# INTRODUCTION

# TO

# STATISTICAL ENGINEERING

*S.J. Morrison*

# Contents

# Foreword

Jim Morrison is to be congratulated on producing this very important book. It used to be thought that to make the nearest thing to a perfect car, it was necessary for each component to be produced to satisfy the very narrowest specifications. This is the philosophy that produced the Rolls Royce. Unfortunately the car was not only exceptional in reliability but also exceptional in cost.

It is remarkable that as far back as 1957 Jim Morrison came up with a very different and important concept. This was to use in engineering design the concept of transmission of error. With this approach, it became clear that to produce low error transmission in the characteristics of an assembly certain components had to satisfy very tight specifications and these were expensive to achieve. However other components that had much less effect on the performance of the assembly could have much wider and less expensive specifications. He showed us how to find out which components must have very narrow specifications and which could be much less narrow. By spending money where it would do the most good, it was possible to produce a car at a moderate cost whose performance and reliability were extremely high. Morrison's concept can be applied in all areas of engineering design. His concept has had profound effects. Those companies that ignore it do so at their peril.

History is full of examples where the origin of an important concept was not known or was ignored until a much later time. This has been true in the case of robust design described above. Sometimes not only has the originator of the idea been forgotten but the essentials which he developed have been misapplied. In particular, Jim had pointed out the importance of knowing, at least approximately, the variances of the components in order to determine the variance of the assembly. In later versions of this concept such matters have been given far too little attention. We are particularly grateful, therefore, for this book in which Jim describes these techniques with clarity and accuracy.

**George E.P.Box FRS**

Emeritus Professor, University of Wisconsin, USA

Honorary Member of the American Society for Quality

Inaugural holder of the George Box Medal for outstanding
    contributions to business and industrial statistics awarded by the
    European Network of Business and Industrial Statisticians (ENBIS)

# Preface

This introductory text on statistical engineering is written by an engineer for an engineering readership. It is hoped it will appeal both to practising engineers and to students, and (indeed) to school leavers contemplating engineering as a career. It may also be useful to managers who are concerned with the quality of manufactured products.

In spite of all the effort over centuries to achieve absolute precision, engineering is still (and probably always will be) beset by variability which is manifest in many different ways – properties of raw materials, the environment, measurement error, process variability, etc.

Statisticians, too, are beset by variability. If variability did not exist their branch of mathematics would (probably) never have come into existence. Variability is their focal point. They have developed powerful analytical techniques which can be of enormous benefit to society in general and to specialists in other branches of science and technology in particular.

Engineers using statistical methods need not concern themselves with profound issues of statistical inference or the subtleties of statistical mathematics. They require only familiarity with relevant statistical methods, an understanding of how they work and how to use them safely without running into danger. Some familiarity with statistical terminology is also desirable so that they can communicate with statisticians when the need arises. That is what this book is all about.

The sequence of topics is not linked in any way to the theoretical development of mathematical statistics. The text begins with a non-mathematical examination of the nature of variability in engineering data, followed by an explanation of some basic statistical methods for dealing with variability. It then follows the pursuit of variability reduction in manufacturing industry starting with production, followed by engineering design, then research and development. Finally, measurement, statistical computing, and quality management are dealt with as background topics. Although it is convenient to use manufacturing industry as a vehicle for demonstrating the use of statistical methods it must be emphasised they are widely applicable in other branches of engineering.

Statistical methods provide the only satisfactory way of dealing with the variability that exists in every engineering situation. The buck doesn't stop at ground level. The responsibility for dealing with variability is carried by engineers and managers at all levels right up to chief executive. The engineer who is lacking in statistical skill is less than competent to handle variability. For that reason statistical engineering should be a continuing professional development interest for practising engineers irrespective of seniority.

Engineering students must recognise that statistical skills will be important to them in their future careers no matter what branch of technology they enter or how high they set their sights. As fee-payers they are entitled to look critically at their academic curricula. If they graduate in an academic establishment at which no provision is made for teaching the elements of statistical engineering they will find themselves later in life competing on unequal terms with statistically competent engineers who are better equipped to deal with the reality of the world.

There is a message here, too, for school leavers who are considering a professional career in engineering. They should enquire carefully about the curriculum of any engineering degree course they are thinking of entering. If there is no evidence of statistical engineering content they should pass it by and look at the next on their list before committing themselves.

