

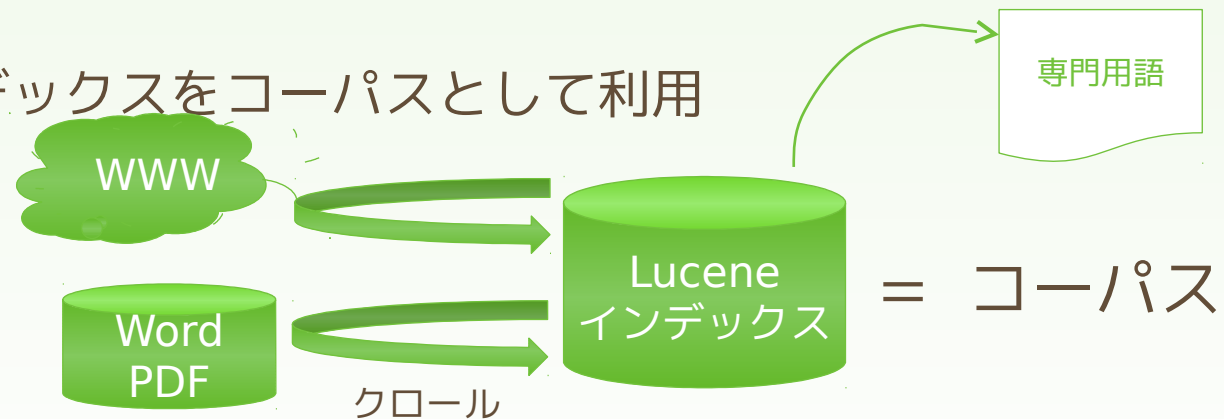
# Luceneインデックスからの 専門用語抽出

関口宏司@ロンウイット

# 背景／目的

- ・ 専門用語抽出
  - ・ サジェスチョン
  - ・ クラスタリング
  - ・ MLT
  - ・ もしかして検索

- ・ Luceneインデックスをコーパスとして利用



# 専門用語抽出

- 単名詞
- 複合名詞
- ターム性
  - ある言語単位（複合名詞など）の持つ専門分野固有の概念への関連性の強さ
  - 専門文書を書いた専門家の概念に直結していると考えられる
- ユニット性
  - ある言語単位がコーパス中で安定して使用される度合い

手順

# 単名詞バイグラムと頻度



# 実際の例（頻度と異なり数）

	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

[トライグラム      統計](2)  
[トライグラム      抽出](1)

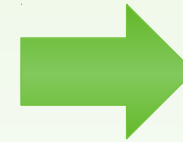
#RDN(トライグラム)=2

## 連接頻度

$$\#LN(N) = \sum_{i=1}^n (\#L_i)$$
$$\#RN(N) = \sum_{i=1}^m (\#R_i)$$

# 実際の例（連接頻度）

	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



# 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： $CN = N_1 N_2 \dots N_L$

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： $FR(N)$

# 実際の例（CNスコア、連接頻度の場合）

	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	

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



[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

$$LR(\text{トライグラム}) = \sqrt{(5 + 1)(3 + 1)} \approx 4.90$$

出現頻度を考慮した重み付け

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

$$FLR(\text{トライグラム}) = 3 \times \sqrt{(5+1)(3+1)} \approx 14.70$$

# G-value

$$C_{value}(CN) = (length(CN) - 1) \times \left\{ n(CN) - \frac{t(CN)}{c(CN)} \right\}$$

CN: 複合名詞  
length(CN): CNの長さ（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

## MC value

$$MC_{value}(CN) = length(CN) \times \left\{ n(CN) - \frac{t(CN)}{c(CN)} \right\}$$

$$MC_{value}(\text{トライグラム}) = 1 \times \left\{ 10 - \frac{7}{6} \right\} \approx 8.83$$

※ 参考文献[1]とカウント結果が異なる

CN: 複合名詞  
length(CN): CNの長さ（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

