

# Unsupervised Word Sense Disambiguation Rivaling Supervised Methods (David Yarowsky)

Kokou Valentin

29.11.2010

Proseminar Word Sense Disambiguation

# Word Sense Disambiguation

- Vorkommen mehrdeutiger Worte
- Bestimmung aller möglichen Bedeutungen
- Auswahl der richtigen Bedeutung

# Überwachte Lernverfahren

- Klassifikator für jedes Wort
- Sehr hohe Präzision
- Sehr zeitaufwendig („teuer“)
- Benötigt großen annotierten Korpus

# Unüberwachtes Lernverfahren

- konkurriert mit Performanz überw. Lernverfahren
- „günstiger“ als Überwachte Lernverfahren
  - benötigt sehr kleinen annotierten Korpus
- basierend auf 2 Eigenschaften natürlicher Sprache
  - One Sense per Collocation (OSPC)
  - One Sense per Discourse (OSPD)
- selbst korrigierend

# One Sense per Collocation

- eine Bedeutung innerhalb einer Kollokation
- enge semantische Beziehung
- enger Zusammenhalt des Sachverhalts

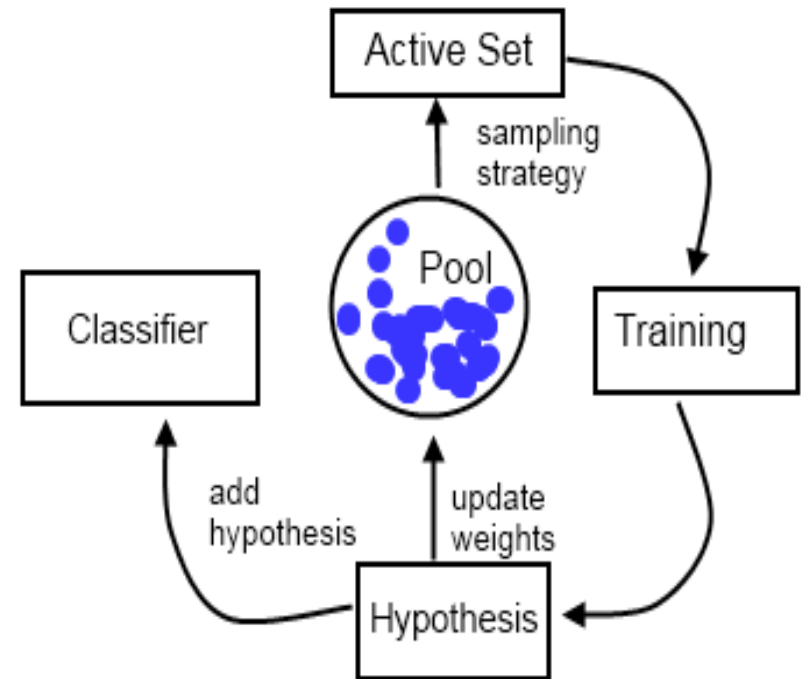
# One Sense per Discourse

- Tendenz nur eine Bedeutung pro Diskurs/Dokument
- Beeinflusst Klassifikationen

Wort	Bedeutung	Genauigkeit	Anwendbar
plant	living/factory	99,8%	72.8%
bass	fish/music	100%	58.8%
palm	tree/hand	99.8%	35.5%

# Bootstrapping

- Iterativer Prozess
- Vermehrung Startmenge
- Reduzierung der Restmenge



# Schritt 1

- Identifiziere alle Beispiele im großen Korpus
- Speichere den Kontext als Zeilen
- Anfangs in nicht annotiertes Trainingsset

# Beispiele für Trainingsset

<b>Sense</b>	Trainings Beispiele
?	... company said the <i>plant</i> is still operating ....
?	...Computer disk drive <i>plant</i> located...
?	... divide life into <i>plant</i> and animal kingdom
?	Although thousands of <i>plant</i> and animal species...

