

### Figure 5.1. MODELLING THE SHAPE OF DATA DISTRIBUTIONS

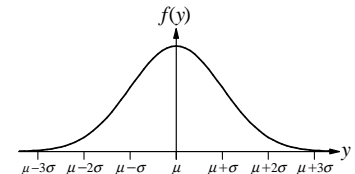
Figures 4.1 and 4.8 of these Course Materials dealt with numerical measures of the *location* and the *variation* of a data set; we now turn to a model (or idealization) for one *shape* of a data distribution.

#### 1. Normal Distributions

Although a large number of shapes is *possible* for the distribution of a data set, in practice it is observed (by examining suitable histograms, for example) that many distributions fall into one of a limited number of categories. A common shape has a central peak with a roughly symmetrical falling away on either side – for instance, see Figures 2.8b (paper thickness results) and 2.9b (coin weights). A *mathematical model* (or *idealization*) of this shape is the *normal (probability) distribution*, whose *probability density function* (or *p.d.f.*) is equation (5.1.1) at the right above. The graph of this function is shaped like the cross-section of an inverted bell – in many texts, it is described as a *symmetrical bell-shaped curve*. The equation has two *parameters*,  $\mu$  and  $\sigma$ , which (independently) determine the position of its centre ( $\mu$ ) and the width of its peak ( $\sigma$  – the *larger*  $\sigma$ , the *wider* the peak). The parameter  $\mu$  is called the *mean* of the distribution and  $\sigma$  is its (probabilistic) *standard deviation*. The value of  $\mu$  is the centre of the graph – the  $y$ -value of its highest point; the value of  $\sigma$  can be roughly assessed by eye because it is the distance from the mean to either point of inflection. As the diagram indicates, the part of the normal p.d.f. we typically *see* covers an interval of about  $3\sigma$  either side of  $\mu$  or about  $6\sigma$  in total.

$$f(y) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2}\left[\frac{y-\mu}{\sigma}\right]^2} ; \quad -\infty < y < \infty \quad \text{-----(5.1.1)}$$

$$Y \sim N(\mu, \sigma) \quad \text{-----(5.1.2)}$$



Probabilities are *estimated* in practice by *proportions*; because proportions are represented by bar *areas* in histograms, and because we model data histograms by probability distributions, it is *area* under a p.d.f. that represents probability.

#### 2. Variates and Random Variables

A *variate* is a characteristic associated with each element of a population – variates may be *categorical* (like colour or sex or marital status) but our concern here is with variates (like length or weight) that take *numerical* values. In data-based investigating, variate values are usually measured for a *sample* of elements selected from an appropriate study population; if the sample contains  $n$  elements, we denote these (response) variate values by subscripted lower-case Roman letters  $y_1, y_2, y_3, \dots, y_n$  (or, more compactly,  $y_j, j=1, 2, \dots, n$ ).

When we use a probability distribution to model the shape of the distribution of a variate, the variable in the equation of the *model* is called a *random variable*; for example,  $y$  in the normal p.d.f. (5.1.1) above is (the value of) a normal random variable. We are familiar with the term ‘variable’ from algebra and calculus; the addition of the adjective *random* for a variable in a probability distribution can serve to remind us that, by the act of modelling, we have asserted there is a probability associated with each value the variable takes on; these probabilities come, in part, from the *method of selecting* the elements which comprise the sample and which yield the data. It is the task of the *Plan* for an investigation (proper selecting *and* proper measuring, etc.) to provide a reasonable basis for regarding variate values as values of a random variable.

As summarized in Ttable 5.1.1 at the right, the usual notation convention is to use a (subscripted) *lower-case italic* letter  $y$  (or  $y_j$ ) for the value(s) of a random variable and an *upper-case italic* letter  $Y$  (or  $Y_j$ ) for the random variable(s). When the random variable  $Y$  has a *normal* distribution with parameters  $\mu$  and  $\sigma$ , we write symbolically  $Y \sim N(\mu, \sigma)$ , shown as equation (5.1.2) above. Letters for random variables are usually chosen from near the *end* of the alphabet, subject to the caveat, to assist problem solving, of using letters whose identity is easy to remember – for instance,  $W$  for weight,  $L$  for length,  $T$  for time; letters near the *beginning* of the alphabet are usually kept for *events* (see Part 7 of the Course Materials).

