

英語教師のための コーパス入門

第3回

コーパスからの 頻度リストを活用する



明海大学 助教授
投野由紀夫

はじめに

さて、いよいよ自分の手元にあるコーパス・データを利用して何ができるのかを具体的に見ていくことにしたい。なお、今回から基本的にコーパスの処理は Oxford University Press から発売されている WordSmith Tools または Michael Barlow 氏の開発した MonoConc Pro というプログラムを用いて行う。WordSmith の入手方法およびインストール方法に関しては、作者の Mike Scott のページ (<http://www.liv.ac.uk/~ms2928/index.htm>) を参考にしていただきたい。デモ版がダウンロードできるので、まずはそれをもとにどのような処理ができるか、今回の内容をもとにチェックしてみるとよいだろう。MonoConc Pro は Barlow 氏の webpage を参考にしていただきたい。

コーパスから頻度リストを作成する

コーパスから得られる最も有益な情報の1つは、「単語リスト」だ。そのサンプル全体がどのような語彙によって成り立っているか、というデータ全体の鳥瞰図を「単語リスト」は示してくれる。早速、WordSmith を使って単語リストを作成してみる手順を紹介しよう。

■ WordSmith で単語リストを作成する

- 1) Tools Controller から [Tools] - [Wordlist] を選択。
- 2) Wordlist の画面で左隅の緑の丸のアイコンをクリックする。
- 3) Getting started の画面で [Choose texts now] を選択。
- 4) ここで任意のファイルを選択する。
- 5) Getting started の画面に戻ったら、[Make a wordlist now] のボタンをクリック。
- 6) 自動的にワードリストが作成される。

図1が WordSmith によって作成された British National Corpus の書き言葉部分の単語リストである。

Rank	Word	Count
1	the	5,789,039
2	of	2,937,837
3	and	2,430,343
4	to	2,432,815
5	in	2,032,599
6	a	1,850,041
7	that	999,267
8	is	869,648
9	for	833,322
10	with	823,258
11	he	772,591
12	it	672,732
13	you	631,073
14	we	630,060
15	she	606,182
16	as	573,485
17	him	514,122
18	her	506,912
19	me	491,975
20	us	419,652
21	my	419,595
22	your	411,957
23	our	408,080
24	its	406,087
25	myself	396,598
26	yourself	395,383

図1：BNC written corpus の单語リスト

異なり語数 (types)、異なり語の総語数に占める割合 (type/token ratio) およびその標準化された値 (Standardized type/token)、平均単語長 (Ave. word length)、文の数 (sentences)、平均文長 (Sent. length)、段落数 (paragraphs)、段落の長さ (Para. length)などの詳細な情報がコーパスごとに瞬時に得られるので、大変便利だ。

Rank	Word	Count
1	the	550,561,493
2	of	90,748,630
3	and	87,324
4	to	0.42
5	Standardized Type	44.04
6	Ave. Word Length	4.69
7	Sentences	4,698,431
8	Ave. Sentence Length	10.05
9	Paragraphs	1,573,265
10	Ave. Paragraph Length	57.55
11	Ave. Headline Length	10.53
12	Ave. Headline Length	0

図2：コーパス統計情報のファイル

図1の单語リストは頻度順のリストであるが、WordSmith ではこれ以外にアルファベット順のリストと、コーパス全体の総語数などの統計値をまとめたファイルを作成してくれる。図2はその統計情報をまとめたファイルである。

は、どのような单語がより頻繁に用いられるかをデータによって直接確かめることによって、学習語彙の選定に役に立てる、ということである。「学習語彙の選定などは、文科省とか教科書の著者がすることだ」と思うかもしれない。しかし、コーパスを用いることで、中高の教員が自らの手で選択したコーパス・データをもとに語彙表を作成し、重要度に関する明確なイメージを持ったり、教材の配列や選択を考えたりということが可能な時代になってきている。実際に試してみれば、それほどむずかしいことではなく、誰もが利用できるのだ。

たとえば、Brown Corpus と LOB Corpus のデータが手元にあれば、100万語のアメリカ英語とイギリス英語の語彙リストが瞬時に作成できる。これを合成して両者に共通な基礎語彙リストを作れば、英米の違いを加味したリストを自分なりにアレンジできる(注:合成するためには2つのコーパスを同時にWordSmithなどの検索ソフトにかけねばよい)。これにBritish National Corpus 1億語のデータを加えて、3つの標準コーパスに載っている上位2000語くらいを抜き出せば、それだけでかなり一般的な基本語リストができ上がる。