This book introduces a broad range of statistical methods that are relevant to engineering. These are presented with the minimum of mathematics and the maximum of explanation. Where statistical jargon is used the words and phrases are printed in italics at the first entry so that the meaning will be self-evident from the context. The object is to build bridges of understanding between the professional disciplines of engineering and statistics.

To assist the readers who may wish to take the subject further than a basic introduction (particularly in areas of research) extensive reference lists are provided at each chapter end. In addition four appendices offer guidance for further study. A fifth appendix accommodates statistical tables.

# Acknowledgements

# 1

_____

# Nature of Variability


There is no engineering product so simple that only one source of variability affects its dimensions or properties. Take two examples of products which are relatively simple in their physical appearance – high carbon steel wire and milk bottles.

The tensile strength of steel wire depends on numerous factors: the carbon content of the ingot from which rods were made in the rolling mill; the temperature of the heat treatment furnace through which the rods were passed; the rate of passage through the furnace; the temperature of the quenching bath; the ambient temperature in the heat treatment shop; the number of dies through which the rods were drawn to finished wire size; the rate of drawing; the ambient temperature in the wire mill, etc. Variability in any of these factors is likely to generate variability in tensile strength.

One of the hazards of a milkman's life is the possibility of being stopped in the street by a weights and measures inspector. Milk bottles are filled to a predetermined level on automatic machines. The capacity at that level is determined by the external profile of the bottle and by its wall thickness. The bottles are made on multi-head automatic machines by dropping gobs of molten glass into metal moulds (one at each work station), piercing them hollow, then inflating them with compressed air until they fill the moulds. The external profile can be affected by different settings at each work station, by mould differences, by fluctuations in air pressure, by sagging after release from the moulds, and by malfunctioning of the automatic timing gear which controls the various functions. The wall thickness is determined by the setting of the gob feeder and this in its turn is affected by the viscosity of the glass, the forehearth temperature, the action of the shears which cut off successive gobs from the continuous flow of the feeder. Variability in any of these process factors may contribute to variability in the volumetric capacity of bottles in continuous production.

It must be assumed that most engineering products which are infinitely more complex than steel wire or milk bottles will be equally susceptible to a multitude of factors located in raw materials, components, processes and the environment which are capable of affecting the properties and

dimensions of a finished product.  It is therefore important for engineers to have an understanding of the way in which random combinations of independent sources can affect the variability of a finished product.

This can be demonstrated with random combinations of the variables R, Y & B in Table 1.1.  These single-digit numbers in the range 0-9 were generated by throwing twenty-faced icosahedron coloured dice (red, yellow and blue) with the numbers zero to nine engraved twice on each die.  The dice were invented in the 1950/60 period by Mr Yasushi ISHIDA and patented by Tokyo-Shibaura Electric Company.  They were marketed and distributed by the Japanese Standards Association for demonstrating the principles of statistical quality control.  In the discussion that follows the data in Table 1.1 will be used to demonstrate some of the phenomena of variability that are encountered in engineering data without resort to the mathematics of probability theory.  It is hoped this will help the reader to understand the relevance of statistical methods to be described later.

**Table 1.1** Dice Scores

| R | Y | B | R+Y+B | Mean | Range | R×Y |
|---|---|---|-------|------|-------|-----|
| 0 | 6 | 5 | 11 ) | | | 0 |
| 0 | 8 | 9 | 17 ) | | | 0 |
| 4 | 6 | 5 | 15 ) | 13.8 | 6 | 24 |
| 7 | 0 | 6 | 13 ) | | | 0 |
| 9 | 4 | 0 | 13 ) | | | 36 |
| 1 | 9 | 4 | 14 ) | | | 9 |
| 7 | 0 | 3 | 10 ) | | | 0 |
| 7 | 3 | 6 | 16 ) | 12.2 | 9 | 21 |
| 2 | 4 | 1 | 7 ) | | | 8 |
| 1 | 9 | 4 | 14 ) | | | 9 |
| | | | (continued for one hundred trials) | | | |

One hundred trials were conducted, but only the first ten are recorded in the table.  Readers who are not convinced that the trials are properly reported are at liberty to conduct their own time-consuming experiments.  Also recorded in the table are the sums R+Y+B, and the products RxY, along with the *mean* and the *range* of groups of five.  In statistical terms, the mean of a set of data is the sum of the individuals divided by the number of individuals.  The range is the difference between the largest and smallest individuals.