Data	Roman	Model	Italic
Variate value	$y$ or $y_j$	Random variable value	$y$ or $y_j$
		Random variable	$Y$ or $Y_j$

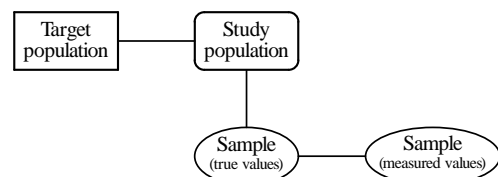
Letters for random variables are usually chosen from near the *end* of the alphabet, subject to the caveat, to assist problem solving, of using letters whose identity is easy to remember – for instance,  $W$  for weight,  $L$  for length,  $T$  for time; letters near the *beginning* of the alphabet are usually kept for *events* (see Part 7 of the Course Materials).

The characteristic of a random variable commonly of interest is a *probability*; for instance, the probability the random variable  $Y$  is greater than 4 [*i.e.*, symbolically,  $\Pr(Y > 4)$ ], or the probability the random variable  $Z$  takes the value  $z$  [*i.e.*,  $\Pr(Z = z)$ ].

#### 3. Attributes and Model Parameters: Estimating

The focus of Section 2 above is on *individual* elements and their variates, whose values we model as the values of random variables. In this Section 3, we focus on three *groups* of elements:

- \* the *target population*: the group of elements to which the investigator(s) want Answer(s) to the Question(s) to apply;
- \* the *study population*: a group of elements *available* to an investigation;
- \* the *sample*: the group of elements selected from the study population and *actually used* in an investigation.



(continued overleaf)

A schema showing these three groups, including the distinction between true and measured variate values, is given overleaf on page 5.3 at the lower right; the vertical line in the middle of this schema reminds us the sample is a *subset* of the study population.

Associated with these three *groups* of elements are:

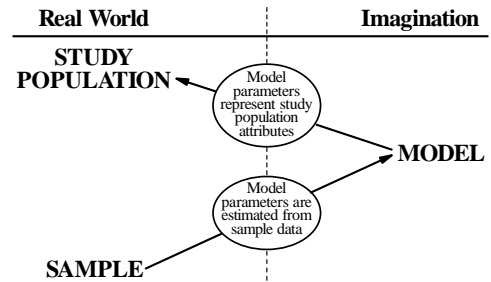
- \* *attributes*: quantities defined as a function of the response (and, perhaps, explanatory) variates over the *group* of elements. Familiar (simple) attributes are averages, proportions, medians and standard deviations. The importance of attributes is that:
  - Answer(s) to Questions(s) are usually phrased in terms of attribute values, as illustrated by the newspaper articles reprinted in Part 4 of the Course Materials;
  - five of our six categories of *error* are defined in terms of attributes.

A (probability) *model parameter* is a constant (usually denoted by a *Greek* letter) in a (probability) model that *represents* a study population *attribute*. Model-based methods of analysis in statistics use data from a sample to *estimate* values of model parameters which then represent reasonable values for study population attributes and, hence, for Answer(s) to Question(s) of interest. We distinguish:

- \* a *point* estimate: a *single* value for an estimate; **AND**:
- \* an *interval* estimate: an *interval* of values for an estimate.

For a sample of (response) variate values where the normal distribution is an appropriate model, the mean  $\mu$  is estimated by the sample average  $\bar{y}$  and  $\sigma$  is estimated by the sample standard deviation  $s$  – these are both *point* estimates. We can think graphically of the process of estimating  $\mu$  by  $\bar{y}$  and  $\sigma$  by  $s$  as approximating the histogram of a data set by the normal p.d.f. that has the same ‘centre’ and the same ‘width’ as the histogram.

