

Neue Verfahren und Evaluierungsmaße für Anfragesegmentierungen

Anna Beyer

Bauhaus-Universität Weimar
anna.beyer@uni-weimar.de

Bachelorverteidigung
3. Februar 2012

- 1 Motivation
- 2 Anfragesegmentierung
- 3 Segmentierungsverfahren
- 4 Evaluierung
- 5 Zusammenfassung



1 Milliarde Suchanfragen täglich

Quelle: <http://www.google.com/insidesearch/underthehood.html>

The Google logo is displayed in its characteristic multi-colored font (blue, red, yellow, blue, green, red).

tokio hotel bill

Was sucht der Nutzer?

Motivation

tokio hotel bill

料金領收書 RECEIPT		簿目 第 號 Book No.	
客名 客名 客人	姓名 姓名 姓 名	册 號 册 號 Number of Guest	2
Mr. <u>上</u>	女中名 Name of Waitress		
	名 氏	保 保 Person in Charge	大 八 郎
下記之通り領收候也 Received The followin. G With Thanks			
茶 代 茶代 Meal & Tea	9 0 5	價 價 Yen	
花 代 花代 For The Flowers of Waitress			
花 代 外 / 料 金 (料 理 酒 其 他) For The liquid & Food.			
品 名 Description	金 額 Amount	品 名 Description	金 額 Amount
宿 代 AT Night		酒 代 (米) Beer (米)	
酒 酒		飲 料 (洋) Cola (洋)	
茶 代 茶代		香 煙 (洋) Cigarette (洋)	3 2
		雜 費 (洋) Misc (洋)	
計 計		計 計	7 8 0
預 交 料 金 預交料金 Y.A. 金		預 交 料 金 預交料金 Y.A. 金	2 3 2
尚 欠 料 金 尚欠料金 Amount 尚欠料金 Amount			
8 第 〇 七 四 號			
営業場所 横浜市西区本町二丁目 日丁 番地 五 牧 入 部 Honmaku, Yokohama, Japan.			
領收者氏名 又一名稱 17.8.25 HOTEL YOKOHAMA 印 留 印 領 收 者 印 留 印 領 收 者 印 留 印 STAMP.			

Quelle: <http://ahoy.tk-jk.net/MoreImages4/WillyHotelBill.jpg>

Anfragesegmentierung

Definition

Einteilung einer Suchanfrage in Sinneinheiten

Syntax

Kennzeichnung der Segmente durch Anführungszeichen

z.B. "tokio hotel" "bill"

Vorteil

Verbesserte Retrieval Performance

[Li et al., SIGIR 2011]

Situation

Segmentierung bei weniger als 1,12 % der Suchanfragen

[White und Morris, SIGIR 2007]

Situation

Segmentierung bei weniger als 1,12 % der Suchanfragen

[White und Morris, SIGIR 2007]

Lösung

Automatische Segmentierung vor der Suche

Segmentierungsverfahren

Problemdefinition

- ▶ Gegeben Anfrage mit n Worten
- ▶ Gesucht “Beste” Segmentierung aus
 2^{n-1} möglichen Segmentierungen

Beispiel

Anfrage new york times

Segmentierungen new | york | times
 new | york times
 new york | times
 new york times

Verwandte Arbeiten

PMI-basiert	[Risvik et al., WWW 2003] [Jones et al., WWW 2006] [Huang et al., WWW 2010]
Überwachtes Lernen	[Bergsma and Wang, EMNLP-CoNLL 2007] [Bendersky et al., SIGIR 2009]
Häufigkeitsbasiert	[Zhang et al., ACL-IJCNLP 2009] [Hagen et al., SIGIR 2010] [Mishra et al., WWW 2011] [Li et al., SIGIR 2011]
Wikipedia-basiert	[Tan and Peng, WWW 2008] [Hagen et al., WWW 2011]
Retrieval feedback	[Brenes et al., CERI 2010] [Bendersky et al., CIKM 2010] [Bendersky et al., ACL 2011]

Idee

Berechnung einer Punktzahl für jede mögliche Segmentierung

Beispiel

new york times

Idee

Berechnung einer Punktzahl für jede mögliche Segmentierung

Beispiel

Mögliche Segmente

Segment s
new york
york times
new york times

Idee

Berechnung einer Punktzahl für jede mögliche Segmentierung

Beispiel

Häufige Phrasen = bessere Segmente

Segment s	$freq(s)$
new york	165,36
york times	17,60
new york times	17,55

Idee

Berechnung einer Punktzahl für jede mögliche Segmentierung

Beispiel

Wikipedia-Titel = hochwertige Segmente

Segment s	$freq(s)$	Wiki
new york	165,36	✓
york times	17,60	-
new york times	17,55	✓

