



Research Work and Statistical Tools

(A special reference of data based research)

Dr. R. B. Singh

Associate Professor - Department of Statistics

D. N. College, Meerut

rbs.aparna@gmail.com

ABSTRACT: Statistics is a wide subject useful in almost all disciplines especially in Research studies. Each and every researcher should have some knowledge in Statistics and must use statistical tools in his or her research, one should know about the importance of statistical tools and how to use them in their research or survey. The quality assurance of the work must be dealt with: the statistical operations necessary to control and verify the analytical procedures as well as the resulting data making mistakes in analytical work is unavoidable. This is the reason why a multitude of different statistical tools is, required some of them simple, some complicated, and often very specific for certain purposes. In analytical work, the most important common operation is the comparison of data, or sets of data, to quantify accuracy (bias) and precision. Fortunately, with a few simple convenient statistical tools most of the information needed in regular laboratory work can be obtained: the "t-test", the "F-test", and regression analysis. Clearly, statistics are a tool, not an aim. Simple inspection of data, without statistical treatment, by an experienced and dedicated analyst may be just as useful as statistical figures on the desk of the disinterested. The value of statistics lies with organizing and simplifying data, to permit some objective estimate showing that an analysis is under control or that a change has occurred. Equally important is that the results of these statistical procedures are recorded and can be retrieved. The key is to sift through the overwhelming volume of data available to organizations and businesses and correctly interpret its implications. But to sort through all this information, you need the right statistical data analysis tools. Hence in this paper, I have made an attempt to give a brief report or study on Statistical tools used in research studies.

KEYWORDS: Quantify accuracy, Analytical procedures, Quality assurance, Data analysis tools.

1. INTRODUCTION

The subject Statistics is widely used in almost all fields like Biology, Botany, Economics, Sociology, Commerce, Medicine, Education, Physics, Chemistry, Bio-Technology, Psychology, Zoology etc.. While doing research in the above fields, the researchers should have some awareness in using the statistical tools which helps them in drawing rigorous and good conclusions. The most well known Statistical tools are the mean, the arithmetical average of numbers, median and mode, Range, dispersion, standard



deviation, inter quartile range, coefficient of variation, etc. There are also software packages like SAS and SPSS which are useful in interpreting the results for large sample size.

The Statistical analysis depends on the objective of the study. The objective of a survey is to obtain information about the situation of the population study. The first Statistical task is therefore is to do a descriptive analysis of variables. In this analysis it is necessary to present results obtained for each type of variable. For qualitative and dichotomous variables, results must be presented as frequencies and percentages. For quantitative variables, the presentation is as means and deviations. After this analysis, you can access the association between variables and predictive analysis based on multiple regression models. You can also use software packages like SPSS, EPIInfo, STATA, Minitab, Open Epi, Graph pad and many others depending on your usage and familiarity with the software. You should also start looking at the distributions of age, gender, race and any measures of socio-economic status that you have (income, education level, access to medical care). These distributions will help to inform your analysis in terms of possible age- adjustment, weighting and another analytical tools available to address issues of bias and non representative samples.

Survey analysis is one of the most commonly used research methods, scholars, market researchers and organization of all sizes use surveys to measure public opinion. Researchers use a wide range of statistical methods to analyze survey data. They do this using statistical software packages that are designed for research professionals. Popular programs include SAS, SPSS and STATA. However, many forms of survey data analysis can be done with a spread sheet program such as EXCEL, which is part of Microsoft's popular office package. EXCEL and other spreadsheet programs are user-friendly and excellent for entering, coding and storing survey data.

2. METHODS

2.1. Context Chart:

This display method is used to understand the context of the data found. When building thematic frames, the data included in each frame must be connected by context to be useful. Once the context chart is complete, partial analysis (partial analysis is often used to validate variables or themes) or interim analysis (interim analysis is finding an early direction or theme in the data) can be performed on the data findings. By using the context chart the researcher shows the interrelationship of the data while keeping the research questions in mind.

