

第3章 汎用コーパス概観

3.1 はじめに

コンピュータ技術の発達に伴って、これまで数多くのコーパスが編纂されてきました。その数は近年急速に増加し、コーパスサイズもますます巨大化しています。Bookmarks for Corpus-based Linguists (<http://tiny.cc/corpora>) では、現在どのようなコーパスが提供されているかの情報一覧を見ることができますが、その数や種類の多さには驚かされます。著作権やその他の理由のため、研究者個人あるいは特定の研究機関でのみ使用が許されるコーパスもありますが、誰でも無料または有料で利用できるように広く一般のユーザーに公開されているコーパスもあります。この章では様々な研究目的に広く対応できるように編纂された「汎用コーパス」について概説します。それぞれのコーパスの特徴を知ることによって目的に応じて使い分けることができます。

3.2 Brown 系コーパス

Brown コーパスは、1961 年当時のアメリカ英語の書き言葉の実態を反映したコーパスを作ることを目指して編纂された、世界初の電子化テキスト集です。様々なテキストジャンルからサンプルが収集され、総計 100 万語のコーパスが作成されました。これがいわゆる「第 1 世代コーパス」と呼ばれるコーパスの原型になり、その後このデザインに倣って、1961 年のイギリス英語版のコーパス Lancaster-Oslo/Bergen Corpus (LOB), Brown/LOB の 30 年後の出版資料を用いた「クローン」とも呼ばれる Freiburg Brown (Frown) と Freiburg LOB (FLOB), 2005-2007 年のア

メリカ英語、イギリス英語を収集した AmE06, BE06 など様々なコーパスが作られてきました。これにより、同じ年代のアメリカ英語とイギリス英語の比較や、年代間の言語変化の調査が可能となりました。

Sketch Engine (3.6 参照) では Brown コーパスが無料で使えます。図 1 は Word Sketch という機能を使って、文法関係単位で動詞 make がどのような語と結びつきやすいか、すなわちコロケーションを表示したものです。例えば [object] という欄には実際に make の目的語として使われる語が頻度の高い順に並べられています。実用的な用例を選ぶのに役に立ちます。

[サイズ] 各 100 万語

[英語の種類] アメリカ英語: Brown, Frown, AmE06 / イギリス英語:

LOB, FLOB, BE06

[利用形態]

The ICAME Corpus Collection (CD-ROM) → 料金: NOK3,500, <http://icame.uib.no/newcd.htm> から注文。* AmE06 と BE06 は含まない。

1,700 万語を超える 20 のコーパスを収録。

Sketch Engine (<https://www.sketchengine.co.uk/>) → 料金: £52.00 (年間),

* 様々な言語の 100 以上のコーパスを提供。Brown の利用は無料。

make (verb) Brown freq = 2322 (1975.0 per million)				
object 1370 5.7	subject 461 2.5	unary_rels	pp_adj_comp 200 44.6	pp_durant 9 7.5
effort 43 8.76	people 5 6.3	reflexive 21 6.3	clear 17 10.1	year 4 5.81
use 25 8.61	man 10 6.29		difficult 13 9.89	
mistake 17 8.61		pre_object 344 9.7	possible 19 9.89	wh_comp 4 0.6
decision 19 8.53	modifier 220 0.4	him 52 9.88	impossible 6 9.19	when 4 8.59
statement 19 8.44	sure 8 10.01	them 40 9.69	easy 7 8.95	
difference 18 8.35	first 7 9.08	me 22 9.36	happy 5 8.89	
contribution 12 8.02	only 14 8.82	it 151 9.18	necessary 6 8.52	
appearance 12 8.01	always 8 8.49	itself 6 8.52	available 4 7.63	
attempt 11 7.89	often 6 8.48	us 9 8.46		
sense 14 7.89	also 10 8.15	her 10 8.4	adj_comp 123 3.2	
payment 11 7.85	even 5 8.02	himself 7 8.25	sure 21 10.88	
money 12 7.77	not 36 7.94	you 24 7.82	pursuant 5 9.85	
remark 10 7.77	already 4 7.89	one 4 7.59	explicit 4 9.72	

