Measuring and Modelling Lexical Diversity across Languages

Christian Bentz & Douwe Kiela cb696@cam.ac.uk



Outline

Background

- Definition, measures
- Qualitative results: lexical diversity space, lexical diversity scales

Statistical Modelling

- Sociolinguistic factors, language families and regions
- Quantitative results: simple linear model, mixed-effects model

Future Directions

• Potential repercussions on word order flexibility





Definition

Lexical diversity

The distribution of word forms used to encode a constant information content





Frequency distributions

Order types (word forms delimited by white spaces) according to their token frequencies (Zipf,1932,1949)



Frequency distributions

Order types (word forms delimited by white spaces) according to their token frequencies (Zipf,1932,1949)





Frequency distributions

Order types (word forms delimited by white spaces) according to their token frequencies (Zipf,1932,1949)



Zipf-Mandelbrot's law
(Zipf,1949; Mandelbrot
1953)
$$f(r_i) = \frac{C}{\beta + r_i^{\alpha}},$$
$$C > 0,$$
$$\alpha > 0,$$
$$\beta > -1,$$
$$i = 1, 2, \dots, n$$







Shannon entropy (Shannon & Weaver, 1949)

$$H = -K \sum_{i=1}^{k} p_i \times \log_2(p_i)$$
$$p_i : \frac{\text{frequency of } w_i}{\text{total number of tokens}}$$



Shannon entropy (Shannon & Weaver, 1949)

$$H = -K \sum_{i=1}^{k} p_i \times \log_2(p_i)$$
$$p_i : \frac{\text{frequency of } w_i}{\text{total number of tokens}}$$





Shannon entropy (Shannon & Weaver, 1949)

$$H = -K \sum_{i=1}^{k} p_i \times \log_2(p_i)$$
$$p_i : \frac{\text{frequency of } w_i}{\text{total number of tokens}}$$











Linguistic Studies

ZM-law

Zipf (1932, 1935, 1949), Ha et al. (2006), Baroni (2009), Popescu et al. (2008,2009, 2010), Baixeries et al. (2013), Bentz et al. (2014)

Linguistic Studies

ZM-law

Zipf (1932, 1935, 1949), Ha et al. (2006), Baroni (2009), Popescu et al. (2008,2009, 2010), Baixeries et al. (2013), Bentz et al. (2014)

• Shannon H

Gries 2012 \rightarrow Possible connection with language learning?



Linguistic Studies

ZM-law

Zipf (1932, 1935, 1949), Ha et al. (2006), Baroni (2009), Popescu et al. (2008,2009, 2010), Baixeries et al. (2013), Bentz et al. (2014)

• Shannon H

Gries 2012 \rightarrow Possible connection with language learning?

Type-Token Ratio

Tweedie & Baayen (1998), Baayen (2001)



Linguistic Studies

ZM-law

Zipf (1932, 1935, 1949), Ha et al. (2006), Baroni (2009), Popescu et al. (2008,2009, 2010), Baixeries et al. (2013), Bentz et al. (2014)

• Shannon H

Gries 2012 \rightarrow Possible connection with language learning?

Type-Token Ratio

Tweedie & Baayen (1998), Baayen (2001)



LDT in Diachrony: Old English and Modern English

Bentz, Kiela, Hill & Buttery (2014) Zipf's law and the grammar of languages. A quantitative study of Old and Modern English parallel texts. *Corpus Linguistics and Linguistic Theory.* aop.



LDT in Diachrony: Old English and Modern English

Bentz, Kiela, Hill & Buttery (2014) Zipf's law and the grammar of languages. A quantitative study of Old and Modern English parallel texts. *Corpus Linguistics and Linguistic Theory.* aop.



What drove this change in frequencies?

Inflections: OE land, lande, landes \rightarrow ModE land, in the land, of the land



What drove this change in frequencies?

Inflections: OE land, lande, landes \rightarrow ModE land, in the land, of the land



Christian Bentz & Douwe Kiela cb696@cam.ac.uk — Lexical diversity 11/24



Inflections in German and English

Lemmatization: yields frequency distributions without inflected types for English and German



Inflections in German and English

Lemmatization: yields frequency distributions without inflected types for English and German



Inflections and frequencies

Having inflections in a language automatically creates lower frequency items

- \rightarrow higher entropy (higher uncertainty)
- \rightarrow higher type-token ratios
- \rightarrow lower ZM-parameters



Synchronic analyses: LDT across languages of today

Parallel texts

- Universal Declaration of Human Rights (UDHR) in ~400 languages
- Parallel Bible Corpus (~800 languages)
- Europarl Corpus (21 languages)









CAMBRIDGE ASSESSMENT

Lexical Diversity Space

Calculating Zipf-Mandelbrot parameters for 359 languages (words delimited by white spaces)

Christian Bentz & Douwe Kiela cb696@cam.ac.uk — Lexical diversity 15/24



Lexical Diversity Space



Calculating Zipf-Mandelbrot parameters for 359 languages (words delimited by white spaces)



Lexical Diversity Scales

Calculating Shannon entropy H and Type-Token Ratios for 359 languages

Christian Bentz & Douwe Kiela cb696@cam.ac.uk — Lexical diversity 16/24



Lexical Diversity Scales

Calculating Shannon entropy H and Type-Token Ratios for 359 languages





Statistical Modelling of Lexical Diversity

Is it possible to predict LDTs of languages?

