

# Semantic Argument Classification

28. Januar 2015

Julian Baumann, Kevin Decker, Maximilian Müller-Eberstein

Institut für Computerlinguistik  
Universität Heidelberg



UNIVERSITÄT  
HEIDELBERG  
ZUKUNFT  
SEIT 1386



# Gliederung

Semantic Argument  
Classification

Julian Baumann, Kevin  
Decker, Maximilian  
Müller-Eberstein

Problemstellung

Daten

Problemstellung

Umsetzung

Features

Featureextraktion

Schwierigkeiten

Experimente

Ausblick

Literatur

Problemstellung

Daten

Problemstellung

Umsetzung

Features

Featureextraktion

Schwierigkeiten

Experimente

Ausblick

Literatur

Referenzen



# Semantic Argument Classification

## Semantic Argument Classification

Julian Baumann, Kevin  
Decker, Maximilian  
Müller-Eberstein

### Problemstellung

2

#### Daten

Problemstellung

#### Umsetzung

Features

Featureextraktion

Schwierigkeiten

#### Experimente

#### Ausblick

#### Literatur

#### Referenzen

## Was ist Semantic Argument Classification?

- ▶ Zuweisung bestimmter Rollen in einem Satz ⇒ „Wer tut wem was an?“
- ▶ It operates stores mostly in Iowa and Nebraska
- ▶ [Arg0 *It*][Pred *operates*][Arg1 *stores*][ArgLoc *mostly in Iowa and Nebraska*]



# Daten

## Semantic Argument Classification

Julian Baumann, Kevin  
Decker, Maximilian  
Müller-Eberstein

### Problemstellung

### Daten

3

Problemstellung

### Umsetzung

Features

Featureextraktion

Schwierigkeiten

### Experimente

Ausblick

Literatur

Referenzen

- ▶ NLTK
- ▶ PropBank



## Semantic Argument Classification

Julian Baumann, Kevin Decker, Maximilian Müller-Eberstein

### Problemstellung

#### Daten

4

Problemstellung

#### Umsetzung

Features

Featureextraktion

Schwierigkeiten

#### Experimente

Ausblick

Literatur

Referenzen

- ▶ versucht generalisierte Argumente zu verwenden → Parser
- ▶ Argumentrollen sind für jedes Verb in Frames organisiert → weniger spezifisch
- ▶ ARG0 = Proto-Agent
- ▶ ARG1 = Proto-Patient
- ▶ ARG2-ARG5 = Argumente mit steigender Intensität



## Semantic Argument Classification

Julian Baumann, Kevin Decker, Maximilian Müller-Eberstein

### Problemstellung

#### Daten

5

Problemstellung

#### Umsetzung

Features

Featureextraktion

Schwierigkeiten

#### Experimente

Ausblick

Literatur

Referenzen

- ▶ [ARG0 *He*][Predicate *wrote*][ARG1 *a book*]
- ▶ [ARG0 *He*][Predicate *wrote*][ARG2 *about them*]
- ▶ [ARG0 *He*][Predicate *wrote*][ARG3 *a book for his children*]
- ▶ → verschieden Rollen



## Semantic Argument Classification

Julian Baumann, Kevin  
Decker, Maximilian  
Müller-Eberstein

### Problemstellung

### Daten

6

Problemstellung

Umsetzung

Features

Featureextraktion

Schwierigkeiten

Experimente

Ausblick

Literatur

Referenzen

- ▶ Subkorpus aus WSJ und Brown Corpus, bestehend aus ungefähr 1.000.000 Wörtern
- ▶ 112.917 Prädikat-Argument Strukturen annotiert nach PropBank-Annotationsschema
- ▶ 292.975 Instanzen
- ▶ wsj/00/wsj\_0001.mrg 1 10 gold publish.01 p—a 10:0-rel 11:0-ARG0



# Penn Treebank

## Semantic Argument Classification

Julian Baumann, Kevin  
Decker, Maximilian  
Müller-Eberstein

## Problemstellung

## Daten

Problemstellung

## Umsetzung

Features

Featureextraktion

Schwierigkeiten

## Experimente

## Ausblick

## Literatur

## Referenzen

7

```
((S
  (NP-MNR-SBJ
    (NP (DT The) (NN way) )
    (SBAR
      (WHADVP-1 (-NONE- 0) )
      (S
        (NP-SBJ (NNP MacArthur) )
        (VP (VBD said)
          (NP (PRP his) (NN line) )
          (ADVP-MNR (-NONE- *T*-1) ))))
    (: - -)
```

