

RELATIONSHIPS IN STATISTICS: Association and Causation – Reality and Illusion

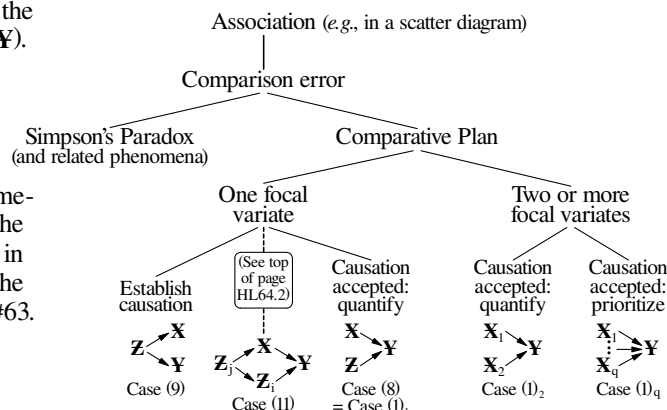
1. Data Visualization

A precept of investigating relationships in statistics is *looking at the data*; that is, *visualizing* the relationship – for example, the course STAT 442 has the title *Data Visualization*. A caveat to this precept is that deciding *how* to visualize the data and how to *interpret* the display(s) must be informed by adequate statistical understanding – seeing may be illusory. In this Highlight #60, we deal in turn with the three main types of Question which arise when investigating statistical relationships.

1. Is there an *association* of \mathbf{X} and \mathbf{Y} and, if so, what is its *form* and/or *magnitude* and/or *direction*?
2. What is the *reason* for an association – e.g., can we *establish* that \mathbf{X} *causes* \mathbf{Y} ?
3. *Accepting* a relationship as causal: what is its *form* – e.g., what is the *effect* of \mathbf{X} on (the average of) \mathbf{Y} ?
which explanatory variate is the *most important* cause of (variation in) \mathbf{Y} ?

This ‘natural’ order of such Questions may differ appreciably from the one in which they are discussed in an introductory course; also, the emphasis in introductory discussion is usually on the relationship of *one* focal variate (\mathbf{X}) to *one* response variate (\mathbf{Y}).

The main issue in our discussion of reality and illusion in data-based investigating of statistical relationships is the limitation imposed on Answers by comparison error arising from confounding by lurking variate(s) [e.g., \mathbf{Z}]. A framework for our present discussion is shown in the schema at the right; it starts with an observed association (e.g., of \mathbf{X} and \mathbf{Y} in the simplest case) [or with *lack* of (expected) association]. The five causal structures are from Statistical Highlights #59 and #63.



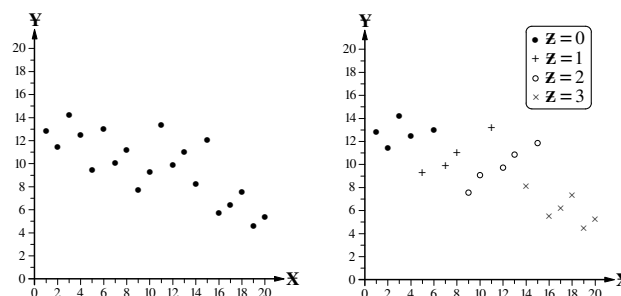
2. Questions of the First Type

One way to investigate association is with a **scatter diagram**

– in its simplest form, Cartesian axes with dots (or other symbols), the coordinates of whose centre are the \mathbf{X} and \mathbf{Y} values of each bivariate observation. However, when examining such diagrams, it is easy to overlook the limitation on an Answer about the \mathbf{X} - \mathbf{Y} relationship imposed by different points on the scatter diagram having *differing* values of a lurking variate \mathbf{Z} . This matter is illustrated by the two versions of the *same* scatter diagram at the right below:

- * in the left-hand version in which \mathbf{Z} values are *ignored*, we see an \mathbf{X} - \mathbf{Y} relationship that could reasonably be modelled by a straight line with a *negative* slope.
- * in the right-hand version, where different symbols for the points denote four different values of some (non-focal) explanatory variate \mathbf{Z} , the straight-line \mathbf{X} - \mathbf{Y} relationship can have a slope which is (close to) zero (when \mathbf{Z} is 0), positive (when \mathbf{Z} is 1 or 2) or negative (when \mathbf{Z} is 3).

