

# Linguistic Annotation of Computer-Mediated Communication (not only) an Explorative Analysis

Kay-Michael Würzner

Alexander Geyken

Lothar Lemnitzer

Bryan Jurish

University of Potsdam

Berlin-Brandenburgische Akademie of Sciences and Humanities

# The Big Picture

- DERIC (Beißwenger et al. 2012)
  - Reference corpus for CMC
  - Representative and balanced
  - Joint project of University of Dortmund and BBAW

**Acquisition** Collecting data (in progress)

**Mark-Up** Standard format for different sources (i.e. TEI, Lemnitzer and Ermakova)

**Annotation** Adding linguistic information

## What we understand as *linguistic annotation*:

**Tokenizing** detection of word and sentence boundaries

**Lemmatization** identification of a word's base form and possible part of speech (PoS)

**PoS-Tagging** computation of a word's most probable PoS in context

→ Requirement for search engine indexing

# Tokenizing – Task

- Segmentation of continuous text into words (or *tokens*) and sentences
- Classification of certain classes of tokens
  - Abbreviations
  - Numbers
  - Special characters
  - Foreign alphabets
- Normalization of hyphenations

## Tokenizing – Challenges

- Reduced or inconsistent punctuation

Hallo haben sie fragen zum BAföG

- Special vocabulary

Wenn Ihre Antwort komplett ist, tippen Sie bitte am Ende ein  
\*E\*

- Sequences of special characters

????????? was das

- Addressing of chat partners

asset1, ...; keller > > ...; baloo: ...; @günni ... etc.

# Tokenizing – Approaches

- Rule-based
  - Set of regular expressions defining tokens
  - Construction of deterministic finite-state automaton
  - Usually done by a *scanner generator* (i.e. flex, re2c)

```
[A-z]+          {return WORD;}  
[0-9]+          {return INT;}  
. / [ \n\t] {return SB;}
```

- very fast but highly error-prone

## Tokenizing – Approaches II

- Statistical
  - Detection of likely token or sentence boundaries using context information
  - Kiss and Strunk (2006)
    - \* Unsupervised; trained on raw text
    - \* Initial, minimal rule-based step
    - \* Disambiguates · using an indepency assumption:

$$P(\cdot|w) = p = P(\cdot|\neg w)$$

- \* Additional features as word length and the number of internal periods to refine decision
- \* Used in OpenNLP

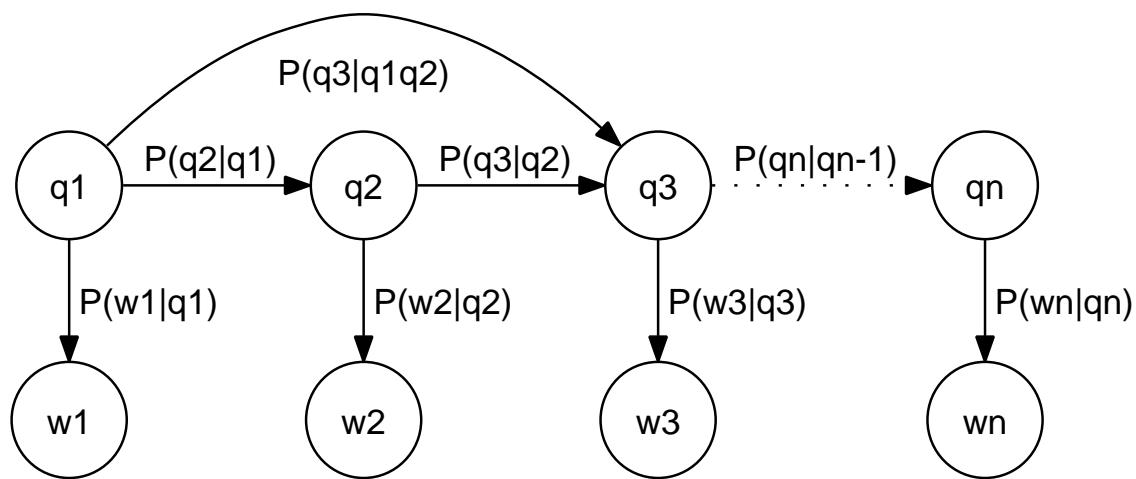
# Tokenizing – Approaches III

- Statistical
  - Laboreiro et al. (2010)
    - \* Especially dedicated to micro blogging
    - \* Supervised; trained on pre-tokenized text
    - \* Initial, minimal rule-based step
    - \* Training of a support vector machine classifier
    - \* Very promising results but **no end-of-sentence detection**

