



Language Evolution WiSe 2023/2024

Lecture 9: Quantitative Linguistics

21/11/2023, Christian Bentz



Overview

Introduction

Section 1: Word Frequency Distributions

- Lexical Typology

- Morphological Typology

- Writing Systems

- Real World Salience

Section 2: Simple Measures

- Type-Token-Ratios (TTR)

- Repetition Rates

Section 3: Quantitative Laws

- Zipf's Law of Word Frequencies

- Zipf's Law of Abbreviation

- Menzerath-Altmann Law

- Heaps' Law

Summary

References



Introduction



What is unique about human language?



“If a Martian scientist [...] received from Earth the broadcast of an extensive speech [...] what criteria would [...] determine whether the reception represented the effect of an animate process on Earth, or merely the latest thunderstorm on Earth?

[...]

It seems that the only criteria would be the arrangement of occurrences of the elements [...]: the arrangement of the occurrences would be neither of **rigidly fixed regularity** [...] nor yet a **completely random scattering** of the same.”

Zipf (1965). The psycho-biology of language, p. 187.

Introduction

Section 1: Word Frequency Distributions

Section 2: Simple Measures

Section 3: Quantitative Laws

Summary

References



Competing Definitions of *Language*

- ▶ **Formal Language Theory**
- ▶ **Faculty of Language**
 - ▶ Recursion
 - ▶ Rich Language Faculty (Narrow Sense)
- ▶ **Minimalism**
 - ▶ Strong Minimalist Thesis
 - ▶ Minimalist Layers Hypothesis
- ▶ **Usage-Based Grammar**

Introduction

Section 1: Word
Frequency
Distributions

Section 2: Simple
Measures

Section 3:
Quantitative Laws

Summary

References



Definition (Usage-Based)

“While all linguists are likely to agree that *grammar is the cognitive organization of language*, a **usage-based theorist** would make the more specific proposal that grammar is the cognitive organization of one’s **experience with language.**”

Bybee (2006). From usage to grammar: The mind’s response to repetition.



Introduction

Section 1: Word
Frequency
Distributions

Section 2: Simple
Measures

Section 3:
Quantitative Laws

Summary

References



Definition (Usage-Based)

From the **usage-based** perspective **language** is ultimately a **mapping** from phonetic shapes (or hand shapes in sign language, or graphemes in writing) to semantic or pragmatic context. The strength of this mapping is determined by the frequency of co-occurrence.

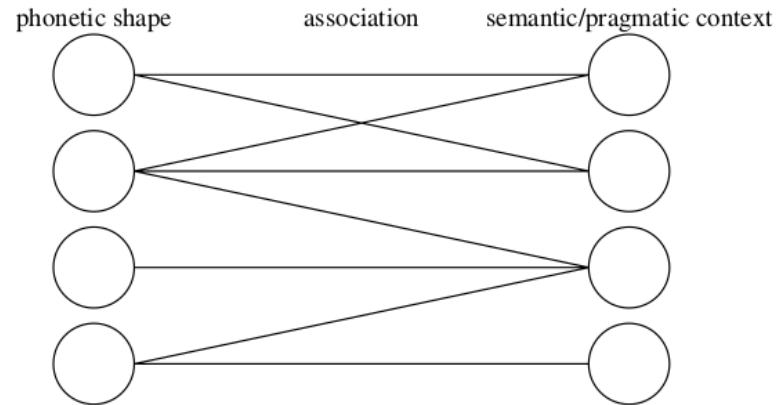


FIGURE 3. Variable associations of form and meaning in a linguistic sign.

Bybee (2006). From usage to grammar: The mind's response to repetition.

Introduction

Section 1: Word Frequency Distributions

Section 2: Simple Measures

Section 3: Quantitative Laws

Summary

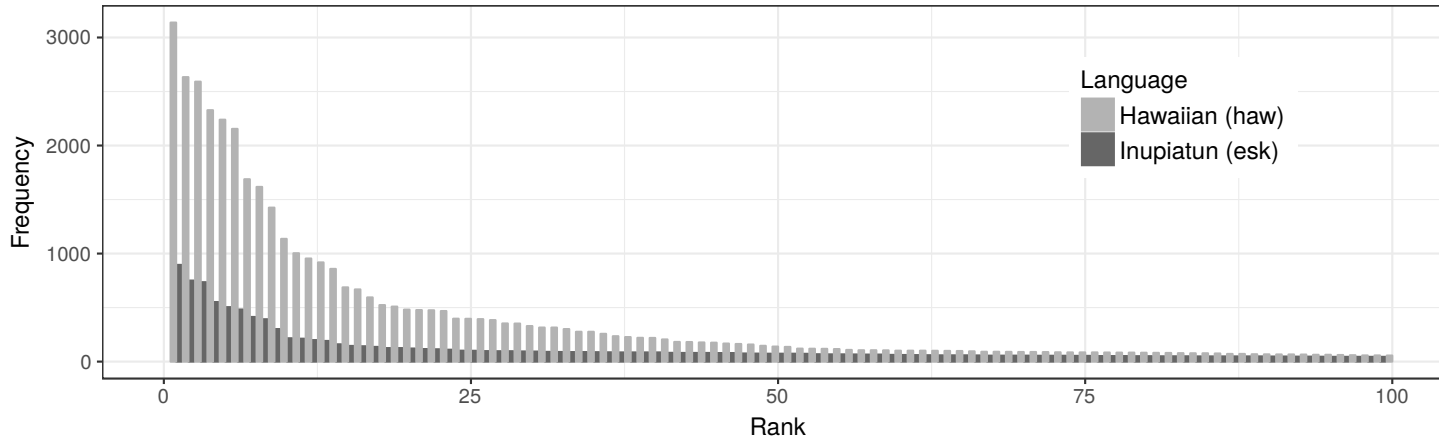
References



Section 1: Word Frequency Distributions



Word Frequency Distributions



Introduction

Section 1: Word Frequency Distributions

Section 2: Simple Measures

Section 3: Quantitative Laws

Summary

References

Hawaiian (haw)

40001001 O ke kuauhau na ka hanauna o Iesu Kristo , ka mamō a Davida , ka mamō a Aberahama.

40001002 Na Aberahama o Isaaka ; na Isaaka o Iakoba ; na Iakoba o Iuda a me kona poe hoahanau;

[...]

Iñupiatun (esk)

40001001 Uvva ukua aglanjich sivullianjñ Jesus Christ-ŋum , kinguvianjupluni David-miñ Abraham-miñ!u .

40001002 Abraham aapagigaa Isaac-ŋum , Isaac-li aapagigaa Jacob-ŋum , Jacob-li aapagigaa Judah-ŋum aniqataiñ!u .

[...]