8

The *frequency distributions* are as follows

| R, Y&B | Frequency | R+Y+B | Frequency |
|--------|-----------|-------|-----------|
| 0 | 30 | 0,1 | 0 |
| 1 | 38 | 2,3 | 1 |
| 2 | 20 | 4,5 | 2 |
| 3 | 38 | 6,7 | 7 |
| 4 | 29 | 8,9 | 12 |
| 5 | 31 | 10,11 | 15 |
| 6 | 29 | 12,13 | 24 |
| 7 | 32 | 14,15 | 17 |
| 8 | 21 | 16,17 | 4 |
| 9 | 32 | 18,19 | 9 |
|   |    | 20,21 | 5 |
|   |    | 22,23 | 3 |
|   |    | 24,25 | 1 |
|   |    | 26,27 | 0 |

These can be represented graphically in Figs.1.1 & 1.2.



Fig. 1.1  Individual dice scores

In a perfect world one might expect Fig.1.1 to display thirty scores in each of the ten categories 0-9, but the bar chart (or *histogram,* to use a statistical term) shows some degree of irregularity.  If bias was suspected it would be necessary to run a much more extensive series of trials to show whether the dice were loaded in favour of scores 1 and 3 at the expense of scores 2 and 8.  In the absence of such evidence it can be assumed that the scoring conforms to a rectangular distribution and that the irregularity is no more than is commonly encountered in real life collections of data.

Fig. 1.2 R+Y+B dice scores

In sharp contrast, the bar chart for the sum of the three colours (Fig. 1.2) shows an entirely different pattern of distribution. There is a marked central tendency around a mean score of 13.5 which is easy to explain. All possible combinations of scores on the three dice are equally likely. There are many different combinations yielding totals of 10,11,12,13,14 or 15, but very few which can yield extreme values of 0,1,2,3 or 24,25,26,27. In fact there is only one combination 0+0+0 which could yield 0 and only one other combination 9+9+9 which could yield 27, and neither occurred in this relatively small set of trials.

Symmetrical bell-shaped distributions exhibiting a central tendency are commonplace in engineering data. It is not unreasonable to argue these are indicative of random combination of independent factors contributing to the variability of the data and to suggest that analytical statistical methods might be used to identify and control them.

However, it must not be assumed that other patterns of distribution will not occur in engineering data. The distribution of products of red and yellow scores, RxY is highly *skewed* (i.e. asymmetric) as shown in Fig.1.3.



Fig. 1.3 R×Y dice scores

10

Skewed distributions do occur in engineering when the effect of a contributory factor is one-sided. For example, in a thermionic valve electrons are emitted from the heated cathode and are attracted by a positive voltage on the anode. They have to pass through the grid (a helix of fine wire) to which a negative voltage is applied to control the current. Any lack of uniformity in the grid helix can only increase, not reduce, the an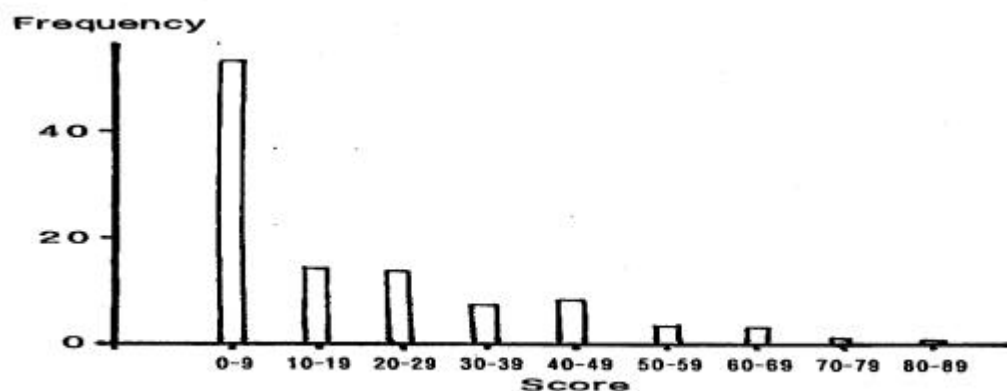ode current. Again, in a cylindrical mechanical product zero ovality is the ultimate degree of perfection. Any finite degree of ovality is positive if it is regarded as the excess of the major diameter over the minor without regard to orientation. In such circumstances skewed data distributions are inevitable.