23





# Features

## Semantic Argument Classification

Julian Baumann, Kevin Decker, Maximilian Müller-Eberstein

Problemstellung

Daten

Problemstellung

Umsetzung

Features

Featureextraktion

Schwierigkeiten

Experimente

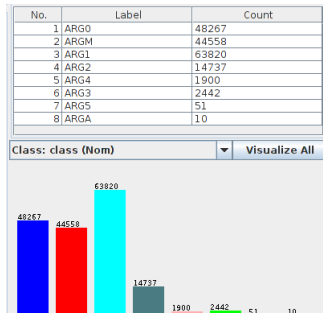
Ausblick

Literatur

Referenzen

8

- ▶ Predicate
- ▶ Path
- ▶ Phrase Type
- ▶ Position
- ▶ Voice





# Predicate

## Semantic Argument Classification

Julian Baumann, Kevin  
Decker, Maximilian  
Müller-Eberstein

Problemstellung

Daten

Problemstellung

Umsetzung

Features

Featureextraktion

Schwierigkeiten

Experimente

Ausblick

Literatur

Referenzen

9

- ▶ lemmatisierte Prädikat
- ▶ 3966 distinct

23



# Path

## Semantic Argument Classification

Julian Baumann, Kevin Decker, Maximilian Müller-Eberstein

## Problemstellung

## Daten

## Problemstellung

## Umsetzung

## Features

## Featureextraktion

## Schwierigkeiten

## Experimente

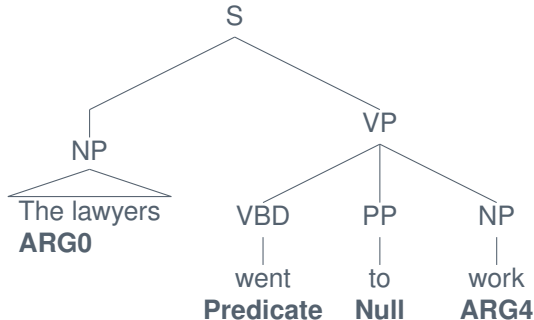
## Ausblick

## Literatur

## Referenzen

10

- ▶ beschreibt Pfad zwischen ARG und Predicate
- ▶ vereinfacht z.B. NP-SBJ  $\rightarrow$  NP
- ▶ extrahiert über Lowest Common Ancestor
- ▶ beispielsweise: NP $\uparrow$ S $\downarrow$ VP $\downarrow$ VBD
- ▶ 41737 distinct



23



# Phrase Type

## Semantic Argument Classification

Julian Baumann, Kevin Decker, Maximilian Müller-Eberstein

- beschreibt die Kategorie des Argument
- z.B: NP, MD, PP, SBAR
- 65 distinct feature values

Problemstellung

Daten

Problemstellung

Umsetzung

Features

Featureextraktion

Schwierigkeiten

Experimente

Ausblick

Literatur

Referenzen

11

No.	Label	Count
1	NP	97599
2	MD	6458
3	PP	29042
4	NN	1311
5	ADVP	10080
6	S	9957
7	-NONE-	1596
8	VBG	141
9	SBAR	10444

Class: class (Nom) Visualize All



# Position

## Semantic Argument Classification

Julian Baumann, Kevin Decker, Maximilian Müller-Eberstein

## Problemstellung

## Daten

Problemstellung

## Umsetzung

### Features

Featureextraktion

Schwierigkeiten

## Experimente

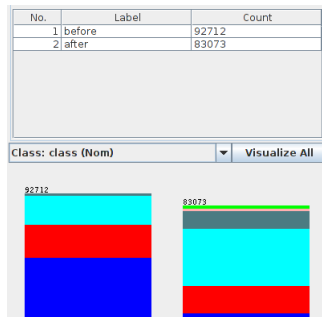
## Ausblick

## Literatur

## Referenzen

12

- Beschreibt, ob das Argument vor oder nach dem Prädikat steht
- Binäres Feature



23



# Voice

## Semantic Argument Classification

Julian Baumann, Kevin Decker, Maximilian Müller-Eberstein

## Problemstellung

## Daten

## Problemstellung

## Umsetzung

## Features

## Featureextraktion

## Schwierigkeiten

## Experimente

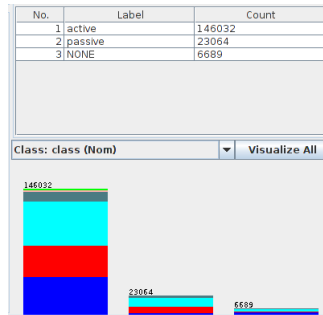
## Ausblick

## Literatur

## Referenzen

13

- ▶ gibt an, ob das Prädikat aktiv oder passiv ist
- ▶ größtenteils annotiert
- ▶ 3 distinct feature values: active, passive, unknown