# Schritt 2

- Für jede mögliche Bedeutung des Wortes
  - Finde kleine Anzahl an Trainingsbeispielen
  - mögliche Untermenge manuell taggen



<b>Sense</b>	<b>Training Examples</b>
<b>A</b>	... divide life into <i>plant</i> and animal kingdom
<b>A</b>	....Although thousands of <i>plant</i> and animal species...
<b>?</b>	... molecules found in <i>plant</i> and animal tissue
<b>?</b>	...Nissan car and truck <i>plant</i> in Japan is...
<b>B</b>	...Computer disk drive <i>plant</i> located...
<b>B</b>	... company said the <i>plant</i> is still operating ....

# Schritt 3a

- Trainiere Klassifikationsalgorithmus auf Seed Set (Seed=Keim)
- Decision-list Algorithmus findet Kollokationen

LOGL	Kollokationen	Sense
8.10	<i>Plant life</i> →	<b>A</b>
7.58	Manufacturing plant →	<b>B</b>
7.39	life(within $\pm$ 2-10 words) →	<b>A</b>
7.20	manufacturing(in $\pm$ 2-10 words) →	<b>B</b>



# Schritt 3c

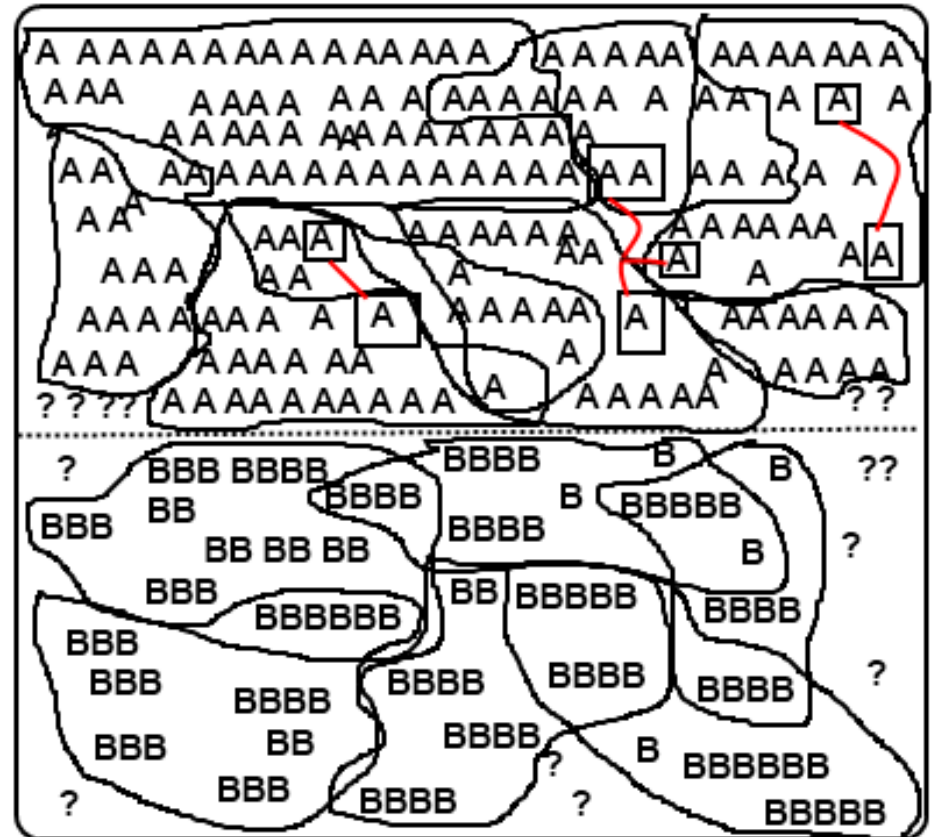
- OSPD Bedingung zum filtern und erhöhen
  - Mehrere Instanzen eines Wortes im Diskurs gleicher Sense Tag
  - Erweitern des Tag zu allen Beispielen im Diskurs
  - Bedingt durch relative Anzahl und WK
  - Erhöhung kann Verbindungen zu neuen Kollokationen bilden
  - Korrigiert fehlerhafte annotierte Beispiele

Change in tag	Diskurs Nummer	Training Samples (Von gleichem Diskurs)
A-> A	742	... the existence of <i>plant</i> and animal life...
A->A	742	....classified as either <i>plant or animal</i> ...
?->A	742	...Although bacterial and <i>plant</i> cells are enclosed
?->A	348	...an aspect of <i>plant life, for example</i> ...
A->A	525	...slight within Arctic <i>plant</i> species...
B->A	525	... are protected by <i>plant</i> parts remaining from ....