2.2. Checklist Matrix:

This display method will determine whether the data is visible or useful as a variable (a variable is an object used for comparison such as “apples” and “oranges”) in the analysis of the qualitative data. The components of the data are broken up by thematic points and placed in labeled columns, rows and point guided rubrics (eg: strong, sketchy, adequate) within the



matrix). The thematic points are then examined for usefulness as a variable according to the numeric strength of the point- guided rubric.

2.3. Pattern-Coded Analysis Table:

This table is created with rows labeled with themes and columns labeled by coded patterns. Pattern coding is a way to add further distinction to a variable-oriented analysis of the data. Often referred to as a cross-case analysis table, the researcher can, at a glance at the rows, render a preliminary analysis of the data collected just by noting which cell the pattern-coded data fills under certain thematic rows.

2.4. Decision-tree modeling:

This method is a chart structured from one central directive. It often resembles a tree with branches. For ex- the central directive may be whether to buy a contract. From the directive two decision boxes are created. Pro & Con. After taking a survey, the researcher creates a branch from the Pro and Con boxes, allowing for a third branch for the undecided. Because the data was collected subjectively/qualitatively, the researcher will have coded the responses earlier by context to determine by pattern if they fail under pro or con. In this display the researcher will write those patterned responses in boxes resembling twigs growing from the appropriate branch to analyze the findings.

Besides this there are some most popular basic methods of analyzing survey data which include frequency distributions and descriptive statistics. Frequency distribution tell you how many people answered a survey question a certain way. Descriptive statistics help describe a set of data through descriptive measures, such as means and standard deviations. Beyond basic techniques, there are more complex analytical methods used in survey research. Researchers may use factor analysis to examine the correlations among different survey questions with the intent of creating index measures for deeper analysis. There are regression techniques to examine how particular variables of interest affect a particular outcome.

2.5. Parametric and non parametric tests:

Choosing the right test to compare measurements is a bit tricky, as you must choose between two families of tests parametric and non-parametric. Many statistical tests are based upon the assumption that the data are sampled from a Gaussian distribution. These tests are referred to as parametric tests. Commonly used parametric tests are listed in the first column of the table and include t test and analysis of variance. Tests that do not make assumption about the probability distribution are referred to as Non parametric tests. All commonly used non parametric tests rank the outcome variable from low to high and then analyze the ranks. These tests are listed in the second column of the table and include the Gottschalk, L. A. Wilcoxon, Mann-Whitney test and Kruskal-Wall's tests which are called distribution free tests.



2.6. Mean

The arithmetic mean, more commonly known as “The average” is the sum of a list of numbers divided by the number of items on the list. The mean is useful in determining the overall trend of a data set or providing a rapid snapshot of your data. Another advantage of the mean is that it’s very easy and quick to calculate.

2.7. Standard Deviation

The standard deviation, often represented with the Greek letter sigma, is the measure of a spread of data around the mean. A high standard deviation signifies that data is spread more widely from the mean, where a low standard deviation signals that more data align with the mean. In a portfolio of data analysis methods, the standard deviation is useful for quickly determining dispersion of data points.

2.9. Regression

Regression models the relationships between dependent and explanatory variables, which are usually charted on a scatterplot. The regression line also designates whether those relationships are strong or weak. Regression is commonly taught in high school or college statistics courses with applications for science or business in determining trends over time.

2.10. Sample Size Determination

When measuring a large data set or population, like a workforce, you don’t always need to collect information from every member of that population – a sample does the job just as well. The trick is to determine the right size for a sample to be accurate. Using proportion and standard deviation methods, you are able to accurately determine the right sample size you need to make your data collection statistically significant.

2.11. Hypothesis Testing

Also commonly called *t* testing, hypothesis testing assesses if a certain premise is actually true for your data set or population. In data analysis and statistics, you consider the result of a hypothesis test *statistically significant* if the results couldn’t have happened by random chance. Hypothesis tests are used in everything from science and research to business and economic.