# 実際の抽出例

※ 接続頻度FLR法（Lucene近似版）を利用

- ・ livedoorニュースコーパス（以下同様）、スポーツウォッチ（上位120位まで）

日本 選手 日本代表 試合 放送 サッカー 野球 代表 監督 チーム 番組 五輪 プロ野球 W杯 一戦 一人 番組内 サッカー解説者 サッカー日本代表 出場	世界 プロ ファン 同番組 ロンドン五輪 試合後 解説 スポーツ番組 優勝 野球解説者 サッカー女子日本代表 元日本代表 試合中 プロ野球解説者 週刊アサヒ芸能 相手 サッカー選手 大会 なでしこジャパン 野球選手	関係 サッカーファン 女子 プロ野球選手 日本中 ゴール 自分 野村克也氏 出演 本田圭佑 ネット上 関係者 日本戦 プレー 関連リンク ブログ スポーツ 練習 報道 楽天	女子サッカー 注目 移籍 開催 リーグ 野村氏 試合前 ロンドン五輪出場 五輪出場 獲得 掲示板上 自身 日本人選手 インタビュー 二人 開幕 ボール Jリーグ 球団 スポーツ選手	一日 なでしこ 巨人 決勝 人気 最終的 選手たち 代表戦 世界選手権 前出 女子W杯 所属 日本サッカー オープン戦 問題 コメント 生出演 一年 日本サッカー協会 日本選手権	勝利 代表選手 発表 北海道日本ハムファイターズ 決勝戦 田中将大 サッカー界 番組中 日本シリーズ 一回 五輪代表 最終戦 日本代表監督 予選 ツイッター上 本田 開幕戦 可能性 北京五輪 スポーツ紙
---	--	---	---	--	--

# 実際の抽出例

※ 接続頻度FLR法（Lucene近似版）を利用

## ・トピックニュース（上位140位まで）

ネット掲示板 関連情報 関連記事 記事 ネットユーザー ネット ネット上 番組 日本 放送 批判 ユーザー 韓国 ネット掲示板 同記事 ニュース 情報 報道 問題 掲示板	同番組 女性 コメント 韓国人 番組内 視聴者 日韓 掲示板 韓国語 発言 ツイッター 批判的 写真 情報番組 テレビ 出演 一人 韓流 意見 ブログ	livedoorニュース 関係 番組中 ファン 写真ギャラリー メディア ユーザ livedoor 関係者 韓国メディア ツイッター上 韓国ネット掲示板 発表 公式サイト 韓国ネットユーザー 被害者 ニュースサイト 放送中 芸能界 放送事故	可能性 人気 選手 生活 ネットユーザ AKB 五輪 テレビ番組 ネットユーザーたち 事件 フジテレビ 韓国内 自分 視聴率 自身 男性 メンバー 話題 出演者 殺到	二人 記事内 タレント 人たち 女性自身 芸人 各メディア 女子 韓国ニュース 逮捕 公開 橋下氏 生活保護 怒り コメント欄 否定的 内容 問題視 関連 紹介	芸能 さん 放送後 社会 記事中 子供 ドラマ 外国人 日本人 活動 騒動 映画 被害 事故 ツイート アイドル 自殺 結婚 事務所 世界	独島 韓国名 発売 政府 フジテレビ系 記者 日本中 関連リンク 同番組内 代表 大阪 ラジオ番組 同ブログ CM 番組放送中 共演者 芸能人 韓国人女性 日本代表 日韓戦
--	--	---	--	---	--	---

# スコア計算方法の違いによる抽出比較

接続種類数LR法	接続頻度LR法	接続種類数FLR法	接続頻度FLR法
選手	日本	選手	日本
日本中	日本代表	日本	選手
日本選手	代表	日本代表	日本代表
日本戦	日本代表選手	チーム	試合
女子選手	日本選手	試合	放送
代表選手	代表選手	監督	サッカー
日本代表選手	日本代表戦	サッカー	野球
日本	サッカー日本代表	一人	代表
日本代表戦	日本戦	一戦	監督
代表戦	日本女子代表	野球	チーム
女子代表選手	女子日本代表	代表	番組
選手ら	選手	ファン	五輪
一日	日本中	出場	プロ野球
一日一日	女子野球日本代表	番組	W杯
同選手	代表戦	優勝	一戦
リー選手	五輪代表	世界	一人
中一	日本サッカー	番組内	番組内
一選手	日本代表監督	同番組	サッカー解説者
サッカー選手	野球選手	試合後	サッカー日本代表
野球選手	女子代表選手	プレー	出場

※ FLR法はいずれもLuceneによる近似実装



# 参考文献

[1] 出現頻度と接続頻度に基づく専門用語抽出

中川、湯本、森 自然言語処理 Vol. 10 No. 1 Jan 2003

[2] Lucene 4.0 Javadoc

[http://lucene.apache.org/core/4\\_0\\_0/index.html](http://lucene.apache.org/core/4_0_0/index.html)

[3] livedoor ニュースコーパス

<http://www.rondhuit.com/download.html#ldcc>