さらに品詞タグ (POS tag) を使って語彙リストを作成できれば活用度は飛躍的にアップする。品詞タグ付コーパスは単語1つ1つを品詞解析プログラム (POS tagger) によって品詞タグ付与を施したもので、BNCはその代表的なものだ。

BNCから单語リストを作成することで、たとえば品詞ごとの頻度表を作成できる。こうすることで、

単純な頻度リスト1000語を持っているよりも、そこから基本動詞リストだけを抜き出したり、最もよく用いる形容詞100を抽出して早めに指導する、といったようなことが可能になる。図3はBNCからの頻度リストをエクセルに移し変えて、オートフィルタという機能を用いて品詞ごとのリストを自動的に作成できるようにしたものである。このような工夫を施すことで、单語リストそのものが非常に利用価値の高いものになる。

Rank	Word	Count
1	is	14,063 last
2	other	13,518 other
3	new	11,552 new
4	good	10,655 good
5	old	6,6999 old
6	great	6,4369 great
7	high	5,2703 high
8	small	5,1626 small
9	different	4,9373 different
10	large	4,7185 large
11	local	4,4920 local
12	social	4,1617 social
13	long	4,0492 long
14	important	3,9295 important
15	young	3,7278 young
16	national	3,7231 national
17	possible	3,4178 possible
18	big	3,3200 big
19	right	3,1630 right
20	early	3,1110 early
21	public	3,0883 public
22	only	3,0775 only
23	able	3,0454 able
24	political	3,0366 political

図3：BNC品詞付リストで高頻度形容詞を見る

単語リストをどう利用するのか？

このような单語リストを実際にどのように中高の英語教員が利用できるのか、考えてみることにしよう。单語リストのおもな活用分野として、今回は以下の3つの点を中心に解説をしてみたい：

- (1) 学習語彙の選択・重み付け
- (2) プロダクション用語彙の選定
- (3) 英作文指導と評価の資料

では、早速実際の応用例を見てみることにしよう。

(1) 学習語彙の選択・重み付け

コーパスから单語リストを作ることのメリットの1つ

英語教師のための コーパス入門

(2) プロダクション用語彙の選定

語彙リストが役に立つ第2の例として、英作文とか英語スピーチのようなプロダクション活動の場面を考えてみよう。最近はライティングの授業もコミュニケーションの手段や目的を重視したライティングや、1文ずつの和文英訳ではなく一定のまとまった英文を書かせるような発展的なライティングが志向されている。スピーチでも、ある特定のテーマやトピックに関して、まずライティングで基本的なアイディアをまとめ、発表活動を行ったりすることが多い。この際に問題になるのは、その課題となるテーマに関するトピック語彙が不足して表現できない、という問題だ。多くの場合、教師は英作文のテキストのモデル・バラグラフの中のボキャブラリーで何とかまかなおうとするのだが、実際の作文ではより広範囲の語彙を与えておかなければ、とても対応できない。

こんなときに威力を発揮するのが、自分のクラスで書かせた作文データをコーパス化しておくことである。たとえば、1年目にライティングで書かせたいトピックをシラバスとして組み、1年がかりで自由英作文を書かせるとしよう。それらのトピックごとの作文を夏休みなどの余裕のあるときを見て、コンピュータに入力する。(もちろん、コンピュータ上で書かせる作業をできればこの手間は省ける。)こうしていくつかトピック別の作文データがミニコーパスとして集まつたら、個別に前述のような方法で単語リストを作る。ただ、これだけだと単純なリストなので何がトピック語彙として重要なのかわからない。そこで登場するのがキーワード分析(Keyword Analysis)というWordSmithの機能である。これは簡単に言うと、自分の作ったミニコーパスと基準になる大きなコーパス(たとえばBrownとかBNCなど)の2つを比較して、基準になるコーパスに比べて極端に高頻度に出てくる単語をそのミニコーパスの特徴を示すキーワードとして抽出する方法である。