3. DATA ANALYSIS

Is the process of systematically applying statistical and/or logical techniques to describe and illustrate, condense and recap, and evaluate data? According to Shamoo and Resnik (2003) various analytic procedures “provide a way of drawing inductive inferences from data and distinguishing the signal (the phenomenon of interest) from the noise (statistical fluctuations) present in the data”.

While data analysis in qualitative research can include statistical procedures, many times analysis becomes an ongoing iterative process where data is continuously collected and



analyzed almost simultaneously. Indeed, researchers generally analyze for patterns in observations through the entire data collection phase (Savenye, Robinson, 2004). The form of the analysis is determined by the specific qualitative approach taken (field study, ethnography content analysis, oral history, biography, **unobtrusive** research) and the form of the data (field notes, documents, audiotape, videotape).

An essential component of ensuring data integrity is the accurate and appropriate analysis of research findings. Improper statistical analyses distort scientific findings, mislead casual readers (Shepard, 2002), and may negatively influence the public perception of research. Integrity issues are just as relevant to analysis of non-statistical data as well.

In deciding which test is appropriate to use, it is important to consider the type of variables that you have (i.e., whether your variables are categorical, ordinal or interval and whether they are normally distributed).

3.1. About the hsb data file

A data file called **hsb2**, high school and beyond, this data file contains observations from a sample of high school students with demographic information about the students, such as their gender socio-economic status and ethnic background It also contains a number of scores on standardized tests, including tests of reading, writing , mathematics and social studies.

3.2. One sample t-test

A one sample t-test allows us to test whether a sample mean (of a normally distributed interval variable) significantly differs from a hypothesized value. The mean of the variable for this particular sample of students which is statistically significantly different from the test value . We would conclude that this group of students has a significantly higher mean on the writing test than the given.

3.3. One sample median test

A one sample median test allows us to test whether a sample median differs significantly from a hypothesized value.

3.4. Binomial test

A one sample binomial test allows us to test whether the proportion of successes on a two-level categorical dependent variable significantly differs from a hypothesized value.

3.5. Chi-square goodness of fit

A chi-square goodness of fit test allows us to test whether the observed proportions for a categorical variable differ from hypothesized proportions.

3.6. Wilcoxon-Mann-Whitney test

The Wilcoxon-Mann-Whitney test is a non-parametric analog to the independent samples t-test and can be used when you do not assume that the dependent variable is a normally distributed interval variable (you only assume that the variable is at least ordinal).



3.7. Chi-square test

A chi-square test is used when you want to see if there is a relationship between two categorical variables.

3.8. Fisher's exact test

The Fisher's exact test is used when you want to conduct a chi-square test, but one or more of your cells has an expected frequency of five or less.

3.9. One-way ANOVA

A one-way analysis of variance (ANOVA) is used when you have a categorical independent variable (with two or more categories) and a normally distributed interval dependent variable and you wish to test for differences in the means of the dependent variable broken down by the levels of the independent variable.

3.10. Kruskal Wallis test

The Kruskal Wallis test is used when you have one independent variable with two or more levels and an ordinal dependent variable. In other words, it is the non-parametric version of ANOVA and a generalized form of the Mann-Whitney test method since it permits 2 or more groups.

3.11. Paired t-test

A paired (samples) t-test is used when you have two related observations (i.e. two observations per subject) and you want to see if the means on these two normally distributed interval variables differ from one another.

3.12. Wilcoxon signed rank sum test

The Wilcoxon signed rank sum test is the non-parametric version of a paired samples t-test. You use the Wilcoxon signed rank sum test when you do not wish to assume that the difference between the two variables is interval and normally distributed (but you do assume the difference is ordinal).

3.13. McNemar test

You would perform McNemar's test if you were interested in the marginal frequencies of two binary outcomes. These binary outcomes may be the same outcome variable on matched pairs (like a case-control study) or two outcome variables from a single group.

