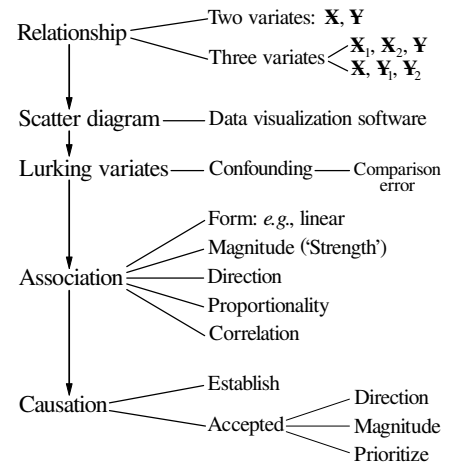University of Waterloo      W. H. Cherry

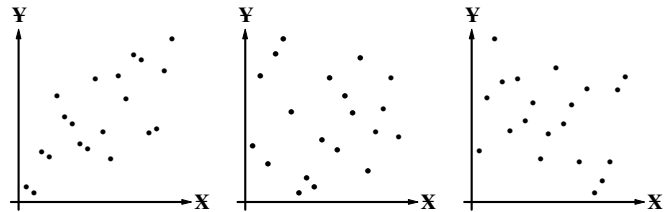# RELATIONSHIPS IN STATISTICS:  Association – Statistical Issues

A **relationship** in statistics arises from the following sequence of happenings.

− We observe that the value of a *response* variate $\mathbf{Y}$ *changes* (*i.e.*, shows *variation*) over the elements or units of a group, such as a target population, a study population, a respondent population or a sample.

    ○ It is implicit that there are one or more *cause(s)* of these changes (*i.e.*, of this variation) in $\mathbf{Y}$.

− We wish to account for these changes (*i.e.*, for this variation) – we introduce the idea of an *explanatory* variate $\mathbf{X}$.

− We look for *association* between the values of $\mathbf{Y}$ and $\mathbf{X}$ (*e.g.*, using a scatter diagram – see below) – a relationship is the *connection* (if any) between *changes* in $\mathbf{X}$ and *changes* in $\mathbf{Y}$ (or in the *average* of $\mathbf{Y}$).

    ○ If (data establish that) $\mathbf{Y}$ remains *un*changed while $\mathbf{X}$ changes (or *vice versa*), there is *no* $\mathbf{X}$-$\mathbf{Y}$ relationship, an idea of *un*connectedness captured by one sense of the word **independent** (but see Note 1 overleaf at the bottom of page HL58.2).

     **+** We should recognize the distinction between the 'behavioural unconnectedness' of *independence* and the 'spatial separateness' captured by *disjoint*, as in 'disjoint events'.

A framework of terminology for discussing data-based investigating of statistical relationships is given in the schema above.

   ∗ A **scatter diagram** is a Cartesion plot with a response variate (or estimated residual) on the vertical axis, an explanatory variate on the horizontal axis.

    − A scatter diagram – a graphical attribute – is a useful way to *look* at data for an $\mathbf{X}$-$\mathbf{Y}$ relationship.  Each element (or unit) appears as a dot (or other appropriate symbol) located at the co-ordinates determined by its $\mathbf{X}$ and $\mathbf{Y}$ values; three examples are shown at the right.

The foregoing description of a relationship in statistics refers to the *association* of $\mathbf{Y}$ and $\mathbf{X}$;  this Statistical Highlight #58 defines association in statistics.  In doing so, it provides background for discussing association between (or among) explanatory variates, and association between them and the response variate, in Statistical Highlight #62.

   ∗ **Association:** if a scatter diagram shows a clustering of its points about, say, a line with positive slope (*i.e.*, we see that, as $\mathbf{X}$ increases, $\mathbf{Y}$ also tends to increase), we say $\mathbf{X}$ and $\mathbf{Y}$ show a (positive) *association*;  there is *moderate* positive association of $\mathbf{X}$ and $\mathbf{Y}$ in the left-hand diagram of the three scatter diagrams at the right above.  The right-hand diagram shows *weak negative* association and the middle diagram shows *little or no* association.