■ WordSmith でキーワードを抽出する

- 1) Tools Controller から [Tools]-[Wordlist] を選択し、それぞれ2つのコーパスについて単語リストを作成する。
- 2) 続いて、[Tools]-[Keyword]を開く。
- 3) 2種類の単語リストを選ぶウインドウが現れるので、左側に生徒の英語コーパス、右側に基準となるBNCなどのコーパスの頻度リストのファイルを選択する。
- 4) [OK]を押すと、双方のリストを比較した結果が出力される。

さて、実際のアウトプットは図4のようになる。

The screenshot shows a Microsoft Excel-like spreadsheet titled 'Keyness' with columns labeled: IDN, Document, WORD, FREQUENCY, DOCUMENTS, WEIGHT, KEYNESS, and PERCENTAGE. The data consists of two rows of words and their metrics. The first row is for '朝食' (breakfast), and the second row is for 'rice'. The 'KEYNESS' column values are notably higher for these words compared to others like 'eat' and 'bread'.

IDN	Document	WORD	FREQUENCY	DOCUMENTS	WEIGHT	KEYNESS	PERCENTAGE
3,153	朝食	朝食	76,514,132	0,57	12,524,3	8,574,9	
3,152	朝食	rice	4,166	0,57	8,467,8		
3,153	朝食	rice	1,932	1,637			
3,154	朝食	eat	774	2,653	5,816	7,633,1	
3,155	朝食	bread	665	2,653	3,256	6,957,0	
3,156	朝食	bread	634	2,657	16,348	4,750,3	
3,157	朝食	eat	591	1,638	4,345	3,561,7	
3,158	朝食	bread	226	0,72	11	3,527,0	
3,159	朝食	rice	1,123	5,529	20,932	9,45	2,681,4
3,160	朝食	rice	427	1,52	17,685	0,92	2,802,3
3,161	朝食	rice	215	0,72	1,230	2,225,2	
3,162	朝食	rice	269	1,52	55,220	0,92	2,129,6
3,163	朝食	rice	516	1,638	119,685	0,12	1,680,1
3,164	朝食	rice	152	0,52	1,759	1,371,7	
3,165	朝食	rice	265	0,33	35,665	0,04	1,073,3
3,166	朝食	rice	442	1,47	169,710	0,21	980,8
3,167	朝食	rice	57	0,12	0	0,11	
3,168	朝食	rice	131	0,42	4,919	0,14	683,5
3,169	朝食	rice	570	1,52	395,393	0,44	787,7
3,170	朝食	rice	285	0,22	98,094	0,11	725,7
3,171	朝食	rice	158	0,52	18,471	0,02	718,1
3,172	朝食	rice	254	0,33	60,485	0,09	683,9
3,173	朝食	rice	112	0,37	7,043	0,02	642,9
3,174	朝食	rice	293	0,35	129,018	0,14	620,8
3,175	朝食	rice	35	0,12	0	0,02	599,8
3,176	朝食	rice	98	0,32	6,282	0,09	588,8
3,177	朝食	rice	182	0,62	50,677	0,06	532,0
3,178	朝食	rice	134	0,24	31,110	0,03	436,0
3,179	朝食	rice	27	0,29	0	0,02	431,8
3,180	朝食	rice	25	0,09	0	0,02	399,8

図4 : Keyword 分析の結果画面

Keynessという部分のスコアが高ければ高いほど、生徒の書いた作文中で一般の英文に比べて特に高頻度に出てきている単語と言える。図4の例では「朝食」に関する作文なので、リストを見ると朝食に関して表現したい語彙(breakfast, rice, eat, bread)が上位に来ていることがわかるだろう。

より精密にするためには、このリストからいわゆる機能語と基本語のリストを除いてしまう(これはstop wordといって普通の単語リストから上位の機能語・一般語をピックアップしてファイルにしておき、フィルタに使うわけだ)。そうすると、トピック語彙が鮮明に見えてくる、というわけだ。

この方法を用いて、1年に5~6本程度の作文の課題をもとにトピック語彙を抽出できれば、翌年の授業に使える。今度は、トピック語彙のリストを与えて書かせてもいいし、それらのトピック語彙のリストとコーパスを両方用意しておき、実際にコーパスで実例を検索させながら指導してもよい。いずれにしても、もっと具体的に「書きたい表現」が何かを科学的に知ることができる方法である。