3.14. One-way repeated measures ANOVA

You would perform a one-way repeated measures analysis of variance if you had one categorical independent variable and a normally distributed interval dependent variable that was repeated at least twice for each subject. This is the equivalent of the paired samples t-test, but allows for two or more levels of the categorical variable. This tests whether the mean of the dependent variable differs by the categorical variable.



3.15. Repeated measures logistic regression

If you have a binary outcome measured repeatedly for each subject and you wish to run a logistic regression that accounts for the effect of these multiple measures from each subjects, you can perform a repeated measures logistic regression.

3.16. Factorial ANOVA

A factorial ANOVA has two or more categorical independent variables (either with or without the interactions) and a single normally distributed interval dependent variable.

3.17. Friedman test

You perform a Friedman test when you have one within-subjects independent variable with two or more levels and a dependent variable that is not interval and normally distributed (but at least ordinal. The null hypothesis in this test is that the distribution of the ranks of each type of score (i.e., reading, writing and math) are the same.

3.18. Ordered logistic regression

Ordered logistic regression is used when the dependent variable is ordered, but not continuous. We do not generally recommend categorizing a continuous variable in this way; we are simply creating a variable to use for this example. The results indicate that the overall model is statistically significant ($p < .0000$) Resnik, D. (2000) as are each of the predictor variables ($p < .000$). There are two cut points for this model because there are three levels of the outcome variable.

One of the assumptions underlying ordinal logistic (and ordinal probity) regression is that the relationship between each pair of outcome groups is the same. In other words, ordinal logistic regression assumes that the coefficients that describe the relationship between, say, the lowest versus all higher categories of the response variable are the same as those that describe the relationship between the next lowest category and all higher categories, etc. This is called the proportional odds assumption or the parallel regression assumption. Because the relationship between all pairs of groups is the same, there is only one set of coefficients (only one model). If this was not the case, we would need different models (such as a generalized ordered logit model) to describe the relationship between each pair of outcome groups. To test this assumption, we can use either the **o model** command (**find it o model**)

3.19. Factorial logistic regression

A factorial logistic regression is used when you have two or more categorical independent variables but a dichotomous dependent variable.

3.20. Correlation

A correlation is useful when you want to see the linear relationship between two (or more) normally distributed interval variables. Although it is assumed that the variables are interval and normally distributed, we can include dummy variables when performing correlations.



3.21. Simple linear regression

Simple linear regression allows us to look at the linear relationship between one normally distributed interval predictor and one normally distributed interval outcome variable.

3.22. Non-parametric correlation

A Spearman correlation is used when one or both of the variables are not assumed to be normally distributed and interval (but are assumed to be ordinal).

3.23. Simple logistic regression

Logistic regression assumes that the outcome variable is binary (i.e., coded as 0 and 1). As with OLS regression, the predictor variables must be either dichotomous or continuous.

3.25. Multiple regression

Multiple regressions is very similar to simple regression, except that in multiple regression you have more than one predictor variable in the equation.

3.26. Analysis of covariance

Analysis of covariance is like ANOVA, except in addition to the categorical predictors you also have continuous predictors as well.

3.27. Multiple logistic regression

Multiple logistic regressions are like simple logistic regression, except that there are two or more predictors. The predictors can be interval variables or dummy variables, but cannot be categorical variables. If you have categorical predictors, they should be coded into one or more dummy variables.

3.28. Discriminant analysis

Discriminant analysis is used when you have one or more normally distributed interval independent variables and a categorical dependent variable. It is a multivariate technique that considers the latent dimensions in the independent variables for predicting group membership in the categorical dependent variable.

3.29. One-way MANOVA

MANOVA (multivariate analysis of variance) is like ANOVA, except that there are two or more dependent variables. In a one-way MANOVA, there is one categorical independent variable and two or more dependent variables.

3.30. Multivariate multiple regression

Multivariate multiple regression is used when you have two or more dependent variables that are to be predicted from two or more predictor variables.

