### 6. Applications of the normal distribution

We have, in the previous Section, 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  $\mu$ . The height of the bell is controlled by the value of  $\sigma$ . See Figure 1.



Figure 1

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

(a) 
$$P(120 \le X \le 155);$$
  
(b)  $P(X \ge 185)$ 

#### Solution

The transformation used in this problem is  $Z = \frac{X - \mu}{\sigma} = \frac{X - 151}{15}$  (a)

$$P(120 \le X \le 155) = P(\frac{120 - 151}{15} \le Z \le \frac{155 - 151}{15})$$
  
=  $P(-2.07 \le Z \le 0.27)$   
=  $0.4808 + 0.1064$   
=  $0.5872$ 

#### Solution (contd.)

(b)

$$P(X \ge 185) = P(Z \ge \frac{185 - 151}{15})$$
  
=  $P(Z \ge 2.27)$   
=  $0.5 - 0.4884$   
=  $0.0116$ 

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 that 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 cm. and 21 cm.

 $1804.0 = (320.0 > Z > 320.0) = (\frac{02-12}{0.1} > Z > \frac{02-01}{0.1}) = (12 > X > 01) = (12 > 12) =$ 

3000.0 = 8604.0 - 3.0 = (3.2 
$$\leq Z$$
) $T = (\frac{02 - 4Z}{3.1}) \leq Z$  si besu noitemroisment ed  
T  $(\frac{02 - 4Z}{3.1}) \leq Z$ ) $T = (4.2 \leq X)T$ 

pue

**Example** 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? (There are many words in the statement of the problem; we must read them carefully to extract the necessary information.)

#### Solution

Let X be the diameter of a piston ring. Then we write  $X \sim N(45, (0.05)^2)$ . The transformation is  $Z = \frac{X-\mu}{\sigma} = \frac{X-45}{0.05}$ . The upper limit of acceptability is  $x_2 = 45.05$  so that  $z_2 = \frac{45.05-45}{0.05} = 1$ . The lower limit of acceptability is  $x_1 = 44.95$  so that  $z_1 = \frac{44.95-45}{0.05} = -1$ . The range of 'acceptable' Z values is therefore -1 to 1. See Figure 2.



Figure 2

Using the symmetry of the curves

$$P(-1 < Z < 1) = 2 \times P(0 < Z < 1)$$
  
= 2 × 0.3413  
= 0.6826.

Thus the proportion of unacceptable items is 1 - 0.6826 = 0.3174, or 31.74%.

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

#### Solution

Now  $\sigma = 0.025$  so that:

$$z_2 = \frac{45.05 - 45}{0.025} = 2$$
 and  $z_1 = -2$ 

Then  $P(-2 < Z < 2) = 2 \times P(0 < Z < 2) = 2 \times 0.4772 = 0.9544$ . Hence the proportion of unacceptable items is reduced to 1 - 0.9544 = 0.0456 or 4.56%.

We observe that less of the area under the curve now lies outside the interval (44.95, 45.05).



Figure 3



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  $100\pm0.5$  ohms. What percentage of gauges are unacceptable?

First, state the upper and lower limits of acceptable resistance and find the Z-values which correspond.



Using a suitable sketch, calculate the probability that  $z_1 < Z < z_2$ .

Your solution





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.

Your solution



Now use the Table to find  $z_2$ , and hence write down the value of  $z_1$ .



Rewrite the formula  $Z = \frac{X-\mu}{\sigma}$  to make  $\sigma$  the subject. Put in values for  $z_2$ ,  $x_2$  and  $\mu$  hence evaluate  $\sigma$ .

#### Your solution

(.q.b 2) 
$$\partial 1.0 = \frac{001 - \delta.001}{1.6} = \frac{\mu - X}{Z} = 0$$