

教育用語彙選定における特徴語抽出及び Wordplot の利用

中條清美¹ 内山将夫² 中村隆宏³ 西垣知佳子⁴

1 日本大学 2 情報通信研究機構 3 小学館 4 千葉大学

1 背景

英語の教授・学習に使える時間は限られている。英語の語彙学習においてはすべての語を習得することは不可能で、また習得する必要もない。そこで、限られた時間を有効に活用するためには無駄を省いた語彙精選が必要となる。

近年、語彙指導の分野では、効率的な学習・指導には「基本語彙」と「専門語彙」という2種類の語彙表で対応すると効果的であるという考え方がある。例えば、最初に2,000-3,000語程度の基本語彙を学習した後に、学習者の学問領域や職域に合わせて目標分野を限定し、その特定分野の専門語彙を指導することで学習効率を高めるといったものである¹。

そこで、語彙教材の開発に、基本語彙と専門語彙という2種類の教育用語彙の選定作業が不可欠となる。基本語彙は「どの分野にも広く出現する語彙」であり、実際に書かれたり話されたテキストを大量に集めた British National Corpus (BNC)² のような一般分野コーパスから高頻度語を選定するという手法が有効である。

一方、ビジネスや科学技術といった「特定分野に特徴的な語彙」である専門語彙の選定に、基本語彙と同様の手法を用いると、中・低頻度の「専門性の高い語」が選定されにくいという問題がある。専門語彙の場合、当該分野の頻度順リストを見るだけでは不十分であり、一般分

野の語彙リストのような「特徴のない語彙」と対照させて、両者の出現頻度数に大きな差のある「特徴のある語彙」を抽出する方法が有効である³。

さらに、学習者の英語習熟度は一様ではない。そのため教育用の専門語彙選定には、初級・中級・上級といった学習者の習熟度レベルに合致した段階的な専門語彙表が求められる。

以上のような状況に鑑みて、我々は、専門語彙をなるべく簡便で客観的に、しかも学習者の語彙レベルに応じて選定するツールの開発を行ってきた。そのために、1) 9種の統計指標を用いて特徴語⁴を抽出する方法¹⁾、さらに語彙の相対的な専門度を視覚的に把握しながら特徴語を抽出する方法として、2) 単語散布図 (Wordplot) を提案した²⁾。

本稿では、これら2種類の手法を BNC の Commerce 分野と Corpus of Professional English (CPE)⁵ に適用して、ビジネス分野および科学技術分野の特徴語を抽出し、その抽出結果の有効性を検討した。

2 研究方法

2.1 特徴語抽出に使用した言語資料

特徴語の抽出には、以下の専門分野および一般

¹ 例えば、Nation (2000:187).

² <http://www.natcorp.ox.ac.uk/World/HTML/thebib.html>

³ 例えば、竹蓋幸生 (1981) 『コンピューターの見た現代英語』東京: エデュカ出版

⁴ 本稿では、以下では、頻度分布における何らかの特徴に基づいて抽出された語彙を「特徴語」と定義する。

⁵ <http://www.perc21.org/>

分野コーパスにおける単語の出現頻度付き語彙リストを用いた⁶。

1) 専門分野の言語資料

① BNCのCommerce分野726万語より作成した頻度100以上の2,973語(延べ5,883,249語)

② CPE科学技術コーパス2,000万語より作成した頻度100以上4,467語(延べ13,527,303語)

2) 一般分野の言語資料

BNC全体より作成した頻度100以上の13,956語(延べ86,008,037語)^[3]

2.2 特徴語の検討に使用した言語資料

1) 英語母語話者の語彙習得学年リスト⁷

2) 基本語彙リストと専門語彙リスト⁸

2.3 9種の統計指標と特徴語リスト

使用した統計指標は、頻度(Freq)、ダイス係数(Dice)、コサイン(Cosine)、補完類似度(CSM)、対数尤度比(LLR)、カイ二乗値(Chi2)、イエーツの補正公式(Yates)、自己相互情報量(PMI)、マクニマーのテスト(McNemar)の9種である⁹。

⁶ 数詞、固有名詞、略語等は人手で除外した。これらの語は、多目的の言語使用に対応する教育用語彙表には不要と考えられる。我々は、外国語としての英語教育における学習語彙は、現代英語を代表する1億語のBNCにおける頻度100以上の13,956語の範囲内の語彙で十分であると考えている。そこで、専門分野の語彙リストからBNC13,956語に含まれていない語を除外した。

⁷ 学年の決定には1~3年はBasic Elementary Reading Vocabularies (Harris & Jacobson, 1972), 4~16年生はThe Living Word Vocabulary (Dale & O'Rourke, 1981)を参照した。両資料に未収録の語は17年生とした。

⁸ 詳細は3.2参照。

⁹ これらの指標はパラメタa,b,c,d(a:専門リストに単語Xが出現した回数, b:BNCリストに単語Xが出現した回数, c:専門リストの延べ語数-a, d:BNCリストの延べ語数-b)によって計算される。各指標の詳細と定義式は文献^[14]を参照されたい。