# Tokenizing – Our Approach

- Tokenizing as PoS tagging
- moot (Jurish 2003)
  - 2nd order Hidden Markov Model
  - *Observations* words; *states* tags
  - Viterbi optimization . . .
- Minimal rule set
- Aggregation of (pre-)tokens
- Small tagset

## Tokenizing – Our Approach II



$$\tau(w_{i\dots n}) = \arg \max_{q_{1\dots n} \in T^n} P(q_{1\dots n} | w_{1\dots n})$$

$$P(q_{1\dots n} | w_{1\dots n}) = \prod_{i=1}^n P(w_i | q_i) P(q_i | q_{i-2} q_{i-1})$$

# Tokenizing – Our Approach III

- Classification of pre-tokens wrt.
  - Class  $C = \{alpha, alpha-stopword, numeric, eos, dot, comma, quote, other\}$
  - Spelling  $S = \{upper, lower, caps, *\}$
  - Length  $L = \{1, \leq 3, \leq 5, long\}$
  - Abbreviation  $A = \{known, unknown\}$
- Classification of token state wrt.
  - Beginning of token  $bot \in \mathbb{B}$
  - Beginning of sentence  $bos \in \mathbb{B}$

# Tokenizing – Our Approach IV

- Text  $T \in K \times S \times L \times A$  (minus „impossible“ combinations)  
→ *Aggregation* of tokens
- Tag  $\in T \times \mathbb{B} \times \mathbb{B}$  (minus „impossible“ combinations)  
→ Tagset!

## Tokenizing – Example

Das	cls:alpha—sw_cas:upper_abbr:uk_len:3_bos:1_bot:1
Unternehmen	cls:alpha_cas:upper_abbr:uk_len:long_bos:0_bot:1
verkauft	cls:alpha_cas:lower_abbr:uk_len:long_bos:0_bot:1
er	cls:alpha—sw_cas:lower_abbr:uk_len:3_bos:0_bot:1
1984	cls:num_cas:*_abbr:uk_len:5_bos:0_bot:1
fuer	cls:alpha—sw_cas:lower_abbr:uk_len:3_bos:0_bot:1
2,5	cls:num_cas:*_abbr:uk_len:3_bos:0_bot:1
Milliarden	cls:alpha_cas:upper_abbr:uk_len:long_bos:0_bot:1
.	cls:dot_cas:*_abbr:uk_len:1_bos:0_bot:1

# Tokenizing – Evaluation

- Dortmund chat corpus (Storrer and Beißwenger)
  - Base corpus containing 140,000 posts ( $\approx 10^6$  tokens)
  - Freely available release corpus ( $\approx 500,000$  tokens)
  - Chats from different contexts: university, media,  
„Plauder-Chats“
- Manually tokenized subset (ca. 100,000 tokens)
- Includes different chat scenarios (1 : 1, 1 : n, m : n)

## Tokenizing – Evaluation II

	manual		automatic	
+ Items	105,279	(100.0 %)	103,290	(100.0 %)
– Match	100,337	(95.3 %)	100,337	(97.1 %)
– NoMatch	4,942	(4.7 %)	2,953	(2.9 %)
<hr/>				
+ Items	39,747	(100.0 %)	39,158	(100.0 %)
– Match	38,515	(96.9 %)	38,515	(98.4 %)
– NoMatch	1,232	(3.1 %)	643	(1.6 %)

- Still losing 34 % of the unmarked sentence boundaries
- Training on Tiger: 92 % correct for CMC and 99 % for Tiger

# Morphological Analysis – Task

- Categorisation of words regarding their *possible* PoS
- Mapping of words to their base form (i.e. *lemma*)
- Due to word formation processes, infinite number of possible words

## Morphological Analysis – Challenges

- Writing errors

dadruch

- Missing whitespace

3,9Millionen

- Dialect-specific words

zwüsched, zwüschađ, zwüschet

- Complex bounded compounds

26/26-Lösung, 2-Fach-BA

# Morphological Analysis – Approach

- Finite-state morphology
  - Large lexicon (word, category, inflection class)
  - Productive rules to mimic word formation processes
  - Implemented with (weighted) finite-state transducers
  - Weights as complexity estimate of the word formation
    - NN · NN → NN ⟨10⟩
  - Beesley and Karttunen (2003), Schmid et al. (2004), Geyken and Hanneforth (2006)

## Morphological Analysis – Example

```
[1]> apply "Telekommunikation"
```

```
Telekommunikation[NN SemClass=abstr Gender=fem Number=sg Case=*)
```

```
tele/K+Kommunikation[NN SemClass=abstr Gender=fem Number=sg  
Case=*) <5>
```

```
Tele/N#Kommunikation[NN SemClass=abstr Gender=fem Number=sg  
Case=*) <10>
```

...

```
Telekom/PN#Muni/N#Kat/N#lon[NN SemClass=k_g_dingnat  
Gender=neut Number=sg Case=nom_acc_dat] <30>
```

...