Fortunately statistical methods are available which are not confined exclusively to data that conforms to a symmetrical distribution. When skewed distributions are encountered in engineering data they can often be handled more easily by making a logarithmic transformation of the data.

The data in Table 1.1 can be used to demonstrate relationships between *samples* and *populations*. This is a matter of considerable importance to engineers who often have to draw valid conclusions from quite small samples of data. For example, in the early stages of development of a new product it is necessary to check measurements of a few prototypes to determine whether the population will be on target and whether the (unavoidable) spread of variability will lie comfortably within specification tolerance limits. In this instance the prototype data can be treated as a sample from a population that does not yet exist, yet a prediction has to be made.

This situation is simulated in the third and fourth columns of Table 1.1 by taking the mean value and range of R+Y+B scores in successive groups of five trials. This resulted in the following twenty mean values, not one of which coincided with the mean of the original set of R+Y+B scores (12.9). The nearest was 13.2, but the extreme examples were 10.2 and 15.2. Clearly, there were many instances in which the sample mean would not have given a good estimate of the population mean.

| 13.8 | 12.2 | 13.4 | 11.2 | 12.2 | 11.2 | 12.0 | 14.0 | 10.8 | 15.0 |
| 13.8 | 14.4 | 15.2 | 12.2 | 14.0 | 12.2 | 13.2 | 14.2 | 13.2 | 10.2 |

The range of R+Y+B scores over each group of five trials gave the following results.

| 6 | 9 | 13 | 16 | 10 | 15 | 4 | 9 | 8 | 12 | 15 | 19 | 12 | 5 | 13 | 7 | 9 | 17 | 13 | 7 |

If the range is taken as a crude measure of overall variability (as many development engineers have been known to do in the past when writing specification tolerances) it is clear that not even the highest value (19) recorded in this set of trials would embrace the span of the distribution shown in Fig.1.2. Most of the others would fall very far short of this requirement.

The relatively small sets of data used by engineers at the development stage of a new product can be regarded as samples from a population which will exist when full scale production starts. The discrepancies in mean value and variability which can exist between a sample, and the population from which it is drawn, identify a serious hazard along the road from design, through development to production of manufactured products. It is to be hoped that the straightforward demonstration of the risks given above will alert engineers to the dangers and persuade them to listen more carefully to the advice of statisticians, or (better still) develop some statistical skill on their own account.

So, if range is not to be regarded as a satisfactory measure of overall variability what else can we do? Consider the following small set of data:

$$\boxed{16 \quad 18 \quad 16 \quad 10 \quad 14}$$

The location of the data on a scale of measurement can be identified by calculating the mean value.

$$(16+18+16+10+14)/5 = 74/5 = 14.8$$

The *deviates* of the individuals from the mean are

$$
\begin{aligned}
16.0 - 14.8 &= \phantom{-}1.2 \\
18.0 - 14.8 &= \phantom{-}3.2 \\
16.0 - 14.8 &= \phantom{-}1.2 \\
10.0 - 14.8 &= -4.8 \\
14.0 - 14.8 &= -0.8
\end{aligned}
$$

The sum of these deviates, taking account of positive and negative signs, will be zero. Suppose we square them before adding them together?

$$1.2^2 + 3.2^2 + 1.2^2 + (-4.8)^2 + (-0.8)^2 = 1.44 + 10.24 + 1.44 + 23.04 + 0.64 = 36.80$$

This *sum of squares* is a powerful overall measure of variability which gives equal weight to all of the individuals, not just the extreme values. It does, however, respond to the size of the data. If data from the same source had ten values the sum of squares would be (roughly) twice as large.

This can be overcome by dividing the sum of squares by the number of individuals to give a *mean square:*

$$\frac{36.80}{5} = 7.36$$

In some situations the divisor should be one less than the number of individuals, but more of that later in Section 2.2!

