## Some Preliminary Thoughts on Statistics and Background Information on SPSS (Part 2)
by H.P.L. Molloy and T Newfields

*This article highlights more basic statistical concepts for SPSS users. After pointing out how counting, ordering, and measuring differ, the procedure for making histograms with SPSS is described. This segment concludes by noting how parametric and non-parametric distributions contrast.*

*The previous article, online at http://jalt.org/test/mn_1.htm, mentioned five ways to increase the statistical power of a study while mentioning some basic statistical concepts. The final article, which will appear in the next issue, clarifies four more basic statistical concepts then consider various types of curve distributions. The article will conclude by explaining what statistical outliers are and how to interpret them, and then finally suggest a few self-study questions.*

### 3. Counting, ordering, and measuring

As we mentioned in the last article, statistical procedures are concerned with numbers. No matter what we are studying, we have got to convert the information into to numbers. There are three ways to do this: counting, ordering, and measuring. The numbers we work with are called variables.

We can count how many things belong to certain categories: men vs. women; high school students planning to continue to university, those planning to work after graduation, and those undecided; extroverts vs. introverts from the Myers-Briggs Type Indicator test. If we count 45 extroverts and 38 introverts, we are using nominal variables. Nominal variables are category variables: something is either in a given category or is not in that category.

We can also rank or order things. We can ask students to rank, in terms of politeness, five different ways of making requests in English. We can also ask language users to rate the acceptability of various correct, incorrect, and questionable sentences or ask people to arrange five size adjectives in order of size. All of these procedures are creating ordinal variables. Ordinal variables are rank variables: they show some ordering in what we are studying: we may find that "Would you mind lending me a bit of money," is considered more polite than "Can you give me some money," which in turn is considered more polite than "Give me some money." It is important to note that we have no idea how much more polite one thing is than another. We only know the rank.

We can also measure things. We can give students the TOEFL® test and measure their English proficiency; we can time students and measure how quickly they can find the correct answers to a reading comprehension test; we can check how many words they recognize in a list of 500 commonly used words. With measurement, we know, for example, how much more able one student is than the next. We know how many more minutes one student takes to find reading comprehension answers than the next. We know how many more words one knows than the next. When we measure, we use some sort of scale that shows us how far apart two numbers are. The scales, or rulers, in these examples are the scores of all previous TOEFL test takers, time, and an ad hoc 0-500 vocabulary word list.

There are two kinds of measurement. When it is not possible to get a zero score, the measurement is interval measurement. Interval measurement lets us know how far apart things are but do not tell us how much of something someone has; only how much more or less he or she has than someone else.

When it is possible to get a zero score, the measurement is ratio measurement. Not only do we know how far apart things are, but we also know how much of something someone has. With ratio measurement, and with no other kind, it is possible to say that person A, who scored 20, has twice as much of something than person B, who scored 10. Height, weight, and income are examples of ratio measurements. It is possible to conceive of ratio measurement in our field, but it is not often used. Consider the problem of linguistic ability, for example: though we might safely say that someone's ability in Basque is zero, when can we say someone has twice as much ability as that person?

## 4. Distributions

Statistical distributions refer to the shapes that graphs of numbers being studied take when placed in a chart. The lowest scores are places on the left side of the *x*-axis (which is horizontal), and the highest scores on the right side of this axis. On the *y*-axis (which is vertical) the numbers used in a study appear.

To make a distribution of your data, count the number of scores of one level (15, for example). If you find six of them, put a mark at the junction of 15 on the *x*-axis and 6 on the *y*-axis. Do that until you've counted all the numbers you have. This is called a histogram.

Distributions are important for two reasons. First, they can tell easily some things about your data. If a histogram distribution looks weird, there are reasons for that. If it looks flat or too narrow or too spread out, you may have problems. Second, and more importantly, many statistical procedures are based on specific distributions. Your distribution has to look like one of these for you to be able to use that particular statistical test.

As we mentioned, statistics shows us the chances that the numbers we are studying are random. This is often done, roughly, by comparing the properties of your numbers with those of several well-studied distributions. Some of the best known of these distributions are the chi-square distribution, the t-distribution, the F distribution, and the *z*, or normal, distribution.

### *4.1 Histograms in SPSS*

**Menu**
1.   Choose the "Analyze" menu.
2.   Choose the "Descriptive statistics" menu.
3.   Choose the "Frequencies" submenu.
4.   On the left side of the dialogue box, highlight the variables you are interested in.
5.   Click on the arrow between the two white boxes.
6.   Click on "Charts."
7.   Click on "Histograms" and "With normal curve."
8.   Click on "Continue."
9.   Click on "OK" or (optionally), click on "Statistics."
10.  (Optional) click on "Std. deviation," "Variance," "Minimum," "Maximum," "S.E. mean," "Mean," "Median," "Mode," "Skewness," and "Kurtosis."
11.  Click on "Continue."
12.  Click on "OK."

**Syntax**

To do the same thing using syntax, type this:

```
FREQUENCIES
VARIABLES=[variable 1] [variable 2] [variable 3] [variable n]  /FORMAT=NOTABLE
/STATISTICS=STDDEV VARIANCE MINIMUM MAXIMUM SEMEAN MEAN MEDIAN MODE
SKEWNESS SESKEW KURTOSIS SEKURT
/HISTOGRAM NORMAL
/ORDER ANALYSIS
```

Here, replace "[variable 1]" (and so on) with the name(s) of the variables you're interested in. The names you should use are those are the tops of the columns in SPSS.

### 5. Parametric vs. nonparametric statistics

Parametric statistics are those that work by comparing your numbers with one of the standard distributions. (Note that the chi-square distribution is not a parametric distribution.) This is why parametric statistics have so many requirements about sample or *N* sizes, skew, kurtosis, and the like: they only work properly if you are working with numbers that have distributions close enough to one of the standard, well-studied distributions. This is why parametric statistics are considered more powerful and more desirable than nonparametric statistics: the assumption that the data more or less match standardized distributions enables us to compare across studies.

Many nonparametric statistics do not depend on your having a distribution of a certain shape. Instead, they depend on permutations, possible combinations of things. Figuring possible permutations very quickly becomes computationally tedious and, surprisingly quickly, comes to involve numbers that are presently beyond our ability to compute, even with supercomputers. Most nonparametric statistics end up using *z*, or normal, distributions with larger numbers, which is why nonparametric statistics are usually associated with small *N* or sample sizes.

Does the difference between parametric and nonparametric statistics matter? No, not really. Each gives us the same thing: the odds that our data are not random.

### 6. Summary

This article has briefly mentioned three ways of arranging data: counting, ordering, and measuring. It has also pointed out some of the limitations and strengths of parametric statistics. The concluding article, to appear this autumn, will mention how to calculate some descriptive statistics using SPSS and then consider various types of curve distributions.

**HTML**: http://jalt.org/test/mn_2.htm         /         **PDF**: http://jalt.org/test/PDF/Molloy-Newfields2.pdf