NOTES: 1. When looking at a scatter diagram of bi-variate data to assess an \mathbf{X} - \mathbf{Y} relationship, we recognize that experience *outside* statistics with diagrams involving Cartesian axes provides poor preparation for statistics – it is difficult for later statistical training to overcome a mindset (*unconcerned* 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.



2. Looking at multivariate data to try to detect *patterns* which answer Questions about relationships can be aided by statistical software that shows a point cloud in three dimensions on a computer screen with options like:
 - rotating the point cloud in real time; ● using colour to distinguish subsets of the points;
 - linking points (e.g., by using colour) across scatter diagrams which show point clouds for different subsets of the variates (illustrated in Program 10, *Multidimensional Data Analysis in Against All Odds: Inside Statistics*.)
3. The foregoing discussion and scatter diagrams in this Highlight #60 draw attention to the distinction between *conditioning* on \mathbf{Z} and *ignoring* \mathbf{Z} when investigating relationships.
 - * Conditioning is *subdividing*, as discussed at the upper left of page HL51.3 in Statistical Highlight #51.
 - * A **marginal** (probability) distribution, referred to in Section 6 on page HL51.4 in Highlight #51 and illustrated in Tables HL51.10 to HL51.13, is an example of ‘ignoring’ the variate which is absent from the marginal distribution – for instance, in Table HL50.11, \mathbf{X}_2 is absent, in Table HL51.12, \mathbf{X}_1 is absent, and in Table HL51.13, \mathbf{Y} is absent.

NOTES: 3. The right-hand scatter diagram overleaf on page HL60.1 shows the marginal distribution of \mathbf{X} and \mathbf{Y} if we think of the \mathbf{Z} direction as coming vertically up from the page. With the \mathbf{Z} values as given at the upper right of the right-hand version of the diagram and thinking of the page as the plane $\mathbf{Z}=0$, the first five points of the cloud would lie *on* the page; the remaining 14 points would then lie progressively further *above* the page in groups as one moves to the right across the diagram. This discussion reminds us that a marginal distribution is a *projection* – we *see* the marginal distribution of \mathbf{X} and \mathbf{Y} if we look vertically *down* on the diagram (*i.e.*, we look along the \mathbf{Z} axis) to project the three-dimensional point cloud on to the plane of the page.

- It is interesting to speculate on the extent to which the ideas of conditioning and marginalizing (or projecting) provide a basis for understanding the ways in which mathematical models *approximate* reality, bearing in mind a maxim of the late Dr. George E.P. Box, a respected U.S. statistician: *All model are wrong, some are useful.*

4. Statistical Highlight #51 about Simpson's Paradox and related phenomena shows how interpreting *tabulated* data and their visual representation must be informed by adequate statistical understanding, including formulating a *clear* Question and the role of lurking variates.

- Anscombe's four regression data sets in Statistical Highlight #53 illustrate the statistical hazards of working with numerical data summaries without also *looking* at the data.

3. Questions of the Second Type

For the second type, we distinguish four reasons ('cases') for *association* of variates \mathbf{X} and \mathbf{Y} in the presence of a possible confounder (or lurking variate) \mathbf{Z} ; the four relevant cases from Statistical Highlight #59 are shown symbolically at the right, where an arrow denotes causation:

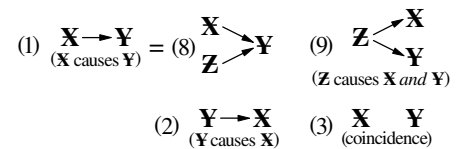
* \mathbf{X} causes \mathbf{Y} (in the presence of confounder \mathbf{Z});

* \mathbf{Y} causes \mathbf{X} ;

* \mathbf{Z} causes \mathbf{X} and \mathbf{Y} – that is, \mathbf{Z} is a cause of \mathbf{X} and \mathbf{Y} (so-called **common cause**);

* coincidence [which often means that both \mathbf{X} and \mathbf{Y} are associated with *time* –

i.e., coincidence is often case (9) in which \mathbf{Z} is time (whatever, if anything, 'causation' by time means)].



