



Tips and Tools for Sociologists

Frances M. Shaver and Bill Reimer

Table 1: Summarizing Univariate Distributions - Bill Reimer and Frances Shaver

September 9, 2002

Summary measures appropriate to all levels of measurement:

proportion: $P = \frac{f}{N}$ percent: $\% = \frac{f}{N} \times 100$

ratio: $\frac{f_1}{f_2}$

rate: $\frac{f \text{ of actual cases}}{f \text{ potential cases}}$

% rate of change: $\frac{f_{time2} - f_{time1}}{f_{time1}} \times 100$

frequency distributions
form of the distribution

measures of central tendency: provide a summary of the location of a set of scores

measures of dispersion (variation): provide a summary of the spread in a set of scores

Measures of central tendency and measures of dispersion are appropriate to particular levels of measurement:

Level of Measurement Some Appropriate Statistics

	central tendency	dispersion
nominal:	<p>mode</p> <ul style="list-style-type: none"> the category with the highest frequency 	<p>variation ratio</p> $v = 1 - \frac{f_{modal}}{N}$ <ul style="list-style-type: none"> varies between 0 & 1 the smaller the variation the more accurate the mode is as a summary measure
ordinal:	<p>median</p> <ul style="list-style-type: none"> reflects the tendency of the <u>whole distribution</u> to be either high or low represents the middle of a distribution 	<p>range</p> <p>decile range</p> $d = d_q - d_1$ <p>interquartile range</p> $Q = Q_3 - Q_1$ <ul style="list-style-type: none"> may be used as indices of variation only for comparisons on the <u>same</u> set of ranks
interval/ratio:	<p>mean</p> $\bar{X} = \frac{\sum_{i=1}^N X_i}{N}$ <ul style="list-style-type: none"> represents the balance point in a distribution 	<p>standard deviation</p> $s = \sqrt{\frac{\sum_{i=1}^N (X_i - \bar{X})^2}{N}}$ <ul style="list-style-type: none"> as variation increases, s increases

Table 2: Some Appropriate Statistics for Describing Association Between Two Variables

Level of measurement of the dependent variable	Level of measurement of the independent variable		
	nominal	ordinal	interval/ratio
nominal	<ul style="list-style-type: none"> • Guttman's lambda (λ): Freeman 1965:71 • Phi (ϕ): A&Z 1958:138 • Pearson's C: A&Z 1958:141 • Yule's Q: Johnson 1988:138 	<ul style="list-style-type: none"> • Coefficient of differentiation/theta (θ): Freeman 1965:108 • Guttman's lambda (λ): Freeman 1965:71 [but some data is lost] 	<ul style="list-style-type: none"> • correlation ratio/eta (η) • Guttman's lambda (λ): Freeman 1965:71 • Gamma (γ) (no ties): A&Z 1958:152; Freeman 1965:79; L&M 1980:235
ordinal	<ul style="list-style-type: none"> • Coefficient of differentiation/theta (θ): Freeman 1965:108 	<ul style="list-style-type: none"> • Gamma (γ) (no ties): A&Z 1958:152; Freeman 1965:79; L&M 1980:235 • Somers' d_{yx} (ties on y): L&M 1980:235 • Tau-b (τ_b) (ties on x and y): L&M 1980:236 • Spearman's rho (r_s): Siegal 1956:202 	<ul style="list-style-type: none"> • Coefficient of Multiserial Correlation (M): Freeman 1965:131
interval/ratio	<ul style="list-style-type: none"> • Correlation ratio/eta (η): Freeman 1965:120; A&Z 1958:155 	<ul style="list-style-type: none"> • Coefficient of Multiserial Correlation (M): Freeman 1965:131 	<ul style="list-style-type: none"> • Pearson's r: A&Z 1958:113; Freeman 1965:89

Guttman's lambda (λ) like r is a ratio which provides an index of the degree in reduction in error (PRE) (i.e. (reduction in error)/(original error)). The only assumption underlying the use of this statistic is the level of measurement. λ , γ , and r are all symmetrical, that is, values taken by each variable are guessed on the basis of knowledge of values taken by the other. λ_a is asymmetrical, i.e. values taken by one value are guessed from values taken from the other, but the reverse is not true. λ and λ_a vary between 0 and +1.

Phi (ϕ) and **Pearson's C** are measures of association between 2 nominal scales that are based on the notion of statistical independence (i.e. observed versus expected frequencies if no association is expected). ϕ is appropriate only if the data has been dichotomized (i.e. a 2 x 2 table). ϕ^2 and C^2 are measures of how much difference there is in the table as a whole, between observed and expected frequencies. In other words, a proportional reduction in error (PRE) interpretation can be used with ϕ^2 and C^2 . ϕ and C are measures of the extent of association or relation between 2

attributes. ϕ varies between 0 and 1 but C does not achieve a value of 1 even for perfect association. It does, however equal 0 if there is no association.

Yule's Q (Q) is a measure of association for ordinal variables (or nominal variables if the sign is dropped). It is only appropriate for 2 x 2 tables. It is identical to gamma for a 2 x 2 table. Varies from -1 to +1 for ordinal variables and from 0 to +1 for nominal variables.

Coefficient of differentiation. Theta (θ) is a measure of the difference between the proportion of comparisons in which members of one class predominate and the proportion in which members of another class predominate. Like λ_a it is an asymmetric measure of association. It varies between 0 and 1.

Coefficient of Multiserial Correlation (M) is an adaptation of Pearson's r to describe the degree of association between one ordinal and one interval scale. We must assume that there is a linear relationship between X and Y and that a normal distribution constitutes a reasonable image of what the ordinal variable might look like if we were able to measure it with greater precision. It varies between -1 and +1.