Summing squares to measure variability is the foundation on which statistical analysis is built. In modern usage 'statistics' implies much more than simply recording events. In the hands of a competent engineer statistical analysis is a powerful tool which should not be neglected. Now read on!

**2**

_____

# Basic Statistical Methods

The absence of a sound statistical element in an engineering degree is a serious weakness.  A course in quality assurance embracing techniques of applied statistics along with principles of operations management would be appropriate (Morrison, 1997).  The necessary basic statistical methods are presented in this chapter.  The elements of operations management are presented in Chapter 7.

## 2.1 Variance

Engineers wishing to make extensive use of statistical methods must first come to terms with _statistical variance_.  There appears to be common ground between engineering and statistics because it is probably true to say that in the early development of their subject statisticians borrowed the concept of _moments_ from mechanics.  The moment of the first order is used to determine the location of a set of data on the scale of measurement in which the individual values were recorded and the moment of the second order is used to measure their dispersal (i.e. the variability).

If $n$ individual values $x_i$ in a set of data are represented by the symbols $x_1$, $x_2$, $x_3$ ... $x_i$ ...$x_n$ then the mean $\overline{x}$ and the variance $V(x)$ are given by

$$\overline{x} = \frac{1}{n}\sum x_i$$

$$V(x) = \frac{1}{n}\sum (x_i - \overline{x})^2$$

Consider a data set of five values 1,8,8,9,6.  The mean and variance can be calculated as follows:

$$\sum x_i = 1+8+8+9+6 = 32$$

$$\therefore \overline{x} = \frac{32}{5} = 6.4$$

$$\sum (x_i - \overline{x})^2 = (1-6.4)^2 + (8-6.4)^2 + (8-6.4)^2 + (9-6.4)^2 + (6-6.4)^2 = 41.20$$

$$\therefore V(x) = \frac{41.20}{5} = 8.24$$

When calculating the sum of squares of the *deviates* of the individuals about the mean it is often more convenient to use the algebraic identity

$$\sum(x_i - \bar{x})^2 = \sum x_i^2 - (\sum x_i)^2 / n$$

The quantity $\sum x_i^2$ is sometimes referred to as the *crude sum of squares*. $(\sum x_i)^2 / n$ is then the *correction factor* and $\sum(x_i - \bar{x})^2$ is the *adjusted sum of squares*. Applying this procedure to the set of data above gives the same result as before:

$$\sum x_i^2 = 1^2 + 8^2 + 8^2 + 9^2 + 6^2 = 246$$

$$(\sum x_i)^2 / n = 32^2 / 5 = 204.80$$

$$\sum(x_i - \bar{x})^2 = 246.00 - 204.80 = 41.20$$

It is not very convenient to have the mean and the variability expressed in different units of measurement, such as 'miles per hour' and 'miles per hour squared'. To overcome this difficulty the square root of variance is termed the *standard deviation, sigma*.

$$\acute{o} = \sqrt{V(x)}$$

For the set of data considered above:

$$\acute{o} = \sqrt{V(x)} = \sqrt{8.24} = 2.87$$

The set of five values 1,8,8,9,6 can now be summarised in statistical terms as having a mean value $\bar{x}$ = 6.4 and a standard deviation $\acute{o}$ = 2.87. Note that the standard deviation (or the variance) is a powerful measure of variability taking account of every individual in the data set, not just the extreme values 1 & 9. In this respect it is superior to the range (i.e. the difference between the largest and smallest values) which is often used by engineers as a crude measure of variability. Moreover, no prior assumption is made about the shape of the parent distribution. The entire data set is taken just as it stands.

Engineers will appreciate that the statistical mean is analogous to a centre of gravity and the statistical standard deviation is analogous to the radius of gyration of a rotating mass.

It should be pointed out that statisticians sometimes use two other central values besides the mean when discussing a set of data. The *median* is the mid-point of the data when the individuals are arranged in order of

magnitude. The *mode* is the most commonly occurring value. The mean, the median and the mode sometimes coincide exactly but this is not an invariable rule.

## 2.2 Divisor 'n' or 'n-1'?

Engineers using a hand-held calculator for statistical calculations may be perplexed to find two keys labelled 'σxn' and 'σxn-1' (or corresponding sub-routines in computer software). Which one should be used?

