

重要複合語の自動判別方式の評価

類似例尺度を重要複合語収集というタスクに応用し、その有用性を検証した。類似例尺度がある閾値を越えるものを探索対象から除くというヒューリスティクスを導入することにより、作業がどの程度効率化され、また副作用としてどの程度の重要複合語のロスが見込まれるかを実験で評価した。

1 実験方法

今回の実験には日経新聞(CD-ROM版)97年[4]の約20万記事のタイトルと本文を利用し、形態素解析器ANIMA[3]を施して単語+品詞の列に分解し、品詞が名詞または未知語で構成されるすべてのバイグラムとその頻度をカウントした。形態素総数は44M個、(名詞/未知語)バイグラムの総数は6.9M個、種類数は965,923種類であった。

[第1段階:対数尤度による順位付け]最初に対象バイグラムを簡略化版の対数尤度をスコアとして順位付けを行った。その定義は以下で与えられる。

$$LL(AB) = fr(AB) \cdot \log \frac{fr(AB) \cdot N}{fr(A) \cdot fr(B)}$$

主尺度として対数尤度(の簡略版)を用いることにした理由はDailleらによる比較実験で相互情報量に対する優位性が示されていたからである[1, 2](1994)。上位10万個を取り、それらに対して類似例尺度を計算した。計算に要した時間はCompaq Alpha server(300MHz)で3時間弱であった。

[ヒューリスティクスによるふるい落とし]次に以下に示すようなヒューリスティクスによりふるい落としを行い、さらに両方の構成語が2文字以上という条件を付けて残ったものの上位1000個(対数尤度による順位)を実験対象として用いた。

構成語が2文字以上という条件を付けた理由は、1文字語を含む場合には形態素解析誤りの場合が多数含まれていて今回の実験の目的には適さないと判断したためである。

否定性因子としては、強汎用語因子(左右2種類)、部分性因子(左右2種類)、分解表現への換言性因子(補う助詞の種類で5種類)の合計9種類を用いた。以下各々の否定性因子について説明する。本実験とは直接には関係しないので詳しい説明を省略する。

2 予備実験

初めに予備実験として類似例尺度が陽性のものと陰性のものを取り、それらが主観的な重要性判断とどのように関係するかを調べた。

類似例因子の陰陽の判定基準となる閾値を0.25に設定した。その場合、上記のようにして選んだ1000個のバイグラムには類似例因子が陽性となるもの(類似例尺度が0.25以上)が159種類、また陰性のものが841種類ということになる。

