

Coreference Resolution

Maschinelles Lernverfahren zur Koreferenz-Auflösung

Nikolina Koleva

UdS

12. Mai 2011

Womit befasst sich die Koreferenz-Auflösung?

Referieren zwei sprachliche Ausdrücke auf das selbe Individuum?

Referieren zwei sprachliche Ausdrücke auf das selbe Individuum?

- **Bezugswort**: worauf referiert wird
- **Anapher**: der referierende Ausdruck

Referieren zwei sprachliche Ausdrücke auf das selbe Individuum?

- **Bezugswort**: worauf referiert wird
- **Anapher**: der referierende Ausdruck

"Alice spielt draußen. Sie genießt das schöne Wetter."

Referieren zwei sprachliche Ausdrücke auf das selbe Individuum?

- **Bezugswort**: worauf referiert wird
- **Anapher**: der referierende Ausdruck

"Alice spielt draußen. Sie genießt das schöne Wetter."

Koreferenzkette: **Alice-Sie**

Warum ist Koreferenz-Auflösung wichtig?

Warum ist Koreferenz-Auflösung wichtig?

- 1 zur Diskursanalyse

Warum ist Koreferenz-Auflösung wichtig?

- 1 zur Diskursanalyse
- 2 allgemein zum Sprachverstehen

Warum ist Koreferenz-Auflösung wichtig?

- 1 zur Diskursanalyse
- 2 allgemein zum Sprachverstehen
- 3 wichtige Teilaufgabe in NLP-Systeme
→ insbesondere in Informationsextraktion(IE) Systeme

Der Ansatz von Soon et al.

- 1 Bestimmung der Markables (potenzielle Bezugswörter oder Anapher)

- 1 Bestimmung der Markables (potenzielle Bezugswörter oder Anapher)
- 2 Festlegung ihrer Attributvektoren

- 1 Bestimmung der Markables (potenzielle Bezugswörter oder Anapher)
- 2 Festlegung ihrer Attributvektoren
- 3 Generieren der Trainingsinstanzen

- 1 Bestimmung der Markables (potenzielle Bezugswörter oder Anapher)
- 2 Festlegung ihrer Attributvektoren
- 3 Generieren der Trainingsinstanzen
- 4 Trainieren des Modells

- 1 Bestimmung der Markables (potenzielle Bezugswörter oder Anapher)
- 2 Festlegung ihrer Attributvektoren
- 3 Generieren der Trainingsinstanzen
- 4 Trainieren des Modells
- 5 Bilden der Koreferenzketten für das Testdokument

- sprachlicher Ausdruck, auf den sich andere sprachliche Ausdrücke beziehen können

- sprachlicher Ausdruck, auf den sich andere sprachliche Ausdrücke beziehen können
- oft: Markable Grenzen \neq Token Grenzen

- sprachlicher Ausdruck, auf den sich andere sprachliche Ausdrücke beziehen können
- oft: Markable Grenzen \neq Token Grenzen
-

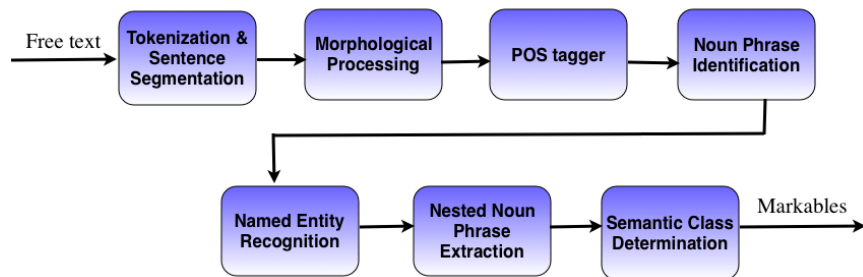
Separately, Clinton transition officials said that Frank Newman, 50, vice chairman and chief financial officer of BankAmerica Corp. , is expected to be nominated as assistant Treasury secretary for domestic finance.