All mathematical models are idealizations and all are products of the intellect and the imagination. It may be helpful to think of the model as a *link* between the *sample* and the *study population*; a pictorial representation of this idea is shown at the right above.



The matters discussed in this Section 3 (and in the Appendix on the last side of the Figure, page 5.6) look ahead to the use of probability models in statistics, which is pursued in Part 6; our *immediate* concern in Part 5 is to become familiar with the properties of continuous probability models.

**NOTE:** 1. To maintain the distinction between the real world (represented by the data) and the model, we use different words – ‘average’ and ‘mean’ – for their measures of location; unfortunately, we do not have this option for the two measures of variation, which are both called ‘standard deviation’. In the early stages of learning statistics, it is helpful to, at least in our minds, add the respective adjectives ‘data’ and ‘probabilistic’ to distinguish the two uses of standard deviation. This terminology is summarized in Table 5.1.2 at the right.

Attribute	Real World	Model
Location	Average	Mean
Variation	(Data) standard deviation	(Probabilistic) standard deviation

**4. Using the Normal Distribution**

To be able to *use* the normal distribution, we need two skills.

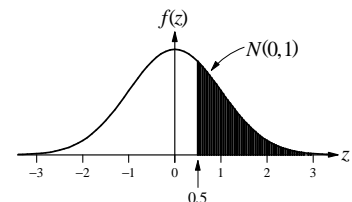
- \* First, we learn to look up a *table* of the standard normal distribution [denoted  $N(0, 1)$ , whose p.d.f. is given as equation (5.1.3) at the right]; such a table furnishes us with *areas* which we use as *probabilities* (see Figure 5.4 of the Course Materials, and also pages T-2, T-3, the last line of page T-11 and the front flyleaf of the text).
 
$$f(z) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}z^2} \quad ; \quad -\infty < z < \infty \quad \text{-----(5.1.3)}$$
- \* Second, we learn about *standardizing* so we can obtain probabilities for the  $N(0, 1)$  distribution; if  $Y \sim N(\mu, \sigma)$ , the result (5.1.4) at the right above – in which we *subtract the mean* and *divide by the standard deviation* – holds; we then can write equation (5.1.5), in which we have standardized to convert a probability for the  $N(\mu, \sigma)$  distribution to the equivalent probability for the  $N(0, 1)$  distribution. The letter  $Z$  is commonly used for a random variable with the standard normal distribution.
 
$$\frac{Y - \mu}{\sigma} \sim N(0, 1) \quad \text{-----(5.1.4)}$$

$$\Pr(Y > y) = \Pr\left(\frac{Y - \mu}{\sigma} > \frac{y - \mu}{\sigma}\right) = \Pr\left[N(0, 1) > \frac{y - \mu}{\sigma}\right] \quad \text{-----(5.1.5)}$$

A numerical illustration of equation (5.1.5) is: if  $Y \sim N(8, 6)$  and  $Z \sim N(0, 1)$ , then:

$$\Pr(Y > 11) = \Pr\left(\frac{Y - \mu}{\sigma} > \frac{11 - 8}{6}\right) = \Pr[N(0, 1) > 0.5] = \Pr[Z > 0.5];$$

this probability is represented by the area under the standard normal p.d.f. to the *right* of 0.5, shown with solid shading in the diagram at the right. We learn in Figure 5.3 how to look up this probability from the  $N(0, 1)$  table in Figure 5.4.



(continued)

**Figure 5.1. MODELLING THE SHAPE OF DATA DISTRIBUTIONS (continued 1)****5. Benefits of Using a Normal Model**

Three benefits arise from using the normal distribution.

