

# Berkson's paradox

**Berkson's paradox** also known as **Berkson's bias** or **Berkson's fallacy** is a result in conditional probability and statistics which is often found to be counterintuitive, and hence a veridical paradox. It is a complicating factor arising in statistical tests of proportions. Specifically, it arises when there is an ascertainment bias inherent in a study design. The effect is related to the explaining away phenomenon in Bayesian networks.

The most common example of Berkson's paradox is a false observation of a *negative* correlation between two positive traits, i.e., that members of a population which have some positive trait tend to lack a second. Berkson's paradox occurs when this observation appears true when in reality the two properties are unrelated—or even *positively* correlated—because members of the population where both are absent are not equally observed. For example, a person may observe from their experience that fast food restaurants in their area which serve good hamburgers tend to serve bad fries and vice versa; but because they would likely not eat anywhere where **both** were bad, they fail to allow for the large number of restaurants in this category which would weaken or even flip the correlation.

It is often described in the fields of medical statistics or biostatistics, as in the original description of the problem by Joseph Berkson.

## Contents

- Statement**
  - Explanation
- Examples**
- See also**
- References**
- External links**

## Statement

The result is that two independent events become conditionally dependent (negatively dependent) given that at least one of them occurs. Symbolically:

If 



0
<
P
(
A
)
<
1
,
0
<
P
(
B
)
<
1
,
and
P
(
A
|
B
)
=
P
(
A
)
,
then
P
(
A
|
B
,
A
∪
B
)
<
P
(
A
|
A
∪
B
)
.


{\displaystyle 0 < P(A) < 1, 0 < P(B) < 1, and P(A|B) = P(A), then P(A|B,A\cup B) < P(A|A\cup B) .}

- Event *A* and event *B* may or may not occur
- P*(*A*|*B*), a conditional probability, is the probability of observing event *A* given that *B* is true.
- Explanation: Event *A* and *B* are independent of each other
- P*(*A*|*B*,*A*∪*B*) is the probability of observing event *A* given that *B* and (*A* or *B*) occurs. This can also be written as *P*(*A*|*B*∩(*A*∪*B*))
- Explanation: The probability of *A* given both *B* and (*A* or *B*) is smaller than the probability of *A* given (*A* or *B*)

In other words, given two independent events, if you consider only outcomes where at least one occurs, then they become negatively dependent, as shown above.

### Explanation

The cause is that the *conditional* probability of event *A* occurring, *given* that it or *B* occurs, is inflated: it is higher than the *unconditional* probability, because we have *excluded* cases where *neither* occur.

*P*(*A*|*A*∪*B*) > *P*(*A*)
conditional probability inflated relative to unconditional

One can see this in tabular form as follows: the yellow regions are the outcomes where at least one event occurs (and ~**A** means "not **A**").

	<b>A</b>	<b>~A</b>	
<b>B</b>	<b>A &amp; B</b>	<b>~A &amp; B</b>	
<b>~B</b>	<b>A &amp; ~B</b>	<b>~A &amp; ~B</b>	

For instance, if one has a sample of 100, and both *A* and *B* occur independently half the time ( *P*(*A*) = *P*(*B*) = 1/2 ), one obtains:

	<b>A</b>	<b>~A</b>	
<b>B</b>	<b>25</b>	<b>25</b>	
<b>~B</b>	<b>25</b>	<b>25</b>	

So in 75 outcomes, either *A* or *B* occurs, of which 50 have *A* occurring. By comparing the conditional probability of *A* to the unconditional probability of *A*:

*P*(*A*|*A*∪*B*) = 50/75 = 2/3 > *P*(*A*) = 50/100 = 1/2

We see that the probability of *A* is higher (2/3) in the subset of outcomes where (*A* or *B*) occurs, than in the overall population (1/2). On the other hand, the probability of *A* given both *B* and (*A* or *B*) is simply the unconditional probability of *A*, *P*(*A*), since *A* is independent of *B*. In the numerical example, we have conditioned on being in the top row:

	<b>A</b>	<b>~A</b>	
<b>B</b>	<b>25</b>	<b>25</b>	
<b>~B</b>	25	25	

Here the probability of *A* is 25/50 = 1/2.

Berkson's paradox arises because the conditional probability of *A* given *B* *within the three-cell subset* equals the conditional probability in the overall population, but the unconditional probability within the subset is inflated relative to the unconditional probability in the overall population, hence, within the subset, the presence of *B* decreases the conditional probability of *A* (back to its overall unconditional probability):

*P*(*A*|*B*,*A*∪*B*) = *P*(*A*|*B*) = *P*(*A*)
*P*(*A*|*A*∪*B*) > *P*(*A*)

## Examples

Berkson's original illustration involves a retrospective study examining a risk factor for a disease in a statistical sample from a hospital in-patient population. Because samples are taken from a hospital in-patient population, rather than from the general public, this can result in a spurious negative association between the disease and the risk factor. For example, if the risk factor is diabetes and the disease is cholecystitis, a hospital patient *without* diabetes is *more* likely to have cholecystitis than a member of the general population, since the patient must have had some non-diabetes (possibly cholecystitis-causing) reason to enter the hospital in the first place. That result will be obtained regardless of whether there is any association between diabetes and cholecystitis in the general population.