Bestimmung der Markables

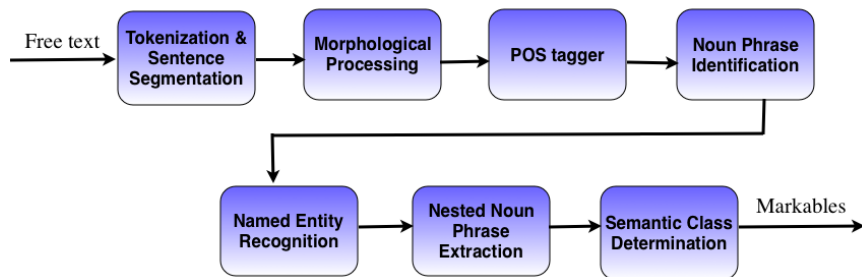
- mittels einer Pipeline von NLP module

- mittels einer Pipeline von NLP module

- mittels einer Pipeline von NLP module



- mittels einer Pipeline von NLP module



- getestet gegen 100 Trainingsdokumente → 85 % Übereinstimmung

12 Attribute der Markables

12 Attribute der Markables

i ist das potenzielle Bezugswort; **j** ist die Anapher

12 Attribute der Markables

i ist das potenzielle Bezugswort; **j** ist die Anapher

- *DIST*: Distanz zwischen i und j (0,1,2,3...)

12 Attribute der Markables

i ist das potenzielle Bezugswort; **j** ist die Anapher

- *DIST*: Distanz zwischen i und j (0,1,2,3...)
- *I_PRONOUN*: ist i ein Pronomen (true, false)

12 Attribute der Markables

i ist das potenzielle Bezugswort; **j** ist die Anapher

- *DIST*: Distanz zwischen i und j (0,1,2,3...)
- *I_PRONOUN*: ist i ein Pronomen (true, false)
- *J_PRONOUN*: ist j ein Pronomen (true, false)

12 Attribute der Markables

i ist das potenzielle Bezugswort; *j* ist die Anapher

- *DIST*: Distanz zwischen *i* und *j* (0,1,2,3...)
- *I_PRONOUN*: ist *i* ein Pronomen (true, false)
- *J_PRONOUN*: ist *j* ein Pronomen (true, false)
- *STR_MATCH*: stimmen die Strings *i* und *j* überein (true, false)
→ *a boy* und *this boy*

12 Attribute der Markables

i ist das potenzielle Bezugswort; **j** ist die Anapher

- *DIST*: Distanz zwischen i und j (0,1,2,3...)
- *I_PRONOUN*: ist i ein Pronomen (true, false)
- *J_PRONOUN*: ist j ein Pronomen (true, false)
- *STR_MATCH*: stimmen die Strings i und j überein (true, false)
→ *a boy* und *this boy*
- *DEF_NP*: ist j definite Nominalphrase (true, false)
→ *the car/girl* etc.

12 Attribute der Markables

i ist das potenzielle Bezugswort; **j** ist die Anapher

- *DIST*: Distanz zwischen i und j (0,1,2,3...)
- *I_PRONOUN*: ist i ein Pronomen (true, false)
- *J_PRONOUN*: ist j ein Pronomen (true, false)
- *STR_MATCH*: stimmen die Strings i und j überein(true, false)
→ *a boy* und *this boy*
- *DEF_NP*: ist j definite Nominalphrase (true, false)
→ *the car/girl* etc.
- *DEM_NP*: ist j demonstrative Nominalphrase (true, false)
→ *this car/girl* etc.