以上の統計指標を用いて、Commerce分野およびCPEリストの各語の出現状況を、BNCでの当該語の出現状況と比較してその語の特徴度を示す指標値を求めた。指標値にしたがって降順にソートし、上位500語から、ビジネス分野と科学技術分野の特徴語リストを9種ずつ作成した。

2.4 Wordplotと特徴語リスト

専門分野の英語語彙を一般分野の英語語彙に対照させて可視化する方法として、単語散布図を利用するWordplotを提案した。この単語散布図では、横軸が一般分野における単語の出現頻度(Fglobal)と割合(Pglobal)を表し、縦軸が専門分野における単語の出現頻度(Flocal)と割合(Plocal)を表す。

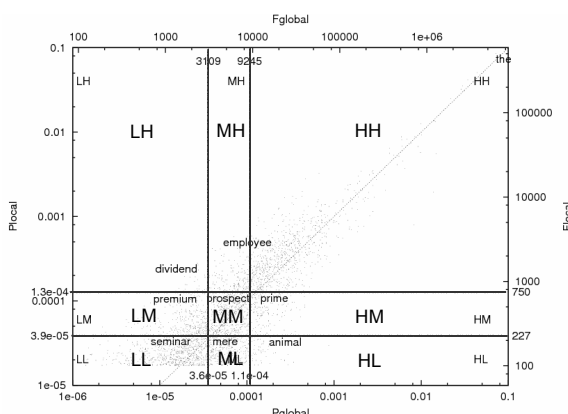


図1 Commerce分野のWordplot(単語散布図)

図1は、横軸にBNCにおける出現頻度を、縦軸にCommerce分野における出現頻度をとり、単語数によって各軸を3等分した9領域よりなるWordplot(単語散布図)である。左上の領域には、一般分野での出現頻度が低く、専門分野での出現頻度が高い単語が位置する。Wordplotの9領域に属する語から、Commerce分野とCPEについて9種ずつ特徴語リストを得た。

3 結果

ビジネス分野と科学技術分野の特徴語を2種類の手法によって抽出し、抽出結果の有効性を、学年レベル、基本語彙と専門語彙の含有率の観点から検討した。両分野の特徴語抽出結果にはほぼ同様の傾向が見られたので、以下では紙幅の関係で、Commerce分野について報告する。

3.1 学年レベル

英語母語話者の語彙習得学年リストを用いて、Commerce分野の特徴語を容易に理解できるのは、米国の何年生ぐらいなのかを調査した結果を表1に示した。表1の左半分は、9種の統計指標による特徴語上位500語の平均学年を指標ごとに示す。表1の右半分は、Wordplotの9領域に配分された特徴語の平均学年を領域ごとに示した。

表1 Commerce分野の特徴語の平均学年

統計指標	対象語数	平均学年	Wordplot 領域	対象語数	平均学年
Freq	500	3.3	HH	771	3.5
Dice	500	3.3	HM	164	3.1
Cosine	500	4.3	HL	56	2.4
CSM	500	5.0	MH	191	6.2
LLR	500	6.1	MM	506	5.8
Chi2	500	6.5	ML	293	4.8
Yates	500	6.5	LH	29	8.2
PMI	500	9.0	LM	320	8.9
McNemar	500	10.2	LL	641	9.0

表1左では、統計指標によって異なる学年レベルの語が抽出されていることが分かる。平均すると、Freq, Dice, Cosine, CSMの特徴語は母語話者の3~5年生レベル、LLR, Chi2, Yatesは6年生レベル、PMI, McNemarは9~10年生レベルに相当する語を抽出した。

表1右に示したWordplotも領域ごとに異なる学年レベルの語を抽出している。Wordplotを縦に3区分した右列にあたるHH, HM, HL(図1の区分参照)は母語話者の2~3年生レベル、中列にあたるMH, MM, MLは4~6年生レベ

ル、左列にあたるLH, LM, LLは8~9年生レベルであった。

以上の結果から、9種の統計指標とWordplotの両手法ともレベル別の語彙を適切に選定し分けており、統計指標ではFreqからCSMは初級用、LLRからYatesは中級用、PMIとMcNemarは上級用に、WordplotではHH, HM, HLは初級用、MH, MM, MLは中級用、LH, LM, LLは上級用の語彙選定に適することが判明した。

3.2 基本語彙と専門語彙の含有率

各特徴語リストにどの程度、基本語彙と専門語彙が含まれるかの割合を調査した。本稿では基本語彙リストとして次の4種を用いた。① Function words(文法的な役割を持つ機能語)¹⁰、② Basic vocabulary (General Service List¹¹のように学習者が最初に学ぶべき基礎的な語彙)、③ Textbook vocabulary(日本人学習者が中学・高校教科書を通して習得する語彙)¹²、④ Academic vocabulary(大学生が一般教養科目を理解するために必要な基本語彙)¹³。そして、専門語彙の例としてBusiness dictionaryの見出し語リストをBusiness vocabularyとして用いた¹⁴。

図2には、9種の統計指標によって抽出された特徴語上位500語に4種類の基本語彙がどの程度含まれるか、また専門語彙であるBusiness dictionaryの見出し語との一致度を調べた結果を示した¹⁵。図3にはWordplotの9領域に分布した特徴語についての結果を領域ごとに示した。

¹⁰ Nation^[5]のFunction words(pp.430-431)を使用した。

¹¹ West, M. (1953). *A General Service List of English Words*. Longman.