To clarify this it is necessary to consider the relationship between a sample and the larger population from which it was drawn. To make a distinction between population and sample the symbols used for mean and standard deviation will be $\overline{X}$ and $\acute{o}$ for the population, $\overline{x}$ and $s$ for the sample. It will be assumed that the purpose of calculating $\overline{x}$ and $s$ from the sample data will be to estimate the unknown parameters $\overline{X}$ and $\acute{o}$ of the population.

With the knowledge that the 1,8,8,9,6 data set in the previous section originated as five throws of an unbiased twenty-faced die capable of generating a rectangular distribution of numbers in the range zero to nine the data can be treated as a random sample from an infinite population with a mean value $\overline{X}$ = 4.5. The true sum of squares about the population mean is therefore:

$$\sum(x_i - \overline{X})^2 = (1-4.5)^2 + (8-4.5)^2 + (8-4.5)^2 + (9-4.5)^2 + (6-4.5)^2 = 59.25$$

Obviously this is greater than the sum of squares previously calculated about the sample mean:

$$\sum(x_i - \overline{x})^2 = 41.20$$

As has already been seen in the previous chapter it is quite unusual for a sample mean to coincide exactly with the population mean. In the general case, when the population mean is not known in advance, the sum of squares about the sample mean will underestimate the true sum of squares and the sample standard deviation $s$ will underestimate the population standard deviation $\acute{o}$. This can be compensated by using the divisor 'n-1' when calculating the variance. This is not just a fudge – there is sound mathematical reasoning to show this gives the best estimate of the population standard deviation $\acute{o}$.

Hence the rule: If the purpose is simply to calculate the variance of a set of data, use the divisor 'n', but if the purpose is to estimate the standard deviation of the population from which the sample may have been drawn, use the divisor 'n-1'. Obviously the difference is neither here nor there in large data sets but there is a considerable difference with small samples which quite often crop up in engineering (for example, at the prototype development stage of a new product).

So, when estimating the variance of a population from sample data we use the expression

$$s^2 = \frac{1}{n-1}\sum(x_i - \overline{x})^2$$

Computed in this way $s^2$ is an *unbiased estimator* of the population variance $V(x)$. The quantity *n*-1 is referred to as the *degrees of freedom* associated with the estimate. The sum of the *n* deviations $(x_i - \overline{x})$ is zero by virtue of the definition of the mean. If values are assigned to *n*-1 individuals the remaining one is already determined. In many forms of statistical analysis the degrees of freedom are identified by the symbol **u**.

## 2.3 Covariance and Correlation

Engineers will sometimes encounter *bivariate* data in which two variables such as $x$ and $y$ appear to be correlated. *Statistical covariance* (cov) measures the degree of association, using sums of products in place of sums of squares:

$$Cov\ (x, y) = \frac{1}{n}\sum(x_i - \overline{x})\ (y_i - \overline{y})$$

where $y_i$ is the individual value of $y$ associated with an individual $x_i$.

As in the case of sums of squares, there is a useful algebraic identity for simplifying the calculation of *sums of products*:

$$\sum(x_i - \overline{x})\ (y_i - \overline{y}) = \sum x_i y_i - \frac{1}{n}\sum x_i \sum y_i$$

A graphical interpretation of covariance is given in Fig.2.1 where individuals are plotted on an *X,Y* co-ordinate field.
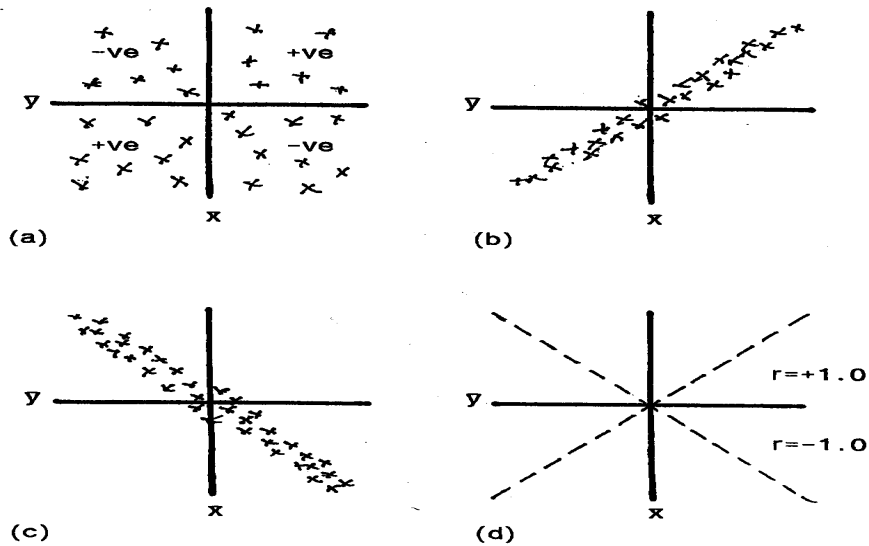
Fig. 2.1  Covariance and correlation