図 1 Brown コーパスにおける make のコロケーション

3.3 British National Corpus

1980年代以降、コンピュータ技術の飛躍的な発達により、Brown系コーパスより大きな、いわゆる「第2世代コーパス」と呼ばれる大規模コーパスが編纂されるようになりました。その中で今日もっとも多くのサイトで無料提供されているのがBritish National Corpus (BNC)です。BNCは、1億語を有するコーパスで、1990年代のイギリス英語の実態を反映する均衡コーパスを目指して作成されました。第1世代のコーパスと比べると、規模は100倍になりました。

BNCには10%の話し言葉(1,000万語)が含まれています。その半分は、講義、説教、ラジオ番組などの公共のイベントからの録音、残りの半分は個人の日々の会話を録音したものです。残りの90%を占める書き言葉の部分は、25%のフィクションと75%のノンフィクションから構成されています。コーパスサイズが大きいことや、構成や出典情報に十分な配慮がされていることから、BNCは現在多くの言語学者に使用されています。

CORPUS.BYU.EDUではBNCを語、句、品詞レベルで検索することが可能で、ジャンル指定もできます。また、文法構造やコロケーションの検索、ジャンル間の比較も可能です。このサイトで提供されているその他のコーパス(後述のCOCAやTIMEなど)の検索結果と比較できるというのも魅力です。CORPUS.BYU.EDUの操作方法是第4章で、利用方法については第5章で詳しく述べます。図2はBYU-BNCで動詞banを検索し、ジャンル別の頻度を示したものです。圧倒的多数が新聞で使用されているのが分かります。

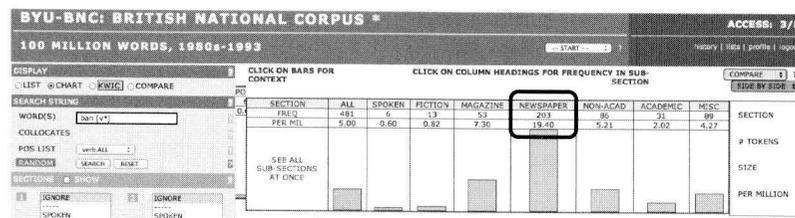


図2 BNCにおけるban(動詞)のジャンル別分布

BYU-BNC以外にBNCはBNCwebやIntelliTextやPhrases in English (PIE)のサイトでも使用できます。それぞれ使い勝手が違うので用途に合わせて使い分けるとよいでしょう。BNCwebは直感的に使用しやすいインターフェースで、非常に多機能です。共起関係を表すさまざまな統計値が用意されているので、特にコロケーションを調べる際に非常に有効です。BNCwebの操作方法是第6章で詳しく述べます。

IntelliTextはLeeds大学のCentre for Translation Studiesによって開発されたコーパス検索用のインターフェースで、BNCの他にも様々な言語のコーパスを検索することができます。検索結果(View Results)の表示機能では、通常のソート機能に加え、Kellyタブを開くとCEFRレベル(欧州評議会が提唱するヨーロッパ言語共通参照枠に基づく5技能6段階の習熟度レベル)別に単語を色分けして表示することができます(CEFRについては17.2.1参照)。図3はenvironmentを検索した結果に中級以上の語に色づけをして表示しています。テキストの難易度の判断に使ったり、レベルに合った語彙を選び、空欄補充問題などを作るのに役立ちます。

PIEではn-gramを検索することができます。「n-gram」とは、あるテキストから切り出した連続したn個の語の並びのことで、その連続した並びの頻度情報などから、意味のまとまりをもったフレーズを抽出してく

