

How to count words?

In this lecture we want to count words and for this we have to ask ourselves what a word actually is? We will learn different methods to compare words and get an insight into the linguistic sub-discipline of lexicography and morphology. We will put this knowledge into a transducer that will enable us to normalize texts and gather statistics about words. Finally, we discuss how our solution is transferable to other languages, such as Chinese.

Text:

In this lecture **we** want to count words and for this **we** have to ask ourselves: what **are** words actually? **We** will learn different methods to compare words and get an insight into the linguistic sub-discipline of lexicography and morphology. **We** will put this knowledge into a transducer that will enable us to normalize texts and gather statistics about words. Finally, **we** discuss how our solution **is** transferable to other languages, such as Chinese.

- Should „**We**“ and „**we**“ count as the same word?
- Should „**is**“ and „**are**“ be considered equal?
- ...

Language:

„I do uh **main-** mainly business data processing.“

„Seuss’s **cat** in the hat is different from other **cats!**“

- „uh“: should we also count speech disfluencies?
- „**main-**“ How to count fragments?
- What about plural –s?

Fuzzy String Matching

Technique of finding strings that match a pattern approximately

https://en.wikipedia.org/wiki/Approximate_string_matching

- **Optical Character Recognition (OCR) errors:**

- **Spelling Errors:**

- upper / lower casing,
- Typing errors,
- ...

- **Phonetically ambiguous words: e.g. “to”, “too”, “two”**

- **Pronunciation complicated or transcription unclear:**

- “Supercalifragilisticexpialidocious”

Pronunciation (IPA): /,su:pər,kæli,frædʒɪ,lɪstɪk,ɛkspi,æli'doʊʃəs/

- Proper names: „Maier“, „Meier“, „Mayr“

Wä, g'il'mēsē 'wīlg'
laē äx'ēdxēs gālay

↓ OCR

ITä, g'il_mēsē \$wīlg_
laē ä_r_ēdvēs gālay

Example: Spelling Errors



Gierafe

Gieraffe

Girafe

Girafhe



Which version is „close“ to the correct *german* version (Giraffe)?

Example: Spelling Errors



Giraffe

Correct spelling with 7 characters

Gierafe

Error?

Example: Spelling Errors



Giraffe

Gierafe

Correct spelling with 7 characters

1 insertion („e“)

1 deletion („f“)

2 errors $2/7 = 0.286$

Example: Spelling Errors



Giraffe

Correct spelling with 7 characters

Gierafe

1 insertion („e“)

1 deletion („f“)

2 errors $2/7 = 0.286$

Gieraffe

1 insertion („e“)

1 error $1/7 = 0.143$

Example: Spelling Errors



Giraffe

Correct spelling with 7 characters

Gierafe

1 insertion („e“)

1 deletion („f“)

2 errors $2/7 = 0.286$

Gieraffe

1 insertion („e“)

1 error $1/7 = 0.143$

Girafe

1 too few („f“)

1 error $1/7 = 0.143$

Example: Spelling Errors



Giraffe

Correct spelling with 7 characters

Gierafe

1 insertion („e“)

1 deletion („f“)

2 errors $2/7 = 0.286$

Gieraffe

1 insertion („e“)

1 error $1/7 = 0.143$

Girafe

1 too few („f“)

1 error $1/7 = 0.143$

Girafhe

1 substitution
(„h“ instead of „f“)

1 error $1/7 = 0.143$

Example: Spelling Errors



Giraffe

Gierafe

Gieraffe

Girafe

Girafhe

Correct spelling with 7 characters

1 insertion („e“)

1 deletion („f“)

2 errors

$2/7 = 0.286$

1 insertion („e“)

1 error

$1/7 = 0.143$

1 too few („f“)

1 error

$1/7 = 0.143$

1 substitution

(„h“ instead of „f“)

1 error

$1/7 = 0.143$

„Edit distance“
or
„Levenshtein-Distance“

WER

How to search for similar strings?



Let (U, d) be a metric space, i.e. U be our „universe of objects“ and $d: U \times U \rightarrow \mathbb{R}^+$ a distance metric satisfying

- $d(x, y) = 0 \Leftrightarrow x = y$
- $d(x, y) = d(y, x)$
- $d(x, z) \leq d(x, y) + d(y, z)$

Let (U, d) be a metric space, i.e. U be our „universe of objects“ and $d: U \times U \rightarrow \mathbb{R}^+$ a distance metric satisfying

- $d(x, y) = 0 \Leftrightarrow x = y$
- $d(x, y) = d(y, x)$
- $d(x, z) \leq d(x, y) + d(y, z)$

Idea

Given a new query $q \in U$ and a maximum distance k , retrieve all strings in our vocabulary $V \subset U$ with a distance at most k from q , i.e.