A Question about the *actual* reason for an observed \mathbf{X} - \mathbf{Y} association can be answered by a process of *elimination*:

- coincidence [case (3)] requires extra-statistical knowledge to rule it out;

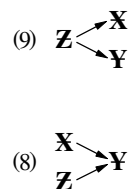
– in a universe with an almost uncountable number of variates changing in value over time, coincidental associations are likely *numerous* but may largely go *unnoticed*;

- identifying correctly which is the response and which the explanatory variate can rule out case (2);

- the remaining two reasons are *direct causation* [case (1) = case (8)] and a *common cause* \mathbf{Z} [case (9)] – an investigation with a comparative Plan that yields acceptable limitation on the Answer imposed by comparison error (due to type 2b confounding – see Statistical Highlight #3) can then be used to try to rule out case (9); if successful, this leaves (direct) causation [case (1)] as the (likely) reason for the association.

– The Plan that manages confounding by a common cause for \mathbf{X} and \mathbf{Y} [case (9)] will also manage possible confounder(s) \mathbf{Z} (or \mathbf{Z}_i) changing as focal variate \mathbf{X} changes [case (8)] in a way that makes *acceptable* the limitation imposed by comparison error (due to type 2a confounding); the common theme of managing type 2 confounding [cases (9) and (8)] is holding \mathbf{Z} *fixed* as \mathbf{X} changes in value.

+ Case (8) is, of course, case (1) with confounder \mathbf{Z} shown explicitly – \mathbf{Y} is a **common response** to \mathbf{X} and \mathbf{Z} .



For Questions of the second type, comparison error results in an Answer that *misidentifies* the cause of \mathbf{Y} – for instance:

- an Answer which says \mathbf{X} is a cause of \mathbf{Y} when (in reality) it is not [case (9)],
- an Answer which says \mathbf{X} is *not* a cause of \mathbf{Y} when (in reality) it is [case (8)].

Diagrams (4) to (8) on the facing page HL60.3 and page HL60.4 illustrate these matters in more detail. A real-world example occurred in *The Globe and Mail* on March 7, 2015, pages M1 and M5, in the article: **\$42 AN HOUR WHY CANADA'S YOUNG ACADEMICS ARE ON THE PICKET LINES.**

The excerpt (from page M5) relevant to this discussion is given at the right. The three variates are:

\mathbf{X} : proportion of the Canadian population with a PhD,

\mathbf{Y} : Canada's level of productivity and innovation,

\mathbf{Z} : Canada's level of development.

The assumption is *that more graduate students equal increased productivity and innovation.* BUT:

It could well be that countries at high levels of development produce both innovation and PhDs.

Thus, there is a need to establish whether the applicable causal structure is our case (1) [= case (8)] or case (9).

In addition to the *causal* Question, there are also statisti-

(1) $\mathbf{X} \rightarrow \mathbf{Y}$

(8) $\mathbf{X} \rightarrow \mathbf{Y}$ and $\mathbf{Z} \rightarrow \mathbf{Y}$

(9) $\mathbf{Z} \rightarrow \mathbf{X}$ and $\mathbf{Z} \rightarrow \mathbf{Y}$

Do we have too many PhDs?

Yet, if graduate programs have grown, it has been because governments have waved money at university administrations, operating under the assumption that more graduate students equal increased productivity and innovation.

In 2011, Ontario announced that it would fund an additional 6,000 master's and doctoral spots even as its undergraduate per-student funding, adjusted for inflation, has been declining for decades.

Turns out that the current instructions for building a PhD student may not translate into the economic gains that governments assume.

