

『コーパス日本語学ガイドブック』の概要

(同書の抜粋)

表紙	1
はしがき	2
目次	3
第1部本編の各記事のとびら	4~8
奥付	9

コーパス日本語学ガイドブック

田野村忠温・服部匡・杉本武・石井正彦

特定領域研究「日本語コーパス」日本語学班

2007

はしがき

平成18～22年度文部科学省科学研究費補助金特定領域研究「代表性を有する大規模日本語書き言葉コーパスの構築：21世紀の日本語研究の基盤整備」（領域代表者：前川喜久雄氏、領域略称：「日本語コーパス」）を1つの柱とする国立国語研究所の日本語書き言葉コーパス構築5か年計画が現在進められている。これにより、2011年には待望の現代日本語の大規模な均衡コーパスをついに誰もが利用できるようになる。

電子媒体の言語資料の特性を生かす形でコーパスを日本語研究に利用するにはそれなりの知識や技術が必要となる。本書は、特定領域研究「日本語コーパス」の計画研究班の1つである日本語学班のメンバー4人がこれまでの経験を踏まえ、日本語研究にコーパスを役立てようとする方の参考になりそうな初歩的な情報を提供するものである。どこからでも読める形にしてあるので、興味のあるところから読んで実際にパソコン上で試していただければ幸いである。

第1部本編では、日本語研究者が自ら独自の目的・方法でコーパスを利用できるようになるための方法や手順などを紹介する。第2部各種解説編には、日本語コーパスを利用するうえで必要となる一般的な知識や既成のソフトウェアに関する情報などを収めた。コーパスを研究に生かしたいがどこから手を付けてよいか分からないという方にとって本書がコーパス処理の世界への道案内となることを願っている。

本書の各章が直接の処理対象とするのはいわゆるプレーンテキストである。一般に利用可能な付加情報付きコーパスがない現状では好むと好まざるとにかかわらず選択の余地がないわけであるが、プレーンテキストには、豊富に入手でき、扱いが容易で、特定の文法解釈に制約されることがないという利点がある。プレーンテキストという形態でのコーパスは、付加情報付きのコーパスが作られその利用が普及した後も日本語研究の精密化に重要な役割を果たし続けることであろう。

なお、本書では現在日本語研究者に最も広く使われていると思われるOSであるWindows XP日本語版の使用を想定している。もっとも、本書の内容の多くはWindowsのバージョンには依存しないはずである。また、日本語テキストの文字コードはShift-JISであることを前提としている。

本書および添付CDの内容に関して更新や訂正の情報があれば特定領域研究「日本語コーパス」のWebサイトにある日本語学班のページ(http://www.tokuteicorpus.jp/g_jpling/)に掲載していく予定である。不審な点などがあればまずはそちらをご覧くださいませければ幸いである。

田野村 忠温

2007年9月

目次

はしがき	1
目次	3

第1部 本編

日本語 KWIC ソフトウェア KWIC (田野村忠温)	7
AWK プログラミング入門 (田野村忠温)	31
Perl・sortf・秀丸エディタを用いた用法分析の一例 (服部匡)	87
perl によるテキスト処理 (杉本武)	97
Excel で語彙表を作成する (石井正彦)	109

第2部 各種解説編

コマンドプロンプト (田野村忠温)	151
ファイル名拡張子 (服部匡)	159
正規表現・文字コード (田野村忠温)	161
エディタを用いた文字列検索——EmEditor の場合—— (田野村忠温)	169
エディタを用いた文字列検索——秀丸エディタの場合—— (服部匡)	173
KWIC 索引生成ソフトウェアの紹介 (服部匡)	177
利用可能なコーパスの紹介 (服部匡)	181
コーパス日本語学研究文献目録	185

日本語KWICソフトウェアKWIC

田野村 忠温

コーパス利用の普及のためには簡単に使えるコーパス処理ソフトウェアの存在が不可欠だと思われるが、残念ながら日本語に関しては研究者が手持ちの電子テキストをそのまま使って簡便にKWIC索引を生成することのできるソフトウェアがない。

そうした状況にかんがみ、このたび日本語KWICソフトウェアを作成してみた。ここではその機能・用法と、検索結果の利用に関する簡単なヒントについて述べる。

内容

- 1 概要
 - 2 インストールと使用の準備
 - 3 用法(1) 簡易検索
 - 4 用法(2) 連続検索
 - 5 定義ファイルの詳細
 - 6 補足・注意事項
- 補説A デスクトップからの実行—RunKWIC
補説B 検索結果のエディタでの利用

AWK プログラミング入門

田野村 忠温

テキスト処理に適したプログラミング言語として AWK と Perl、ほかには Python や Ruby が広く使われている。ここでは AWK とソートプログラム `sortf` を使ってテキストから文字列を検索したり KWIC 索引を生成したりする方法を紹介する。例題では英語のテキストを処理対象とするが、日本語のテキストを処理する場合にも基本的な考え方はそのまま通用する。