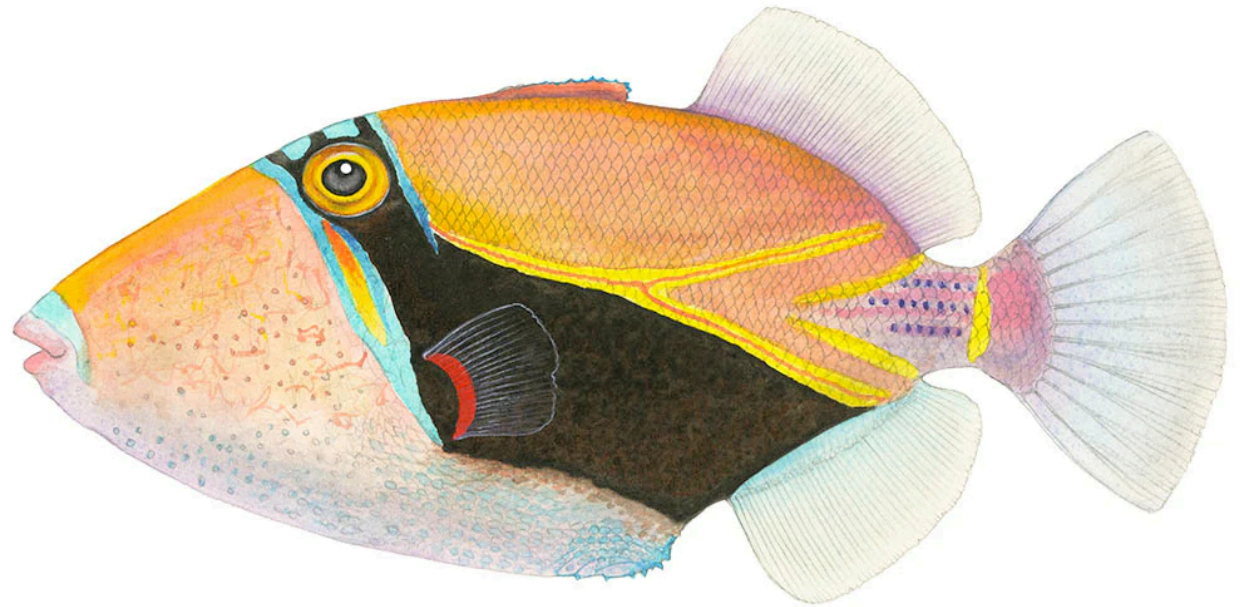
Mayer and Cysouw (2014). A massively parallel Bible corpus.



Differences in the Lexicon

“Centuries before there were marine biologists and scientific methods for classifying fish and other marine life, Pacific Islanders were passing on their accumulated knowledge about the behavior of each of hundreds of varieties of fish.”

Nettle and Romaine
(2000). *Vanishing voices*,
p. 56.



humuhumunukunukuapua'a
“Triggerfish with a snout like a pig”



Differences in Morphological Typology

- (1) **Hawaiian** (haw, PBC 41006018)

A ua olelo aku o loane ia ia [...]

Then PERF say to SUBJ Johan he.DAT [...]

“Then Johan said to him [...]”

- (2) **Turkish** (tur, PBC 41006004)

Ýsa da on-lar-a [...] *de-di*

Jesus also 3P-PL-DAT [...] say-3SG.PERF

“Jesus also said to them [...]”

- (3) **Iñupiatun** (esk, PBC 41006004)

Aglaan Jesus-ngum itna-ġ-ni-ġai [...]

But Jesus-ERG this-say-report-3S.to.3PL

“But Jesus said to them (it is reported) [...]”

Bentz (2018). Adaptive languages: An information-theoretic account of language diversity.

Introduction

Section 1: Word
Frequency
Distributions

Section 2: Simple
Measures

Section 3:
Quantitative Laws

Summary

References



Side Note: Distributional Typology

Typologists are moving away from **aggregated** and **paradigmatic** perspectives on morphology and rather towards a **distributional** view.

Introduction

Section 1: Word
Frequency
Distributions

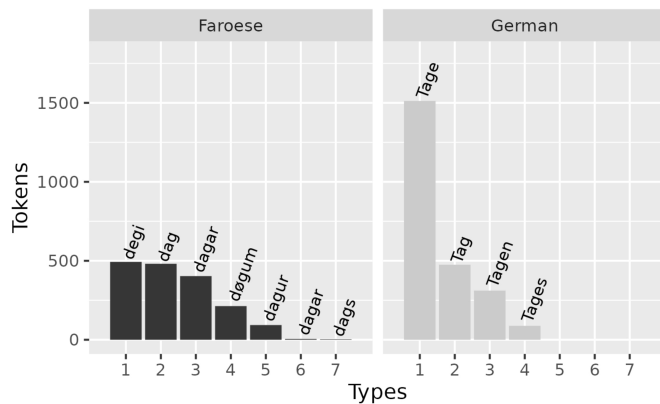
Section 2: Simple
Measures

Section 3:
Quantitative Laws

Summary

References

a) Distributional

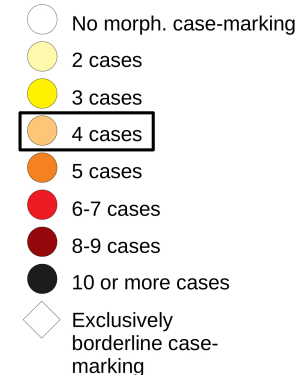


b) Paradigmatic

Singular		Plural	
Nom.	<i>dagur</i>	Nom.	<i>dagar</i>
Gen.	<i>dags</i>	Gen.	<i>daga</i>
Dat.	<i>degi</i>	Dat.	<i>døgum</i>
Acc.	<i>dag</i>	Acc.	<i>dagar</i>

Singular		Plural	
Nom.	<i>Tag</i>	Nom.	<i>Tage</i>
Gen.	<i>Tages</i>	Gen.	<i>Tage</i>
Dat.	<i>Tage</i>	Dat.	<i>Tagen</i>
Acc.	<i>Tag</i>	Acc.	<i>Tage</i>

c) Aggregated (WALS)



Bentz and Verkerk (forthcoming). Sociolinguistic typology: complexification and simplification.

Wälchli (2012). Indirect measurement in morphological typology.



Differences in Writing Systems

Introduction

Section 1: Word
Frequency
Distributions

Section 2: Simple
Measures

Section 3:
Quantitative Laws

Summary

References

ISO 15924	Script	Writ. Sys.	No.	Example*
Latn	Latin	Alphabet	38	And they remembered his words ,
GreK	Greek	Alphabet	1	Και ενεθυμηθησαν τους λογους αυτου .
Deva	Devanagari	Abugida	1	तब उस की बातें उन को स्मरण आई ।
Geor	Georgian	Alphabet	1	და მთავისებრი სიტყუანი მისნი .
Hang	Hangul	Alphabet	1	저희가 예수의 말씀을 기억하고
Mymr	Burmese	Abugida	1	မိန့်တော်မူခဲ့သောစကားများကိုပြန်သတိရ၍ ။ -
Cyrl	Cyrillic	Alphabet	1	И они вспомнили эти слова Его .

*Verse number 42024008 of the New Testament.