Questions of statistical interest about an association are:

− what is its **form**? – for example, can the trend be modelled by a *straight line* (*i.e.*, is the trend *linear*)?

− what is its **magnitude**? – for linear association, what is the magnitude of the *slope*?

− what is its **direction**? – for linear association, is the slope *positive* or *negative*?

    **+** **Proportionality** refers to a straight-line $\mathbf{X}$-$\mathbf{Y}$ association *through the origin.*

    **+** The sign of the direction (positive or negative) of a linear association is *also* the sign of correlation, but the connection between the *magnitudes* of slope and correlation is more complicated – see Section 8 in Statistical Highlight #66.

− **Correlation:** a numerical measure of *tightness of clustering* of the points on a scatter diagram about a straight line – historically, correlation is denoted r (c would have been a better choice) and its values lie in the interval [−1, 1]; the respective correlations are about +0.7, 0 and −0.25 for the three scatter diagrams at the right above.

    **+** If the points of a scatter diagram lie *on* a straight line with positive slope, r = +1;

    **+** if the points of a scatter diagram lie *on* a straight line with negative slope, r = −1;

    **+** if the points of a scatter diagram are haphazardly spread over its rectangular area, r is zero or close to it.

Correlation is discussed and illustrated in detail in Statistical Highlight #66.

The background discussion for comparison error [*e.g.*, at the beginning of Section 2 on the overleaf side (page HL57.2) in Statistical Highlight #57] refers to a group of elements (or units) with a lurking variate ($\mathbf{Z}$) whose distribution of values differs, over the elements (or units) of the group, for different values of the focal variate $\mathbf{X}$.  A consequence of this behaviour of $\mathbf{Z}$ is that the values of $\mathbf{X}$ and $\mathbf{Z}$ are *associated*, as illustrated in the scatter diagrams given overleaf on page HL58.2, for respondent populations with 4 or 9 elements and $\mathbf{Z}$ values (shown beside the points) like 0, 1, 2 and 3. [*Distinct* $\mathbf{Z}$ values for *all* population elements, as in diagrams (1) and (5), is rare in real populations.]

2006-06-20

○ In diagram (1) at the right, the element with $\mathbf{Z}=2$ when $\mathbf{X}=0$ has $\mathbf{Z}=1$ when $\mathbf{X}=1$; thus, the change in the average of $\mathbf{Y}$ (indicated by a short horizontal line) from 2.6 to 3.6, as $\mathbf{X}$ changes from 0 to 1, no longer reflects *only* the effect of changing $\mathbf{X}$; a limitation is therefore imposed on the Answer about the $\mathbf{X}$-$\mathbf{Y}$ relationship by comparison error due to the behaviour of $\mathbf{Z}$ not being taken into account (or due to confounding by $\mathbf{Z}$).

+ Because $\mathbf{Z}$ changes with $\mathbf{X}$, there is a (weak) $\mathbf{X}$-$\mathbf{Z}$ association, quantified by a correlation of about $-0.11$ over the eight $(\mathbf{X}, \mathbf{Z})$ values; by contrast, when $\mathbf{Z}$ does *not* change with $\mathbf{X}$, the $\mathbf{X}$-$\mathbf{Z}$ correlation is *zero*; this is the case in diagrams (6), (7) and (8) at the right below diagram (1) – these three diagrams are used in Statistical Highlight #62 to illustrate the formal definition of what *we* mean by causation in statistics (the short horizontal lines again show the average of $\mathbf{Y}$ for each $\mathbf{X}$ value).

An extension of the illustration in diagram (1) is to the case of repeated values involving *more than two* $\mathbf{X}$ values.

○ In diagram (2), if $\mathbf{Z}$ has the *same* value (say 1) for all nine units whose $\mathbf{X}$ and $\mathbf{Y}$ values yield this scatter diagram, there is *no* $\mathbf{X}$-$\mathbf{Y}$ relationship in the sense that the $\mathbf{X}$-$\mathbf{Y}$ correlation is zero.

