

第5章: 係り受け解析

日本語Wikipedia (<https://ja.wikipedia.org/>) の「人工知能」 (<https://ja.wikipedia.org/wiki/%E4%BA%BA%E5%B7%A5%E7%9F%A5%E8%83%BD>)」に関する記事からテキスト部分を抜き出したファイルが `ai.ja.zip` に収録されている。この文章を `CaboCha` (<https://taku910.github.io/cabocha/>) や `KNP` (<http://nlp.ist.i.kyoto-u.ac.jp/index.php?KNP>) 等のツールを利用して係り受け解析を行い、その結果を `ai.ja.txt.parsed` というファイルに保存せよ。このファイルを読み込み、以下の問に対応するプログラムを実装せよ。

40. 係り受け解析結果の読み込み（形態素）

形態素を表すクラス `Morph` を実装せよ。このクラスは表層形（`surface`）、基本形（`base`）、品詞（`pos`）、品詞細分類1（`pos1`）をメンバ変数に持つこととする。さらに、係り受け解析の結果（`ai.ja.txt.parsed`）を読み込み、各文を `Morph` オブジェクトのリストとして表現し、冒頭の説明文の形態素列を表示せよ。

41. 係り受け解析結果の読み込み（文節・係り受け）

40に加えて、文節を表すクラス `Chunk` を実装せよ。このクラスは形態素（`Morph` オブジェクト）のリスト（`morphs`）、係り先文節インデックス番号（`dst`）、係り元文節インデックス番号のリスト（`srcs`）をメンバ変数に持つこととする。さらに、入力テキストの係り受け解析結果を読み込み、1文を `Chunk` オブジェクトのリストとして表現し、冒頭の説明文の文節の文字列と係り先を表示せよ。本章の残りの問題では、ここで作ったプログラムを活用せよ。

42. 係り元と係り先の文節の表示

係り元の文節と係り先の文節のテキストをタブ区切り形式ですべて抽出せよ。ただし、句読点などの記号は出力しないようにせよ。

43. 名詞を含む文節が動詞を含む文節に係るものを抽出

名詞を含む文節が、動詞を含む文節に係るとき、これらをタブ区切り形式で抽出せよ。ただし、句読点などの記号は出力しないようにせよ。

44. 係り受け木の可視化

与えられた文の係り受け木を有向グラフとして可視化せよ。可視化には、[Graphviz](http://www.graphviz.org/) (<http://www.graphviz.org/>)等を用いるとよい。

45. 動詞の格パターンの抽出

今回用いている文章をコーパスと見なし、日本語の述語が取りうる格を調査したい。動詞を述語、動詞に係っている文節の助詞を格と考え、述語と格をタブ区切り形式で出力せよ。ただし、出力は以下の仕様を満たすようにせよ。

- 動詞を含む文節において、最左の動詞の基本形を述語とする
- 述語に係る助詞を格とする
- 述語に係る助詞（文節）が複数あるときは、すべての助詞をスペース区切りで辞書順に並べる

「ジョン・マッカーシーはAIに関する最初の会議で人工知能という用語を作り出した。」という例文を考える。この文は「作り出す」という1つの動詞を含み、「作り出す」に係る文節は「ジョン・マッカーシーは」、「会議で」、「用語を」であると解析された場合は、次のような出力になるはずである。

作り出す で は を

このプログラムの出力をファイルに保存し、以下の事項をUNIXコマンドを用いて確認せよ。

- コーパス中で頻出する述語と格パターンの組み合わせ
- 「行う」「なる」「与える」という動詞の格パターン（コーパス中で出現頻度の高い順に並べよ）

46. 動詞の格フレーム情報の抽出

45のプログラムを改変し、述語と格パターンに続けて項（述語に係っている文節そのもの）をタブ区切り形式で出力せよ。45の仕様に加えて、以下の仕様を満たすようにせよ。

- 項は述語に係っている文節の単語列とする（末尾の助詞を取り除く必要はない）
- 述語に係る文節が複数あるときは、助詞と同一の基準・順序でスペース区切りで並べる

「ジョン・マッカーシーはAIに関する最初の会議で人工知能という用語を作り出した。」という例文を考える。この文は「作り出す」という1つの動詞を含み、「作り出す」に係る文節は「ジョン・マッカーシーは」、「会議で」、「用語を」であると解析された場合は、次のような出力になるはずである。

