

# Statistical Design Considerations for Lexicon-Based Information Extraction By Mireya Diaz, PhD

- Large proportion of information extraction (IE) via natural language processing (NLP) is based on lexicons
- Design considerations for a useful lexicon are: Lexicon size, document's word length, prevalence of nonstop words, unique tokens (lemmas) distribution, words' length distribution

Probability that sample will contain one or more instances of the token of interest follows a Poisson distribution<sup>1,2</sup>:

 $p(x > 0) = 1 - e^{-\lambda}$  $\Lambda$ : mean rate of token's occurrence

For each token:  $\lambda_i = f_i/(N/n)$   $f_i$ : frequency of token *i* in the corpus

N: size of corpus n: size of sample

Juckett<sup>2</sup> worked further the Poisson arrival of tokens within a sample focusing in the capture probability for each token given its length providing the following formulae:

P: aggregate capture probability across all token sizes  $b_j$ : overall frequency of tokens of size j $W_s, W_i$ : smallest and largest token sizes chosen

$$P_j = \sum_{i=1}^{m_j} a_{i,j} (1 - e^{-\lambda}) \quad \text{with} \quad \sum_{i=1}^{m_j} a_{i,j} = 1 \quad \text{for each value of } j$$

Example: for 1 token/10 documents and a desired capture probability of 0.95 it would be required 30 documents to extract the token

### Department of Biomedical Sciences, Division of Epidemiology and Biostatistics

Example: would like to build a corpus for words in the toxic substances and disease registry Glossary from www.atsdr.cdc.gov contains 166 terms (75 unigrams -20 are acronyms, 52 bigrams, 26 trigrams, 13 n-grams). Frequency of these terms in Pubmed documents were obtained (displayed in table below under terms and average probability columns). For comparison purposes: frequency of words from the Thorndike and Lorge English list as extracted by EllegaRd<sup>1</sup>.

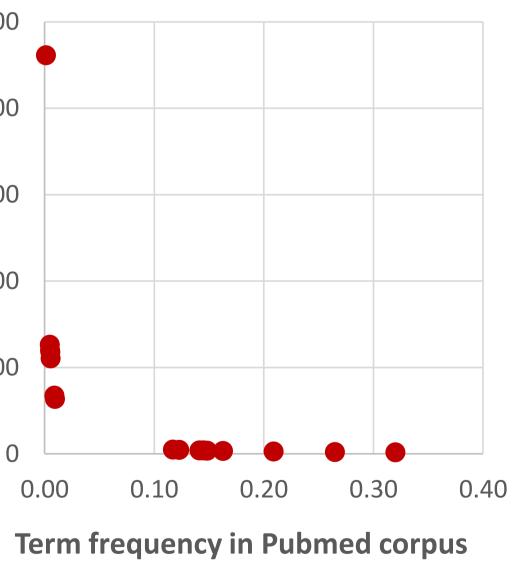
Freq level	Terms	Average probability	Words <sup>1</sup>	Average probability <sup>1</sup>	(	Corpu their
<b>10</b> <sup>1</sup>	1	$1.8 \ge 10^{-1}$	-	_		
<b>10</b> <sup>2</sup>	24	$3.8 \ge 10^{-2}$	12	$1.33 \ge 10^{-2}$		25000
<b>10</b> <sup>3</sup>	29	3.5 x 10 <sup>-3</sup>	109	2.67 x 10 <sup>-3</sup>		
104	35	$4.0 \ge 10^{-4}$	969	$2.58 \ge 10^{-4}$		20000
<b>10</b> <sup>5</sup>	37	3.9 x 10 <sup>-5</sup>	4,810	$3.08 \ge 10^{-5}$	<b>S</b> )	
106	19	<b>2.8</b> x 10 <sup>-6</sup>	11,100	3.47 x 10 <sup>-6</sup>	ents)	15000
107	11	$4.5 \ge 10^{-7}$	16,000	3.71 x 10 <sup>-7</sup>		
<b>10</b> <sup>8</sup>	9	6.0 x 10 <sup>-8</sup>	20,000	3.95 x 10 <sup>-8</sup>	docum	10000

- Over 23K documents need to be extracted to build the desired • corpus if all 166 terms are sought.
- If 5 terms with the smallest frequency are excluded, then only 6,329 documents would be necessary to obtain the other terms with a frequency of 0.047% or higher, with a 0.95 capture probability.

### References

- 1. EllegåRd A. Estimating vocabulary size. WORD 1960;16:2:219-244.
- 2. Juckett D. A method for determining the number of documents needed for a gold standard corpus. J Biomed Informatics 2012;45:460-470.

## ous size for terms depending on eir frequency and word length



size

Corpus