+ This *lack* of $\mathbf{X}$-$\mathbf{Y}$ relationship is also reflected by the slope of *zero* for the straight line (shown dashed) which summarizes the trend in the points of the scatter diagram.

+ When interpreting a scatter diagram like (2), it is easy to confuse *explicit* knowledge that there is the same $\mathbf{Z}$ value among the elements, with *assuming* this to be the case by *ignoring* the elements' $\mathbf{Z}$ value(s) – see Statistical Highlight #29.

+ In diagrams (6) to (8) above, reminiscent of an *experimental* Plan with *two* values of the focal variate, we can accommodate *different* values of the potential confounder $\mathbf{Z}$ among the elements; by contrast, in diagram (2) above, reminiscent of an *observational* Plan, the elements must have the *same* $\mathbf{Z}$ value to meet the requirement for $\mathbf{Z}$ to remain fixed to avoid the limitation imposed on an Answer about an $\mathbf{X}$-$\mathbf{Y}$ relationship by comparison error due to this lurking variate.

[Experimental and observational Plans are discussed in Sections 4 and 5 on pages HL63.3 to HL63.6 in Statistical Highlight #63 (and also in Statistical Highlight #9).]

○ Diagram (3) is visually the *same* as diagram (2) but the $\mathbf{Z}$ values *change* with $\mathbf{X}$ – the association of $\mathbf{X}$ and $\mathbf{Z}$ can be quantified as a correlation of about $+0.7$; as indicated by the dashed lines, there is now a (strong) *positive* $\mathbf{X}$-$\mathbf{Y}$ association among points for which $\mathbf{Z}$ values are held fixed (*i.e.*, for points with the *same* $\mathbf{Z}$ value).

○ In diagram (4), again visually the same as diagrams (2) and (3), a *different* distribution of the *same* set of $\mathbf{Z}$ values as in diagram (3) yields a (strong) *negative* $\mathbf{X}$-$\mathbf{Y}$ association – the $\mathbf{X}$-$\mathbf{Z}$ correlation is again about $+0.7$.

+ In diagrams (3) and (4), the $\mathbf{X}$-$\mathbf{Y}$ relationship is the *same* for the three values of $\mathbf{Z}$; the matter of *different* $\mathbf{X}$-$\mathbf{Y}$ relationships for different $\mathbf{Z}$ values is pursued in Statistical Highlight #29 (and also in Statistical Highlight #60).

+ Like diagram (1), diagrams (3) and (4) illustrate, in a broader context, the limitation imposed on an Answer about an $\mathbf{X}$-$\mathbf{Y}$ relationship by comparison error, when the elements' $\mathbf{Z}$ values do not remain fixed (are not the same) as $\mathbf{X}$ changes, and this behaviour is *not* taken into account (*e.g.*, when interpreting an $\mathbf{X}$-$\mathbf{Y}$ scatter diagram).

A special case is when $\mathbf{Z}$ changes with $\mathbf{X}$ but in such a way that their values have *zero* correlation; an illustration is shown at the right in diagram (5), which is adapted from diagram (6) above. In such a situation, *despite* the confounding, it *is* possible (under an assumption of *additive* effects) to estimate the effect of $\mathbf{X}$ on the average of $\mathbf{Y}$.

● This idea is exploited in Design of Experiments (DOE) when investigating a relationship with *two or more* focal variates – see Notes 3 to 6 (starting at the botom of page HL68.1) in Statistical Highlight #68.

**NOTE:** 1. The discussion above shows that, when looking at a scatter diagram of bivariate data to assess an $\mathbf{X}$-$\mathbf{Y}$ relationship, experience *out*side statistics with diagrams involving Cartesian axes provides poor preparation for statistics – it is difficult for later statistical training to overcome a mindset (*un*concerned with lurking variates) that arises from more formative earlier experience with such diagrams, starting in elementary school, with on-going exposure in the media, and continuing up to post-secondary-level courses, including calculus and algebra.

● As discussed in Statistical Highlights #60, #51 and #29, lurking variates can affect the interpretation of the presence, *or* the absence, of an observed association (in even the simplest situation) involving $\mathbf{X}$ and $\mathbf{Y}$.