# Featureextraktion

## Semantic Argument Classification

Julian Baumann, Kevin  
Decker, Maximilian  
Müller-Eberstein

### Problemstellung

### Daten

Problemstellung

### Umsetzung

Features

**Featureextraktion**

Schwierigkeiten

### Experimente

### Ausblick

### Literatur

### Referenzen

14

```
featureList = [...] # zu extrahierende Features
extArgList = []
for pbInstance in pbInstances :
    for pbArg in pbInstance.arguments :
        for feature in featureList :
            extArgList.append(extFeature(feature, pbArg,
pbInstance))
# write features to file in ARFF
```

23



# Featureextraktion

## Semantic Argument Classification

Julian Baumann, Kevin Decker, Maximilian Müller-Eberstein

## Problemstellung

## Daten

Problemstellung

## Umsetzung

Features

## Featureextraktion

Schwierigkeiten

## Experimente

## Ausblick

## Literatur

## Referenzen

15

wsj/00/wsj\_0041.mrg 38 14 gold talk.01 vn-3a 0:1-ARGM-ADV 12:1-ARG0 14:0-rel 15:1-ARG1-about

ARGInstanceBuilder

```
(S  
  (PP-LOC  
    (IN Against)  
    (NP  
      (NP (DT a) (NN shot))  
      (PP (IN of) (NP (NNP Monticello)))  
      (VP  
        (VBN superimposed)  
        (NP (-NONE- *))  
        (PP-CLR (IN on) (NP (DT an) (JJ American) (NN flag))))))  
  (, ,)  
  (NP-SBJ (DT an) (NN announcer))  
  (VP  
    (VBZ talks)  
    (PP-CLR  
      (IN about)  
      (NP  
        (NP (DT the) (`` ``) (JJ strong) (NN tradition))
```

ARGInstance processed features

ARFFDocument attributes, data





# ARFF

## Semantic Argument Classification

Julian Baumann, Kevin Decker, Maximilian Müller-Eberstein

### Problemstellung

### Daten

Problemstellung

### Umsetzung

Features

### Featureextraktion

Schwierigkeiten

### Experimente

### Ausblick

### Literatur

### Referenzen

16

@relation SAC\_All

@attribute predicate {join,publish,name,use, make, cause, ...}

@attribute phraseType {NP, MD, PP, NN, ADVP, S, ...}

@attribute position {before, after}

@attribute path {NP^S!VP!VP, MD^VP^S!VP!VP,...}

@attribute voice {active, passive, NONE}

@attribute class {ARG0, ARGM, ARGa, ARG1, ...}

@data

join, NP, before, NP^S!VP!VP, active, ARG0

join, MD, before, MD^VP^S!VP!VP, active, ARGM

join, NP, after, NP^VP^VP^S!VP!VP, active, ARG1

join, PP, after, PP^VP^VP^S!VP!VP, active, ARGM

join, NP, after, VP^VP^S!VP!VP, active, ARGM

23



# Schwierigkeiten

## Semantic Argument Classification

Julian Baumann, Kevin  
Decker, Maximilian  
Müller-Eberstein

### Problemstellung

### Daten

### Problemstellung

### Umsetzung

### Features

### Featureextraktion

### Schwierigkeiten

### Experimente

### Ausblick

### Literatur

### Referenzen

17

- ▶ PropBankChain- und PropBankSplitTreePointer
- ▶ Verwendung einer externen PennTreeBank
- ▶ einige Feature (bsp. path) nehmen sehr viele Werte an
- ▶ Speicherbedarf der Algorithmen J48 und SVM