Since we are considering the product of the $x$ and $y$ deviates from their means it is appropriate to use an origin at the centroid of the data with axes representing the deviates $(x_i - \bar{x})$ and $(y_i - \bar{y})$.  These divide the field into four quadrants, upper right, upper left, lower left, lower right.  The products will be positive in the upper right and lower left quadrants. They will be negative in the upper left and lower right quadrants.

If there is no association between the $x$ and $y$ variates, as in Fig.2.1(a), the positive and negative products will cancel out.  If there is a strong association then either the positive products will predominate, as in Fig.2.1(b), or the negative products, as in Fig.2.1(c).

A dimensionless *correlation coefficient* $r$ can be used to measure the degree of association:

$$r = \frac{\Sigma(x_i - \bar{x})\ (y_i - \bar{y})}{\sqrt{\Sigma(x_i - \bar{x})^2\ \Sigma(y_i - \bar{y})^2}}$$

If there is a perfect association between the $x$ and $y$ variates as in Fig.2.1(d) the square of the sum of products will be numerically equal to the product of the sums of squares and the correlation coefficient will be unity.  It will be positive for a rising gradient and negative for a falling gradient depending on whether the products of the $x$ and $y$ deviates are positive or negative.

In the case of Fig.2.1(a) the correlation coefficient will be zero.  In Fig.2.1(b) and Fig.2.1(c) the correlation coefficient will have intermediate values within the range $\pm 1.0$.

## 2.4 Normal Distribution

Symmetric bell-shaped distributions of the type shown in Chapter 1, Fig.1.2 can be modelled in statistical terms using the so-called *Normal Distribution* (sometimes referred to as the Gaussian distribution after the celebrated German mathematician).

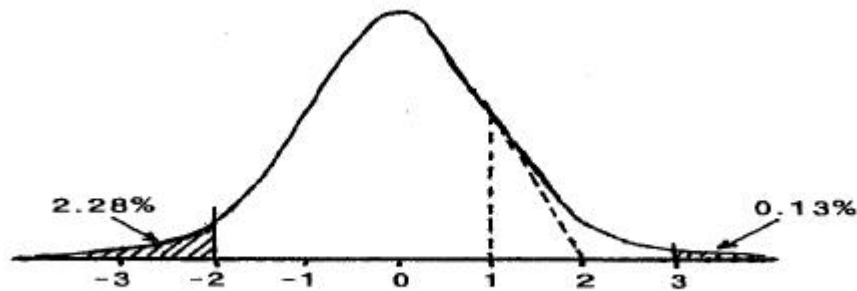The profile of this distribution is shown in Fig.2.2.



Fig 2.2 Normal distribution: $\bar{x} = 0$, $\acute{o} = 1.0$

As seen in this diagram the mean is zero, the standard deviation is unity, and the tails of the distribution extend to three standard deviations (and beyond, to infinity). The equation for the normal frequency curve is

$$\boldsymbol{f}(x) \;=\; \frac{1}{\sqrt{2\partial}} e - \frac{1}{2}x^2$$

The area to the left of the *ordinate* at $x$ is given by

$$\ddot{O}(x) \;=\; \frac{1}{\sqrt{2\partial}} \int\limits_{-\infty}^{x} e^{-\frac{1}{2}t^2} dt$$

Extensive tables of the values of $\ddot{O}(x)$ and $\ddot{o}(x)$ are to be found elsewhere in the standard statistical literature but it is the area to the right of the ordinate that is of special interest to production engineers because it can be used to quantify the proportion of rejects falling outside specification tolerance limits. Values of the variate $x$ corresponding to specific percentages of outliers *P*% are tabulated for convenient reference in Appendix E.

