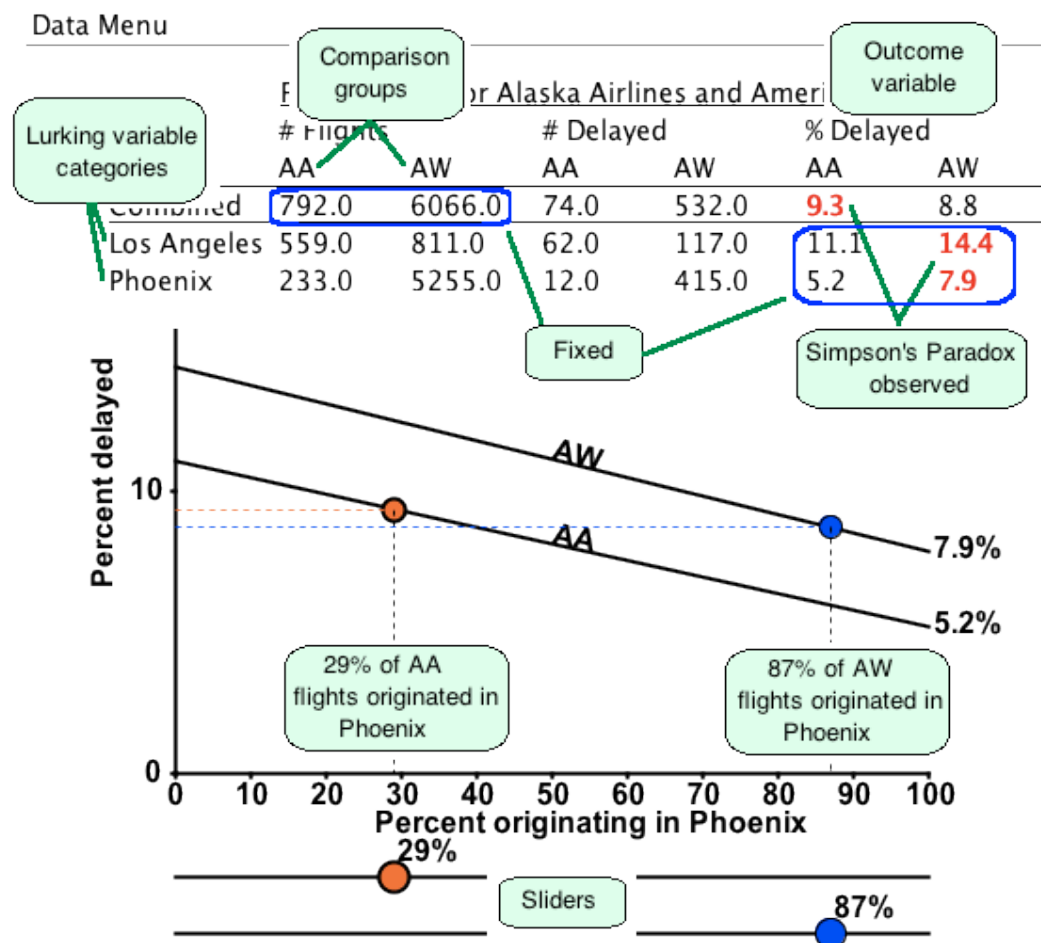Simpson's Paradox

**Overview**: Simpson's Paradox is the name given to the phenomenon in which relationships observed between groups reverse when the groups are divided into subgroups based on a lurking variable.

**Example**: Select the dataset, 'Airlines' from the data menu.

<u>The Table:</u> The first line of the table shows a count of Alaskan Airlines (AA) flights (792), a count of America West (AW) flights (6066), the number of each of these that was delayed (74 for AA and 532 for AW) and finally the percentage delayed for each airline. From this row, we see that Alaskan Airlines had a greater percentage of delayed flights than did America West Airlines (9.3% versus 8.8%).



In the second and third rows of the table, the data for each airline has been divided into subgroups based on the place of origin of the flight (e.g. 559 of the 792 AA flights originated in Los Angeles and 233 originated in Phoenix). For these subgroups, we see

that a greater percentage of the AW flights originating in Los Angeles were delayed (14.4% versus 11.1 % for AA) and a greater percentage of the AW flights originating in Phoenix were delayed (7.9% versus 5.2% for AA).

Whereas AA had a greater number of delayed flights when the data were combined, the relationship is reversed when we divide the data into subgroups based on the lurking variable 'Place of origination'. Why this occurs can be seen in the plot.

The Plot: For each of the comparison groups (AA and AW) the plot shows the percentage of delayed flights as a function of the percentage of flights originating from Phoenix.  Colored dots on the lines indicate the actual percentages of flights that originated in Phoenix for each of the comparison groups.  We see that 87% of AW flights originated in Phoenix whereas only 29% of AA flights began there.  Since Phoenix flights were less likely to be delayed than Los Angeles flights (5.2% versus 11.1% for AA and 7.9% versus 14.4% for AW) and the vast majority of AW flights began there, AW's overall percentage of delayed flights is lower than that of AA for which most of the flights originated in Los Angeles.

The Sliders: The sliders allow the user to adjust the percentage of flights originating in Phoenix for each of the two airlines and to see how this affects the observed relationships.  As a slider is adjusted, a circle on the corresponding line (circle color the same as slider dot color) moves. Dashed lines from the circles to the axes highlight the relationship between the variable values.

As the percentage of flights originating in Phoenix is changed, the data in the table is updated to reflect this. The 'combined' counts for each category are fixed as are the percentages of delayed flights for each airline from each origination point.
When the Simpson's Paradox is observed, i.e. the airline with the greater percentage delayed for the 'combined' data is different from the airline with the greater percentage delayed for each of the origination points, the percentages in the table are highlighted in red, otherwise they are green.