

PRINCIPLES OF MEDICAL STATISTICS

III.—PRESENTATION OF STATISTICS

IN dealing with a series of observations the first object must be to express them in some simple form which will permit, directly or by means of further calculations, conclusions to be drawn. The publication, for instance, of a long series of case results is not particularly helpful (beyond providing material for interested persons to work upon), for it is impossible to detect from the unsorted mass of raw material relationships between the various factors at issue. The worker must first consider the questions which he believes the material is capable of answering, and then determine the form of presentation which brings out the true answers most clearly. For instance, let us suppose the worker has amassed a series of after-histories of patients treated for gastric ulcer and wishes to assess the value of the treatment given, using as a measure the amount of incapacitating illness suffered in subsequent years. There will be various factors, the influence of which it will be of interest to observe. Is the age or sex of the patient material to the upshot? Division of the data must be made into these categories and tables constructed to show how much subsequent illness was in fact suffered by each of these groups. Is the after-history affected by the type of treatment? A further tabulation is necessary to explore this point. And so on. The initial step must be to divide the observations into a relatively small series of groups, those in each group being considered alike in that characteristic for the purpose in hand. To take another example, in the Table showing the fatality-rate from scarlet fever of hospital cases, children within each year of age up to 10 and in each five-year group from 10 to 20 are considered alike with respect to ages.

Table showing Hospital Cases of Scarlet Fever,
1905-14

Fatality-rate at ages

Age in years.	Number of cases.	Number of deaths.	Fatality-rate (%).
0-	46	18	39.1
1-	383	43	11.2
2-	881	50	5.7
3-	1169	60	5.1
4-	1372	36	2.6
5-	1403	24	1.7
6-	1271	22	1.7
7-	986	21	2.1
8-	864	6	0.7
9-	673	5	0.7
10-	1965	14	0.7
15-20	513	3	0.6

The fatality, or case-mortality, rate is the proportion of patients with a particular disease who die.

It is, of course, possible that by this grouping we are concealing differences. The fatality-rate at 0-6 months may differ from the fatality-rate at 6-12 months, at 12-18 months it may differ from the rate at 18-24 months. To answer that question further subdivision—if the number of cases justifies it—would be necessary. In its present form (accepting

the figures of hospital cases at their face value) the grouping states that fatality declines nearly steadily with age, a conclusion which it would be impossible to draw from the 11,526 original unsorted and ungrouped records. The construction of a *frequency distribution* is the first desideratum—i.e., a table showing the frequency with which there are present individuals with some defined characteristic or characteristics.¹

As a general rule the distribution should be drawn up on a fine basis—i.e., with a considerable number of groups, for if this basis prove too fine owing to the numbers of observations being few, it is possible to double or treble the group-interval by combining the groups. If, on the other hand, the original grouping is made too broad, the subdivision of the groups is impossible without re-tabulating much of the material.

Statistical Tables

To return to the original table, this may be used in illustration of certain basic principles in the presentation of statistical data.

(i) The contents of the table as a whole and the items in each separate column should be clearly and fully defined. For lack of sufficient headings, or even any headings at all, many published tables are quite unintelligible to the reader without a search for clues in the text (and not always then). For instance, if the heading given on the left of the table were merely "age," it would not be clear whether the groups refer to years or months of life. The unit of measurement must be included.

(ii) If the table includes rates, the base on which they are measured must be clearly stated—e.g., death-rate per cent., or per thousand, or per million, as the case may be (a very common omission in published tables). To know that the fatality-rate is "20" is not helpful unless we know whether it is 20 in 100 patients who die (1 in 5) or 20 in 1000 (1 in 50).

(iii) Whenever possible the frequency distributions should be given in full. These are the basic data from which conclusions are being drawn and their presentation allows the reader to check the validity of the author's arguments. The publication merely of certain values descriptive of the frequency distribution—e.g., the arithmetic mean or average, severely handicaps other workers. For instance, the information that the mean age at death of patients with cancer of the lung is 54.8 years and with cancer of the stomach is 62.1 years is of very limited value in the absence of any knowledge of the distribution of ages at death in the two classes.

(iv) Rates or proportions should not be given alone without any information as to the numbers of observations upon which they are based. In presenting experimental data, and indeed nearly all statistical data, this is a fundamental rule (which, however, is constantly broken). For example, the fatality-rate from small-pox in England and Wales (ratio of registered deaths to notified cases) was 42.9 per cent. in 1917, while in the following year, 1918, it was only 3.2 per cent. This impressive difference becomes less convincing of a real change in virulence when we note that in 1917 there were but 7 cases notified, of whom 3 died, and in 1918 only 63 of whom 2 died.