## Morphological Analysis – Evaluation

- No gold standard for standard or non-standard text available
- Alternatively, evaluation of coverage of TAGH (Geyken and Hanneforth 2006)

	Token	Types
<i>known</i>	84,599	10,346
<i>unknown</i>	9,372	5,459
Coverage %	88.92	47.24

- Recognition rate of TAGH for modern newspaper text: 99 %

## Morphological Analysis – Improvements

- Error tolerant morphological analysis: Jurish (2011)
  - Developed for historical texts
  - Uses phonological similarity and edit distance

	Token	Types
<i>known</i>	84,599	10,346
<i>unknown</i>	9,372	5,459
<i>fixed</i>	7,183	3,984
Coverage %	97.41	85.74

- Historical model applied to test corpus without adjustments

# Morphological Analysis – Top 10 Improvements

Type	Frequency	Type	Frequency
<i>hi</i>	98	<i>gibts</i>	33
<i>nich</i>	51	<i>kinder</i>	32
<i>nix</i>	46	<i>musik</i>	31
<i>leute</i>	42	<i>z.B.</i>	29
<i>adelheid</i>	40	<i>tach</i>	27

# Part-of-Speech Tagging – Task and Approaches

- Disambiguation of possible PoS of a word in context
- (partly) rule-based (Brill 1995) vs. fully statistical (Brants 2000)
- Supervised (trained on manually tagged text)
- Training materials available in the form of tree banks (leaves and pre-terminals)

# Part-of-Speech Tagging – Challenges

- Mimicry of spoken language
- Ellipses
- Single word sentences

Geeehts?

Ich!!!

Ein sehr starkes Schlafmittel...

# Part-of-Speech Tagging – Approach

- moot (Jurish, 2003)
- Trained on Tiger
- 2nd order *Hidden Markov Modell*: *observations* words, *states* PoS tags
- Viterbi maximization over PoS tag sequences for a sentence
- Optional restriction to categories assigned by the morphology

# Part-of-Speech Tagging – Example (Jurish 2003)

Input:	Linda	wird	die	Mannschaft	verstärken	.
<b>Morphological Analysis:</b>	$\left\{ \begin{array}{l} \text{NE.} \textit{first}, \\ \text{NE.} \textit{last} \end{array} \right\}$	$\left\{ \begin{array}{l} \text{VAFIN.} \textit{3rd.sg.pres}, \\ \text{VVFIN.} \textit{3rd.sg.pres} \end{array} \right\}$	$\left\{ \begin{array}{l} \text{ART.} \textit{sg.nom.fem}, \\ \vdots \\ \text{PDS.} \textit{nom.sg.fem}, \\ \vdots \\ \text{PRELS.} \textit{acc.pl} \end{array} \right\}$	$\left\{ \begin{array}{l} \text{NN.} \textit{masc.sg.nom}, \\ \vdots \\ \text{NN.} \textit{fem.sg.*} \end{array} \right\}$	$\left\{ \begin{array}{l} \text{VFIN.} \textit{1st.pl.pres}, \\ \vdots \\ \text{VVINF} \end{array} \right\}$	$\{ \$\cdot \}$
<b>Tag Extraction:</b>	{NE}	$\left\{ \begin{array}{l} \text{VVFIN}, \\ \text{VAFIN} \end{array} \right\}$	$\left\{ \begin{array}{l} \text{ART}, \\ \text{PDS}, \\ \text{PRELS} \end{array} \right\}$	{NN}	$\left\{ \begin{array}{l} \text{VFIN}, \\ \text{VVINF} \end{array} \right\}$	$\{ \$ \}$
<b>Disambiguation:</b>	NE	VAFIN	ART	NN	VVINF	\$.

## Part-of-Speech Tagging – Evaluation

- $\approx 10\%$  of the test corpus manually annotated
- Not (yet) enough to train a model but to evaluate existing tools
- moot vs. TreeTagger

	moot	moot <sub>tagh</sub>	moot <sub>fixed</sub>	TreeTagger
<i>match</i>	7,424	6,999	7,948	7,673
<i>nomatch</i>	2,908	3,330	2,386	2,661
Correctness %	71.9	67.8	76.9	74.3

- `tree-tagger -cap-heuristics -sgml -eos-tag "<eos>"`

# Summary

- Promising results for tokenizing
- Morphological analysis yet to be evaluated
- PoS tagging disastrous: better models and pre-processing needed

# Many thanks for your attention

Big shout to Gabriella Pein, Maria Ermakova, Michael Beißwenger,  
Julian Heister und Frank Wiegand