output all $x^* \in V: d(x^*, q) \leq k$

- There are different edit distances for string sequences
- Not all edit distances satisfy the symmetry relation $d(x, y) = d(y, x)$ of a distance metric

https://en.wikipedia.org/wiki/Edit_distance

- **Three types of errors:**

- $I := \# \text{Insertions}$ („too much“)
- $D := \# \text{Deletions}$ („too few“)
- $S := \# \text{Substitutions}$ („confusion“)
- $N := \# \text{SymbolsOfCorrectString}$

- **Above metrics on word level => Word Error Rate**

$$WER = \frac{S + D + I}{N}$$

Input

```
X[1..M], Y[1..N]
```

```
// 1-indexed, of length m and n respectively
```

Initialization

```
d[0..M, 0..N] := zeros()
```

```
For all i: d[i,0] := i
```

```
For all j: d[0,j] := j
```



```
// set all elements in 0-indexed array to zero
```

Recurrence Relation

```
For j from 1 to N:
```

```
  For i from 1 to M:
```

$$d[i, j] := \min \begin{cases} d[i-1, j] + 1 & \text{// deletion} \\ d[i, j-1] + 1 & \text{// insertion} \\ d[i-1, j-1] + \begin{cases} 2; & \text{if } X[i] \neq Y[j] \\ 0; & \text{if } X[i] = Y[j] \end{cases} & \text{// substitution} \end{cases}$$

Termination:

```
d[N,M] is the distance
```

■ Fuzzy String Match:

Grapheme Sequence

TO
TOO
TWO

Phoneme Sequence

T UW1
T UW1
T UW1

} works.

Robert
Rupert

R AA1 B ER0 T
R UW1 P ER0 T } Does not work

Robert => Hash: R163
Rupert => Hash: R163

} Wie findet man diesen Hash?

■ Robert C. Russell and Margaret King Odell

■ Patented in 1918:

1. Retain the first letter of the name drop all other occurrences of a, e, i, o, u, y, h, w.
2. Replace consonants with digits as follows (after the first letter):
 1. b, f, p, v \rightarrow 1
 2. c, g, j, k, q, s, x, z \rightarrow 2
 3. d, t \rightarrow 3
 4. l \rightarrow 4
 5. m, n \rightarrow 5
 6. r \rightarrow 6
3. If two or more letters with the same number are adjacent in the original name (before step 1), only retain the first letter; also two letters with the same number separated by 'h' or 'w' are coded as a single number, whereas such letters separated by a vowel are coded twice. This rule also applies to the first letter.
4. If you have too few letters in your word that you can't assign three numbers, append with zeros until there are three numbers. If you have four or more numbers, retain only the first three.

1. **Text translated in tokens: Word segmentation**
2. **Normalisation: gather comparability**
 - Normalizing
 - Upper- and lower-casing
 - Morphology
 - Lemmatization/stemming
3. **Sentence Segmentation**

Tokens vs. Types

Distinguish two ways of talking about words

Token vs. Typen?



1 Kranich/Type



10 Kraniche/Tokens



1 Individuum or
„identity“



Token vs. Typen?



Beispiel: „HELLO“

#Tokens: 5

#Types: 4 (here: E, O, H, L,)

1 Kranich/Type



10 Kraniche/Tokens



1 Individuum or
„identity“



Token vs. Typen?



Beispiel: „HELLO“

#Tokens: 5

#Types: 4 (hier: E, O, H, L,)

Beispiel: „There are cars.“

#Tokens: 3

#Types: 3 (there, are, cars)



cars = car?
are = were = be = is?