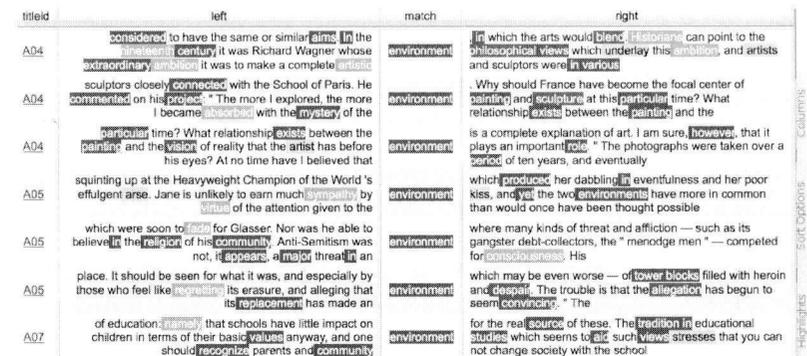


図3 語彙レベル別検索結果

3.5 Corpus of Contemporary American English (COCA)

COCAは上述のCORPUS.BYU.EDUで2008年に公開されたコーパスです。使いやすいインターフェースを持つ優れたツールで、様々な英語の調査が可能です。1990年以降1年ごとに2,000万語のセクションで構成され、現在も毎年2,000万語ずつ拡張し続け、2013年8月現在4億6,000万語を有します。無料で提供されるコーパスの中で最大、かつアメリカ英語コーパスとしては唯一の大規模均衡コーパスと言えるでしょう。データはサイズの等しい5つのジャンル(話し言葉, フィクション, 大衆雑誌, 新聞, 学術誌)に分けられています。

図6ではglobal warmingのジャンル別、年代別頻度を示しています。主に雑誌で用いられ、2005年以降頻度が飛躍的に伸びているのが分かります。また、話し言葉はラジオやテレビ番組の台本のない会話を文字起こしたもので、BNC中の人前で話したスピーチなどの部分とは十分比較可能です。英米語間での話し言葉の比較に使うことができます。COCAの操作方法と活用事例については第4章と第5章を参照して下さい。

CORPUS.BYU.EDUには、1億語を超えるコーパスが複数提供されており、これらのコーパス間での比較が可能なので、言語調査に役立つツールとしてますます広い支持を得ているサイトと言えるでしょう。

[サイズ] 4億6,000万語

[英語の種類] アメリカ英語

[利用形態] COCA (<http://corpus.byu.edu/coca/>) → 無料(登録については4.2)

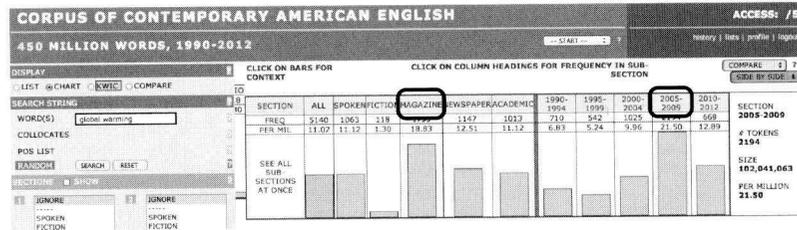


図6 COCAにおけるglobal warmingのchart表示

3.6 Sketch Engine

Sketch Engineはこれまで紹介してきたBNCやCOCAなどと違いコーパスそのものの名前ではありません。数多くのコーパスへのアクセスや多機能な分析ツール、コーパス作成支援などのサービスを提供する有料サイトです。変種を含め50以上の言語の、様々な目的に応じて作られた100以上のコーパスにアクセスすることができます。その中には前述のBrown系コーパスやBNCも含まれます。ここで使用できるコーパスの多くは巨大で、最大のenTenTen12はBNCの約100倍の100億語を突破しました。このような巨大コーパスの特徴は、すべてのテキストがコンピュータのウェブ自動巡回により特定のドメインから集められ、自動でタグ付けやレマ化がされていることです。このような膨大なデータから検索をして用例を抽出できるので、低頻度の語句の用例を見つけることも可能です。

Sketch Engineは非常に多機能で、Word Sketch機能では、文法関係単位で検索語と他の語との結びつきの相性の良し悪し、すなわちコロケーションを表示することができます。それ以外にも、Thesaurus機能では、対象語と同じように用いられる語を検出し、類義語の可能性のある語として提示してくれます。Sketch-Diff機能では2語間の用法の違いが色付けして表示されます。図7(次ページ)は日本語の「心」の訳語として使われる2つの英単語heartとmindがどのように異なる使われ方をしているのかを表現しています。これら2語が主語のとき、どのような動詞が使われているのか(subject_ofの欄)、どのような形容詞が補語になっているのか(adj_subject_ofの欄)を表しています。一方の語で特徴的に見られる共起語が、他方の語とは全く結合していないことが見て取れ、両語は異なる環境で使われることが分かります。