表 1. 類似例尺度の閾値を0.27とし、それを越える（類似例因子が+）のパイグラムとそれ以下のもの各30個に重要度の主観判断を下し、主観判断別にまとめた結果

主観判断	類似例因子	A-B	類似例尺度	対数尤度	主観判断	類似例因子	A-B	類似例尺度	対数尤度
+1	-	付加-価値	0.11	10001	0	-	信用-取引	0.19	2472
+1	-	医療-保険	0.24	8520	0	-	経常-収益	0.22	2128
+1	-	自己-破産	0.13	4865	0	-	水道-メーター	0.11	1851
+1	-	数値-目標	0.19	4364	0	-	商品-ファンド	0.19	1691
+1	-	衛星-放送	0.21	4310	0	-	財政-支出	0.23	1539
+1	-	カラ-出張	0.12	3594	0	-	実効-税率	0.16	1446
+1	-	在宅-介護	0.16	2968	0	-	航空-貨物	0.23	1431
+1	-	使い-勝手	0.11	2133	0	-	宇宙-飛行士	0.07	1260
+1	-	コスモ-石油	0.13	1237	0	-	パネル-討論	0.14	1120
+1	-	所得-隠し	0.13	944	0	-	下値-不安	0.13	1119
+1	+	第一-勸銀	0.29	883	0	-	中華-料理	0.15	1119
-1	-	ク リ ン ト	0.24	6469	0	-	緊急-融資	0.23	1094
		ン-大統領			0	-	土砂-崩れ	0.05	893
-1	-	全会-一致	0.03	2096	0	+	増収-増益	0.30	7050
-1	-	和平-交渉	0.23	1309	0	+	為替-相場	0.28	5469
-1	-	特例-措置	0.19	1141	0	+	部品-メーカー	0.32	5254
-1	-	シ ア ヌ ー	0.08	1114	0	+	衆院-本会議	0.30	2260
		ク-国王			0	+	金融-収支	0.27	1847
-1	-	休業-日数	0.19	1003	0	+	電機-メーカー	0.27	1815
-1	+	医療-機関	0.31	5566	0	+	債券-先物	0.31	1684
-1	+	同日-午前	0.29	4209	0	+	家電-メーカー	0.32	1665
-1	+	ポイント-改	0.29	2422	0	+	当期-利益	0.36	1600
		善			0	+	回収-作業	0.31	1532
-1	+	処理-能力	0.31	2079	0	+	給与-振り込み	0.30	1342
-1	+	メーカー-各	0.30	1972	0	+	行革-本部	0.29	1161
		社			0	+	調整-局面	0.28	1021
-1	+	ドル-前後	0.37	1701	0	+	株式-含み益	0.31	919
-1	+	財政-状況	0.28	1458	0	+	買い-意欲	0.31	900
-1	+	流通-業界	0.38	1352	0	+	自民党-執行部	0.30	894
-1	+	前期-推定	0.27	1278	0	+	税率-引き下げ	0.36	887
-1	+	推進-会議	0.34	1092	e	-	共産-党大会	0.21	1357
-1	+	除去-作業	0.28	1017					
-1	+	チェルノムイ	0.30	953					
		ルジン-首相							

陰陽各30個ずつをランダムに取り、それらがどちらであるか分からないように混ぜ合わせて並べ、筆者の判断で重要複合語として取りたいか、捨てたいか、どちらでもよいかで+1, -1, 0を降った。少しでも迷った場合はすべて0としたので0の数が約半数になった。

表1は主観判断別に結果をまとめたものである。第2コラムの+/-は類似例因子の陰陽を示す。

主観的に重要と判定したもの11個の内、10個は類似例因子が陰性であり、類似例因子がプラスで重要と判定されたのは「第一/勧銀」一つであった。この結果は類似例因子がプラスのものを捨ててしまってもロスが少ないこと、すなわち類似例因子を否定的なフィルターとして利用できそうなことを示している。

3 作業量削減効果と重要語損失率の関係

類似例因子の閾値を小さく設定すれば、それだけ類似例因子がポジティブになりやすくなり、探索対象から除かれるバイグラムの個数が大きくなる。またその副作用として多くの重要な複合語が失われる可能性もある。

図1は類似例因子の閾値 θ を変化させた場合の作業量削減率と重要複合語損失率の推移を示したグラフである。

このグラフで重要なことはグラフの最初の部分がフラットになっていることである。すなわちある程度までは作業量削減しても損失が非常に軽微に押えられることをしめしている。例えばこの場合のフラットになっている右端は閾値を0.23に設定した場合で、34%の削減率で損失が0%であることを示している。しかしもう少し安全に見積もって閾値を0.25に設定した場合でも、26%の削減率を達成している。

4 結言

4.1 成果のまとめ

コーパスから抽出した多数の日本語単語バイグラムを対象とし、統計的なスコアによる順位付けの後、さらに重要な複合語の濃度を高める手段として、種々のヒューリスティクスに加えて顕著な類似例のある複合語は平凡である可能性が高いので除くという方法を提案した

新聞から抽出した単語バイグラムのセットを対象とした実験によりその効果を確認した。統計的スコア上位1000位を対象とした場合、重要複合語をほとんど損失することなく25~30%の作業量を削減できた。