- *NUMBER*: kongruieren i und j bezüglich Numerus (true, false)

12 Attribute der Markables

- *NUMBER*: kongruieren i und j bezüglich Numerus (true, false)
- *SEMCLASS*: is-a Hierarchy

- *NUMBER*: kongruieren i und j bezüglich Numerus (true, false)
- *SEMCLASS*: is-a Hierarchy
 - ▶ wenn $SEMCLASS(i \text{ bzw. } j) \subset SEMCLASS(j \text{ bzw. } i)$
oder $SEMCLASS(i) = SEMCLASS(j) \rightarrow \text{true}$

The top model, Mrs. Claudia Schiffer has three children.
top model \in "person"; Mrs. Claudia Schiffer \in "female"
"person" \subset "female"

- *NUMBER*: kongruieren i und j bezüglich Numerus (true, false)
- *SEMCLASS*: is-a Hierarchy
 - ▶ wenn $SEMCLASS(i \text{ bzw. } j) \subset SEMCLASS(j \text{ bzw. } i)$
oder $SEMCLASS(i) = SEMCLASS(j) \rightarrow \text{true}$

The top model, Mrs. Claudia Schiffer has three children.
top model \in "person"; Mrs. Claudia Schiffer \in "female"
"person" \subset "female"

- ▶ $SEMCLASS(i) \neq SEMCLASS(j) \rightarrow \text{false}$

- *NUMBER*: kongruieren *i* und *j* bezüglich Numerus (true, false)
- *SEMCLASS*: is-a Hierarchy
 - ▶ wenn $SEMCLASS(i \text{ bzw. } j) \subset SEMCLASS(j \text{ bzw. } i)$
oder $SEMCLASS(i) = SEMCLASS(j) \rightarrow \text{true}$

The top model, Mrs. Claudia Schiffer has three children.

top model \in "person"; Mrs. Claudia Schiffer \in "female"
"person" \subset "female"

- ▶ $SEMCLASS(i) \neq SEMCLASS(j) \rightarrow \text{false}$
- ▶ $SEMCLASS(i) = \text{unknown}, SEMCLASS(j) = \text{unknown}$
und $STR_MATCH = \text{false} \rightarrow \text{unknown}$

- *GENDER*: kongruieren i und j bezüglich Genus (true, false, unknown)

- *GENDER*: kongruieren i und j bezüglich Genus (true, false, unknown)
 - ▶ *SEMCLASS*(i und j)=unknown → *GENDER*(i und j)=unknown

12 Attribute der Markables

- *GENDER*: kongruieren i und j bezüglich Genus (true, false, unknown)
 - ▶ *SEMCLASS*(i und j)=unknown → *GENDER*(i und j)=unknown
- *PROPER_NAME*: sind i und j Eigennamen(true, false)

- *GENDER*: kongruieren i und j bezüglich Genus (true, false, unknown)
 - ▶ *SEMCLASS*(i und j)=unknown → *GENDER*(i und j)=unknown
- *PROPER_NAME*: sind i und j Eigennamen(true, false)
- *ALIAS*: sind i und j **Named Entities**, die auf das selbe Entität referieren (true, false)

- *GENDER*: kongruieren i und j bezüglich Genus (true, false, unknown)
 - ▶ *SEMCLASS*(i und j)=unknown → *GENDER*(i und j)=unknown
- *PROPER_NAME*: sind i und j Eigennamen(true, false)
- *ALIAS*: sind i und j **Named Entities**, die auf das selbe Entität referieren (true, false)
 - ▶ **Person** *Mr. Simpson* und *Bent Simpson*

- **GENDER**: kongruieren i und j bezüglich Genus (true, false, unknown)
 - ▶ $SEMCLASS(i \text{ und } j) = \text{unknown} \rightarrow GENDER(i \text{ und } j) = \text{unknown}$
- **PROPER_NAME**: sind i und j Eigennamen (true, false)
- **ALIAS**: sind i und j **Named Entities**, die auf das selbe Entität referieren (true, false)
 - ▶ **Person** *Mr. Simpson* und *Bent Simpson*
 - ▶ **Datum** *01-08* und *Jan. 8*

- **GENDER**: kongruieren i und j bezüglich Genus (true, false, unknown)
 - ▶ $SEMCLASS(i \text{ und } j) = \text{unknown} \rightarrow GENDER(i \text{ und } j) = \text{unknown}$
- **PROPER_NAME**: sind i und j Eigennamen (true, false)
- **ALIAS**: sind i und j **Named Entities**, die auf das selbe Entität referieren (true, false)
 - ▶ **Person** *Mr. Simpson* und *Bent Simpson*
 - ▶ **Datum** *01-08* und *Jan. 8*
 - ▶ **Organization** *IBM* und *International Business Machines Corp.*