¹² *New Horizon English Course 1, 2, 3* と *Unicorn English Course I, II, Reading* より作成した3,245語(延べ語数38,937語)

¹³ <http://www.vuw.ac.nz/lals/research/awl/awlinfo.html>

¹⁴ *Longman Business English Dictionary* (Pearson Education Limited, 2000)の見出し語4,565語

¹⁵ FreqとDice, Chi2とYatesは同じ値であったので、図2では同軸上に示した。

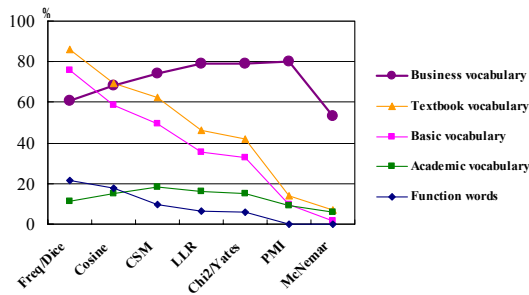


図2 統計指標による特徴語の基本・専門語彙含有率

図2より、LLR、Chi2、Yates、PMIの4指標の上位500語の約8割がBusiness dictionaryの見出し語と一致する。一方、FreqからMcNemarにかけて基本語彙の含有率は段階的に減少している傾向が見られることから、PMIに専門性の高い語彙が抽出されていることが推定される。実際、PMIの特徴語上位には、lading, buyout, arbitrage, subcontractor, liquidity, volatilityなど専門的な経済用語が順位付けられている。

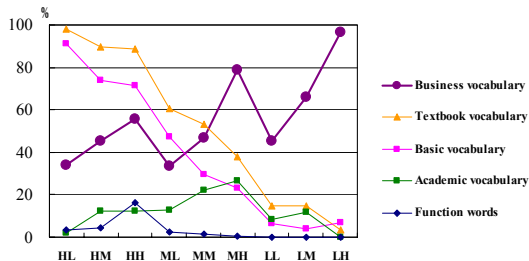


図3 Wordplotによる特徴語の基本・専門語彙含有率

図3に見られるように、LHの領域語のほぼすべてがBusiness dictionaryの見出し語と一致する。LHにはdividend, equity, discount, monetaryやtakeover, audit, mergerのような経済用語が含まれる。次いでMH、LMの順に一致度が高い。MHに含まれる語は、asset, credit, stock, cash, loan, debtやemployee, consumer, investor, buyer, shareholderのようなビジネス用語である。LMでは専門語彙の含有率は少し下がるが、premium, aggregate, lease, stake, innovation,

monopoly, incentiveなどの経済英語が含まれる。これらの領域では基本語彙の含有率は低く、高い精度で専門語彙を取得できることがわかる。

3.1の結果と併せて考察すると、統計指標のうち中級用にはLLR、Chi2、Yatesが、上級用にはPMIが高い精度で専門語彙を選定し、Wordplotでは中級用にはMH、上級用にLH、LMの領域に入る語を選定すれば高い精度で専門語彙を選定できることが判明した。

4 おわりに

本稿では、ビジネス分野と科学技術分野の専門語彙を対象として統計指標とWordplotを使用して特徴語を抽出し、その専門語彙抽出の有効性を検証した。本稿で検討した2種類の選定手法を適切に選択し利用することで、学習者に応じて、初級・中級・上級といった段階的な専門語彙を効率的に選定できると考える。

参考文献

- [1] Chujo, K. and Utiyama, M. (2006). Selecting level-specific specialized vocabulary using statistical measures. *System*, 34 (2), 255-269.
- [2] Utiyama, M. and Chujo, K. (2007). Linking word distribution to technical vocabulary. *Journal of the College of Industrial Technology, Nihon University*, 40 (in press) .
- [3] Chujo, K. (2004). Measuring vocabulary levels of English textbooks and tests using a BNC lemmatised high frequency word list. In: J. Nakamura, N. Inoue, and T. Tabata (Eds.), *English Corpora under Japanese Eyes* (pp. 231-249). Amsterdam: Rodopi.
- [4] 内山将夫・中條清美・山本英子・井佐原均 (2004). 「英語教育のための分野特徴単語の選定尺度の比較」『自然言語処理』11 (3): 165-197.
- [5] Nation, I. S. P., (2001). *Learning Vocabulary in Another Language*. Cambridge: Cambridge University Press