Many researchers familiar with traditional multivariate analysis may not recognize the tests above. They do not see Wilks' Lambda, Pillai's Trace or the Hotelling-Lawley Trace statistics, the statistics with which they are familiar. It is possible to obtain these statistics



using the **mvtest** command written by David E. Moore of the University of Cincinnati.

3.31. Canonical correlation

Canonical correlation is a multivariate technique used to examine the relationship between two groups of variables. For each set of variables, it creates latent variables and looks at the relationships among the latent variables. It assumes that all variables in the model are interval and normally distributed. Stata requires that each of the two groups of variables be enclosed in parentheses. There need not be an equal number of variables in the two groups.

The output above shows the linear combinations corresponding to the first canonical correlation. At the bottom of the output are the two canonical correlations. Because the output from the **canon** command is lengthy, we will use the **cantest** command to obtain the eigenvalues, F-tests and associated p-values that we want. Note that you do not have to specify a model with either the **canon** or the **cantest** commands if they are issued after the **canon** command.

3.32. Factor analysis

Factor analysis is a form of exploratory multivariate analysis that is used to either reduce the number of variables in a model or to detect relationships among variables. All variables involved in the factor analysis need to be continuous and are assumed to be normally distributed. The goal of the analysis is to try to identify factors which underlie the variables. There may be fewer factors than variables, but there may not be more factors than variables. For our example, let's suppose that we think that there are some common factors underlying the various test scores. We will first use the principal components method of extraction (by using the **pc** option) and then the principal components factor method of extraction (by using the **pcf** option). This parallels the output produced by SAS and SPSS.

4. RESULTS AND DISCUSSION

4.1. Quantitative and qualitative data :

In advanced studies, a researcher may approach his topics quantitatively, qualitatively or with the use of a mixed methodology. When opting for a qualitative approach, researchers have several options in analyzing the data. The use of matrices, charts, tables and other visual displays are common tools used. With visual displays, the researchers can pare down the often abundant subjective data that has been gathered and determine what will be useful variables in his qualitative data analysis.

One way educational researchers work to overcome the challenge of repeatability is to distinguish, in their reports, between repeatable practices and the non repeatable results that emerged from those practices. Quantitative research can demonstrate rigor by including a wide variety of numerical and statistical data Schroder, K.E., Carey, M.P., Venable, P.A. (2003), while the rigor of qualitative research is harder to demonstrate because it often



involves the qualitative analysis of qualitative data. For example in literary studies, researchers apply interpretive models to texts such as poems or novels. A literary researcher can apply a wide variety of interpretation models and can apply a single interpretive model in multiple ways to a variety of texts. Therefore it is difficult to generate a unifying set of criteria for determining whether that researcher's work is truly rigorous. When the researcher is applying qualitative models of analysis to qualitative or numerical data, the research process can be long and tedious because the researcher must carefully pore over the data in detail while crafting the analysis. For example to write a comprehensive historical account, a historian must examine hundreds of primary historical records and secondary historical accounts. Even after spending all his time and energy examining records and accounts, the historian has no guarantee that it covered everything. One way to compensate for the time-consuming problem of qualitative research is to promote qualitative research projects, such as writing historical accounts, as team based or collaborative.

After collection of data, the selection of statistical test is more important. To select the right test, two questions arise, What kind of data have you collected ? and what is your goal ? Accordingly you have to select the statistical test.

4.2. Limitations to qualitative research:

Qualitative Research is a broad term that refers to research methods most commonly used in fields such as Sociology, anthropology, ethnography and other human and social sciences. The strongest objection to qualitative research is that the quality of the research depends too greatly on the individual researcher (Silverman, S., Manson, M. (2013). Because the researcher designs the type of questions, he or she can in adherently influence the results due to her own personal beliefs.