- **GENDER**: kongruieren i und j bezüglich Genus (true, false, unknown)
 - ▶ $SEMCLASS(i \text{ und } j) = \text{unknown} \rightarrow GENDER(i \text{ und } j) = \text{unknown}$
- **PROPER_NAME**: sind i und j Eigennamen (true, false)
- **ALIAS**: sind i und j **Named Entities**, die auf das selbe Entität referieren (true, false)
 - ▶ **Person** *Mr. Simpson* und *Bent Simpson*
 - ▶ **Datum** *01-08* und *Jan. 8*
 - ▶ **Organization** *IBM* und *International Business Machines Corp.*
- **APPOSITIVE**: ist j Apposition von i (true, false)

Separately, Clinton transition officials said that **Frank Newman**, 50, **vice chairman** and chief financial officer of BankAmerica Corp. , is expected to be nominated as assistant Treasury secretary for domestic finance.

Beispiel für einen Attributvektor

Attribut	Attributwert	Kommentar
DIST	0	i und j sind im selben Satz
I_PRONOUN	false	i ist kein Pronomen
J_PRONOUN	false	j ist kein Pronomen
STR_MATCH	false	i und j stimmen nicht überein
DEF_NP	false	j ist keine definite NP
DEM_NP	false	j ist keine demonstrative NP
NUMBER	true	i und j sind beide im Singular
SEMCLASS	true	i und j sind beide Personen
GENDER	true	i und j sind beide maskulin
PROPER_NAME	false	nur i ist Eigenname
ALIAS	false	j ist keinen Alias von i
APPOSITIVE	true	j ist Apposition von i

Table: Attributvektor für das Markables-Paar (**i**=Frank Newmann, **j**=vice chairman)

... ((*Eastern Airlines*)⁵_{a2} *executives*)⁶ notified ((*union*)⁷_{e1} *leaders*)⁸ that
(*the carrier*)⁹ wishes to discuss (*selective* ((*wage*)¹⁰_{c2} *reductions*)_{d2})¹¹ on
(*Feb. 3*)¹²_{b2}. (((*Union*)¹³_{e2} *representatives*)¹⁴ *who could be reached*)_{f1} said
(*they*)¹⁵_{f2} hadn't decided whether (*they*)¹⁶_{f3} would respond. By proposing
(*a meeting date*)¹⁷_{b3}, (*Eastern*)¹⁸_{a3} moved (*one step*)¹⁹ closer toward
reopening (*current high – cost contrast agreements*)²⁰ with
((*its*)²¹_{a4} *unions*)²²_{e3}.

... ((*Eastern Airlines*)⁵_{a2} executives)⁶ notified ((*union*)⁷_{e1} leaders)⁸ that
(*the carrier*)⁹ wishes to discuss (*selective* ((*wage*)¹⁰_{c2} reductions)^{d2})¹¹ on
(*Feb. 3*)¹²_{b2}. (((*Union*)¹³_{e2} representatives)¹⁴ who could be reached)^{f1} said
(*they*)¹⁵_{f2} hadn't decided whether (*they*)¹⁶_{f3} would respond. By proposing
(*a meeting date*)¹⁷_{b3}, (*Eastern*)¹⁸_{a3} moved (*one step*)¹⁹ closer toward
reopening (*current high – cost contrast agreements*)²⁰ with
((*its*)²¹_{a4} unions)²²_{e3}.