Gutierrez-Vasques, Bentz, and Samardžić (2023). Languages through the looking glass of BPE compression.



Differences in Translation

Korean (*kor*)

- (5) 저희가 예수의 말씀을 기억하고
jeo-hui=ga yei-su=ui mal-sseum=eul gi-eog=ha-go
1P-PL=SUBJ Jesus=POSS speak.HON=OBJ remember=COM
Literal translation: “And we remember Jesus’ speech.”
English verse: “And they remembered his words.”

Gutierrez-Vasques, Bentz, and Samardžić (2023). Languages through the looking glass of BPE compression.

Introduction

Section 1: Word
Frequency
Distributions

Section 2: Simple
Measures

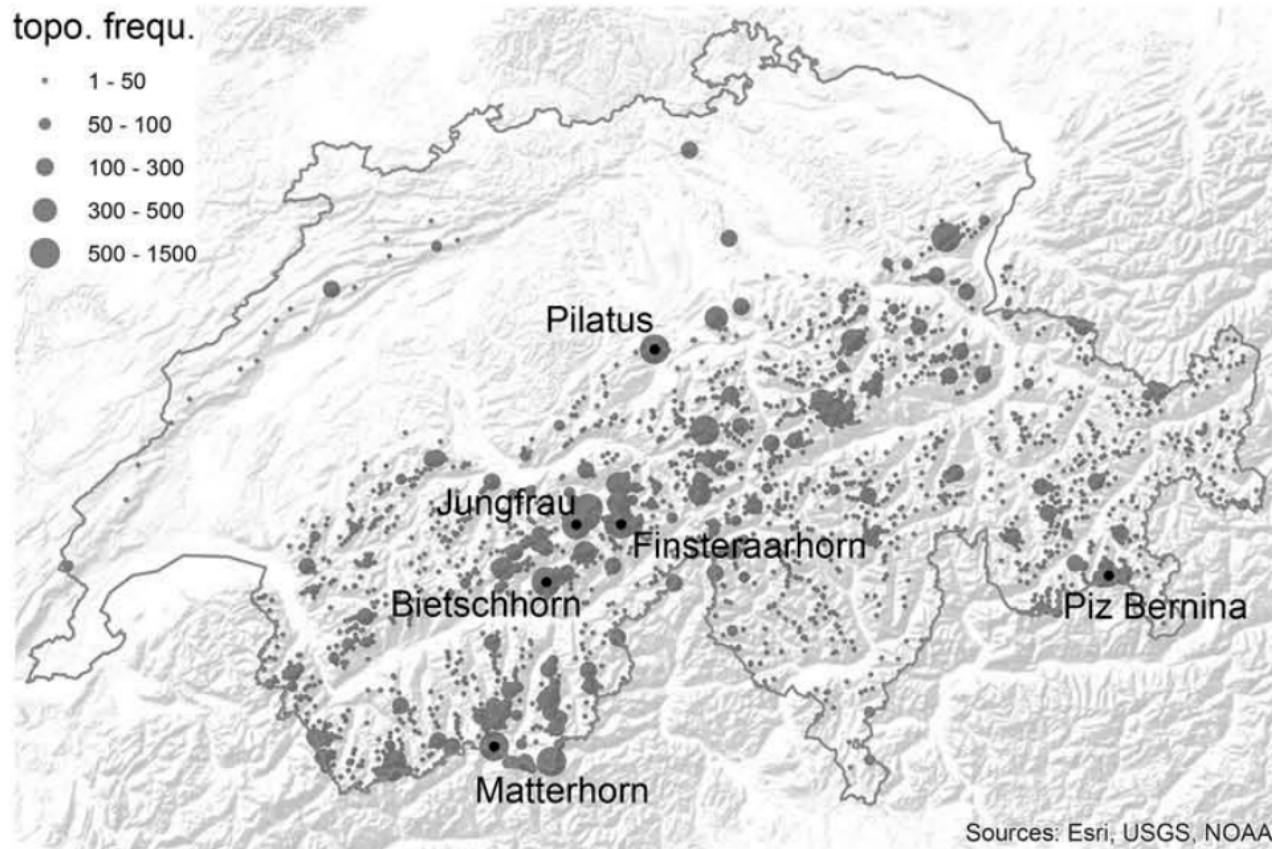
Section 3:
Quantitative Laws

Summary

References



Differences in the “Real World”



Introduction

Section 1: Word
Frequency
Distributions

Section 2: Simple
Measures

Section 3:
Quantitative Laws

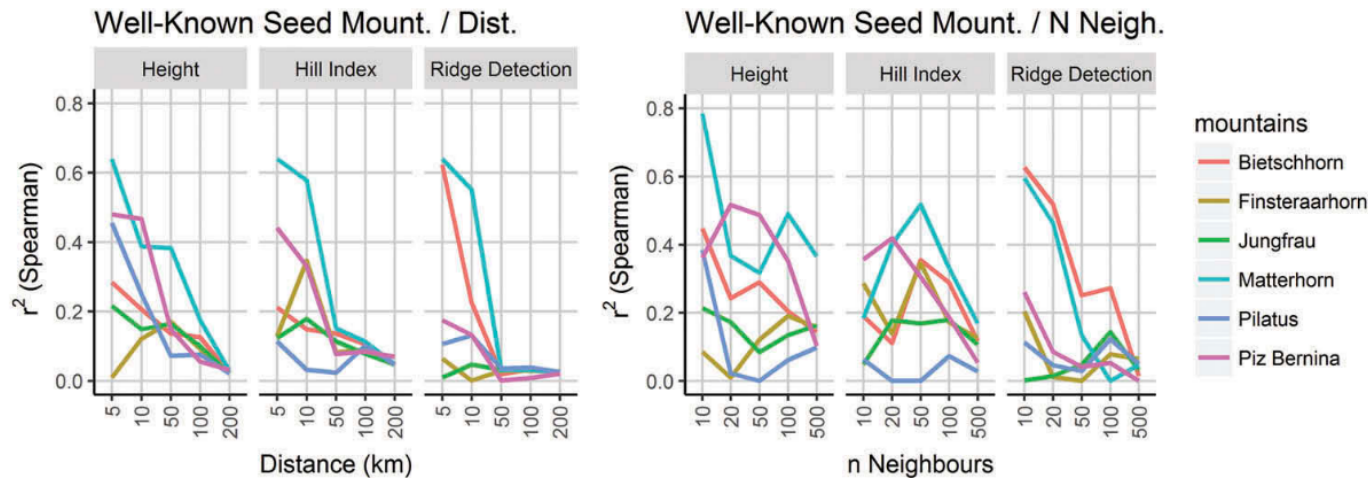
Summary

References

Derungs & Samardžić (2017). Are prominent mountains frequently mentioned in text?



Differences in the “Real World”



Introduction

Section 1: Word Frequency Distributions

Section 2: Simple Measures

Section 3: Quantitative Laws

Summary

References

Figure 5. The relation toponym frequency: spatial measure tested for different spatial extents and a set of well-known seed mountains.