"It is an open question whether PhDs are necessary to produce innovation and productivity, or can you do that with MAs," said Daniel Munro, an analyst at the Confer-

RELATIONSHIPS IN STATISTICS: Association and Causation – Reality and Illusion (continued 1)

cally challenging matters of *measuring* a country's productivity and innovation (number of patents?) and level of development.

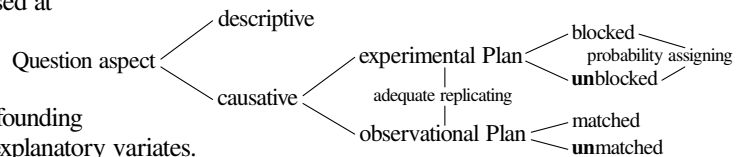
ence Board of Canada who will be releasing a report he co-authored on PhDs in the Canadian labour market this spring.

Dr. Munro says it is true that countries with a higher percentage of PhDs have higher level of productivity and innovation, and tend to produce more patents. However, the relationship between very advanced education and innovation is uncertain. It could well be that countries at high levels of development produce both innovation and PhDs.

4. Questions of the Third Type

For the third type of Question, provided that the \mathbf{X} - \mathbf{Y} relationship is causal [case (1) on the facing page HL60.2], acceptable limitation imposed by comparison error (due to type 2a confounding) can usually be achieved by a comparative experimental (but *not* an observational) Plan developed and executed as discussed at length and illustrated in Statistical Highlights #57 to #70. The schema at the right, from page HL63.2

in Statistical Highlight #63, reminds us of Plan components available to investigators to manage possible confounding by known and by unknown and unmeasured *non-focal* explanatory variates.



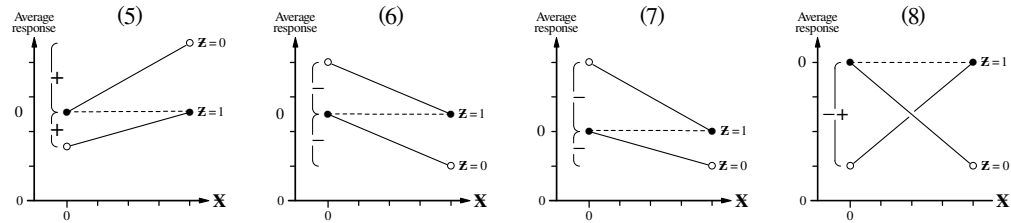
For Questions of the second and third types, which involve an inference from association to causation, the limitation imposed by comparison error (due to type 2 confounding – see Statistical Highlight #3) on Answers can be summarized in two precepts:

- * **Partial or complete false positive:** when *coincidence* can be ruled out as a reason, an \mathbf{X} - \mathbf{Y} association indicates *causation* of \mathbf{Y} but *not necessarily* (only or even in part) by \mathbf{X} .
- * **False negative:** association of \mathbf{X} and \mathbf{Y} may be *absent* even when they have a *causal* relationship.

These precepts formalize the ideas in Note 1 in Statistical Highlight #59 and are illustrated in more detail in eight diagrams below and overleaf on page HL60.4, which show how the *average* of \mathbf{Y} changes with \mathbf{X} in the presence of a (binary) confounder \mathbf{Z} ; except for diagram (4) [case (9) on the facing page HL60.2], the causal connections are those in case (8). In these eight diagrams:

- the circles (○) represent response variate averages that do *not* occur under the Plan because these averages would require the elements' \mathbf{Z} values to be *different* from their actual values (counterfactuals – see page HL9.5 in Statistical Highlight #9).
 - the dots (●) represent response variate averages for elements whose \mathbf{Z} values mean that these averages *are* observed (data),
 - the solid lines are the *unobserved* 'reality'; ● the dashed lines are the observed 'illusion';
 - the horizontal axis has a *quantitative* scale, rather than showing \mathbf{X} as an indicator variate with values of 0 and 1.
- In the left-hand diagram (1) below, the (positive) effect of \mathbf{X} on (the average of) \mathbf{Y} is represented by the two braces to the *right* of the \mathbf{Y} -axis; however, if there is confounding (\mathbf{Z} changing from 0 to 1 as \mathbf{X} changes), the effect of \mathbf{X} on \mathbf{Y} would be *observed* as the brace to the *left* of the axis. Thus, in diagram (1):
 - the dashed line involving comparison error yields a *wrong* (exaggerated) magnitude for the effect of \mathbf{X} on \mathbf{Y} .
 - the slope of the solid lines showing the relationship of \mathbf{X} to the average of \mathbf{Y} are *unaffected* by the value of \mathbf{Z} – that is, there is no *interaction* of \mathbf{X} and \mathbf{Z} .
 - Diagram (2) is like diagram (1) except there *is* interaction of \mathbf{X} and \mathbf{Z} – the solid lines have *different* slopes.
 - Diagrams (1) and (2) illustrate partial false positives; association of \mathbf{X} and \mathbf{Y} *does* correspond to causation of \mathbf{Y} by \mathbf{X} but confounding by \mathbf{Z} distorts reality (creates illusion).
 - In diagram (3) with a *different* (negative) \mathbf{Z} - \mathbf{Y} relationship, comparison error is *wrong direction* for the effect of \mathbf{X} on \mathbf{Y} .
 - An Answer which, due to comparison error, is a *wrong direction* for an \mathbf{X} - \mathbf{Y} relationship is also a case of *wrong value* (and, usually, *wrong magnitude*) but, as in Simpson's Paradox (see Statistical Highlight #51), the more dramatic (directional) manifestation of comparison error is usually emphasized.
 - The right-hand diagram (4) illustrates (without interaction) a (complete) *false positive* Answer – there is \mathbf{X} - \mathbf{Y} association without \mathbf{X} - \mathbf{Y} causation; this is the situation when \mathbf{Z} is a common cause of \mathbf{X} and \mathbf{Y} [case (9) on the facing page HL60.2].
 - The two solid lines in diagram (4) have *zero* slope so we say \mathbf{X} and \mathbf{Y} are *independent* when $\mathbf{Z}=0$ and when $\mathbf{Z}=1$ – for (unit) change in \mathbf{X} , there is *no* change in (the average of) \mathbf{Y} *provided* \mathbf{Z} is (held) fixed – see also Notes 8 and 7 on page HL51.8 in Highlight #51 and Section 13 on page HL89.18 in Statistical Highlight #89.
 - A false negative Answer, as in diagram (5) overleaf, can occur in other ways; in diagram (6), such an Answer again arises from confounding *without* interaction, and in diagrams (7) and (8) [as in diagram (5)] from *both* confounding *and* interaction.
 - In diagrams (2), (5) and (7), interaction of \mathbf{X} and \mathbf{Z} is *incidental* to the manifestation of comparison error when quantifying the magnitude and/or direction of the effect of \mathbf{X} on \mathbf{Y} ; only in diagram (8) is the interaction *essential* to this manifestation.

+ (Complete)
false negatives
involve the
special case of
exact cancel-
lation of the
effects of \mathbf{X}
and \mathbf{Z} on \mathbf{Y} .



The foregoing discussion of diagrams (1) to (8), involving type 2 confounding, is summarized in Table HL60.1 at the right.

NOTES: 5. For Questions of the third type, concerned with quantifying the effect of \mathbf{X} on \mathbf{Y} , the eight \mathbf{X} -Average response diagrams overleaf on page HL60.3 and above are only *illustrative* of comparison error because its particular manifestation depends on the interplay of several matters affecting diagram appearance, including:

- the magnitude(s) and direction(s) of the slope(s) of the \mathbf{X} - \mathbf{Y} relationships;
- the magnitude of the \mathbf{Z} - \mathbf{Y} relationship, reflected in the *vertical separation* of the solid lines;
- the absence or presence of interaction (of \mathbf{X} and \mathbf{Z} in their effects on \mathbf{Y});
- the *forms* of the \mathbf{X} - \mathbf{Y} and \mathbf{Z} - \mathbf{Y} relationships (e.g., linear or *nonlinear*).