[サイズ] ukWaCは15億語、enTenTen12は100億語など。多数のコーパスを提供。

[英語の種類] 多種。ウェブ上で収集されたデータが多い。

[利用形態] Sketch Engine (<https://www.sketchengine.co.uk/>) → 料金:

heart/mind ukWaC freqs = 216719 / 214661

heart	6.0	4.0	2.0	0	-2.0	-4.0	-6.0	mind
and/or	22956	27656	0.9	1.1				
consciousness	0	144	0.0	5.9				
Body	0	82	0.0	5.5				
body	147	4470	2.7	7.6				
heart	246	2464	4.4	8.2				
emotion	29	222	3.5	6.3				
spirit	224	712	5.0	6.7				
will	70	164	4.0	5.2				
intellect	31	77	4.9	6.0				
brain	374	263	6.4	5.9				
soul	1609	606	8.6	7.1				
conscience	161	54	6.8	5.1				
muscle	136	35	5.4	3.4				
stomach	64	11	5.3	2.7				
mind	2464	332	8.2	4.8				
breathing	43	0	5.0	0.0				
vessel	132	0	5.1	0.0				
spleen	33	0	5.3	0.0				
spade	38	0	5.5	0.0				
gut	50	0	5.5	0.0				
transplant	56	0	5.6	0.0				
artery	65	0	5.8	0.0				
circulation	153	0	6.5	0.0				
kidney	251	0	7.5	0.0				
liver	314	0	7.6	0.0				
lung	852	0	9.0	0.0				
subject_of	21717	23761	1.5	1.8				
boggle	0	470	0.0	9.3				
wander	0	278	0.0	7.5				
numb	0	107	0.0	7.1				
blow	0	278	0.0	6.9				
drift	0	63	0.0	5.7				
reel	0	43	0.0	5.6				
bend	0	51	0.0	5.3				
race	25	92	5.9	5.9				
leap	113	11	6.8	3.3				
stop	364	19	5.6	1.3				
swell	38	0	5.3	0.0				
soar	42	0	5.3	0.0				
flutter	36	0	5.6	0.0				
yearn	37	0	5.7	0.0				
melt	72	0	5.8	0.0				
rend	40	0	5.9	0.0				
wrench	55	0	6.3	0.0				
skip	93	0	6.4	0.0				
thump	69	0	6.6	0.0				
bleed	123	0	6.9	0.0				
pump	180	0	7.1	0.0				
ache	149	0	7.5	0.0				
sink	410	0	7.9	0.0				
pound	233	0	8.1	0.0				
beat	1151	0	8.2	0.0				
adj_subject_of	2988	3251	1.7	2.0				
cluttered	0	9	0.0	5.7				
receptive	0	14	0.0	5.7				
blank	0	33	0.0	5.6				
incapable	0	21	0.0	5.6				
alert	0	18	0.0	5.6				
polluted	0	8	0.0	5.1				
bent	0	11	0.0	5.0				
apt	0	8	0.0	4.6				
sharp	0	38	0.0	4.6				
capable	14	86	2.4	5.1				
pure	69	21	5.4	3.7				
content	33	0	4.6	0.0				
enlarged	8	0	4.8	0.0				
evil	34	0	5.0	0.0				
troubled	11	0	5.0	0.0				
glad	41	0	5.1	0.0				
healthy	113	0	5.2	0.0				
hasty	8	0	5.5	0.0				
faint	22	0	5.9	0.0				
merry	18	0	6.1	0.0				
sore	29	0	6.2	0.0				
steadfast	20	0	7.0	0.0				
thumping	12	0	7.1	0.0				
restless	39	0	7.4	0.0				
deceitful	46	0	8.4	0.0				

図7 Sketch-Diffでheartとmindの類似点・相違点を表示

£52.00 (年間)。ただし Open Corpora として公開されている British Academic Spoken English Corpus (BASE), British Academic Written English Corpus (BAWE), Brown は無料 (3.2 参照)。