Correlation ratio or **eta** (η) like r and λ , η is a proportional reduction in error (PRE) measure. The magnitude of η^2 like r^2 , expresses the proportion of shared variance between X and Y . To put it another way, η^2 is equal to the proportion of variance in Y (the measured variable) which is associated with subclasses in X (the nominal scale). η varies between 0 and 1.

Gamma (γ or G) is a coefficient of association between two sets of ordered observations based on their mutual predictability in terms of the relative number of agreements and inversions in the order of rankings. γ may be computed for data with ties and without ties, however, T_x and T_y are not treated as being lack of evidence for association. γ varies between -1 and +1.

Somers' d_{yx} is like Gamma except that ties on y are treated as evidence of lack of association. It varies between -1 and +1.

Tau-b (τ_b) is like Gamma except that ties on x and y are treated as evidence of lack of association. It varies between -1 and +1 and is for square tables only.

Spearman's rho (ρ_s or r_s) is a symmetric measure of association suitable for ordinal level variables. It compares two sets of ranks. Identical ranks on two variables is treated as a perfect association. It is sensitive to large differences in ranks.

Pearson's r is a ratio which provides a standard index of the degree of reduction in error (i.e. (reduction in error)/(original error)). r^2 tells us the proportion of the variance in the Y observations that is accounted for by variation in X . To determine association using r you must assume that there is a linear relationship between X and Y ; that X and Y are both measured at the interval level; that the conditional distributions of Y for each X are normal and their variances equal. r varies between -1 and +1.

(All measures increase as the association increases. The sign tells you whether the association is positive or negative.)

Table 3: Some Appropriate Statistical Tests for Evaluating the Null Hypothesis

Level of measurement of the dependent variable	Level of measurement of the independent variable		
	nominal	ordinal	interval/ratio
nominal	<ul style="list-style-type: none"> the χ^2 test⁺: Siegal:104; A&Z:256; Freeman:215 Fisher's Exact test⁺: A&Z:264; Siegal:96 	<ul style="list-style-type: none"> the significance of τ^{+2}: A&Z:176; Siegal:220 the significance of γ^{+2}: Freeman:162 the significance of r_s^{+2}: Siegal:210 	<ul style="list-style-type: none"> Fisher's Analysis of Variance (F) T-test (t) [note: $t^2 = F$]
ordinal	<ul style="list-style-type: none"> the median test⁺¹: Siegal:111 the Mann Whitney U⁺¹: Siegal:116 Kolmogrov-Smirnov test⁺¹: A&S:269; Siegal:127 	<ul style="list-style-type: none"> the significance of τ^{+2}: A&Z:176; Siegal:220 the significance of γ^{+2}: Freeman:162 the significance of r_s^{+2}: Siegal:210 	<ul style="list-style-type: none"> test of significance of M
interval/ratio	<ul style="list-style-type: none"> the Z-score test^{*1} the t-test^{*1}: Levin:131 the F-test^{*1}: Levin:164 the randomization test⁺¹: A&Z:272; Siegal:152 	<ul style="list-style-type: none"> the significance of M^{*2}: Freeman:211 	<ul style="list-style-type: none"> the significance of r^{*2}: A&Z:277; Freeman:178

* a parametric test 1. Restricted to a nominal scale containing only 2 subclasses.
 + a non-parametric test 2. See Table 2 for a description of this measure of association.

- All statistical tests are based on the assumption of random sampling (a procedure which provides that every observation has an equal chance of appearing in the sample).
- All statistical tests evaluate the significance of the null Hypothesis (H_0).
- All statistical tests have measurement requirements. They are outlined in Table 3.
- Data measured by either nominal or ordinal scales should be analyzed by the non-parametric methods. Data measured in interval or ratio scales may be analyzed by parametric methods if the assumptions of the parametric statistical model are appropriate.

A **Parametric Statistical Test** is a test whose model specifies certain conditions about the parameters of the population from which the research sample was drawn. Elements of the usual parametric statistical model:

- the observations must be independent
- the observations must be drawn from normally distributed populations
- these populations must have the same variance (or, in special cases, they must have a known ratio of variance)

In the case of analysis of variance (the F test), another condition is added:

- the means of these normal and homoscedastic populations must be linear combinations of effects due to columns and/or rows. That is, the effects must be additive.

These conditions are not ordinarily tested: they are assumed to hold. The meaningfulness of the results of a parametric test depends on the validity of these assumptions. If you cannot meet these assumptions or if the scores under analysis do not meet the requirements of at least an interval scale then non-parametric statistical tests must be used.

A **Non-parametric Statistical Test** is a test whose model does not specify conditions about the parameters of the population from which the sample was drawn. Certain assumptions are associated with most non-parametric statistical tests: i.e. that the observations are independent and that the variable under study has underlying continuity, but these assumptions are fewer and much weaker than those associated with parametric tests.

If all the assumptions of the statistical model are met and if the variables under analysis achieve an interval level of measurement, a parametric statistical test is the most powerful (i.e. the most likely of all tests to reject H_0 when H_0 is false). However, the power of any non-parametric statistical test may be increased by simply increasing the size of N.

References

- Anderson, Theodore R. and Morris Zelditch, Jr.
1958 A Basic Course in Statistics With Sociological Applications New York: Holt, Rinehart and Winston.
- Freeman, Linton C.
1965 Elementary Applied Statistics: for students in the behavioral sciences New York: John Wiley & Sons.
- Johnson, Allan G.
1988 Statistics New York: Harcourt Brace Jovanovich.
- Loether, Herman J. and Donald G. McTavish
1988 Descriptive and Inferential Statistics: an introduction (3rd edition) Boston: Allan and Bacon.
- Siegel, Sidney
1956 Nonparametric Statistics for the Behavioral Sciences New York: McGraw-Hill.