AWK はほかの言語に比べて機能が限られているが、手軽に学べるという利点がある。また、AWK を学んだ経験は、さらに進んで Perl その他の言語を学ぶときにも大いに役に立つであろう。

内容

- AWK のインストールなど
- 学習開始時の準備
- §1 AWK の第 1 歩
- §2 AWK の第 2 歩
- §3 いくつかの応用例
- §4 if 文
- §5 for 文
- §6 配列
- §7 for 文 (その 2)
- §8 文字列の処理
- §9 文字列の処理 (その 2)
- (解説 1) サンプルデータ
- (解説 2) ソートプログラム `sortf`

Perl・sortf・秀丸エディタを用いた用法分析の一例

服部 匡

Perl スクリプトによって簡単な KWIC 索引を生成する。それを利用し、Perl・sortf・秀丸エディタを組み合わせることで語の用例の仮分類や計数を行っていく方法の一例を紹介する。

内容
はじめに
事前の準備
§1 データの準備
§2 分析開始時の準備
§3 KWIC 索引の生成
§4 KWIC 索引のソート
§5 KWIC 索引からの用例仮分類
おわりに

perl によるテキスト処理

杉本 武

コーパスを用いた日本語研究には、テキスト・データが不可欠であるわけであるが、現在では、一般的には、各紙の新聞データや青空文庫のような小説のテキスト・データが利用できる。ただし、これらのデータは、必ずしも日本語研究のために作成されたものではないため、検索などの処理をするためには扱いにくい点があることがある。

ここでは、上のような一般的に入手できるテキストデータから必要な情報を抽出、整形し、扱いやすいデータを作成する perl のスクリプト、また、そのデータを検索するための簡易 KWIC のスクリプトを紹介する。

内容

- 1 データ形式
 - 1.1 1 行の単位
 - 1.2 ヘッダー情報
 - 1.3 文字の正規化
- 2 スクリプトを使用するための準備
- 3 テキスト・データの変換
 - 3.1 プレーンなテキスト・データの場合
 - 3.2 毎日新聞データの場合
- 4 テキストの検索: 簡易 KWIC
 - 5 ライブラリ
 - 5.1 regular.pl
 - 5.2 divsent.pl

Excel で語彙表を作成する

石井 正彦

難しいプログラムや特別なソフトを使わず、手軽に語彙表を作ることはできないか。こうした要望に応えるべく、Excel (Microsoft Office Excel 2003) を使って語彙表を作成する方法を、新聞コーパスを例として、まとめてみた。全体を8回に分け、各回とも、「今日の目標」「解説」「新出の技法」「準備」「作業の手順」という構成にした。なお、Excel による語彙表の作成は、伊藤雅光『計量言語学入門』(大修館書店、2002) も行っているが、伊藤の方法はあらかじめ作成した語彙表のデータを Excel に読み込んでいろいろな作業を行うのに対して、ここでは、Excel 自身に語彙表を作らせるようにしている。

内 容

- § 0 はじめに (フォルダとファイルの準備)
- § 1 新聞コーパスを記事種別 (紙面) ごとに分割する
- § 2 コーパスから語彙調査に必要なレコードだけを取り出す
- § 3 形態素解析プログラムを使って、コーパス (新聞文章) の単位切りと同語異語判別を行う
- § 4 形態素解析の結果を Excel に読み込んで、記事種別ごとの「語彙表」を作る
- § 5 Excel で、すべての記事種別の頻度が一覧できる「(自立語の) 層別語彙表」を作る
- § 6 Excel で、「全紙面 (層別) 語彙表」に、使用率・順位・出現紙面数などの情報を加える
- § 7 「全紙面 (層別) 語彙表」から度数分布表とヒストグラムを作る
- § 8 「全紙面 (層別) 語彙表」から「基幹語彙」「特徴語彙」を取り出す

特定領域研究「日本語コーパス」平成19年度研究成果報告書 (JC-L-07-01)

コーパス日本語学ガイドブック

田野村忠温・服部匡・杉本武・石井正彦

2007(平成19)年9月7日発行

文部科学省科学研究費補助金特定領域研究「日本語コーパス」日本語学班

〒562-8558 大阪府箕面市粟生間谷東 8-1-1 大阪外国語大学 田野村忠温研究室

電話・FAX: 072-730-5177

E-mail: [tanomura\(at\)osaka-gaidai.ac.jp](mailto:tanomura(at)osaka-gaidai.ac.jp) (2007年9月30日まで)

[tanomura\(at\)let.osaka-u.ac.jp](mailto:tanomura(at)let.osaka-u.ac.jp) (2007年10月1日より)