The frequency of occurrence of so-called toponyms (in this case names of famous mountains) in texts is significantly correlated with measures of spatial salience (e.g. height), especially if a text is written in a location close-by.



Real World Salience \leftrightarrow Word Frequencies



.... Matterhorn ... Matterhorn ... Bietschhorn ... Jungfrau ... Matterhorn ...
Pilatus ... Matterhorn ... Finsteraarhorn ... Bietschhorn ... Matterhorn ...
Finsteraarhorn ... Matterhorn ... Matterhorn ... Jungfrau ... Bietschhorn ...
Matterhorn ... Matterhorn ... Bietschhorn

Introduction

Section 1: Word
Frequency
Distributions

Section 2: Simple
Measures

Section 3:
Quantitative Laws

Summary

References



Summary

Factors influencing *Word Frequency Distributions*:

- ▶ Lexicon
- ▶ Morphology
- ▶ Writing Systems
- ▶ Translation/Content
- ▶ Real World Salience
- ▶ etc.

Introduction

Section 1: Word
Frequency
Distributions

Section 2: Simple
Measures

Section 3:
Quantitative Laws

Summary

References

Methodological Question

How can we measure the differences in *Word Frequency Distributions*?



Section 2: Simple Measures



Type-Token Ratio (TTR)

$$TTR = \frac{V}{\sum_{i=1}^V f_i}$$

- ▶ \mathcal{V} : set of unique types (*vocabulary*), e.g. $\mathcal{V} = \{A, a, b, \dots\}$, with $|\mathcal{V}| = V$,
- ▶ V : number of character types,
- ▶ f_i : token frequency of given type.

Example

All human beings are born free
and equal in dignity and rights

char.types	freq	word.types	freq
a	5	All	1
A	1	human	1
b	2	beings	1
d	3	are	1
e	5	born	1
f	1	free	1
g	3	and	2
...

$$TTR = \frac{19}{51} = 0.37$$

$$TTR = \frac{11}{12} = 0.92$$

Introduction

Section 1: Word
Frequency
Distributions

Section 2: Simple
Measures

Section 3:
Quantitative Laws

Summary

References



Repetition Rate (R)

$$R = \frac{r}{\sum_{i=1}^V f_i - 1}$$

- ▶ **r**: number of **adjacent repetitions**,
- ▶ **V**: number of types,
- ▶ **f_i** : token frequency of a given type.

Note: r is the number of *actual repetitions*, and the term $\sum_{i=1}^V f_i - 1$ in the denominator is the number of *possible repetitions* given the token frequencies.

Example

All human beings are born free and equal in dignity and rights

char.types	freq
a	5 - 1 = 4
A	1 - 1 = 0
b	2 - 1 = 1
d	3 - 1 = 2
e	5 - 1 = 4
f	1 - 1 = 0
g	3 - 1 = 2
...	...

word.types	freq
All	1 - 1 = 0
human	1 - 1 = 0
beings	1 - 1 = 0
are	1 - 1 = 0
born	1 - 1 = 0
free	1 - 1 = 0
and	2 - 1 = 1
...	...

$$R = \frac{2}{32} = 0.0625$$

$$R = \frac{0}{1} = 0$$

Sproat (2014). A statistical comparison between written language and nonlinguistic symbol systems.

Introduction

Section 1: Word
Frequency
Distributions

Section 2: Simple
Measures

Section 3:
Quantitative Laws

Summary

References



Exercise

Take the strings of characters below and calculate the TTR and R for characters (white spaces are not counted). What is different/similar for the English character sequence compared to the others? What is the problem with the repetition measure R ?

All human beings
uj kd ro su sv sw sx
GGTAGTTAGGGTCT
N01 N01 N01 ZATU6
SWCCSSSSSSSSSS

Introduction

Section 1: Word
Frequency
Distributions

Section 2: Simple
Measures

Section 3:
Quantitative Laws

Summary

References



Section 3: Quantitative Laws

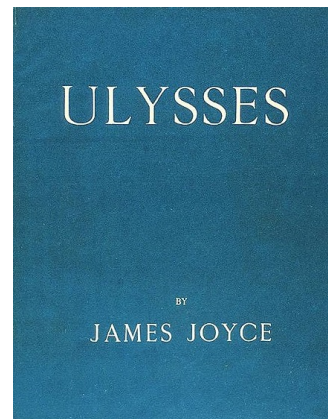


Quantitative Linguistics

Word	Rank	Freq	Char
the	1	12539	3
and	2	9964	3
of	3	7459	2
to	4	7317	2
in	5	3985	2
you	6	3747	3
for	7	3014	3
is	8	2957	2
he	9	2925	2
a	10	2862	1
...
work-then	2742	1	10
world-rulers	2743	1	12
worm	2744	1	4
wormwood	2745	1	8
wounding	2746	1	8
writer	2747	1	6
writers	2748	1	7
zarephath	2749	1	9
zenas	2750	1	5



George Kingsley Zipf



Introduction

Section 1: Word
Frequency
Distributions

Section 2: Simple
Measures

**Section 3:
Quantitative Laws**

Summary

References

Zipf's Law (of Word Frequencies)

Word	Rank	Freq	Char
the	1	12539	3
and	2	9964	3
of	3	7459	2
to	4	7317	2
in	5	3985	2
you	6	3747	3
for	7	3014	3
is	8	2957	2
he	9	2925	2
a	10	2862	1
...
work-then	2742	1	10
world-rulers	2743	1	12
worm	2744	1	4
wormwood	2745	1	8
wounding	2746	1	8
writer	2747	1	6
writers	2748	1	7
zarephath	2749	1	9
zenas	2750	1	5

“[...] we have found a clearcut correlation between the number of different words in the *Ulysses* and the frequency of their usage [...]”:

$$r \times f = C \quad (1)$$

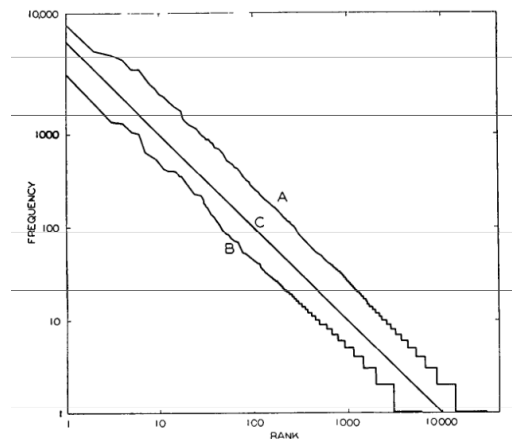


Fig. 2-1. The rank-frequency distribution of words. (A) The James Joyce data; (B) the Eldridge data; (C) ideal curve with slope of negative unity.

Zipf, G. K. (1949). Human behavior and the principle of least effort, p. 24.