- positive Beispiele

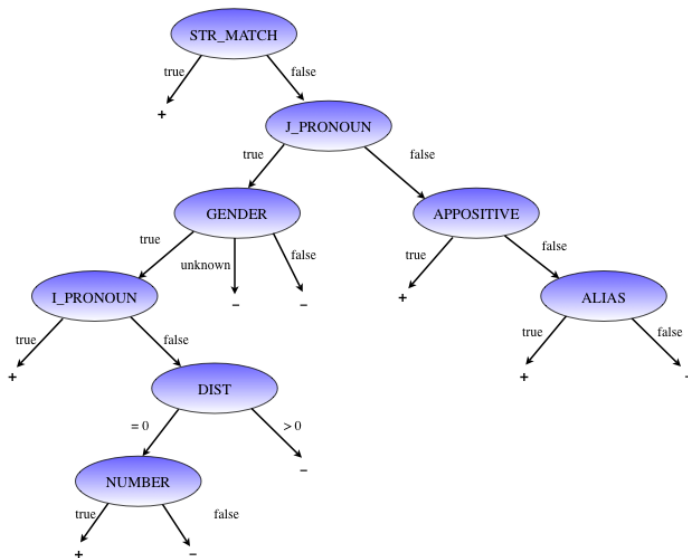
- ▶ ((*union*)⁷, (*union*)¹³)
- ▶ ((*union*)¹³, (*its unions*)²²)

- negative Beispiele

- ▶ ((*the carrier*)⁹, (*union*)¹³)
- ▶ ((*wage*)¹⁰, (*union*)¹³)
- ▶ ((*selective wage reductions*)¹¹, (*union*)¹³)
- ▶ ((*Feb. 3*)¹², (*union*)¹³)

MUC-6 Klassifikator erstellt mit C5

Message Understanding Conference (MUC)



...(Ms. Washington)⁷³ 's candidacy is being championed by (several powerful lawmakers)⁷⁴ including ((her)⁷⁶ boss)⁷⁵, (Chairman John Dingell)⁷⁷ (Dr.(Mich.)⁷⁸) of (the House Energy and Commerce Committee)⁷⁹. (She)⁸⁰ currently is (a counsel)⁸¹ to (the committee)⁸². (Ms. Washington)⁸³ and (Mr. Dingell)⁸⁴ have been considered (allies)⁸⁵ of (the (securities)⁸⁷ exchanges)⁸⁶, while (banks)⁸⁸ and ((futures)⁹⁰ exchanges)⁸⁹ have often fought with (them)⁹¹.

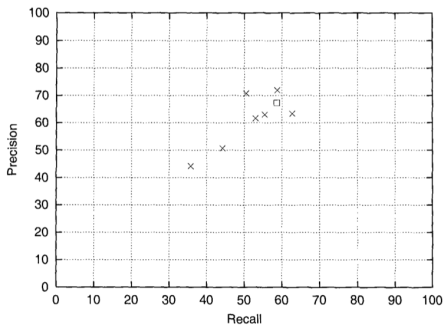
getestete Markables-Paare zur Koreferenzkettenbildung