# Schritt 4

- Trainings Parameter konstant halten
- Algorithmus nähert sich stabiler Restmenge
- Meisten Beispiele zeigen auf multiple Kollokationen mit der selben Bedeutung
- Decision List Alg nutzt nur höchsten Eintrag
- Keine Kombination aller zutreffenden Kollokationen



# Schritt 5

- Klassifikator aus finalelem Trainingsschritt
  - Anwendung auf neue Daten
  - Annotiert original nicht annotierten Korpus
  - Sense Tag und Wahrscheinlichkeiten annotieren
  - Ursprüngliche Seeds nicht länger an Spitze der Liste

# Schritt 5

- Möglich, dass ein Sense A zu einem Sense B wird
  - Bei höherer Kompatibilität zur 2. Klasse
- Störanfälligkeit durch fehlerhafte Ausgangswörter nicht fatal
- Korrigiert, wenn Mehrheit der Ausgangs-WK kohärenten Kollokationsraum bilden

## Final decision list for *plant*

<b>LOGL</b>	<b>Kollokationen</b>	<b>Sense</b>
10.12	<b>plant growth</b> →	<b>A</b>
9.68	<b>Car (within ± k words)</b> →	<b>B</b>
9.64	<b>plant height</b> →	<b>A</b>
9.61	<b>Union (± k words)</b> →	<b>B</b>
9.54	<b>Equipment (within ± k words)</b> →	<b>B</b>
9.51	<b>assembly plant</b> →	<b>B</b>
9.02	<b>Plant species</b> →	<b>A</b>

# Beispiel

... the loss of animal and plant **species** trough extenction...

•Höchste Kollokation im Bsp. Satz „species“ wird genutzt zum Klassifizieren als Sense A (living plant)

# Woher kommen Seed Patterns?

- Nutze Wörter aus Wörterbuch Definitionen
  - Automatisch auslesen von Wörtern mit hoher Frequenz relativ zum gesamten Wörterbuch
- Nutze einzelne abgrenzende Anordnung für jede Klasse
  - Nur diese Anordnung nutzen mit ihrem Kontext
  - Bsp.: flora und factory für plant
- Markiere hervorstechende Korpusanordnungen
  - Wörter mit ungewöhnlich hohem Vorkommen
  - Menschliche Entscheidung über Bedeutung

# Korrektur

- Beispiele aus Restmenge nur solange WK über dem Grenzwert
- Klassifikation kann sich ändern durch neue Beispiele
- Rückkehr zur nicht annotierten Restmenge
- Neuklassifikation durch iterativen Trainingsschritt
- Prozess ist selbstkorrekturfähig

# Rettung

- Trotzdem Möglichkeit der falschen Klassifikation
  - starke Kollokationen Indikator für falsche Klasse
- 2 Techniken im Trainingsalgorithmus sollen dies verhindern
  - 1) taktweise Erhöhung des Kontextrahms
  - 2) zufällige Störung des Grenzwertes der Klassifikation

# OSPD Eigenschaft nutzen

- gute Performanz bei lokalen Kollokationen
  - jedes Token von Zielwort unabhängig behandelt
- Genauigkeit kann verbessert werden
  - Alle Vorkommen in Diskurs haben selbe Bedeutung
- Eigenschaft an 2 Stellen nutzbar
  - nach Schritt 3c (nach jeder Iteration)
  - Ende Schritt 4 nachdem Algo. konvergiert

# OSPD Eigenschaft nutzen

- Fehlerkorrektur
  - Bsp. *plant* erscheint mehrmals in einem Diskurs
  - Tokens vom Algo. durch lokale Kollokationen niedriger Wert
    - Kann durch dominanten Tag des Diskurs überschrieben werden
  - Variablen zur Entscheidung
    - n vollständige Anzahl des Auftretens im Diskurs
    - Anzahl Vorkommen von zugeteilten Mehr-/Minderheiten der Bedeutung
    - Summe von LogL-Verhältnissen