6. The *two* descriptions in the last column of Table HL60.1 above for diagrams (4) to (8) overleaf on page HL60.3 and above remind us that, for quantitative variates, comparison error is [despite the second (more dramatic) descriptions] a *wrong value* (and, usually, a *wrong magnitude*) for the effect of \mathbf{X} on (the average of) \mathbf{Y} . As a consequence, diagrams (4) to (8) illustrate the occurrence of comparison error for Questions of *both* the second and third type.

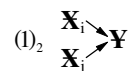
- Similarly, Simpson's Paradox (discussed in Statistical Highlight #51) is, under subdivision on an additional explanatory variate, a *change in value* for a 'treatment' effect, resulting in a (dramatic) change of direction of the 'effect' [analogous to diagram (3) overleaf on page HL60.3].

7. For the (unattainable) ideal [criterion (1) in Statistical Highlight #62] of *all* non-focal explanatory variates remaining fixed when \mathbf{X} changes to make apparent its relationship to \mathbf{Y} , the Answer has *no* limitation imposed by comparison error but, when a possible confounder \mathbf{Z} does *not* remain fixed, comparison error *does* impose a limitation on the Answer about the \mathbf{X} - \mathbf{Y} relationship. These limitations, from the foregoing discussion in this Highlight #60, are summarized in Table HL60.2 below; it reminds us that the nature of the limitation depends on the Question context.

Table HL60.2		Comparison error	Confounding type	Other name
Question context	Question			
Comparing proportions	Does \mathbf{Y} increase or decrease with \mathbf{X} ?	Wrong direction for an \mathbf{X} - \mathbf{Y} association	3	Simpson's Paradox
Establishing causation	Is \mathbf{X} the (or a) <i>cause</i> of \mathbf{Y} ?	Wrong cause identified for \mathbf{Y}	2b	Confounding
Quantifying a treatment effect	What is the effect of \mathbf{X} on \mathbf{Y} ?	Wrong magnitude for effect of \mathbf{X} on \mathbf{Y}	2a	Confounding
		Wrong direction for effect of \mathbf{X} on \mathbf{Y}	2a	Confounding
Improving a process	Is \mathbf{X} the <i>most important</i> cause of \mathbf{Y} ?	Wrong main cause identified for \mathbf{Y}	2a	Confounding
More than one focal variate	What is the effect of each \mathbf{X}_i on \mathbf{Y} ?	Wrong effect of one or more \mathbf{X}_i on \mathbf{Y}	1, 2a	-----

- If Simpson's Paradox (in the first line of Table HL60.2) is observed in a *sample* but the Question involves the corresponding *population*, limitation is imposed on the Answer by *sample* error as well as by *comparison* error.
 - Regardless of whether the investigation involves a (respondent) population census or sample, study, non-response and measurement error also likely need to be managed in the Plan – e.g., see Section 6 on page HL21.4 in Statistical Highlight #21 and also Statistical Highlight #6.

- The second-last line of Table HL60.2 deals with **process improvement** – prioritizing the explanatory variates responsible for variation in \mathbf{Y} ; the relevant causal connections can then be thought of as case (1)₂ [repeated at the right from Statistical Highlight #63]. Prioritizing a set of m explanatory variates \mathbf{X}_1 ($1 = 1, 2, \dots, m$) can be achieved by successively prioritizing *pairs* of variates \mathbf{X}_i and \mathbf{X}_j . The effects of confounding are the *same* as for *quantifying* the effect of \mathbf{X} on \mathbf{Y} and the manifestation of comparison error is identifying the **wrong main cause** of (variation in) \mathbf{Y} .



- Because prioritizing explanatory variates involves *two or more* focal variates, the Plan should make provision for estimating *interaction* effect(s).

+ ['Perfect or type 1 confounding' among some treatment effects arising from using a *fractional* factorial treatment structure may be acceptable in the Question context – see Note 5 on page HL68.2 in Statistical Highlight #68 and also the discussion of type 1 confounding in Statistical Highlight #3.]

- In the *last* line of Table HL60.2, again involving estimating treatment effects for two or more focal variates, the Plan should provide for estimating both *main* and *interaction* effects.