Introduction

Section 1: Word Frequency Distributions

Section 2: Simple Measures

Section 3: Quantitative Laws

Summary

References



Zipf's Law (of Word Frequencies)

Word	Rank	Freq	Char
the	1	12539	3
and	2	9964	3
of	3	7459	2
to	4	7317	2
in	5	3985	2
you	6	3747	3
for	7	3014	3
is	8	2957	2
he	9	2925	2
a	10	2862	1
...
work-then	2742	1	10
world-rulers	2743	1	12
worm	2744	1	4
wormwood	2745	1	8
wounding	2746	1	8
writer	2747	1	6
writers	2748	1	7
zarephath	2749	1	9
zenas	2750	1	5

Another (more common) formulation of the law:

$$f(w) \propto \frac{1}{r^\alpha} \quad (2)$$

The α -parameter is the slope in log-log space (i.e. when both the ranks and frequencies are log transformed). Zipf assumed that $\alpha \sim 1$ holds across languages.

Introduction

Section 1: Word Frequency Distributions

Section 2: Simple Measures

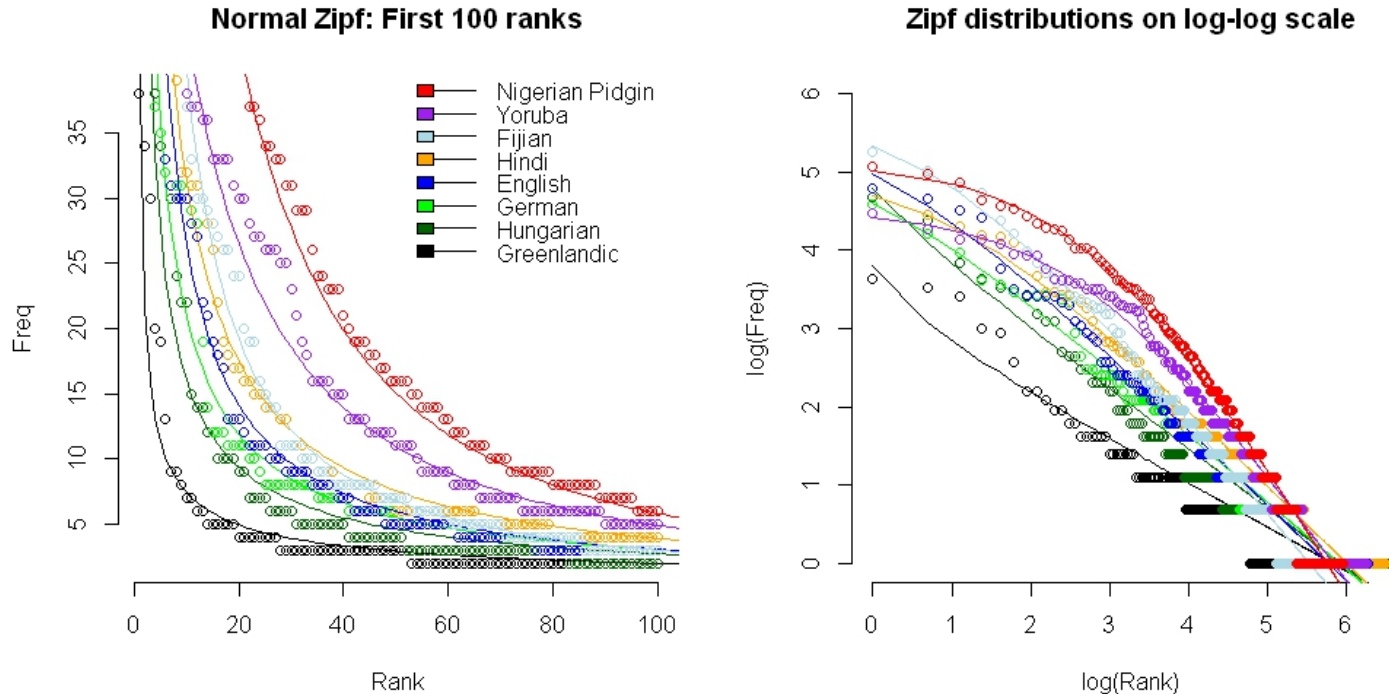
Section 3: Quantitative Laws

Summary

References



Zipf's Law across Languages



Note: This illustrates another version of the law, the *Zipf-Mandelbrot Law*, which has an extra parameter (β) accounting for the deviations from linearity in high frequency items.

Bentz & Kiela (2014). Zipf's law across languages of the world.

Introduction

Section 1: Word Frequency Distributions

Section 2: Simple Measures

Section 3: Quantitative Laws

Summary

References



Zipf's Law of Abbreviation

Word	Rank	Freq	Char
the	1	12539	3
and	2	9964	3
of	3	7459	2
to	4	7317	2
in	5	3985	2
you	6	3747	3
for	7	3014	3
is	8	2957	2
he	9	2925	2
a	10	2862	1
...
work-then	2742	1	10
world-rulers	2743	1	12
worm	2744	1	4
wormwood	2745	1	8
wounding	2746	1	8
writer	2747	1	6
writers	2748	1	7
zarephath	2749	1	9
zenas	2750	1	5

“[...] the magnitude of words tends, on the whole, to stand in an **inverse** (not necessarily proportionate) relationship to the number of occurrences [...]”

Zipf, George K. (1965). The psycho-biology of language: An introduction to dynamic philology, p. 25.

In other words: **more frequent words** tend to be **shorter**.