作り出す で は を 会議で ジョンマッカーシーは 用語を

47. 機能動詞構文のマイニング

動詞のヲ格にサ変接続名詞が入っている場合のみに着目したい。46のプログラムを以下の仕様を満たすように改変せよ。

- 「サ変接続名詞+を（助詞）」で構成される文節が動詞に係る場合のみを対象とする
- 述語は「サ変接続名詞+を+動詞の基本形」とし、文節中に複数の動詞があるときは、最左の動詞を用いる
- 述語に係る助詞（文節）が複数あるときは、すべての助詞をスペース区切りで辞書順に並べる
- 述語に係る文節が複数ある場合は、すべての項をスペース区切りで並べる（助詞の並び順と揃えよ）

例えば「また、自らの経験を元に学習を行う強化学習という手法もある。」という文から、以下の出力が得られるはずである。

学習を行う に を 元に 経験を

48. 名詞から根へのパスの抽出

文中のすべての名詞を含む文節に対し、その文節から構文木の根に至るパスを抽出せよ。ただし、構文木上のパスは以下の仕様を満たすものとする。

- 各文節は（表層形の）形態素列で表現する
- パスの開始文節から終了文節に至るまで、各文節の表現を”->”で連結する

「ジョン・マッカーシーはAIに関する最初の会議で人工知能という用語を作り出した。」という例文を考える。CaboChaを係り受け解析に用いた場合、次のような出力が得られると思われる。

ジョンマッカーシーは -> 作り出した
AIに関する -> 最初の -> 会議で -> 作り出した
最初の -> 会議で -> 作り出した
会議で -> 作り出した
人工知能という -> 用語を -> 作り出した
用語を -> 作り出した

KNPを係り受け解析に用いた場合、次のような出力が得られると思われる。

ジョンマッカーシーは -> 作り出した
 AIに -> 関する -> 会議で -> 作り出した
 会議で -> 作り出した
 人工知能と -> いう -> 用語を -> 作り出した
 用語を -> 作り出した

49. 名詞間の係り受けパスの抽出

文中のすべての名詞句のペアを結ぶ最短係り受けパスを抽出せよ。ただし、名詞句ペアの文節番号が i と j ($i < j$) のとき、係り受けパスは以下の仕様を満たすものとする。

- 問題48と同様に、パスは開始文節から終了文節に至るまでの各文節の表現（表層形の形態素列）を” -> ”で連結して表現する
- 文節 i と j に含まれる名詞句はそれぞれ、 X と Y に置換する

また、係り受けパスの形状は、以下の2通りが考えられる。

- 文節 i から構文木の根に至る経路上に文節 j が存在する場合: 文節 i から文節 j のパスを表示
- 上記以外で、文節 i と文節 j から構文木の根に至る経路上で共通の文節 k で交わる場合: 文節 i から文節 k に至る直前のパスと文節 j から文節 k に至る直前までのパス、文節 k の内容を” | ”で連結して表示

「ジョン・マッカーシーはAIに関する最初の会議で人工知能という用語を作り出した。」という例文を考える。CaboChaを係り受け解析に用いた場合、次のような出力が得られると思われる。

Xは | Yに関する -> 最初の -> 会議で | 作り出した
 Xは | Yの -> 会議で | 作り出した
 Xは | Yで | 作り出した
 Xは | Yという -> 用語を | 作り出した
 Xは | Yを | 作り出した
 Xに関する -> Yの
 Xに関する -> 最初の -> Yで
 Xに関する -> 最初の -> 会議で | Yという -> 用語を | 作り出した
 Xに関する -> 最初の -> 会議で | Yを | 作り出した
 Xの -> Yで
 Xの -> 会議で | Yという -> 用語を | 作り出した
 Xの -> 会議で | Yを | 作り出した
 Xで | Yという -> 用語を | 作り出した
 Xで | Yを | 作り出した
 Xという -> Yを

KNPを係り受け解析に用いた場合、次のような出力が得られると思われる。

Xは | Yに -> 関する -> 会議で | 作り出した。
Xは | Yで | 作り出した。
Xは | Yと -> いう -> 用語を | 作り出した。
Xは | Yを | 作り出した。
Xに -> 関する -> Yで
Xに -> 関する -> 会議で | Yと -> いう -> 用語を | 作り出した。
Xに -> 関する -> 会議で | Yを | 作り出した。
Xで | Yと -> いう -> 用語を | 作り出した。
Xで | Yを | 作り出した。
Xと -> いう -> Yを

🌱 **Updated:** June 7, 2020