An example presented by Jordan Ellenberg: Suppose Alex will only date a man if his niceness plus his handsomeness exceeds some threshold. Then nicer men do not have to be as handsome to qualify for Alex's dating pool. So, *among the men that Alex dates*, Alex may observe that the nicer ones are less handsome on average (and vice versa), even if these traits are uncorrelated in the general population. Note that this does not mean that men in the dating pool compare unfavorably with men in the population. On the contrary, Alex's selection criterion means that Alex has high standards. The average nice man that Alex dates is actually more handsome than the average man in the population (since even among nice men, the ugliest portion of the population is skipped). Berkson's negative correlation is an effect that arises *within* the dating pool: the rude men that Alex dates must have been *even more* handsome to qualify.

As a quantitative example, suppose a collector has 1000 postage stamps, of which 300 are pretty and 100 are rare, with 30 being both pretty and rare. 10% of all his stamps are rare and 10% of his pretty stamps are rare, so prettiness tells nothing about rarity. He puts the 370 stamps which are pretty or rare on display. Just over 27% of the stamps on display are rare (100/370), but still only 10% of the pretty stamps are rare (and 100% of the 70 not-pretty stamps on display are rare). If an observer only considers stamps on display, they will observe a spurious negative relationship between prettiness and rarity as a result of the selection bias (that is, not-prettyness strongly indicates rarity in the display, but not in the total collection).

## See also

- Simpson's paradox

## References

- Berkson, Joseph (June 1946). "Limitations of the Application of Fourfold Table Analysis to Hospital Data" (http://ije.oxfordjournals.org/content/43/2/511.full). *Biometrics Bulletin*. **2** (3): 47–53. doi:10.2307/3002000 (https://doi.org/10.2307%2F3002000). JSTOR 3002000 (https://www.jstor.org/stable/3002000). (The paper is frequently miscited as Berkson, J. (1949) **Biological** Bulletin 2, 47–53.)
- Jordan Ellenberg, "Why are handsome men such jerks? (http://www.slate.com/blogs/how\_not\_to\_be\_wrong/2014/06/03/berkson\_s\_fallacy\_why\_are\_handsome\_men\_such\_jerks.html)"

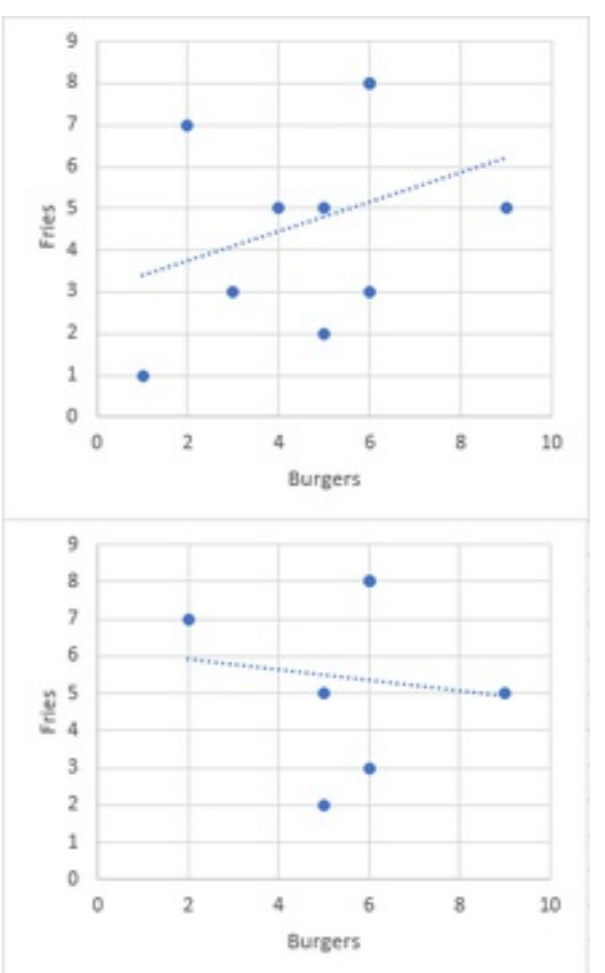
## External links

- Numberphile: Does Hollywood ruin books? (https://www.youtube.com/watch?v=FUD8h9JpEVQ) – An education video on Berkson's paradox in popular culture

Retrieved from "https://en.wikipedia.org/w/index.php?title=Berkson%27s\_paradox&oldid=873685392"

<span></span>	<span></span>
<b>This page was last edited on 14 December 2018, at 13:32 (UTC).</b>	

Text is available under the  Creative Commons Attribution-ShareAlike License; additional terms may apply.
By using this site, you agree to the [Terms of Use](#) and [Privacy Policy](#).
Wikipedia® is a registered trademark of the [Wikimedia Foundation, Inc.](#), a non-profit organization.



An illustration of Berkson's Paradox. The top graph represents the actual distribution, in which a positive correlation between quality of burgers and fries is observed. However, an individual who does not eat at any location where both are bad observes only the distribution on the bottom graph, which appears to show a negative correlation.