Introduction

Section 1: Word Frequency Distributions

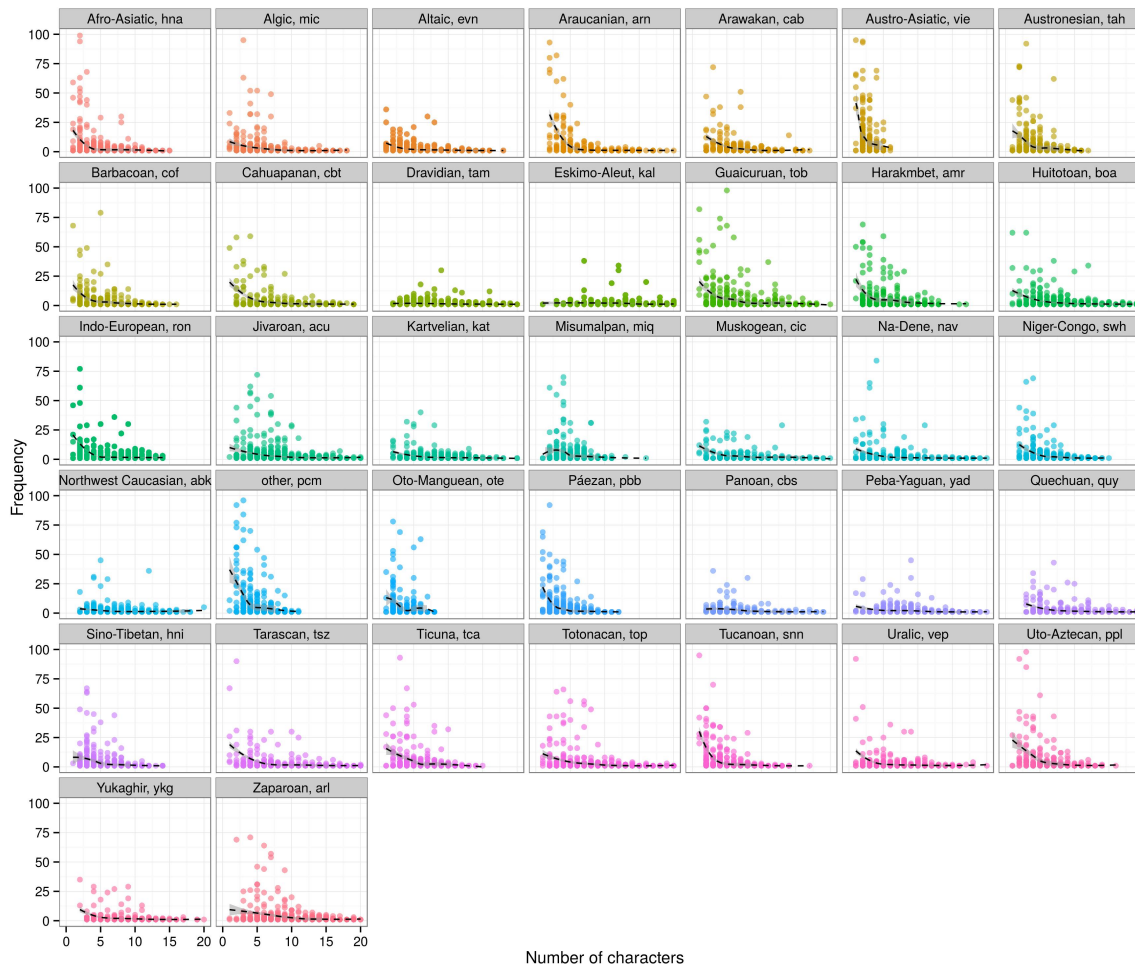
Section 2: Simple Measures

Section 3: Quantitative Laws

Summary

References

Zipf's Law of Abbreviation across Languages



Introduction

Section 1: Word Frequency Distributions

Section 2: Simple Measures

Section 3: Quantitative Laws

Summary

References



Zipf's Law of Abbreviation across Languages

TABLE I

THE CONCORDANCE WITH ZIPF'S LAW OF ABBREVIATION ACROSS 986 LANGUAGES. FOR EACH DATASET, N IS THE NUMBER OF TEXTS OR LANGUAGES, N_{α}^{-} IS THE NUMBER OF NEGATIVE CORRELATIONS BETWEEN WORD FREQUENCY AND WORD LENGTH WITH P-VALUES NOT EXCEEDING α ; N_{α}^{+} IS THE CONVERSE OF N_{α}^{-} FOR POSITIVE CORRELATIONS.

	Texts		Languages	
	PBC	UDHR	PBC	UDHR
N	907	355	801	332
N_1^{-}	907	355	801	332
N_1^{+}	0	0	0	0
$N_{0.05}^{-}$	907	328	801	307
$N_{0.01}^{-}$	907	316	801	296
$N_{0.001}^{-}$	907	283	801	265
$N_{0.0001}^{-}$	907	245	801	230

Bentz & Ferrer-i-Cancho (2016). Zipf's law of abbreviation as a language universal.

Introduction

Section 1: Word
Frequency
Distributions

Section 2: Simple
Measures

Section 3:
Quantitative Laws

Summary

References



Menzerath-Altmann Law



Paul Menzerath (right) with Paulo Lazerda (left) in the Bonn phonetics laboratory.

bə'tsɛr:ɕnuŋ: „,plɑtdiɑ,lɛktɪfə ,ʔɔ:spɾɑ:xə “ 'zɑ:gən ,vɔlən. ɪn 'dɔ:ʧlɑnt ,trɪft ʔe'bən 'nɪçt 'tʃu, vɑs hɛr PɑLMɛr fɪr 'jɑ:pɑn 'fɛstʃɛltə, ʊnt vɑs hɛr GRɑHɑM mu'tɑ:ti:s mu'tɑndi:s (zɔ:!) fɪr 'ʔɑlə spɾɑ:xən 'ʔɑntʃu:nemən fɛrnt': 'nɛmliç, dɑs 'ʔɛ:ɪnfɑx di 'ʔɔ:spɾɑ:xə 'jɛdəs ɪndi'vi:duʊms dər gə'bɪldətən 'klɑsə əl:s 'mʊstər gɛltən dɑ:f. 'dɑs bə,ʃvɛrflə ɪç fɔ'n fɪr grɔ:sbrɪ'tʌniən, əs fɪmt 'zɪçər ,nɪçt fɪr 'frɑ:ŋ:kɛr:ç, fɪr ɪ'tʌ:ljən 'ʔɔ:x nɪçt', ʔʊnt ,gɑn:ʃ ʊnt 'gɑ:r nɪçt fɪr 'dɔ:ʧlɑnt'. 'mʊstərhaftəs 'dɔ:ʧ vɪrt ɪn 'k'ɛrnər 'dɔ:ʧfən 'ʃtɑt ɔ'dər 'lɑntʃɑft gə'spɾɑ:xən; 'ʔɔ:x dər gə'bɪldətə bɛr'li:nər, 'mɪn:çənər, 'vi:nər, 'frɑ:ŋ:kfʊrtər, 'k'œlnər ,ʃpɪçt 'ʔɑ:zɔ: nɪçt ʔɔ'nə 'vɛr:tərəs 'mʊstərhaft'. ʔɔ:f 'k'ɛrnən

'For Germany it is not true what Mr. Palmer noted for Japan and what Mr. Graham seems to assume mutatis mutandis for all languages, i.e. that the pronunciation of any educated individual may serve as a model. I doubt this for Great Britain, it is certainly not true for France, neither for Italy and least of all for Germany. Exemplary German is not spoken in any German city or region; neither the educated man from Berlin nor from Munich, Vienna, Frankfurt or Cologne can be considered to pronounce in an exemplary way'

Braun & Möbius (2022).

Introduction

Section 1: Word
Frequency
Distributions

Section 2: Simple
Measures

Section 3:
Quantitative Laws

Summary

References

Menzerath-Altmann Law

“The greater the whole the smaller its parts”. For example, as the **number of syllables** in words of a language increases, the *relative* number of **phonemes per syllable** decreases.

Note: Menzerath counted the number of words in a German lexicon with a given number of phonemes (“Lautzahl”) and a given number of syllables (“Silbenzahl”). We might expect that the number of phonemes increases linearly with the number of syllables (dashed lines), but this is not the case. Instead, the *relative* (i.e. average) number of phonemes per syllable decreases.