3.7 新聞・雑誌・その他のアーカイブ

英字新聞・雑誌をオンラインで無料で検索・閲覧できるアーカイブもいくつか存在します。語法の確認や適当な用例を探すのに役立ちます。

TIME のサイト (<http://content.time.com/time/archive/>) では、TIME の創刊以来の全記事が検索可能になっています (閲覧は一部無料)。また、上述の CORPUS.BYU.EDU では 1923 年から 2006 年までの Time 誌アーカイブに無料でアクセスすることができます (<http://corpus.byu.edu/time/>)。

New York Times (<http://www.nytimes.com/>) のサイトでは、創刊以来の全記事 (1,300 万記事以上) を検索することができます (1923 年以前と 1986 年以降の記事は毎月 10 件まで無料)。

その他、印刷された書籍などを電子化したテキストアーカイブとして、The Oxford Text Archive (OTA: <http://ota.ahds.ac.uk>), Project Gutenberg (<http://www.gutenberg.org>) があります。これらのテキストをダウンロードし、ファイルに整形を施せばコーパスを自作することができます (第 11 章参照)。また自作コーパスを検索するにはコンコーダンスーが必要です。コンコーダンスーの 1 つである AntConc については第 12 章で詳しく紹介します。

3.8 コーパスとしてのウェブサイトと WebCorp Live

近年注目されている言語調査の資料として、ウェブサイトがあります。Google や Yahoo といった検索エンジンを使用して、インターネット上にある情報を丸ごと言語資料として利用するのです。人の手によって編纂されたどのコーパスよりも膨大なデータを持つウェブサイトは情報の宝庫だと言えると同時に、常に変化し続ける不確定なデータベースだとみることもできます。そのため、使用する際には注意しなければならないことがいくつかあります。

第 1 の問題は、ウェブサイトは常に変わり続けているということです。ある検索エンジンを使って今日検索した結果が明日も同様に得られるとは限らないのです。もう 1 度検索すれば同じ結果にたどり着くと過信せず、必要な情報・データは記録・保存しておくことが重要です。

また、データそのものの信頼性の問題もあります。インターネット上では書き手の情報が常に明らかとは限らないので、英語で書かれたデータをすべて英語の母語話者が書いた用例として扱うことはできません。この問題は検索するサイトの範囲を制限することで、ある程度解決できます。例えば、Google の「検索オプション」を使って検索結果をアメリカ合衆国やイギリスなどに絞り込んだり、ドメインを “ac.uk” や “edu” などの

教育機関用ドメインに限定したりすることで、ある程度のデータの制限ができます(14.2.3 参照)。例えば、日本人英語学習者がよく書く英文“My hometown has a lot of nature.”の妥当性を調べるとします。実際に has a lot of nature という表現を検索エンジンで調べると、ドメインの制限なしで検索すれば、約 11 万 8 千件が該当しますから、「使われる」という判断に至るかもしれません。ところが、“edu”にドメインを絞って検索すると、108 件しかヒットしません。しかも検索結果は、日本人留学生のコメントだったり、英文法の授業の資料で悪い例として取り上げられている場合だったりします。このことから、has a lot of nature は適切な表現ではないと判断することができるでしょう。

Google や Yahoo などの検索エンジンは、本来言語調査用に作られたものではなく商業目的で利用されることを前提としているので、検索表示件数に制約があったり、表示順などが操作されたりして、純粋な語句の検索結果として扱うには注意が必要です。以上の制約や問題点を理解した上で使用すれば、商用の検索エンジンもコーパスでの調査の結果を補足する情報が必要な場合や、非常に低い頻度の言語表現を扱う場合などの用例調査に有効に利用できます。なお Google の検索技法と留意点について詳しくは第 14 章を参照して下さい。

ウェブサイトをコーパスとして利用する際に不便な点としては、KWIC 表示などのコーパスツール機能がないため、検索結果を見渡すに困ることが挙げられます。しかし、例えば WebCorp Live (<http://www.webcorp.org.uk/live/>) で提供されている機能を用いれば、KWIC 表示やコロケーションの表示ができるようになります。WebCorp Live の検索画面では、使用する検索エンジン、言語、使用するサイトなどを選択し、検索語の後処理オプションでは、ソート方法やコロケーション表示の有無、使用するウェブページの期間などを指定できます。その結果、AntConc などのコンコーダンサーで検索したような KWIC 表示やコロケーション情報を得ることができます。WebCorp Live による検索は、まだ辞書にも登録されていないような、新語の使われ方を調べるのに有効です。図 8、9 では

