

Statistics essay: interpreting social data



Interpreting Social Data

The British Household Panel Survey of 1991 measured many opinions, among other things, of the UK population. One of the questions asked was whether the husband should be the primary breadwinner in the household, while the wife stayed at home. Answers to the questions were provided on an ordinal scale, progressing in five ordinances from Strongly disagree to Strongly agree. Results for each ordinance were recorded from male respondents and female respondents. Of survey respondents, 96.75%, or $N = 5325.162$ answered this question of a total survey population of $N = 5500.829$. 3.2%, or $N = 175.667$ of survey respondents did not answer the question. In lay terms, this means approximately 97% of the survey respondents answered the question, while 3% did not.

The study presents ordinal ranking, or ranking in a qualitative manner, of five sets of concordant pairs of variables: the male and female count for those who strongly agree the husband be the primary earner while the wife stays at home, the male and female count for those who agree, the male and female count for those who are neutral, the male and female count for those who disagree, and the male and female count for those who strongly disagree. The sex cross-tabulation presents numeric data for responses for each of the ten variables, arranged in five variable pairs with male and female responses for each variable pair. Data is presented in terms of number of responses for each of the ten variables.

The counts or number of responses for each variable are dependent variables in the data analysis. We know they are dependent variables because first,

they are presented on the y-axis in the chart graphically representing the data. Dependent variables are graphically represented on the y-axis, with independent variables presented on the x-axis. Causally it becomes more difficult to distinguish between dependent and independent variables at first glance. Dependent variables usually change as a result of independent variables. For example, if one were studying the effect of a certain medication on blood sugar in diabetics, the independent variable would be the amount of medication given to the patient. In a test group or cohort of patients, each would be given a set dosage and their blood sugar responses recorded. One patient may respond with a blood sugar reading of 110 when given 20mg of medicine. Another day the patient, again given 20mg of medicine, may respond with a blood sugar reading of 240. The amount of medicine provided to the patient is fixed, or the independent variable. The response of the patient is variable, and believed to be influenced by, or dependent on, the amount of medicine provided. The dependent variable would therefore be the responding blood sugar reading in each patient.

In this survey, independent variables are the five choices of answers available to the survey takers. These five possible responses are presented to each survey respondent, just as the medicine is provided to the patient in the example above. The respondent then chooses his or her reply to the five possible answers, or chooses not to answer the question at all. The amount of those choosing not to answer at all, 3.2%, is considered statistically irrelevant in the analysis of this data. Data related to non-response is not considered from either an independent variable or dependent variable standpoint.

The amount of responses or response count for a given independent variable in the survey is a dependent variable. The response count will change, at least slightly, from survey to survey. This could be a due to change in survey size, response rate or number of those choosing to respond to the statement, or possible minor fluctuation in percentage response for the five answer possibilities. Although the statistical results of the responses should be similar, given a large enough and representative sample for each survey attempt, some variance is likely to occur. The independent - dependent variable relationship in the Husband should earn, wife should stay at home analysis is trickier to get one's mind around than the medical example given above. In the medical example, it is easy to grasp how a medicine could affect blood sugar, and the resulting cause-effect relationship. In this survey, the creation of five answer groups causes the respondents to categorise their opinion into one of the groups, a much more difficult mental construction than more straightforward cause-result examples.

Four examples of dependent variables in these statistics are the number of men who agreed with the statement (525), the number of women who agreed with the statement (520), the number of men who disagreed with the statement (688), and the number of women who disagreed with the statement (997). As described above, we know these are dependent variables because they are caused by the independent variables, the five ordinal answer groups, in the survey.

Overall, empirical data for the results is skewed towards the Disagree / Strongly disagree end of the survey. Three of the independent variables are of particular note. Strongly agree is the lowest response for both men and

women, with Disagree being the highest response for both men and women although according to Gaussian predictions the Not agree/disagree variable should have the highest distribution.

In lay terms, the graphical representation of each of the five possible answers should have looked like a bell-shaped curve. The two independent variables on each end of the chart, Strongly agree and Strongly disagree, should have had a low but approximately equal response. The middle independent variable on the chart, Not agree / disagree, should have been the largest response. This should have produced dependent variables of approximately 935 each for both men and women for the Not agree / disagree variable. Instead, the response for men was 586, or 63% of typical distribution of answers. The response for women was 702, or 75% of the typically distributed answers. The mean, or average, of all responses in this survey is 1065.2, with the mean or average of male responses being 464.6 and the mean or average of female responses being 600.6. Were the responses distributed evenly amongst all five possible answers, these would be the anticipated response counts.

In examining this data, a hypothesis can be put forth that the correlation between the counts on two of the answer possibilities (two of the dependent variables) will be some value other than zero, at least in the population represented by the survey respondents. This hypothesis can be tested using the ordinal symmetric measures produced in the data analysis. As Pilcher describes, when data on two ordinal variables are grouped and given in categorical order, we want to determine whether or not the relative positions of categories on two scales go together' (1990, 98). Three ordinal symmetric

<https://assignbuster.com/statistics-essay-interpreting-social-data/>

measures, Kendall's tau-b, Kendall's tau-c, and Gamma, were therefore calculated to determine if the order of categories on the amount of agreement to the question would help to predict the order of categories on the count or amount of those selecting each ordinal category. The most appropriate measures of association to evaluate this hypothesis are the two Kendall's tau measures. The Kendall tau-c measure allows for tie correction not considered in the Kendall tau-b measure. The results of these measures, value .083 and .102 with approximate T^b of 6.75 indicate there is neither a perfect positive or perfect negative correlation between variables. Results do indicate a low level of prediction and approximation of sampling distribution. The correlation between two of the dependent variables is indeed a value other than zero, proving the hypothesis correct.

Three nominal symmetric measures were also calculated. These showed weak relationship between category and count variables, with a value of only .096 for Phi, Cramer's V, and Contingency Coefficient. These were not used in testing the above hypothesis.

A theory of distribution, Chebyshev's theorem states that the standard of deviation will be increased when data is spread out, and smaller when data is compacted. While the data may or may not present according to the empirical rule (bell-shaped), Chebyshev's theorem contends that defined percentages of the data will always be within a certain number of standard deviations from the mean (Pilcher 1990).

In this example, data is compressed into five possible answer variables. The data does not present according to the empirical rule, but is skewed towards

the disagreement end of the variable scale. However, Chebyshev's theorem does apply relating to the distribution of data according to standard deviation from the mean for nine of the ten dependent variables. The response count of women who Disagree with the statement the Husband should earn, the wife stay at home, was proportionately larger than would be indicated along normal distribution. While the response count for men is also statistically high, it is not beyond the predictions of Chebyshev's theorem. If the survey had been conducted with fewer independent variables, say three ordinances instead of five, the resulting data would be more tightly compacted. If the survey had been conducted with ten ordinances, the data would have been more spread out.

REFERENCES

Pilcher, D., 1990. *Data Analysis for the Helping Professions*. Sage Publications, London.