¹ The choice of groups and the group-intervals and the calculation of averages and other values from such distributions are discussed in numerous text-books of statistical method—e.g., Woods and Russell: *An Introduction to Medical Statistics*. London: P. S. King and Son. 1936.

(Though the low rate of 1918 *may* mark the presence of variola minor.) "It is the essence of science to disclose both the data upon which a conclusion is based and the methods by which the conclusion is attained." By giving only rates or proportions (and by omitting the actual numbers of observations or frequency distributions) we are excluding the basic data. In their absence we can draw no valid conclusion whatever from, say, a comparison of two, or more, percentages. Even when the number of observations is small, 20-25 perhaps, there is no reason why a percentage distribution or a rate should not be calculated but, as discussed later, particular care will have to be exercised in drawing conclusions.

(v) Full particulars of any deliberate exclusions of observations from a collected series must be given, the reasons for and the criteria of exclusion being clearly defined. For example, if it be desired to measure the success of an operation for, say, cancer of the breast, it might, from one aspect, be considered advisable to take as a measure the percentage of patients surviving at the end of 5 years *excluding those who died under the operation itself*—i.e., the question asked is "what is the survival-rate of patients upon whom the operation is successfully carried out?" It is obvious that these figures are not comparable with those of observers who have included the operative mortality. If the exclusion that has been made in the first case is not clearly stated, no one can necessarily deduce that there is a lack of comparability between the records of different observers, and misleading comparisons are likely to be made. Similarly one worker may include among the subsequent deaths only those due to cancer and exclude unrelated deaths—e.g., from accident—while another includes all deaths, irrespective of their cause. Definition of the exclusions will prevent unjust comparisons.

Sometimes exclusions are inevitable—e.g., if in computing a survival-rate some individuals have been lost sight of so that nothing is known of their fate. The number of such individuals must invariably be stated and it must be considered whether the lack of knowledge extends to so many patients as to stultify conclusions. For instance, if 1000 patients were originally observed, 300 are known to be dead at the end of 5 years, 690 are known to be alive, and 10 have been lost sight of, this lack of knowledge cannot appreciably affect the survival-rate. At the best, presuming the 10 are all alive, 70 per cent. survive (700 out of 1000); at the worst, presuming the 10 are all dead, 69 per cent. survive (690 out of 1000). But if 300 are known to be dead, 550 are known to be alive, and 150 have been lost sight of, the upper and lower limits are 70 per cent. surviving (700 out of 1000, presuming the 150 are all alive) and 55 per cent. surviving (550 out of 1000 presuming the 150 are all dead), an appreciable difference. To measure the survival-rate on only those patients whose history is known, or, what comes to the same thing, to divide the 150 into alive and dead according to the proportions of alive and dead in the 850 followed up successfully, is certainly dangerous. The characteristic "lost sight of" *may* be correlated with the characteristics "alive or dead"; in other words, a patient who cannot be traced may be more likely to be dead than a patient who can be traced (or vice versa), in which case the ratio of alive to dead in the untraced cases cannot be the same as the ratio in the traced cases. Calculation of the possible upper rate shows at least the margin of error.

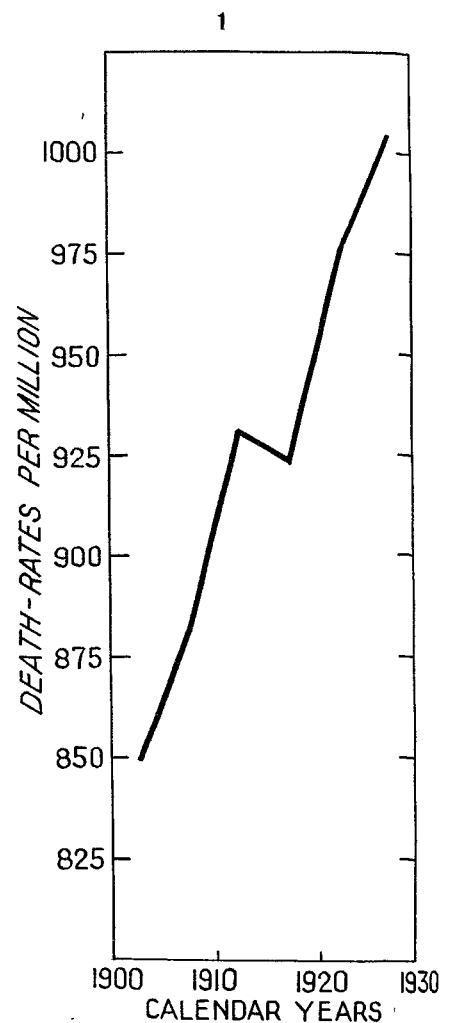
Beyond these few rules it is very difficult, if not

impossible, to lay down laws for the construction of tables. The whole issue is the arrangement of data in a concise and easily read form. In acquiring skill in the construction of tables probably the best way is, as Pearl suggests, to consider critically published tables with such questions as these in mind: "What is the *purpose* of this table? What is it *supposed* to accomplish in the mind of the reader? . . . wherein does its failure of attainment fall?"² Study of the tables published by the professional statistician—e.g., in the Registrar-General's Annual Reports—will materially assist the beginner.

