

目的

- 専門用語の抽出
 - サジェスション
 - メタ情報
 - クラスタリング

専門用語の抽出

- ・ コーパスから固有表現を自動抽出する
- ・ 対象：単名詞とその単名詞からなる複合名詞のみとする
- ・ 固有表現
 - ・ 単名詞
 - これ以上分割できない名詞
 - ・ 複合名詞
 - 専門用語の多くは複合語
 - ・ ターム性
 - ・ ある単語単位(複合名詞など)のもつ専門用語固有概念への関連の強さ
 - ・ 専門用語を書いた専門家の概念に直結していると考えられる
 - ・ tf-idfは表層表現のコーパスでの現れ方を利用した近似表現に過ぎない
 - ターム性を直接的に反映する用語抽出方法が必要
 - ・ ユニット性
 - ・ ある言語単位がコーパス中で使用されている度合い

単名詞の接続統計情報の一般化

- ・ 特定のコーパスを想定したとき、単名詞Nが接続する状況すなわち単名詞バイグラム(2-gram)を以下のように表す

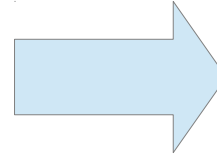
$$\begin{array}{ll} [LN_1 \ N](\#L_1) & [N \ RN_1](\#R_1) \\ [LN_2 \ N](\#L_2) & [N \ RN_2](\#R_2) \\ \vdots & \vdots \\ [LN_n \ N](\#L_n) & [N \ RN_m](\#R_m) \end{array}$$

図 1: 単名詞 N を含む単名詞バイグラムと左右接続単名詞の頻度

- ・ $LN_i(i=1,\dots,n)$ は、単名詞バイアグラム $[LN_i \ N]$ において N の左方に接続する単名詞(n 種類)を表す
- ・ $RN_i(i=1,\dots,m)$ は、単名詞バイアグラム $[N \ RN_i]$ において N の右方に接続する単名詞 RN_i (m 種類)を表す
- ・ $\#L_i(i=1,\dots,n)$ は、 LN_i の頻度
- ・ $\#R_i(i=1,\dots,m)$ は、 RN_i の頻度

実例（頻度と種類）

	トライグラム	統計
	トライグラム	
単語	トライグラム	
クラス	トライグラム	
単語	トライグラム	
	トライグラム	
	トライグラム	
単語	トライグラム	抽出
	トライグラム	統計
	トライグラム	
文字	トライグラム	

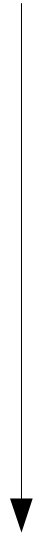


[単語 トライグラム](3)
[クラス トライグラム](1)
[文字 トライグラム](1)
#LDN(トライグラム)=3
#LN(トライグラム)=5

[トライグラム 統計](2)
[トライグラム 抽出](1)
#RDN(トライグラム)=2
#RN(トライグラム)=3

- ・ 接続種類数
 - #LDN(N) : 単名詞バイグラムで単名詞Nの左方にくる単名詞の種類の違い数
 - #RDN(N) : " 右方 "
- ・ 接続頻度
 - #LN(N) : Nの左方に連結して複合名詞を形成する全単名詞の頻度
 - #RN(N) : " 右方 "

専門用語には複合名詞が多い



複合名詞のスコアを定義する必要がある

複合名詞のスコアリング

$$LR(CN) = \left(\prod_{i=1}^L (FL(N_i) + 1)(FR(N_i) + 1) \right)^{\frac{1}{2L}}$$

単名詞 : N_i

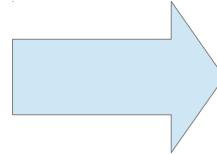
複合名詞 : $CN = N_1N_2...N_L$

単名詞 N の左方スコア関数 : $FL(N)$

単名詞 N の右方スコア関数 : $FR(N)$

実際に計算してみる(CN=トライグラム、L=1)

	トライグラム	統計
	トライグラム	
単語	トライグラム	
クラス	トライグラム	
単語	トライグラム	
	トライグラム	
	トライグラム	抽出
単語	トライグラム	統計
	トライグラム	
文字	トライグラム	



[単語 トライグラム](3)
 [クラス トライグラム](1)
 [文字 トライグラム](1)
 #LDN(トライグラム)=3
 #LN(トライグラム)=5

[トライグラム 統計](2)
 [トライグラム 抽出](1)
 #RDN(トライグラム)=2
 #RN(トライグラム)=3

$$LR(CN) = \left(\prod_{i=1}^L (FL(N_i) + 1)(FR(N_i) + 1) \right)^{\frac{1}{2L}}$$

連接頻度LN,RNを単名詞のスコアとした場合

$$LR(\text{トライグラム}) = \sqrt{(5+1)(3+1)} \approx 4.9$$

候補語の出現頻度を考慮した重み付け

$$FLR(CN) = f(CN) \times LR(CN)$$

$f(CN)$: 候補語CNが単独で出現した頻度

連結頻度をスコアとした場合

$$\begin{aligned} FLR(\text{トライグラム}) &= f(\text{トライグラム}) \times LR(\text{トライグラム}) \\ &= 3 \times \sqrt{(5+1)(3+1)} \\ &= 14.90 \end{aligned}$$

C-Value

比較のために、単名詞バイグラムによらない用語スコアリングとして C-Value(Frantzi and Ananiadou 1996)を考える

$$\text{C-value}(CN) = (\text{length}(CN) - 1) \times (n(CN) - \frac{t(CN)}{c(CN)})$$

CN : 複合名詞¹

$\text{length}(CN)$: CN の長さ (構成単名詞数)

$n(CN)$: コーパスにおける CN の出現回数

$t(CN)$: CN を含むより長い複合名詞の出現回数

$c(CN)$: CN を含むより長い複合名詞の異なり数

$\text{length}(CN)=1$ (CN が単名詞)の場合、C-Valueが0になってしまう

MC-Value (Modified C-Value)

$$\text{MC-value}(CN) = \text{length}(CN) \times (n(CN) - \frac{t(CN)}{c(CN)})$$

○CN = トライグラム の場合

$\text{length}(\text{トライグラム}) = 1$: トライグラムの長さ(構成単名詞数)

$n(\text{トライグラム}) = 10$: コーパスにおけるトライグラムの出現回数

$t(\text{トライグラム}) = 7$: トライグラムを含む、より長い複合名詞の出現回数

$c(\text{トライグラム}) = 6$: トライグラムを含む、より長い複合名詞の異なり数

$$\text{MC-Value}(\text{トライグラム}) = 10 - 7/6 \approx 8.83$$

実際に抽出してみる

- Luceneインデックスをコーパスとして利用



「出現頻度と接続頻度に基づく専門用語抽出」

<http://www.r.dl.itc.u-tokyo.ac.jp/~nakagawa/academic-res/jnlp10-1.pdf>

「Luceneインデックスからの専門用語抽出」

<http://www.slideshare.net/KojiSekiguchi/lucene-terms-extraction>