Idee

Berechnung einer Punktzahl für jede mögliche Segmentierung

Beispiel

Wikipedia-Titel = hochwertige Segmente

Segment s	$freq(s)$	Wiki	$freq'(s)$
new york	165,36	✓	165,36
york times	17,60	-	17,60
new york times	17,55	✓	165,36

Idee

Berechnung einer Punktzahl für jede mögliche Segmentierung

Beispiel

Normalisierung mit Segmentlänge

Segment s	$freq(s)$	Wiki	$freq'(s)$	$weight(s)$
new york	165,36	✓	165,36	330,72
york times	17,60	-	17,60	35,20
new york times	17,55	✓	165,36	496,08

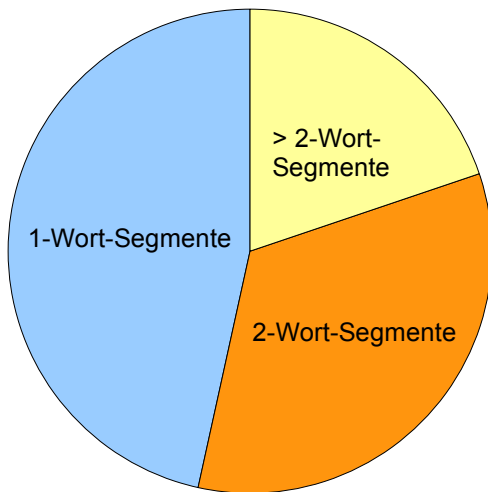
Idee

Berechnung einer Punktzahl für jede mögliche Segmentierung

Beispiel

Punktzahl = Summe Gewichte

Rang	Segmentierung S	$score(S)$
1	new york times	496,08
2	new york times	330,72
3	new york times	35,20
4	new york times	0,00



550 000 Segmentierungen

Idee

Nur Wikipedia-Titel in einer Anfrage segmentieren

Beispiel

Anfrage pictures of new york in the 1920s

Idee

Nur Wikipedia-Titel in einer Anfrage segmentieren

Beispiel

Anfrage pictures of new york in the 1920s

Idee

Nur Wikipedia-Titel in einer Anfrage segmentieren

Beispiel

Anfrage pictures of new york in the 1920s

Segmentierung pictures | of | new york | in | the | 1920s

	Wiki-basiert	WT-Baseline
Nominalanfragen		
Nicht-Nominalanfragen		

Trainingsset ca. 5 000 Anfragen

Performance-Vergleich

	Wiki-basiert	WT-Baseline
Nominalanfragen	0,580	0,497
Nicht-Nominalanfragen		

Trainingsset ca. 5 000 Anfragen

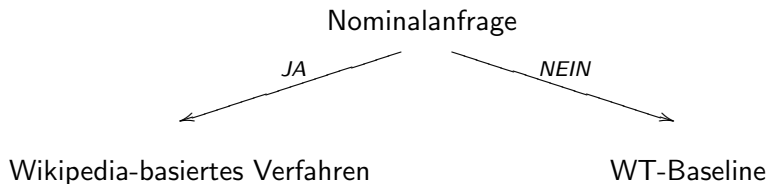
Performance-Vergleich

	Wiki-basiert	WT-Baseline
Nominalanfragen	0,580	0,497
Nicht-Nominalanfragen	0,268	0,492

Trainingsset ca. 5 000 Anfragen

	Wiki-basiert	WT-Baseline
Nominalanfragen	0,580	0,497
Nicht-Nominalanfragen	0,268	0,492

Trainingsset ca. 5 000 Anfragen



Wiki-basiert	WT-Baseline	Hybrid
0,412	0,502	0,535

Testset ca. 50 000 Anfragen

Evaluierung

Grundlage

Manuelle Segmentierungen in Evaluierungskorpora

Beispiel

Stimmen	Segmentierung
2	tokio hotel bill
1	tokio hotel bill

Bisher

Referenz = “Ähnlichste” Segmentierung im Korpus

Beispiel

Vorhersage tokio | hotel bill

Stimmen	Segmentierung
2	tokio hotel bill
1	tokio hotel bill

Bisher

Referenz = “Ähnlichste” Segmentierung im Korpus

Beispiel

Vorhersage `tokio | hotel bill`

	Stimmen	Segmentierung
	2	<code>tokio hotel bill</code>
Referenz \implies	1	<code>tokio hotel bill</code>

