

Inferential Statistics

Harry Ganzeboom

RESMA course “Data Analysis and Report”

Lecture #1

February 3, 2009

IS and DS lingo

- Standard error
- Confidence level
- Confidence interval
- Probability level
- Significance (alpha) level
- Statistically significant
- T-test / t-value
- F-test
- Chi-squared test
- Alpha / beta (type I/II) error
- Mean, median, mode
- Standard deviation
- Percentile scores
- Z-scores
- Association, correlation
- Regression
- Variance, explained variance

IS

- IS is about DS.
- IS adds to descriptive statistics (such as (difference in) means, standard deviation, regression- and correlation coefficients) on sample data, how these statistics would be in the population.
- Note that population statistics are fixed (but unknown), the uncertainty is in the sample.

What if there is no (good) sample?

- IS assumes (simple) random sampling.
- If there is no random sample, IS quantities have no literal interpretation, they are only ‘metaphors’.
 - What could we have said about the population if the observations would have constituted a simple random sample?
- I still find this a very valuable point of view. Others disagree.
- NOTE: If there is no simple random sample, it does not mean that there is no sampling error. You just cannot model it exactly.

Benchmark: simple random sampling (SRS)

- Random: lottery determines inclusion in the sample.
- Simple: there is only one single draw.
- Almost all IS quantities are calculated assuming SRS.
- In practical sampling SRS is hardly ever used.
- Techniques for handling complex random sampling: multi-level analyses, robust estimation, weights.

Random, but not simple

- Systematic sampling with random begin
- Multistage clustered sampling
- Stratified sampling
- PPS sampling
- Evaluation of these sampling schemes requires specialized statistical programs (STATA).

Non-random sampling

- Quota sampling
- Snowball sampling
- Convenience sampling
- Purposive sampling
- IS for these procedures is at best metaphorical

Sampling distributions

- Take a DS quantity: P, M, SD, R, B or whatever.
- Assume a (fixed) population value.
- Draw an infinite number of samples (SRS) of size N and calculate the sample quantity: p, m, sd, r, b or whatever.
- The frequency distribution of all the sample quantities together is called the sampling distribution.
- See Babbie for an excellent introduction.

Sampling distributions

- Have a normal (symmetric, bell-shaped) form
 - If population distribution is normal, or
 - If N is large (> 30)
 - Note that normality arises with larger N in the sample, irrespective of the nature of the distribution in the population.
- Have a t-distribution, when:
 - Sample is small / distribution is very non-normal.
- The (expected) mean of the sampling distribution is the population value of the DS quantity of interest.
- The (expected) standard deviation (SD) of the sampling distribution can be derived mathematically. It always contains an element that resembles $1/\sqrt{N}$: The higher N , the smaller the SD of the sampling distribution.

SE

- The SD of the sampling distribution is called the Standard Error (SE) and is often printed in computer outputs.
- SE denotes the variation of a statistic when many samples (of size N) would be drawn from the population.
- Important SE's: P , M , $M1-M2$, R .
- Most simple SE is that of the pearson correlation: $1/\sqrt{N}$.

Confidence intervals (CI)

- If we know the (normal) shape of the sampling distribution and SE, we can estimate the population value S of the statistic from the sample.
- The best estimate of the statistic S is its value s in the sample.
- The uncertainty in the estimate is formulated as a (95%) confidence interval: $S = s \pm 2*SE$.
- Again: the variation is among the samples, not in the population value. Each sample will give you a different estimate and in the long run 95% of the CI's will contain the population value.

Significance testing

- (Significance) testing is a (binary) test whether a sample statistics could have been produced by a hypothetical population value, called the null-hypothesis (H_0).
- H_0 usually assumes that a population statistic is 0, but this is not necessary.
- Statistical testing is the decision whether the sample statistic lies inside or outside the confidence interval of the population value specified in H_0 .

Confidence level $1-\alpha$

- Confidence levels are chosen by the researcher. α is the chosen probability that the conclusion (population value is in CI) will be wrong.
- This is called the type I error: rejecting H_0 while it is correct.
- If we increase $1-\alpha$, then we will be wrong less often, but we will also be less informative.
- $1-\alpha$ is conventionally very often chosen at 95%. This is arbitrary, but it could never be, say, 50%.

Probability values

- Using the sampling distribution of the H_0 , we can calculate the probability that the sample statistic would arise.
- This probability (p) is calculated by computer programs and printed in the output. The SPSS header for this is “Sig.”.
- If the probability is lower than the significance level α , we reject the H_0 and call the sample result “statistically significant”.
- We will be wrong in $\alpha\%$ of the case and commit an alpha-error (type I error): reject the H_0 , while it is correct.
- We know in advance how often we will make an α -error, but NOT when we make it.

T-test

- We can compare the probability to the significance levels, or alternatively calculate:
- $\text{Effect} / \text{SE} = t$
- If $t > 2$ or $t < -2$, we call the result statistically significant.
- Exact evaluation of t depends upon degrees of freedom, but this matters only for small samples.
- Most computer programs print both p and t .

One- and two-tailed

- Computer programs usually calculate two tailed probability levels: what is the probability that the sample statistic or its negative counterpart would arise if H_0 is true?
- This is so because computer programs cannot know the expected direction of a result.
- However, researchers usually have a one-tailed interest: they have an directional alternative hypothesis in mind.
- One-tailed probabilities are just half the two-tailed ones.
- The decision about one/two tailed is in the mind of the researcher.

β -errors

- A β -error (or type II error) occurs when the H_0 is not true (but some alternative hypotheses H_1), but we accept ('do not reject') the H_0 (i.e. decide that the sample statistic is not significant)
- If we do not have a fixed H_1 , the size of β cannot be calculated. However, we do know some circumstances in which β becomes smaller / larger.
- The ability to avoid β -errors is called the statistical power of a test [*onderscheidingsvermogen*]. Power = $1 - \beta$.

When is β smaller?

- β is smaller:
 - With larger sample size N
 - With more extreme H_1
 - With higher alpha
 - In one-sided problems
 - With better measurement, also higher level of measurement
 - With stronger designs: matching, repeated observations / panel, paired observations, higher explained variance
 - With better samples (SRS or better).

Problems with significance testing

- Significance testing leads to a binary decision (yes/no), but our research often requires nuances.
- Significance level is arbitrarily set (at, say, 5%).
- It all depends on SRS assumptions.
- Statistical significance is not relevance.
- In the end, we know more about H_0 than about H_1 – and we know very little about β .