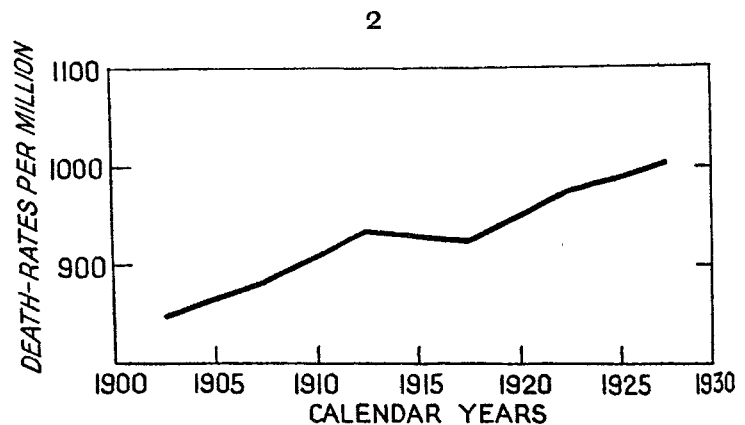
Graphs

Even with the most lucid construction of tables such a method of presentation always gives difficulties to the reader, especially to the non-statistically-minded reader. The presentation of the same material diagrammatically often proves a very considerable aid and has much to commend it if certain basic principles are not forgotten.

(i) The sole object of a diagram is to assist the intelligence to grasp the meaning of a series of figures by means of the eye. If—as is unfortunately often the case—the eye itself is merely confused by a criss-cross of half a dozen, or even a dozen, lines, the sole object is defeated. The criterion must be that the eye can with reasonable ease follow the movement of the various lines on the diagram from point to point and thus observe what is the change in the value of the ordinate (the vertical scale) for a



² Raymond Pearl: Medical Biometry and Statistics. Philadelphia and London: W. B. Saunders Co., Ltd., 2nd ed., 1930.



FIGS. 1 and 2.—Standardised death-rates from cancer in England and Wales in each quinquennium from 1901-05 to 1926-30.

given change in the value of the abscissa (the horizontal scale).

(ii) The second point to bear in mind in constructing *and in reading* graphs is that by the choice of scales the same figures can be made to appear very different to the eye. Figs. 1 and 2 are an example. Both show the same figures—namely, the death-rates (standardised) from cancer in England and Wales in each quinquennium between 1901 and 1930. In Fig. 1 the increase in mortality that has been recorded appears at a cursory glance to be exceedingly rapid and of serious magnitude, while in Fig. 2 a slow and far less impressive rise is suggested. The difference is, of course, due to the difference in the vertical and horizontal scales. In reading graphs, therefore, the scales must be carefully observed and the magnitude of the changes interpreted by a rough translation of the points into actual figures. In drawing graphs undue exaggeration or compression of the scales must be avoided, and it must be considered whether a false impression is conveyed, as quite frequently happens, if the vertical scale does not start at zero but at some point

appreciably above it. Graphs should always be regarded as subsidiary aids to the intelligence and *not* as the evidence of associations or trends. That evidence must be largely drawn from the statistical tables themselves. It follows that graphs should never be a *substitute* for statistical tables. An entirely deaf ear should be turned to such editorial pleading as this: "if we print the graphs would it not be possible to take some of the tables for granted? Having given a sample of the process by which you arrive at the graph is it necessary in each case to reproduce the steps?" The retort to this request is that statistical tables are *not* a step to a diagram, they are the basic data. Without these basic data the reader cannot adequately consider the validity of the author's deductions, and he cannot do any further analysis of the data, if he should wish, without laboriously and inaccurately endeavouring to translate the diagram back into the statistics from which it was originally constructed (and few tasks are more irritating).