Bisher

Referenz = “Ähnlichste” Segmentierung im Korpus

Beispiel

Vorhersage tokio | hotel bill

	Stimmen	Segmentierung
	9	tokio hotel bill
Referenz \implies	1	tokio hotel bill

Bisher

Referenz = “Ähnlichste” Segmentierung im Korpus

Beispiel

Vorhersage tokio | hotel bill

	Stimmen	Segmentierung	
	9	tokio hotel bill	⇒ Faire Referenz
Referenz ⇒	1	tokio hotel bill	

Neu

Erzeugung der Referenz unter Berücksichtigung aller Segmentierungen

Beispiel

Stimmen	Segmentierung
4	tokio hotel bill
3	tokio hotel bill
2	tokio hotel bill
1	tokio hotel bill

Neu

Erzeugung der Referenz unter Berücksichtigung aller Segmentierungen

Beispiel

Stimmen	Segmentierung
4	tokio hotel bill
3	tokio hotel bill
2	tokio hotel bill
1	tokio hotel bill

NB	B	NB	B
tokio	hotel	bill	
6	4	3	7

Neu

Erzeugung der Referenz unter Berücksichtigung aller Segmentierungen

Beispiel

Stimmen	Segmentierung
4	tokio hotel bill
3	tokio hotel bill
2	tokio hotel bill
1	tokio hotel bill

NB	B	NB	B
tokio	hotel	bill	
6	4	3	7



tokio hotel | bill

Beispiel

Referenz	tokio hotel bill	(2 Segmente)
Vorhersage	tokio hotel bill	(3 Segmente)

Beispiel

Referenz	tokio hotel bill	(2 Segmente)
Vorhersage	tokio hotel bill	(3 Segmente)

query 0,00 Vorhersage \neq Referenz

Beispiel

Referenz tokio hotel | bill (2 Segmente)

Vorhersage tokio | hotel | **bill** (**3 Segmente**)

query 0,00

Vorhersage \neq Referenz

segPrec 0,33

1 von 3 der vorhergesagten Segmente korrekt

Beispiel

Referenz tokio hotel | bill (2 Segmente)

Vorhersage tokio | hotel | bill (3 Segmente)

query 0,00

Vorhersage \neq Referenz

segPrec 0,33

1 von 3 der vorhergesagten Segmente korrekt

segRec 0,50

1 von 2 der referenzierten Segmente korrekt

Beispiel

Referenz	tokio hotel bill	(2 Segmente)
Vorhersage	tokio hotel bill	(3 Segmente)

query	0,00	Vorhersage \neq Referenz
segPrec	0,33	1 von 3 der vorhergesagten Segmente korrekt
segRec	0,50	1 von 2 der referenzierten Segmente korrekt
break	0,50	1 von 2 Break-Entscheidungen korrekt

	Wiki-basiert	WT-Baseline	Hybrid
query	0,412	0,502	0,535
segPrec	0,606	0,673	0,693
segRec	0,565	0,708	0,703
break	0,717	0,778	0,785

Testset ca. 50 000 Anfragen

	Wiki-basiert	WT-Baseline	Hybrid
Anfragen/s	3 358	3 797	2 918

Testset ca. 50 000 Anfragen

	Wiki-basiert	WT-Baseline	Hybrid
Anfragen/s	3 358	3 797	2 918

Testset ca. 50 000 Anfragen

1 Milliarde Anfragen/Tag = 12 000 Anfragen/s

Zusammenfassung

Ergebnisse

- ▶ WT-Baseline
- ▶ Hybrides Verfahren
- ▶ Fairere Evaluierungsmaße

Ergebnisse

- ▶ WT-Baseline
- ▶ Hybrides Verfahren
- ▶ Fairere Evaluierungsmaße

Ausblick

- ▶ Optimierung
- ▶ Retrieval Performance

Ergebnisse

- ▶ WT-Baseline
- ▶ Hybrides Verfahren
- ▶ Fairere Evaluierungsmaße

Ausblick

- ▶ Optimierung
- ▶ Retrieval Performance

Vielen Dank