Bentz, Verkerk, Kiela, Hill & Buttery (submitted) Adpative languages: Modelling the co-evolution of population structure and lexical diversity



Statistical Modelling of Lexical Diversity

Is it possible to predict LDTs of languages? Bentz, Verkerk, Kiela, Hill & Buttery (submitted) Adpative languages: Modelling the co-evolution of population structure and lexical diversity

Potential factors

Christian Bentz & Douwe Kiela cb696@cam.ac.uk — Lexical diversity 17/24



Statistical Modelling of Lexical Diversity

Is it possible to predict LDTs of languages? Bentz, Verkerk, Kiela, Hill & Buttery (submitted) Adpative languages: Modelling the co-evolution of population structure and lexical diversity

Potential factors

language families



Statistical Modelling of Lexical Diversity

Is it possible to predict LDTs of languages? Bentz, Verkerk, Kiela, Hill & Buttery (submitted) Adpative languages: Modelling the co-evolution of population structure and lexical diversity

Potential factors

- language families
- language regions



Statistical Modelling of Lexical Diversity

Is it possible to predict LDTs of languages? Bentz, Verkerk, Kiela, Hill & Buttery (submitted) Adpative languages: Modelling the co-evolution of population structure and lexical diversity

Potential factors

- language families
- language regions
- writing systems



Statistical Modelling of Lexical Diversity

Is it possible to predict LDTs of languages? Bentz, Verkerk, Kiela, Hill & Buttery (submitted) Adpative languages: Modelling the co-evolution of population structure and lexical diversity

Potential factors

- language families
- language regions
- writing systems
- non-native influence (McWhorter, 2002, 2007, 2011; Wray & Grace, 2007; Lupyan & Dale 2010, 2012; Trudgill, 2011; Bentz & Winter, 2013)



- LDT measures (ZM-parameters, entropy H, TTR) for 359 language



- LDT measures (ZM-parameters, entropy H, TTR) for 359 language
- L2 ratio information on 231 languages



- LDT measures (ZM-parameters, entropy H, TTR) for 359 language
- L2 ratio information on 231 languages
- \rightarrow Sample: 99 overlapping languages



- LDT measures (ZM-parameters, entropy H, TTR) for 359 language
- L2 ratio information on 231 languages
- \rightarrow Sample: 99 overlapping languages





Mixed-effects regression

Regressing LDT on L2 ratios but controlling for non-independence of datapoints (language families, regions, writing systems) **Results**: α ($\chi^2(1) = 1.38, p > 0.5$), H ($\chi^2(1) = 9.28, p < 0.001$), TTR ($\chi^2(1) = 7.11, p < 0.01$)



CAMBRIDGE ASSESSMEN

Mixed-effects regression

Regressing LDT on L2 ratios but controlling for non-independence of datapoints (language families, regions, writing systems) **Results**: α ($\chi^2(1) = 1.38, p > 0.5$), H ($\chi^2(1) = 9.28, p < 0.001$), TTR ($\chi^2(1) = 7.11, p < 0.01$)



Christian Bentz & Douwe Kiela cb696@cam.ac.uk — Lexical dive 19/24

Questions

• Does that mean non-native speakers ruin languages?





Questions

- Does that mean non-native speakers ruin languages?
- Are there more and less efficient/expressive languages?



Questions

- Does that mean non-native speakers ruin languages?
- Are there more and less efficient/expressive languages?
- Are there simplex and complex languages?



Questions

- Does that mean non-native speakers ruin languages?
- Are there more and less efficient/expressive languages?
- Are there simplex and complex languages?

Suggestions

- Languages are **complex adaptive systems** shaped by the cognitive and conceptual needs of their speakers (Croft, 2000; Ritt, 2004; Christiansen & Chater, 2008; Beckner et al., 2009)
- A lack of lexical diversity might be made up for by encoding of information at a different level (constructions, fixed word order, multi word expressions)



 Permutation entropy: Reflects the mutual dependence of words in ngrams, i.e. word sequences (Zhang et al. 2006; Ramisch et al. 2008)



- Permutation entropy: Reflects the mutual dependence of words in ngrams, i.e. word sequences (Zhang et al. 2006; Ramisch et al. 2008)
- PE = 0 → words are completely dependent (fixed order)



- Permutation entropy: Reflects the mutual dependence of words in ngrams, i.e. word sequences (Zhang et al. 2006; Ramisch et al. 2008)
- PE = 0 → words are completely dependent (fixed order)
- PE = 1 → words are completely independent (free order)



- Permutation entropy: Reflects the mutual dependence of words in ngrams, i.e. word sequences (Zhang et al. 2006; Ramisch et al. 2008)
- PE = 0 → words are completely dependent (fixed order)
- PE = 1 → words are completely independent (free order)





Lexical Diversity (LDT)

• Reflects **information encoding** strategies (inflection, compounding, base vocabulary)



Lexical Diversity (LDT)

- Reflects **information encoding** strategies (inflection, compounding, base vocabulary)
- It can be measured **quantitatively** (ZM-law, Shannon entropy, type-token ratios)



Lexical Diversity (LDT)

- Reflects **information encoding** strategies (inflection, compounding, base vocabulary)
- It can be measured **quantitatively** (ZM-law, Shannon entropy, type-token ratios)
- It is linked to language learning via frequency effects



Lexical Diversity (LDT)

- Reflects **information encoding** strategies (inflection, compounding, base vocabulary)
- It can be measured **quantitatively** (ZM-law, Shannon entropy, type-token ratios)
- It is linked to language learning via frequency effects
- It can be **statistically predicted** across languages (language families, language regions, writing systems, potential non-native influence)



Collaborators



Douwe Kiela



Felix Hill



Andrew Caines



Dimitrios Alikaniotis



Paula Buttery





Christian Bentz & Douwe Kiela cb696@cam.ac.uk — Lexical diversity 23/24

Thank You!

jchris@christianbentz.de¿