		SILBENZAHL										
		z	1	2	3	4	5	6	7	8	9	Summe
Lautzahl	n											
	22										<u>3</u>	3
	21								<u>2</u>	<u>1</u>	<u>1</u>	4
	20							<u>2</u>	<u>3</u>	<u>3</u>	0	8
	19							<u>2</u>	<u>3</u>	0	<u>1</u>	6
	18							<u>5</u>	<u>2</u>	<u>2</u>	<u>1</u>	10
	17						<u>3</u>	<u>10</u>	<u>11</u>	<u>1</u>		25
	16					<u>4</u>	<u>14</u>	<u>20</u>	8	4		50
	15					<u>6</u>	<u>44</u>	<u>29</u>	8			87
	14				<u>1</u>	<u>43</u>	<u>81</u>	<u>50</u>	4			179
	13				<u>4</u>	<u>99</u>	<u>154</u>	<u>52</u>	<u>1</u>			310
12			<u>28</u>	<u>268</u>	<u>200</u>	<u>31</u>					527	
11			<u>157</u>	<u>572</u>	<u>201</u>	<u>10</u>					940	
10		<u>25</u>	<u>566</u>	<u>896</u>	<u>167</u>	<u>3</u>					1657	
9		<u>91</u>	<u>1125</u>	<u>883</u>	<u>43</u>						2142	
8		<u>430</u>	<u>1893</u>	<u>643</u>	<u>12</u>						2978	
7	<u>2</u>	<u>1394</u>	<u>1840</u>	<u>204</u>	<u>1</u>						3441	
6	<u>69</u>	<u>1603</u>	<u>1150</u>	<u>22</u>							2799	
5	<u>444</u>	<u>1492</u>	<u>258</u>								2189	
4	<u>962</u>	<u>1256</u>	<u>7</u>								2225	
3	<u>645</u>	<u>101</u>									746	
2	<u>114</u>	<u>4</u>									118	
1	<u>9</u>										9	
Summe :		2245	6396	<u>6979</u>	3640	920	214	42	11	6	20453	

Vértés (1955).

Introduction

Section 1: Word Frequency Distributions

Section 2: Simple Measures

Section 3: Quantitative Laws

Summary

References

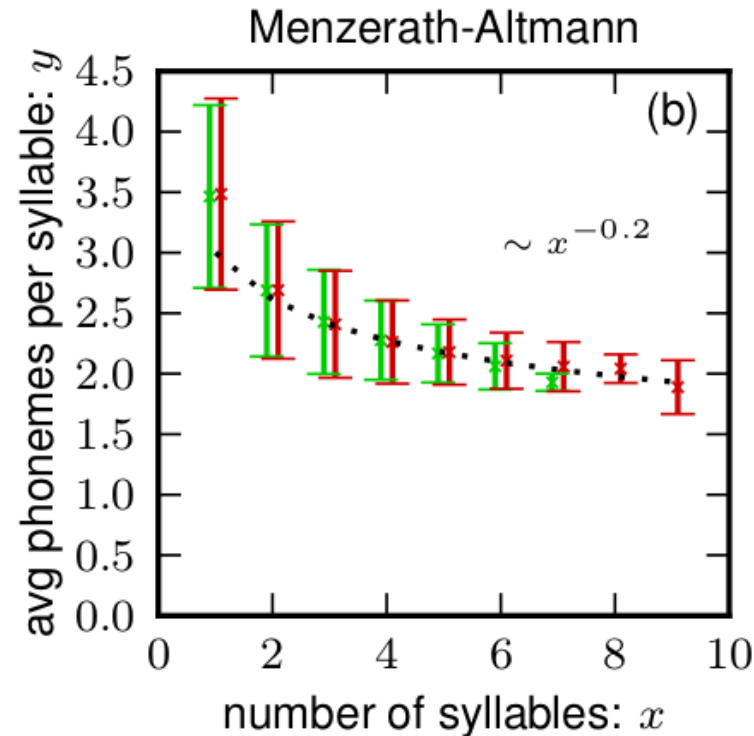
Menzerath-Altmann Law

The modern version of the law is formulated as:

$$y = \alpha x^\beta e^{-\gamma x}, \quad (3)$$

with y representing the size of the parts (e.g. in number of phonemes), x representing the length of the whole (e.g. number of syllables), and α , β , and γ are parameters (different from the ones in other laws of course).

Gerlach & Altmann (2015), p. 3.



green: Moby Dick in English,
red: Wikipedia in English.

Introduction

Section 1: Word
Frequency
Distributions

Section 2: Simple
Measures

Section 3:
Quantitative Laws

Summary

References

Heaps' Law (Herdan's Law)

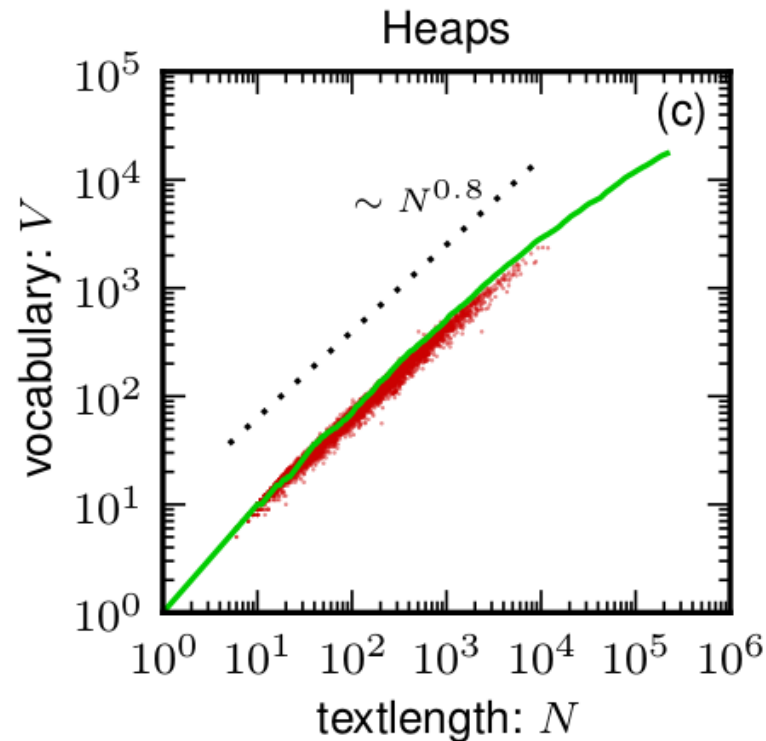
The number of *word types* in the vocabulary (V) grows with the number of *word tokens* in a text (N) according to

$$V \sim N^\alpha, \quad (4)$$

with $0 < \alpha < 1$.

Gerlach & Altmann (2015), p. 3-4.

Note: This does not hold for characters in writing, as there will be a finite number of characters in any writing system. Once this number is reached, V cannot grow further.



green: Moby Dick in English,
 red: Wikipedia in English.