図 8 WebCorp Live の検索画面と後処理オプション

```

208:      in the mechanisms for filtering ideas in crowdsourcing systems. "Right now, most companies still are
245:      Reddit has formally apologized for its role in crowdsourcing that resulted in falsely identifying innocent
351:      Subscription links RBS 2.0 Login / Logout Log in Crowdsourcing, Sociopolitical Legitimacy, and the Boston
364:      "the internet" when those ideas were incorrect. Crowdsourcing can't be effective in a fishbowl environment, as
91:      project will study the impact of an innovative crowdsourcing initiative on Your Paintings, an important new
187:      population being connected to the internet. Crowdsourcing research has focused around two areas (with
262:      any special dietary needs in advance. What is crowdsourcing? Researchers in wide range of academic discipline
83:      Institutional staff interested in launching crowdsourcing activities Representatives of funding and
204:      to learn its benefits and limitations. "Crowdsourcing won't replace R&D departments anytime soon,"
202:      got too far from land. Today the Internet makes crowdsourcing much cheaper and easier to implement and extends
353:      this proves that crowdsourcing was a mistake. Crowdsourcing has failed, and the public should never help the
12:      historical photographs in forgotten boxes. A new crowdsourcing project organized by Stanford researchers is
39:      vdsourcing, either by novel ML methods, or on new crowdsourcing problems. Crowdsourcing for machine learning.
8:      describes the intellectual roots of the idea of crowdsourcing in such concepts as collective intelligence, the
10:      organizations. And he considers the future of crowdsourcing in both theory and practice, describing its
11:      humanities scholars harness the power of crowdsourcing Researchers at Stanford's Center for Spatial and

```

図 9 WebCorp Live での crowdsourcing の検索結果

crowdsourcing という比較的新しい語を検索し、結果を KWIC 表示しています(いくつかの後処理を行っています)。これにより、この語がどのように使用されているか、この非常に新しいビジネス形態が新聞などでどのように取り上げられているのか、知ることができ、時事的な英語の解説に利用できます。

3.9 おわりに

一般に公開されている汎用コーパスを使用例と共に紹介してきました。汎用コーパスのサイズはますます大きくなる傾向にあります、大きいも

のが常によいととは限りません。高頻度の語や一般的な構文を調べる場合にはそれほど大きなコーパスは必要ありませんが、頻度はそれほど高くないが、使用パターンやコロケーションを知りたい語や周辺の構文を調べるときには大規模コーパスが必要です。調べたいことに対して正しいサイズのコーパスを使用することがコーパス検索成功の鍵となるでしょう。

第4章 汎用コーパス COCA の使用法

4.1 はじめに

この章では Corpus of Contemporary American English (COCA) の操作方法とその活用法を説明します。COCA は Brigham Young 大学の Mark Davies によって提供されているサイト, CORPUS.BYU.EDU で公開されている7つのコーパスの1つです。操作方法は共通ですので、COCA の検索方法を学べば他の6つのコーパスも同じように扱うことができます。詳細な説明や検索例は、初期画面の [Help/information/contact] のメニューから見ることができます。特に [GENERAL] の Brief tour (general) は簡潔にまとめられているので、COCA を使って何ができるのかを具体例を見ながら手早く学ぶことができます。

表1 CORPUS.BYU.EDU で利用可能な英語コーパス

コーパス名	語数	国・地域	時期
Global Web-Based English (GloWbE)	1.9 billion	20 countries	2012-13
Corpus of Contemporary American English (COCA)	450 million	American	1990-2012
Corpus of Historical American English (COHA)	400 million	American	1810-2009
TIME Magazine Corpus	100 million	American	1923-2006
Corpus of American Soap Operas	100 million	American	2001-2012
British National Corpus (BYU-BNC)	100 million	British	1980s-1993
Strathy Corpus (Canada)	50 million	Canadian	1970s-2000s