Bezugswort	Anapher	Attributvektor*	Koreferent?
<i>(several powerful lawmakers)</i> ⁷⁴	(her) ⁷⁶	0, t, f, u, f, f, t, f, f, f, f, f	Nein
(Ms. Washington) ⁷³	(her) ⁷⁶	0, t, t, t, f, f, t, f, f, f, f, f	Ja
<i>(the House Energy and Commerce Committee)</i> ⁷⁹	(She) ⁸⁰	1, f, t, f, f, f, t, f, f, f, f, f	Nein
<i>(Mich.)</i> ⁷⁸	(She) ⁸⁰	2, f, t, f, f, f, t, f, f, f, f, f	Nein
<i>(Chairman John Dingell)</i> ⁷⁷	(She) ⁸⁰	3, t, t, f, f, f, t, f, f, f, f, f	Nein
(her) ⁷⁶	(She) ⁸⁰	3, t, t, t, f, f, t, f, f, f, f, t	Ja
<i>(the committee)</i> ⁸²	(Ms. Washington) ⁸³	1, f, t, f, f, f, f, f, f, f, f, f	Nein
<i>(a counsel)</i> ⁸¹	(Ms. Washington) ⁸³	1, f, t, u, f, f, f, f, f, f, f, f	Nein
(She) ⁸⁰	(Ms. Washington) ⁸³	1, t, t, t, f, f, f, f, f, f, f, t	Nein
<i>(the House Energy and Commerce Committee)</i> ⁷⁹	(Ms. Washington) ⁸³	2, f, t, f, t, f, f, f, f, f, f, f	Nein
<i>(Mich.)</i> ⁷⁸	(Ms. Washington) ⁸³	3, f, t, f, t, f, f, f, f, f, f, f	Nein
<i>(Chairman John Dingell)</i> ⁷⁷	(Ms. Washington) ⁸³	4, t, t, f, t, f, f, f, f, f, f, f	Nein
(her) ⁷⁶	(Ms. Washington) ⁸³	4, t, t, t, f, f, f, f, f, f, f, t	Nein
<i>(her boss)</i> ⁷⁵	(Ms. Washington) ⁸³	4, t, f, f, f, f, f, f, f, f, f, f	Nein
<i>(several powerful lawmakers)</i> ⁷⁴	(Ms. Washington) ⁸³	4, t, f, u, f, f, f, f, f, f, f, f	Nein
(Ms. Washington) ⁷³	(Ms. Washington) ⁸³	4, t, t, t, t, t, f, f, f, t, f, f	Ja

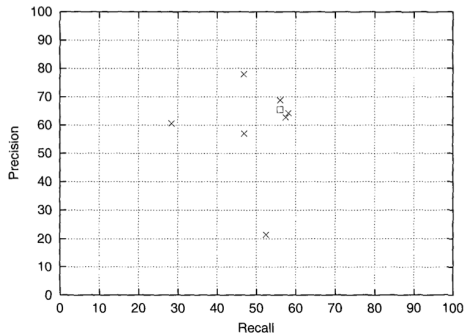
* dist, semclass, number, gender, proper-name, alias, j-pronoun, def-np, dem-np, str-match, appositive, i-pronoun

Evaluation: Precision und Recall

MUC - 6

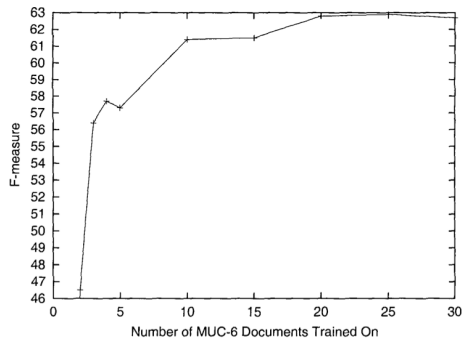


MUC - 7

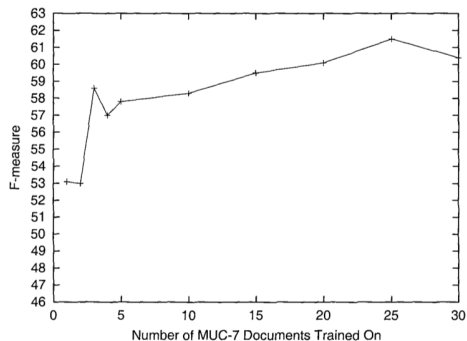


Evaluation: F-Score

MUC - 6



MUC - 7



- inkrementelle Verarbeitung des Eingabetextes

- inkrementelle Verarbeitung des Eingabetextes
- hinreichende Attribute der Markables

- inkrementelle Verarbeitung des Eingabetextes
- hinreichende Attribute der Markables
- Umsetzung auf erhältliche Datensätze

- inkrementelle Verarbeitung des Eingabetextes
- hinreichende Attribute der Markables
- Umsetzung auf erhältliche Datensätze
- echtes Testen und umfassende Evaluation

- [1] Wee Meng Soon, Daniel Chung Yong Lim and Hwee Tou Ng, 2001, "A Machine Learning Approach to Coreference Resolution of Noun Phrases" , *Computational Linguistics*, 27:4, S. 521-544

Danke für Ihre Aufmerksamkeit!