- First, it provides a *data summary* – instead of having to work with a *data set*, we can convey a number (but usually not *all*) of its essential characteristics by saying merely that the data has a *normal distribution with a specified mean and a specified standard deviation*. This can be a significant convenience, particularly in the case of large data sets. [It is understood, of course, that the normality is only *approximate*; also, the location and variation parameters of the normal distribution can take on their values *independently* – that is, the value of one does not influence the value of the other.]
- Second, we can readily *combine* normal distributions; the random variable represented by any *linear combination* of normal random variables has a *normal* distribution, provided the component random variables are *probabilistically independent*. The linear combinations of greatest relevance to statistical methods of data analysis are *sums*, *differences* and *averages*. This matter is taken up in Figure 5.14.
- Third, the distribution of a linear combination (such as a sum or an average) of *non-normal* random variables *tends* to a normal distribution. If there is an *infinite* number of components, the distribution of the combination is *exactly* normal, a theoretical result known as the *Central Limit Theorem*. In *practice*, there are only *finite* combinations, whose distributions therefore exhibit only *approximate* normality (unless the components are themselves normally distributed). How close the approximation is to *exact* normality depends on both the number of components in the linear combination and on the shape of their distribution(s) – the more *symmetrical* it is, the smaller the number of components needed for reasonable approximate normality. This third matter is a central theme of Figure 5.16 (and some later Parts) of the Course Materials.

It is interesting to speculate on the reason(s) behind the mathematical and practical benefits that accrue from normal modelling. It *may* be only coincidence that the mathematics of normal theory works so conveniently in many instances (and *not only* in statistics). Alternatively,  $e = 2.71828\ 18284\ \dots$  is a quantity that, in a sense, *nature* brings to our attention as, for example, the limit as  $n \rightarrow \infty$ , of  $(1 + 1/n)^n$ . It is therefore possible that the wide applicability of normal distribution theory is a reflection of the harmony that results when the mathematics we employ corresponds properly with some aspect(s) of the underlying structure of the physical world. The matter is highlighted in a quotation from Gabriel Lippmann, the French 1908 Nobel Laureate in physics: *Everybody believes in the [normal approximation], the experimenters because they think it is a mathematical theorem, the mathematicians because they think it is an experimental fact*. Paradoxically, *both* beliefs are correct, which reminds us of other seemingly contradictory aspects of experience, like the wave-particle duality of electromagnetic radiation such as light.

**6. Using Random Variables in Probability**

We *use* random variables to describe (or *model*) quantities that take on different values according to chance; one criterion we use to decide which probability model is appropriate in a particular situation is the degree of agreement between the *shape* of the p.d.f. of the model and the shape of the distribution of the data (as assessed from a histogram, for example). The *stochastic* behaviour of response variates in our models arises because of *equiprobable* (or *random*) *selecting* of elements from the relevant population. Some *informal* descriptions of what is meant by the term ‘random variable’ (or its value) are:

- \* a quantity which takes on different real values according to chance;
- \* a numerical outcome of a stochastic phenomenon;
- \* a characteristic which changes from element to element in a sample obtained by equiprobable selecting.

In addition to our concern with distribution *shape* in choosing an appropriate probability model, we usually also need to take account of:

- the *mean* of the distribution (or of the random variable,  $Y$  say), denoted  $\mu_Y$  or  $E(Y)$ ;
- the *standard deviation* of the distribution, denoted  $\sigma_Y$  or  $s.d.(Y)$ .

Formally, a *random variable* is a function which assigns a real number to each point of the sample space ( $S$ ); *i.e.*, a random variable is a function with domain  $S$  and range  $\mathbb{R}$  – it is a mapping from the sample space to the real numbers.

**7. Distributions Other Than the Normal Distribution**

Although the normal distribution is widely used as a model for data distributions, other useful models are:

- the *uniform distribution* (which represents the case of equiprobable values, as in equiprobable selecting, and thus has a p.d.f. which is *rectangular*) – see Figure 5.11;
- the *exponential distribution* (whose p.d.f. is like an exponential decay, and is used to model failure times) – see Figure 5.12;
- the *log normal distribution* (where the *logarithm* of the random variable has a normal distribution; it is used to model biological characteristics which have a lower cut-off at zero but which are not, at least in theory, limited in the *high* values they can take on).

(continued overleaf)

In particular methods of data analysis, as well as the *normal* distribution, we later use of the *t*, *K* and  $\chi^2$  distributions (cf. Figures 6.4, 6.6 and 12.25 of the Course Materials and the text page T-11).