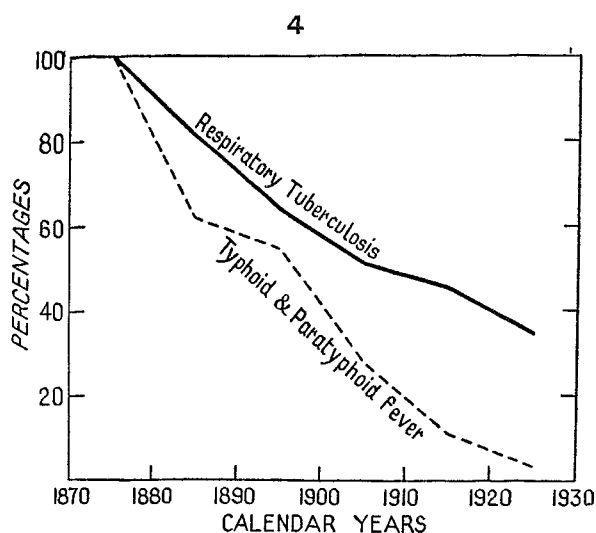
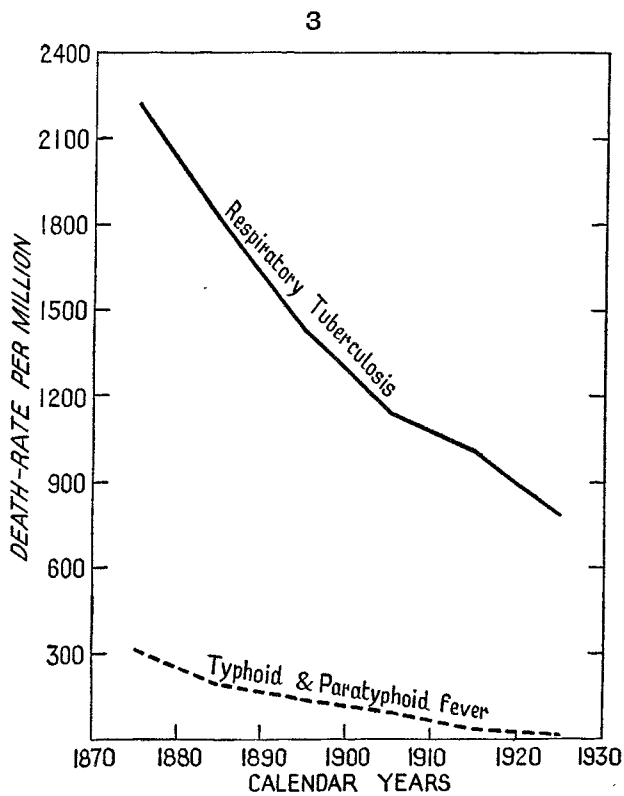
The problem of scale illustrated in Figs. 1 and 2 is also an important factor in the comparison of trend lines. Fig. 3 shows the trend of the death-rates from respiratory tuberculosis and from typhoid fever in England and Wales from 1870 to 1930. Unless the scale of the ordinate (the vertical scale) is carefully considered, the inference drawn from this graph might well be that *relatively* the mortality from respiratory tuberculosis has declined more than the mortality from typhoid fever. Actually the precise reverse is the case—relatively typhoid fever has declined considerably more than respiratory tuberculosis. In 1921–30 the rate from respiratory tuberculosis was 34 per cent. of the rate recorded in 1871–80, while the rate from typhoid fever was but 3·4 per cent. of its earlier level. *Absolutely* respiratory tuberculosis shows the greater improvement (from 2231 deaths per million in 1871–80 to 768 in 1921–30, compared with 321 and 11 for typhoid fever); but *relatively* typhoid fever shows the advantage. If it is the relative degree of improvement that is at issue Fig. 3 is insufficient. For this purpose the rates in each decade may be converted into percentages based upon the rate in the first decade, as is shown in Fig. 4.

It is a *sine qua non* with graphs, as with tables, that they form self-contained units, the contents of which can be grasped without reference to the text. For this purpose inclusive and clearly stated headings must be given, the meaning of the various lines indicated, and a statement made against the ordinate and abscissa of the characteristics to which these scales refer (vide Figs. 1 to 4).

Summary

For the comprehension of a series of figures tabulation is essential; a diagrammatic representation (*in addition to tables but not in place of them*) is often of considerable aid. Both tables and graphs must be entirely self-explanatory without reference to the text. As far as possible the original observations should always be reproduced (in tabulated form showing the actual numbers belonging to each group) and not given only in the form of percentages—i.e., the percentages of the total falling in each group. The exclusion of observations from the tabulated series on any grounds whatever must be stated, the criterion upon which exclusion was determined clearly set out, and usually the number of such exclusions stated. Conclusions should be drawn from graphs only with extreme caution and only after careful consideration of the scales adopted.

A. B. H.



FIGS. 3 and 4.—Standardised death-rates from (a) respiratory tuberculosis, and (b) typhoid and paratyphoid fever in England and Wales. In Fig. 4 the rate from each disease in each decade is expressed as a percentage of the corresponding rate in 1871–80.