# Fehlerkorrektur II

- Wenn Hauptbedeutung größer als Nebenbedeutung
  - Grenzwert in Abhängigkeit von  $n$
  - Nebenbedeutung wird neu markiert
- $n=2$  keine Neuklassifikation da nicht klar, welche dominant
- $n>4$  meisten lokalen Klassifikationen durch dominanten Tag überschrieben
- Im Diskurs ohne dominante Bedeutung werden Tags entfernt

# Tabelle Evaluation

(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
Word	Sense	Samp. Size	% Major Sense	Supvsd Algrtm	Seed Training Options		
					Two Words	Dict. Defn.	Top Colls.
plant	living/ factory	7538	53.1	97.7	97.1	97.3	97.6
space	volume/ outer	5745	50.7	93.3	89.1	92.3	93.5
tank	vehicle/ container	11420	58.2	97.1	94.2	94.6	95.8
bass	fish/ music	1859	57.5	97.8	96.6	97.2	97.7
palm	tree/ hand	1572	56.1	96.5	93.9	94.7	95.8
poach	steal/ boil	585	74.9	97.1	96.6	97.2	97.7
axes	grid/ tools	1344	84.9	95.5	94.0	94.3	94.7
drug	medicine/ narcotic	1380	50.0	93.0	90.4	91.4	92.6
crane	bird/ machine	2145	78.0	96.6	59.6	93.6	94.2
<b>AVG</b>		3936	63.9	96.1	90.6	94.8	95.5

# Tabelle Evaluation

(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)	(10)
Word	Sense	Samp. Size	% Major Sense	Supvsd Algrtm	Seed Training Options			<b>(7)+ OSDP</b>	
					Two Words	Dict. Defn.	Top Colls.	<b>End only</b>	<b>Each Iter.</b>
plant	living/ factory	7538	53.1	97.7	97.1	97.3	97.6	<b>98.3</b>	<b>98.6</b>
space	volume/ outer	5745	50.7	93.3	89.1	92.3	93.5	<b>93.3</b>	<b>93.6</b>
tank	vehicle/ container	11420	58.2	97.1	94.2	94.6	95.8	<b>96.1</b>	<b>96.5</b>
bass	fish/ music	1859	57.5	97.8	96.6	97.2	97.7	<b>98.5</b>	<b>98.8</b>
palm	tree/ hand	1572	56.1	96.5	93.9	94.7	95.8	<b>95.5</b>	<b>95.9</b>
poach	steal/ boil	585	74.9	97.1	96.6	97.2	97.7	<b>98.4</b>	<b>98.5</b>
axes	grid/ tools	1344	84.9	95.5	94.0	94.3	94.7	<b>96.8</b>	<b>97.0</b>
drug	medicine/ narcotic	1380	50.0	93.0	90.4	91.4	92.6	<b>93.9</b>	<b>93.9</b>
crane	bird/ machine	2145	78.0	96.6	59.6	93.6	94.2	<b>95.4</b>	<b>95.5</b>
<b>AVG</b>		3936	63.9	96.1	90.6	94.8	95.5	<b>96.1</b>	<b>96.5</b>

# Schlussfolgerung

- Algorithmus nutzt mehrere Eigenschaften von Sprache aus
  - One Sense per Discourse
  - One Sense per Collocation
- max. Einfluss zur Modellierung großer Vielfältigkeit von Beziehungen zw. Kollokationen
- empfindlicher für größere Spanne an sprachlichen Details
- Fast gleiche Performanz wie überwachtetes Lernverfahren  
95.5% vs. 96.1%
- in Fällen mit OSPD besser 96.5% vs. 96.1%
- großer Sense-tagged Trainingskorpus nicht notwendig

# Zusammenfassung

- Teilweise bessere Performanz (mit/ohne OSPD)
- Gutes Ergebnis bei Homonymen
- Voraussichtlich bei Polyseme weniger gut