Introduction

Section 1: Word Frequency Distributions

Section 2: Simple Measures

Section 3: Quantitative Laws

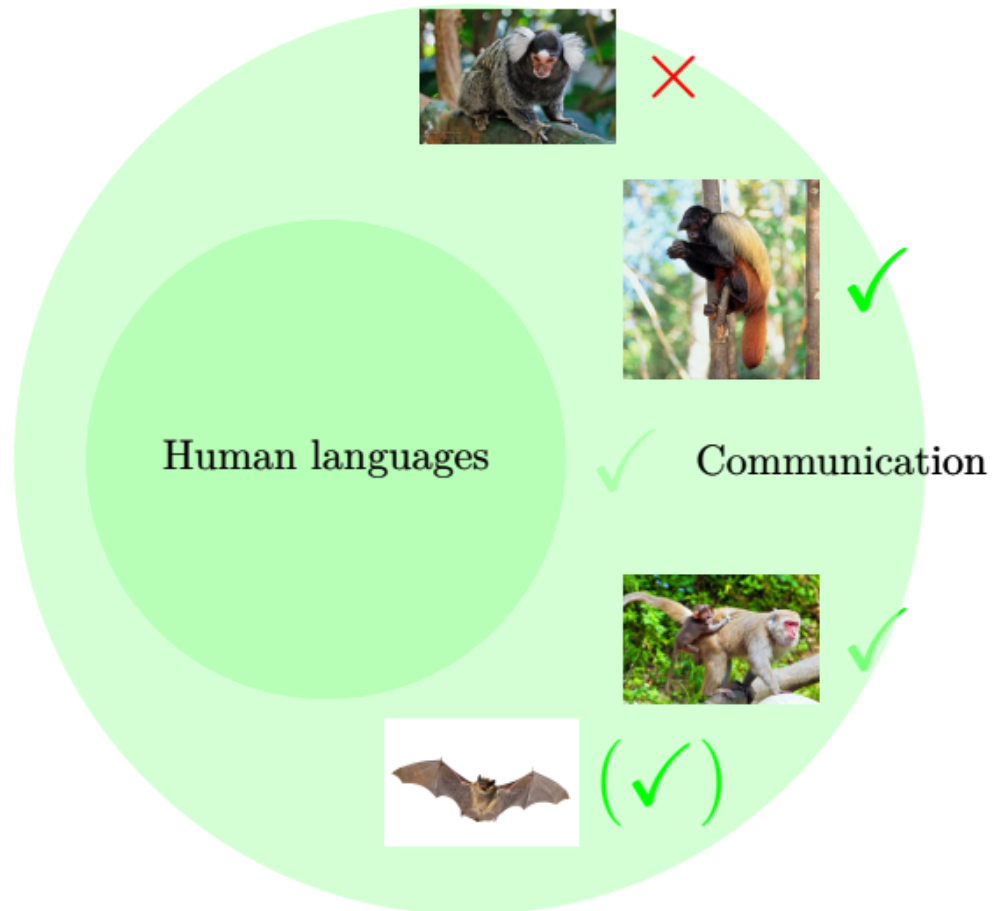
Summary

References



Current Research Question

Are “linguistic” laws also represented in other species?



Introduction

Section 1: Word
Frequency
Distributions

Section 2: Simple
Measures

Section 3:
Quantitative Laws

Summary

References



Summary



Summary

- ▶ **Word frequency distributions** are affected by a host of factors: *lexical typology, morphological typology, real world salience, translation, writing system, etc.*
- ▶ Differences in these distributions can be measured by simple measures such as **type-token-ratios** (TTR).
- ▶ More complicated models include **Zipf's law of word frequencies**.
- ▶ There are further **quantitative laws** found in natural languages: *Zipf's law of abbreviation, Menzerath-Altmann law, Heap's law.*
- ▶ It is a recurrent research question to what extent these are **specific to human language**.

Introduction

Section 1: Word
Frequency
Distributions

Section 2: Simple
Measures

Section 3:
Quantitative Laws

Summary

References



References



References

Bentz, C. (2018). *Adaptive languages: An information-theoretic account of linguistic diversity*. Berlin, Boston: De Gruyter Mouton.

Bentz, C. and Ferrer-i-Cancho, R. (2016). Zipf's law of abbreviation as a language universal. In: Bentz, C., Jäger, G. and Yanovich, I. (eds.) *Proceedings of the Leiden Workshop on Capturing Phylogenetic Algorithms for Linguistics*. University of Tübingen, online publication system.

Bentz, C., and Kiela, D. (2014). Zipf's law across languages of the world: towards a quantitative measure of lexical diversity. In *Evolution of Language: Proceedings of the 10th International Conference (EVLANG10)* (pp. 385-386).

Bentz, C. and Verkerk, A. (forthcoming). Sociolinguistic typology: complexification and simplification.

Braun, A. and Möbius, B. (2022). Armando de Lacerda and his contemporaries: Paul Menzerath. In *Proceedings of Fifth International Workshop on the History of Speech Communication Research (HSCR)*.

Bybee, J. (2006). From usage to grammar: the mind's response to repetition. *Language*.

Derungs, C., and Samardžić, T. (2018). Are prominent mountains frequently mentioned in text? Exploring the spatial expressiveness of text frequency. *International Journal of Geographical Information Science*, 32(5), 856-873.

Gerlach, M. and Altmann, E. (2015). Statistical laws in linguistics. *arXiv*.

Introduction

Section 1: Word
Frequency
Distributions

Section 2: Simple
Measures

Section 3:
Quantitative Laws

Summary

References



Gutierrez-Vasques, X., Bentz, C., and Samardžić, T. (2023). Languages Through the Looking Glass of BPE Compression. *Computational Linguistics*, pp. 1-59.

Nettle, D. and Romaine, S. (2000). *Vanishing voices: The extinction of the world's languages*. Oxford: Oxford University Press.

Mayer, T. and Cysouw, M. (2014). Creating a massively parallel Bible corpus. In *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC'14)*, pages 3158–3163, Reykjavik, Iceland. European Language Resources Association (ELRA).

Sproat, R. (2014). A statistical comparison of written language and nonlinguistic symbol systems. *Language*, 457-481.

Vértes, E. (1955). Review: Die Architektonik des deutschen Wortschatzes. In *Acta Linguistica Academiae Scientiarum Hungaricae*, Vol. 5, No. 3/4, pp. 415-433.

Wälchli, B. (2012). Indirect measurement in morphological typology. In A. Ender, A. Leemann and B. Wälchli (Ed.), *Methods in Contemporary Linguistics* (pp. 69-92). Berlin, Boston: De Gruyter Mouton.

Zipf, G. K. (1949). *Human behavior and the principle of least effort. An introduction to human ecology*. Cambridge: Addison-Wesley Press.

Zipf, G. K. (1965). *The psycho-biology of language. An introduction to dynamic philology*. Cambridge: MIT Press.

Introduction

Section 1: Word
Frequency
Distributions

Section 2: Simple
Measures

Section 3:
Quantitative Laws

Summary

References



Thank You.

Contact:

Faculty of Philosophy

General Linguistics

Dr. Christian Bentz

SFS Keplerstraße 2, Room 168

chris@christianbentz.de

Office hours:

During term: Wednesdays 10-11am

Out of term: arrange via e-mail