**8. Appendix: Response Models**

Equation (5.1.2) shows the usual notation in *probability* for stating the distribution of a random variable, *Y* in this instance. In *statistics*, the *equivalent* expression (5.1.6) at the right, called a *response model*, is more convenient. Its form suggests we think of *two* components making up the random variable *Y* which, under an appropriate Plan for an investigation [including equiprobable selecting (as ‘EPS’ at the end of the model statement (5.1.6) reminds us)], is used to model the values of a response variate of an element:

$$Y \sim N(\mu, \sigma) \quad \text{-----(5.1.2)}$$

$$Y = \mu + R, \quad R \sim N(0, \sigma), \quad \text{EPS} \quad \text{-----(5.1.6)}$$

- \* a *structural component* which (in general) models the effect on the response variate of specific explanatory variate(s);
  - in the case of equation (5.1.6), which is the *simplest* response model, the structural component is merely a constant ( $\mu$ ) and contains no explicit explanatory variates;
- \* a *stochastic component* which models variation about the structural component;
  - in equation (5.1.6), the variation of *Y* about the structural component of the model ( $\mu$  in this instance) is modelled by a *normal* distribution with mean 0 and standard deviation  $\sigma$ .

Using the response model (5.1.6) to find an *interval* estimate for each of the model parameters  $\mu$  and  $\sigma$ , representing the study population *average* response  $\bar{Y}$  and response (data) *standard deviation*  $S$ , is pursued in Part 6 of the Course Materials.

**NOTES:** 2. The foregoing discussion uses the word ‘model’ in two senses:

- the division of *Y* into a structural component and a stochastic component;
- the normal model for the stochastic component.

Only the *second* of these models involves probability.

3. The use of a *probability* distribution (here, the *normal* distribution) to model the stochastic component of a response model illustrates one reason why probability is useful in statistics.

4. Equation (5.1.6) can be extended, as shown in equation (5.1.7) at the right, to model response variate values from a sample of *n* elements selected equiprobably from a study population; like the requirement for equiprobable selecting, the modelling assumption of *probabilistic independence* of the  $Y_j$ ’s has implications for Plan components which address how the data are to be collected.

$$Y_j = \mu + R_j, \quad j=1, 2, \dots, n, \quad R_j \sim N(0, \sigma), \quad \text{independent, EPS} \quad \text{-----(5.1.7)}$$

5. The *formal* definitions of the symbols in equation (5.1.7) are:

$Y_j$  is a random variable whose distribution represents the possible values of the measured response variate for the *j*th element in the sample of *n* elements selected from the study population, if the selecting and measuring processes were to be repeated over and over.

$\mu$  is a model parameter (called the *mean*) which represents the *average* of the measured response variate of the elements of the study population.

$R_j$  is a random variable (called the *residual*) whose distribution represents the possible *differences*, from the structural component of the model, of the measured value of the response variate for the *j*th element in the sample of *n* elements selected from the study population, if the selecting and measuring processes were to be repeated over and over.

$\sigma$  the (probabilistic) *standard deviation* of the normal model for the distribution of the residual, is a model parameter which represents the (data) *standard deviation* of the measured response variate of the elements of the study population; this standard deviation (and, hence,  $\sigma$ ) *quantifies* the *variation* of the measured response variate of the elements of the study population – as this variation increases, so does the study population (data) standard deviation (and, hence, so does  $\sigma$ ).

6. Another way of writing the response model (5.1.6) is equation (5.1.8), where the model for *R* is now *standard* normal.

$$Y = \mu + \sigma R, \quad R \sim N(0, 1), \quad \text{EPS} \quad \text{-----(5.1.8)}$$

**SOURCE:** The Gabrielle Lippmann quotation in Section 5 on the third side of the Figure is taken from Freedman, D., Pisani, R. and R. Purves: *Statistics*. First Edition, W. W. Norton & Company, New York, 1980, page 275.

□ George E.P. Box, a well-known statistician, is quoted as saying: “All models are wrong, some are useful.” Give an explanation, which could be understood by an intelligent but non-technical audience, of what Dr. Box meant.