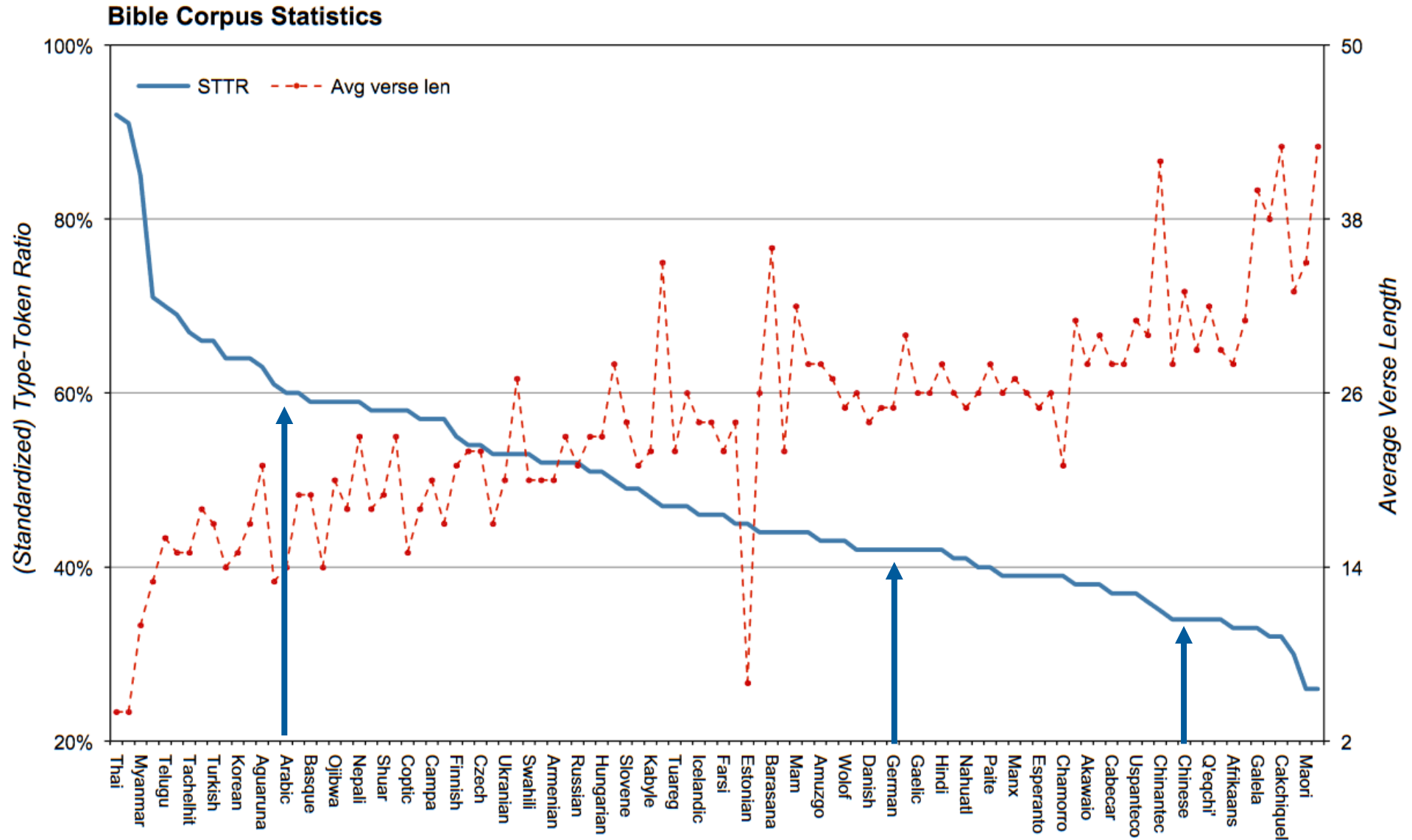
Type: an element of the vocabulary

Token: an instance of that type in running text

■ Church & Gale (1990): $|\text{Typen}| > O(|\text{Tokens}|^{0.5})$

| | $ \text{Tokens} $ | $ \text{Typen} := \text{Vokabular Größe}$ |
|---------------------------------|-------------------|--|
| Switchboard phone conversations | 2 400 000 | 20 000 |
| Shakespeare | 884 000 | 31 000 |
| Google n-gram | 1 Trillionen | 13 000 000 |

Typen-Token-Ration in verschiedenen Sprachen



Tokenization

Defining words

**Segmentation of a text into units on a word level,
aka „words“**

- **For German, English etc: ususally simply words separated by whitespaces**
- **But there are special cases**

„Finland’s capital“

What’re

I’m

isn’t

Hewlett-Packard

State-of-the-art

Lowercase

San Francisco

m.p.h., PhD

Finland, Finlands, Finland’s?

What are

i am

is not

HP, Hewlett Packard

state of the art

lower-case, lowercase, lower case

one token or two?

?

L'ensemble

L, L', Le?

L'ensemble

un ensemble

Lebensversicherungsgesellschaftsangesteller

⇒ Compound splitter required:

- Leben s
- versicherung s
- gesellschaft s
- angesteller

Slang in Japanese:

フォーチュン500社は情報不足のため時間あた\$500K(約6,000万円)
Katakana Hànzì Hiragana Kanji Romaji

Slang in Japanese:

フォーチュン500社は情報不足のため時間あた\$500K(約6,000万円)
Katakana Hànzì Hiragana Kanji Romaji

| | |
|---------------------|------------|
| 那是一句话。 | (Chinese) |
| それは一文です。 | (Japanese) |
| นั่นคือประโยค | (Thai) |
| 그것은 문장입니다. | (Korean) |
| This is a sentence. | (English) |

Segmentation into
words?