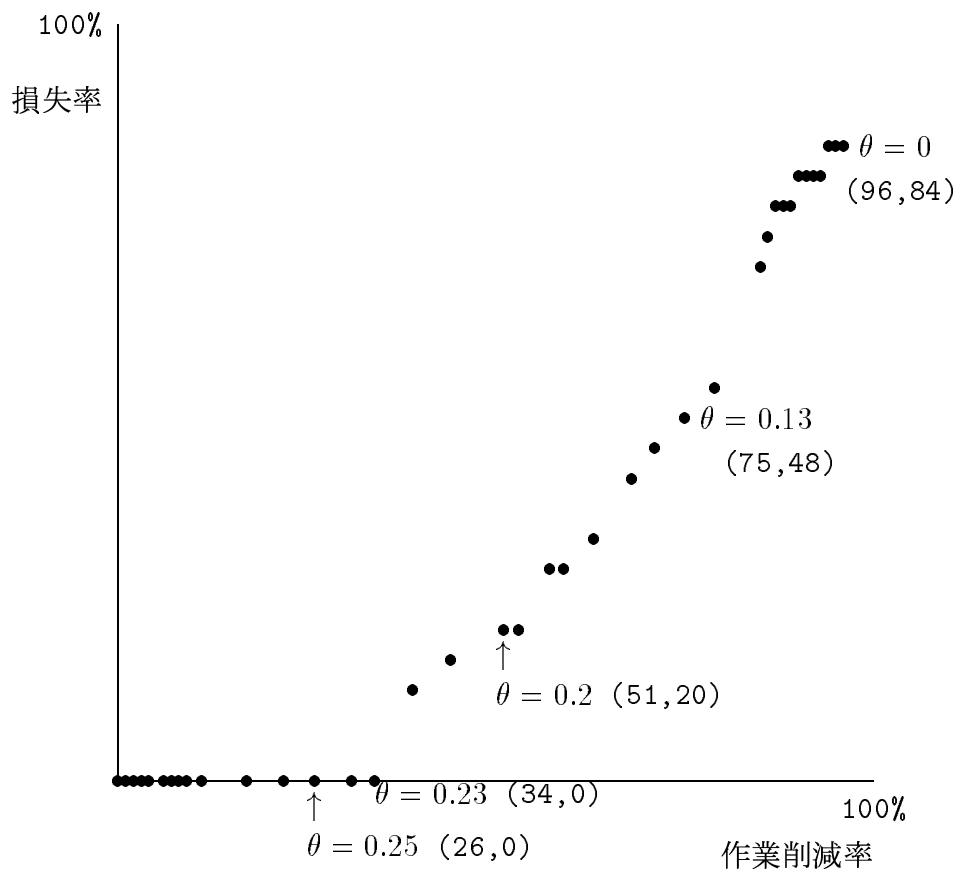


図 1. 類似例因子の閾値(θ)を変えた時の作業量削減効果と重要複合語損失率の推移

4.2 今後の課題

本提案手法は単語類似度の測り方に大きく依存する。従ってその最適化が課題となる。本実験ではバイグラム の右側にくる語としての類似度を計算する場合には左側に伴う語の分布のみを用いた。

しかしながら純粋に単語間の類似度を測る目的であれば、左右片方より両方用いた方が(さらには遠距離の共起関係なども用いた方が)良い結果がでるはずであろう。

ただし、本研究の目的である重要な複合語を抽出するという目的に用いる場合には、どちらが良い結果になるかは明らかとは思われず、実験的に確認する必要がある。

共起分布とシソーラスベースの類似度を用いた場合の比較も興味深い。後者はネットワークタイプの語彙データベース上での単語間距離に基づく類似度であるが、その場合未登録語の影響が不利に働くのではないかと予想される。

5 参考文献

- [1] Béatrice Daille. Study and implementation of combined techniques for automatic extraction of terminology. In *The Balancing Act (Workshop at the 32nd Annual Meeting of the ACL)*, pp. 29–36, New Mexico, 1994.
- [2] Béatrice Daille, Éric Gaussier, and Jean-Marc Langé. Towards automatic extraction of monolingual and bilingual terminology. In *Proceedings of COLING'94*, pp. 515–521, Kyoto, 1994.
- [3] Hirofumi Sakurai and Toru Hisamitsu. A data structure for fast lookup of grammatically connectable word pairs in japanese morphological analysis. In *Proceedings of ICCPOL'99*, pp. 467–471, 1999.
- [4] 日本経済新聞社. 日本経済新聞CD-ROM 1997年版, 1998.