In the following schematic diagram figures extracted from the table show that only a very small fraction (less than 0.2%) of the area under the normal distribution curve lies beyond plus or minus three standard deviations.

| P | 0.2% | 0.1% |
|---|------|------|
| X | 2.88 | 3.09 |

19

In theory the distribution extends to plus/minus infinity but beyond three standard deviations the height of the ordinate becomes vanishingly small. Within three standard deviations of the mean the normal distribution provides a good fit to many of the distributions encountered in engineering data.

In some quarters extrapolations are made beyond *three-sigma* and minutely small probabilities are quoted in parts-per-million, but the practice is dubious and millions of test results would be necessary to validate it. If an increased margin of safety is necessary it is more sensible to specify this in terms of the same scale of measurement as that which was used to record the data, or as a multiple of the standard deviation.

One of the commonest applications of the normal distribution in manufacturing engineering is to predict the proportion of units of product likely to fall outside specification limits. Consider the case of a product whose dimension is intended to meet a specification tolerance of 80±0.3. The process is running slightly above target with a mean of 80.1 and a standard deviation of 0.14.

The tolerance limits can be expressed as multiples of the standard deviation (i.e. *standardised deviates*):

Upper specification limit = (80.3 – 80.1)/0.14 = +1.43
Lower specification limit = (79.7 – 80.1)/0.14 = -2.86

In the following schematic diagram these standardised deviates are inserted between adjacent values extracted from the table of percentage points of the normal distribution in Appendix E. The estimates *P* 7.5% and *P* 0.2% were arrived at by taking note that the value *x*=1.43 is almost exactly midway between the ten and five *percentiles* and the value *x* = 2.86 is closer to the 0.2 percentile than to the 0.5 percentile.

| *P* | 10.0% | *P*≈7.5% | 5.0% | | 0.5% | | *P*≈0.2% | 0.2% |
|---|---|---|---|---|---|---|---|---|
| *X* | 1.28 | 1.43 | 1.64 | | 2.58 | | 2.86 | 2.88 |

In this, and in subsequent schematic diagrams bold arrows pointing to a double-lined box serve to focus readers' attention on the issue under discussion. The approximate equality sign ≈ identifies estimates that have to be determined by interpolating exact figures extracted from the

tables. It is not suggested these schematic diagrams should be constructed on every occasion that reference is made to the tables in Appendix E. They are used here simply to demonstrate the process of visual interpolation.

From the above display it can be seen that the proportion of fall-out is as follows:

$$\begin{array}{ll} \text{Above upper specification limit} & \approx 7.5\% \\ \text{Below lower specification limit} & \approx 0.2\% \\ \hline \text{Total} & \approx 7.7\% \\ \hline\hline \end{array}$$

If a full-dress table of the normal distribution function is used the precise estimate is 7.85%. Does the discrepancy of 0.15% really matter?
If the process was brought back on target the tolerance limits would be at $x = \pm 0.3/0.14 = \pm 2.14$ standard deviations which is not quite midway between the two and one percentiles.

| $P$ 2.0% | $P \approx 1.6\%$ | 1.0% |
|----------|-------------------|------|
| $X$ 2.05 | 2.14 | 2.33 |

The total fall-out would then be $2 \times 1.6 = 3.2\%$, less than half what it had been. This would be advantageous but there would still be work to do to get the variability reduced. The source(s) of variability would have to be identified and brought under closer control. To eliminate fall-out the standard deviation would have to be reduced from 0.14 to 0.10 (one sixth of the overall tolerance). Even then, the process would have to be held strictly on target. If this was not possible a standard deviation less than 0.10 would allow some room for manoeuvre.

Before leaving the Normal distribution it is worth noting that its standard deviation is not just a mathematical abstraction. Fig. 2.2 page 19 shows that the point of inflexion at which the distribution curve changes from concave inwards to concave outwards occurs at one standard deviation and the tangent at that point intersects the base line at two standard deviations. In this way the standard deviation does provide a valid measure of the spread of the distribution, irrespective of the tails which extend to infinity in both directions.