Because qualitative research is so inextricably entwined with the individual researcher, it is extremely challenging for other researchers to repeat qualitative studies. This makes it hard to confirm or deny the results of the original study. For example, in the field of education, one of the challenges of repeating qualitative study is that different elements of the original study can't be repeated, the teachers and students will all be different, as will the school and classroom environment, the methods of teaching and the styles of learning.

4.3. Usage of excel:

Excel, the spread sheet program in Microsoft's popular office Software Package is a powerful application used to manage various types of data. Excel's capabilities, however are not limited to data management. The program Data Analysis tool enables users to analyze data using an array of statistical procedures that range from descriptive measures to rigorous inferential statistics, such as regression and analysis of variance (Smeeton, N., Goda, D. (2003). The data analysis tool is included in all versions of Excel but must be installed by the



user. Fortunately, setting up and using the tool is relatively easy. We can use Data Analysis for Random Number Generation, to test a hypothesis in Excel to Analyze data.

Excel's data analysis capabilities make it possible to conduct some advanced analyses of survey data but not others. However a program known as XL Stat expands the analytical capabilities of Excel. Tools such as SAS and SPSS are designed with research professionals in mind and make a full range of analytical methods possible.

Choosing between parametric and non parametric tests is sometimes easy. You should definitely choose a parametric test if you are sure that your data are sampled from a population that follows a Gaussian distribution (at least approximately). It is not always easy to decide whether a sample comes from a Gaussian population. If you collect many data points (over a hundred or so) you can look at the distribution of data and it will be fairly obvious whether the distribution is approximately bell shaped. (Thompson, B., Noferi, G. 2002). A formal statistical test (Kolmogorov Smirnov test) can be used to test whether the distribution of the data differs significantly from a Gaussian distribution. But the solution depends on sample size. Parametric tests work well with large samples even if the population is non-Gaussian. In other words, Parametric tests are robust to deviate from Gaussian distributions as long as the samples are large. Parametric test is suitable when there are at least two dozen data points in each group.

Non-Parametric tests work well with large samples from Gaussian population. The p Values tend to be a bit too large, but the discrepancy is small. Non parametric tests are only slightly less powerful than parametric tests with large samples. P value is inaccurate for small samples and it tends to be too high.

5. CONCLUSIONS

In this paper, different types of Statistical tools were explained for the purpose of Research and dissertations for different types of fields .So one should have the skill of selecting a statistical tool for their research which renders good conclusions. Still some more information can be given for the researchers for their future research.

REFERENCES

1. Gottschalk, L. A. (1995). Content analysis of verbal behavior: New findings and clinical applications. Hillside, NJ: Lawrence Erlbaum Associates, Inc International Organization of Scientific Research, IOSR.
2. Journal of Statistical Methodology—Elsevier, WWW. journals.elsevier.com/statistical-methodology.
3. Shamoo, A.E., Resnik, B.R. (2003). Responsible Conduct of Research. Oxford University Press.
4. Savenye, Robinson, 2004 Clinical significance of research: A growing concern. Canadian Journal of Nursing Research, 24, 1-4.
5. (Shepard, 2002), Problems in clinical trials go far beyond misconduct. Science. 264(5165): 1538-41.



6. Resnik, D. (2000). Statistics, ethics, and research: an agenda for educations and reform. *Accountability in Research*. 8: 163-88
 7. Schroder, K.E., Carey, M.P., Venable, P.A. (2003). Methodological challenges in research on sexual risk behavior: I. Item content, scaling, and data analytic options. *Ann Behav Med*, 26(2): 76-103.
 8. Silverman, S., Manson, M. (2013). Research on teaching in physical education doctoral dissertations: a detailed investigation of focus, method, and analysis. *Journal of Teaching in Physical Education*, 22(3): 280-297.
 9. Smeeton, N., Goda, D. (2003). Conducting and presenting social work research: some basic statistical considerations. *Br J Soc Work*, 33: 567-573.
 10. Thompson, B., Noferi, G. 2002. Statistical, practical, clinical: How many types of significance
-