„Most common“: Max-Match Segmentation
Research: Neural nets for word segmentation

Max-Match Segmentation

Languages without „obvious“ word boundaries in grapheme sequences

莎拉波娃现在居住在美国东南部的佛罗里达

English: „Sharapova now lives in Florida in the southeast of the United States “

莎拉波娃现在居住在美国东南部的佛罗里达

Longest word in vocabulary? – no.

Vocabulary:

现在的
国东
在美
莎拉波娃
居住
南部
佛罗里达
佛罗里达
居住南部
...

莎拉波娃现在居住在美国东南部的佛罗里达



Longest word in vocabulary? – no.

Vokabular:

现在的
国的
在美国
莎拉波娃
居住
南部
佛罗里达
佛罗里达
居住南部
...

莎拉波娃现在居住在美国东南部的佛罗里达



Longest word in vocabulary? – no.

Vocabulary:

现在的
国的
在美国
莎拉波娃
居住
南部
佛罗里达
佛罗里达
居住南部
...

莎拉波娃现在居住在美国东南部的佛罗里达

莎拉波娃

Longest word in vocabulary? – yes.

莎拉波娃

Vocabulary:

现在的

的

国东

在美

莎拉波娃

居住

南部

佛罗里达

佛

罗

里

达

居

住

南

部

...

莎拉波娃现在居住在美国东南部的佛罗里达



Longest word in vocabulary? – no.

莎拉波娃

Vocabulary:

现在的

的

国东

在美

莎拉波娃

居住

南部

佛罗里达

佛

罗

里

达

居

住

南

部

...

莎拉波娃现在居住在美国东南部的佛罗里达

Longest word in vocabulary? – yes.

莎拉波娃 现在

Vocabulary:

现在

的

国东

在美

莎拉波娃

居住

南部

佛罗里达

佛

罗

里

达

居

住

南

部

...

莎拉波娃现在居住在美国东南部的佛罗里达

Longest word in vocabulary? – no.

莎拉波娃 现在

Vocabulary:

现在的
国的
东部
在美
莎拉波娃
居住
南部
佛罗里达
佛罗里达
居住南部
...

莎拉波娃现在居住在美国东南部的佛罗里达

Longest word in vocabulary? – yes.

莎拉波娃 现在 居住

Vocabulary:

现在的
国东
在美
莎拉波娃
居住
南部
佛罗里达
佛罗里达居住南部
...

莎拉波娃现在居住在美国东南部的佛罗里达

└.....

Longest word in vocabulary? – no.

莎拉波娃 现在 居住

Vocabulary:

现在

的

国东

在美

莎拉波娃

居住

南部

佛罗里达

佛

罗

里

达

居

住

南

部

...

莎拉波娃现在居住在美国东南部的佛罗里达

莎拉波娃 现在 居住 在美 国东 南部 的 佛罗里达

Pinyin: Shā lā bō wá xiànzài jūzhù zài měiguó dōngnán bù de fóluólidá

Thecatinthehat => the cat in the hat

Thetabledownthere => theta bled own there

(correct: the table down there)



Funktioniert nicht für Englisch, Deutsch, ...
Wir häufig mit Grammatiken gelöst.

Segmentierung ist aktives Forschungsfeld in allen Sprachen!

Normalization

Remove noise and other superfluous information, establish comparability.

U.S.A. vs. USA

GM vs. General Motors vs. general motors

Fed vs. fed

US vs. us <= context

Define equivalence classes of terms

Examples: Internet Slang



| Input | Output |
|---------|------------|
| 2moro | tomorrow |
| 2mrrw | tomorrow |
| 2morrow | tomorrow |
| 2mrw | tomorrow |
| tomrw | tomorrow |
| b4 | before |
| otw | On the way |

Examples: Noise



| Input | Output | word stem |
|----------------|---------|-----------|
| ..trouble.. | trouble | troubl |
| trouble< | trouble | troubl |
| trouble! | trouble | troubl |
| <a>trouble | trouble | troubl |
| 1.trouble | trouble | troubl |



We'll get to that in a minute!

- `[!"#$%&'()*+,-./:;<=>?@[\\]^_`{|}~]`
- **Space, line break**
- `<tr>`, `<a>`, `<p>`, ...

Capitalization

iS uPPeR AnD LoWEr CAsiNg ReaLLy IMportant FoR uNDeRStandAbiliTY?

| | |
|---|-------------------------------|
| Sentence start/Sentence case | General syntactic agreements |
| Munich, Audi, United States | Proper names |
| BMW, ICE, US | Abbreviations |
| easyJet A319, WikiWord, WikiCase, PhD, BSc., StGB, GmbH, TzBfG, macOS iPhone, BahnCard, RegionalExpress, InterCityExpress | „Marketing“ |
| I (in English) | Peculiarities of the language |
| ... | |