私が中心に収集しているJEFLL Corpusではこの情報からさらに進んで、「どういう単語が英語でできないか?」という情報をコーパス内に組み込んでいる。この情報があれば、語彙リストとして「生徒が英語でできなかった単語リスト」を作ることが可能だ。それらの詳細な分析はまた別の機会に譲りたいが、そのようなリストの中からもトピック語彙として生徒が

知りたい単語を抽出する貴重なリソースとなるだろう。

このような作文のトピック語彙を先生同士で共有すればよい。できればもとになったコーパス・データもシェアできるともっとよい。そうすることにより、英語教師が準備不足で悩んでいるライティングやスピーチ、ディベートなどの活動にボキャブラリーの面で大いに貢献することができるだろう。

(3) 英作文指導と評価の資料

最後に私がこの2年間行ってきたライティングの授業の話をしよう。昨年明海大学に赴任して初めて担当した Writing I の授業では「一定のまとまつた英文を継続的に書かせる」ということを目標にした。コンピュータ・ルームを利用して、テキストはなし。2週間で1本のトピックに関して自由英作文をさせる。その媒体としては、いわゆる BBS (掲示板) のフリーソフトを自分のサーバーに設置して、そこに書き込みさせるようにした。

これだけではあまりおもしろくない。そこで、全員で BBS に投稿後、1人が最低5名の友人の作文を読んで、それに対して「もっと知りたいこと」をコメントすることを義務付けた。こうすると最低でも5人くらいから自分の作文に対する peer review を得られるので、revision の際の参考になる。

これだけでもまだおもしろくない。そこで、自分の書いている作文の文法の誤りに関する意識を高めるために、私が各作文に対して問題箇所をアステリスク (*) で印をつけて BBS にコメントとしてアップした。添削はしない。考え方である。

学生はこの私の文法の誤り箇所のマークと、5名の友人からの参考意見をもとに再度自分の作文を書き直す。その結果を最終バージョンとして BBS に再度アップする。

今年度はさらにこの方法から工夫をして、最終的にアップした作文をプリントアウトし、1日10回音読させ、1週間後に抜き打ちで皆の前で作文を見ないで口頭で発表するという宿題を課している。これで、書いたものを口で言えるようになるサポートを何とかしようというわけだ。

この方法は BBS というコンピュータ環境を使ってるために、教師の負荷が非常に少ないので大変効果的だ。同じことを紙と鉛筆でやろうとすると、作文のコピーを作ったり、添削の記録をつけたりとい

教員の単純作業の時間が大幅に増えてしまう。熱心な先生はそれをやっていると思うが。

そして最後に登場するのがコーパスである。学期末になると、BBS のデータを個人データとして個別に整理して一人一人の作文コーパスを作る。そして、その個人の作文の評価に単語リストを活用するのだ。学生全員のデータを単語リストにすることで、WordSmith による個人の作文の統計情報が出る。これで、学期中に書いた作文の総語数はどのくらいかがわかる。1つ1つの作文で書く量が増えたか減ったかもわかる。異なり語数を見ればどのくらい語彙が豊富かもわかる。キーワードを抽出すれば、一般の学生の作文に対して、どういう表現を特に使いすぎているか、どういう表現が出てこないか、といった大まかな語彙使用の状況もわかる。これを全体のコーパスを処理した平均データと併記して、一人一人にコードとして返却する。自分がクラス全体でどの程度作文をがんばったかが実際に鮮明に数値化されて、学生は非常に感心するのである。図5はその学生に返却するデータの一例である。

評価項目	平均	学生A
総語数(半年の作文総数)	1,702	2,436
異なり語数	526	687
語彙の豊富さ (全体)	32	28.2
語彙の豊富さ (作文毎)	36	40.3
平均の文の長さ	7.3	8.9
平均の単語の長さ	4	4.27

(このあとに他の学生との語彙比較リストを添付。)

【コメント】クラス平均以上にたくさん書けた。単語もいろいろな表現を使う努力が見られる。

図5：作文評価の返却例

まとめ

コーパスは研究の道具だけでなく、実際の英語指導と結びつけたときにも威力を發揮する。コーパス言語学といっても、それはデータを解析したり、指導法に活かしていくための「方法論」なのであるから、自分の自由な発想で活用すればよい。今回は単語リストにしぶって話をした。次回はもう1つのコーパス検索プログラムの大きな機能であるコンコーダンスをとりあげて、具体的な用例を見たり調べたりすることで英語教員にとってどのような活用が可能か考えてみたい。