In the previous lecture-unit we talked about working with single sets of values; we now move to larger sets of data, where we have more than one piece of information about each of the subjects. What, if anything, changes when you measure two things at the same time? First, you now have to do everything that you did before twice; once for each measure. This includes summarizing each set of values in terms of the center, the spread, and a name for the shape. Nothing new here, just twice as much work. This part is easy. Second, you can and should also try to make some kind of summary statement about the relationship between the measures. This is a bit more complicated.

The main problem that we face when we try to express the relationship between two different measures is much like the issue of "apples and oranges." If the two measures have different units – such as when one is a measure of response time (in milliseconds) and the other is a measure of depression (in Beck Depression Inventory units) – then there is no way to directly compare the two. As I said: apples to oranges.

This is where correlation coefficients come in. One of the beautiful things about correlations is that they can be calculated for any two variables, regardless of what the variables are measures of, or even what classes the variables are. (Yes, calculating correlations involving qualitative variables is a lot more complicated [and won't be discussed in this course], but it can be done.) The reason that correlations can be used with apples and oranges is because they operate in relative terms, instead of absolute terms.

For example, assume that you are interested in the relationship between response time and depression. You measure each of your subjects on both of these things – i.e., you collect bivariate data (bi-variate ... bi = two; variate = variables). But instead of thinking about each subject's response time and depression score in absolute terms, you think of each subject as being relatively fast or slow and relatively depressed or happy. By switching to relative scores, you remove the units. Now you just have two unitless values for each subject.

Technical aside: If you're wondering, the way that this is done is by taking each subject's score, subtracting off the mean score (across all subjects) and then dividing by the standard deviation. Note how this uses information that you would have already calculated in the first step, where you treat each variable on its own. The name for these calculated values is the "standardized score." They go roughly from minus two (for the extreme lowest original scores) to positive two (for the extreme highest). And this is true regardless of what the original values were. Response times, for example, are usually between 300 and 1500 milliseconds, but standardized scores for response time are almost always between -2 and +2. Likewise, typical BDI scores go from near zero to about 25 (for severe depression), but standardized scores are, again, almost always between -2 and +2.

Finally, to get a correlation coefficient, the pairs of standardized scores (across subjects) are combined to provide a single number that tells you whether there is a linear relationship between the two variables. (This is done by calculating the mean product of the standardized scores: i.e., within each subject, multiple the two standardized scores, then add them all up, then divide by N.) Note, however, a very critical word in that sentence: the (plain) correlation coefficient does not measure all possible kinds of relationship between two variables; it only measures the <u>linear</u> relationship. In words: it tells you whether, in general, when one variable goes up, the other goes up, as well (which is a positive correlation); or whether when one variable goes up, the other goes down (which is a negative correlation). If, for example, as one variable

goes up, the other goes up for a while and then starts to come back down again, then this will not show up in the (plain) correlation coefficient, because it isn't a linear (straight-line) relationship. You'd have to use a more sophisticated measure to pick up on a curved relationship like that.

With the above limitation in mind, let's return to the good things about correlations. Not only can you calculate a correlation between any two variables, but any two correlations can be directly compared to each other. This is true for two reasons. First, correlations, themselves, have no units, so it won't be apples and oranges; it will just be two unitless numbers, which, of course, can be directly compared. Second, all correlations, regardless of the original range of values for the two input variables, are always between -1 and +1, with -1 being a perfect, negative, linear relationship and +1 being a perfect, positive, linear relationship at all.)

Cool trick: since both -1 and +1 correlations represent perfect, linear relationships, many people prefer to talk about the square of the correlation, instead. This converts a -1 into a +1, while leaving +1 as +1, such that both kinds of perfect relationship are now coded as +1. We'll talk some more about this in lecture.

One last point on correlations: they are bi-directional. If the correlation between BDI score and response time is +.25 (which is weak but positive relationship, implying that depressed people are a little slower than non-depressed people), then the correlation between response time and BDI score is also +.25 (implying that slow people are little less happy than fast people). This is something that will come back to haunt us later, so try to keep it in mind.