23



# Setup

## Semantic Argument Classification

Julian Baumann, Kevin  
Decker, Maximilian  
Müller-Eberstein

Problemstellung

Daten

Problemstellung

Umsetzung

Features

Featureextraktion

Schwierigkeiten

Experimente

Ausblick

Literatur

Referenzen

18

- ▶ 60% train, 20% dev, 20% test
- ▶ Baseline: ZeroR
- ▶ Naive Bayes, j48 tree, (*libSVM*)
- ▶ bisher: Training auf train, Evaluierung mit dev

23



# Ergebnisse

## Semantic Argument Classification

Julian Baumann, Kevin  
Decker, Maximilian  
Müller-Eberstein

### Problemstellung

### Daten

### Problemstellung

### Umsetzung

### Features

### Featureextraktion

### Schwierigkeiten

### Experimente

### Ausblick

### Literatur

### Referenzen

19

	Precision	Recall	F-Measure
<i>Baseline</i>	<i>0.132</i>	<i>0.364</i>	<i>0.194</i>
Naive Bayes	0.771	0.778	0.770
j48 Tree	0.784	0.786	0.781

23



# Feature Evaluation

## Semantic Argument Classification

Julian Baumann, Kevin  
Decker, Maximilian  
Müller-Eberstein

### Problemstellung

#### Daten

##### Problemstellung

##### Umsetzung

##### Features

##### Featureextraktion

##### Schwierigkeiten

#### Experimente

##### Ausblick

##### Literatur

##### Referenzen

20

	a	b	c	d	e	f	g	h	<- classification
	14497	348	361	272	0	4	0	1	a = ARG0
	291	11394	2143	1009	189	48	0	1	b = ARG1
	3119	707	17064	375	31	27	0	0	c = ARG2
	180	792	1854	2163	29	23	0	0	d = ARG3
	2	217	23	141	379	3	0	0	e = ARG4
	37	289	144	147	170	99	0	0	f = ARG5
	0	13	0	1	0	0	3	0	g = ARG6
	5	0	0	0	0	0	0	0	h = ARG7

23



# Confusion Matrix (Naive Bayes)

Semantic Argument  
Classification

Julian Baumann, Kevin  
Decker, Maximilian  
Müller-Eberstein

Problemstellung

Daten

Problemstellung

Umsetzung

Features

Featureextraktion

Schwierigkeiten

Experimente

Ausblick

Literatur

Referenzen

21

	Precision	Recall	F-Measure	F-Measure Ch
<i>All Features</i>	<i>0.771</i>	<i>0.778</i>	<i>0.770</i>	<i>0</i>
- voice	0.748	0.754	0.745	-0.025
- path	0.778	0.783	0.776	<b>+0.006</b>
- phraseType	0.735	0.747	0.733	-0.037
- position	0.758	0.773	0.757	-0.013
-predicate	0.717	0.732	0.716	<b>-0.54</b>

23



# ToDo

## Semantic Argument Classification

Julian Baumann, Kevin  
Decker, Maximilian  
Müller-Eberstein

### Problemstellung

### Daten

Problemstellung

### Umsetzung

Features

Featureextraktion

Schwierigkeiten

### Experimente

### Ausblick

### Literatur

### Referenzen

- ▶ Path Feature überarbeiten
- ▶ HeadWord Feature implementieren
- ▶ genauere Evaluation
- ▶ *SVM?*
- ▶ Abschlussbericht schreiben

22

23



# Quellen

## Semantic Argument Classification

Julian Baumann, Kevin Decker, Maximilian Müller-Eberstein

Problemstellung

Daten

Problemstellung

Umsetzung

Features

Featureextraktion

Schwierigkeiten

Experimente

Ausblick

Literatur

Referenzen

23

- [1] Omri Abend und Roi Reichart. *Unsupervised Argument Identification for Semantic Role Labeling*.
- [2] Jean Carletta. "Assessing agreement on classification tasks: the kappa statistic". In: *Computational Linguistics* (1996), S. 249–254.
- [3] Daniel Gildea. "Automatic labeling of semantic roles". In: *Computational Linguistics* 28 (2002), S. 245–288.
- [4] Alessandro Moschitti und Cosmin Adrian Bejan. "A Semantic Kernel for Predicate Argument Classification". In: *IN CONLL 2004*. 2004, S. 17–24.
- [5] Sameer Pradhan u. a. *Support Vector Learning for Semantic Argument Classification*. 2005.

23



Vielen Dank für Eure Aufmerksamkeit!  
Noch Fragen?



UNIVERSITÄT  
HEIDELBERG  
ZUKUNFT  
SEIT 1386