

## § 2.2 Proto-Structure

### "Formally Sampling Language"

2.2

- ~~Desirable which samples are possible w/ language~~
- ~~Desirable properties of corpora // as they are currently used~~
  - ~~Size~~
  - ~~Randomness~~
  - ~~Validity~~
  - ~~Limits on sampling method [practical]~~

### - Desirable Properties of a corpus

- Size
- Randomness
- Validity (internal, for external see § 2.1)
- Replicability / Reliability
- Flexibility

Based on current techniques, but not restricted by method

### - Description of current best practice in sample

2.2.1

- Sampling as applied to language, (ie what is applicable)

- A - Population Definition
- B - Selection of a Sampling Frame
- C - Sampling Methods
- D - Determination of size // Power
- E - Post-hoc / Reweighting
- F - Bias Estimation methods

2.2.2

(NB. Perhaps these should be done in-place)

## Comparison of Current Methods to A-F

A - Population Definition	∞, specialist + GP approx.
B - Sampling Frame	↓ why no speech et
C - Sampling Methods	trad. Web to come next section <sup>anal. sample</sup>
D - Size / Power	ling. estimates + Big data findings
E - Post-hoc / Reweighting	
F - Bias Estimation.	- Literature lim ling.

## A - Population Definition (external validity, ~~lower~~ replicability focus)

- Largely linguistic as a problem
  - What was the aim of corpora used today?
  - What is the basis of external validity in §2.1?
  - Can there be improved?
  - Literature on population estimates for languages
  - ~~What is the difference between~~
  - Causal systems, I vs. E language etc w/  $\Delta$

## - Statistical Measures of specific cases (GP is too inf.)

- Methods for pop. estimation (proportions)

~~- Application~~

- Application to language

- Estimates of web size

- " Books/Publications

- " Spoken word

- " All T(ratios)

perhaps this belongs higher up? 2P

↓ add more re. 2.1 linguistic angle.

## B - Sampling Frame (internal validity, ~~reproducibility~~ reliability focus)

- Enumeration or Bounding

- Special-purpose corpora

- General Purpose

- Personal corpora

- ~~Resistant / reserve~~

~~- Also~~

- Subsampling and (ab)use of GP corpora

- Auxiliary Information: Bonus and availability

## C - Sampling Methods

Non Random, compare to current practice

- Snowball Sampling
- web crawler comparison
- Purposive
- compare to GP corpus Building, see notes.

Random, ideal case

- Simple Random
- Focus on problems with  $\infty$  population, practical issues in establishing PC, practical + legal issues with selection
- Stratified Sampling
- How to select strata? (with ling. relevance). Multidimensional strata. Inverse Prob. Weighting. PPS possibilities with web/offline.
- Multistage Sampling
- IPW, how to balance samples for general purpose use?
- Cluster Sampling
- Website-wide samples. Relate to the structure and legality of publishing, the web.
- Adaptive Sampling
- Perhaps this, guided by multistage, offers a good approximation of WC methods as they are currently?

## D - Sample Size // Power

- Criteria for selecting sizes (Linguistic)
  - Common models "should be big enough for  $x$ "
  - Sufficient representation "should contain  $n$   $x$ s"
  - Breadth "should contain  $n$  types of  $x$ "

- Existing corpora - Focus on the rise of big data.

- For multistage designs, how large should the initial sample be?
  - is this small enough to be guided by the user?

sample more.

## E - Reweighting (perhaps should be merged with Bias methods)

- Establishing weights using auxiliary data
  - Possible sources of valid aux data (Focus on possible uses, perhaps too dependent)
- Existing manual resampling in Ling (eg. selecting categories by  $BVC$ )

F - Bias Estimation NB. more of a work in progress than rest.

- Resampling (Focus on generality)
- ~~Comp~~ Comparison to cross data

## § 2.2 notes

- External Validity of Corpora
  - also transferability when qualitative.
- Should 'Big Data' be mentioned more
- Is current practice basically nonrandom. If so, is it valid to base new corpora on the balance of old ones? [no]  $\rightarrow$  AND
- Do we wish to have representative or comparable results?
- Is linguistic data ever simple enough for ratios or regression estimation to be used to guide samples?
- Adaptive sampling looks to be chosen in a way similar to conventional corpus building
  - conditions of interest may be very useful for § 5 - expressing a sampling policy in a repeatable, valid, and useful way
- It would be nice to separate "how to represent a sampling frame" into a separate discussion, but this probably belongs in § 5.