

第4章: 形態素解析

夏目漱石の小説『吾輩は猫である』の文章 (`neko.txt`) をMeCabを使って形態素解析し、その結果を`neko.txt.mecab`というファイルに保存せよ。このファイルを用いて、以下の問に対応するプログラムを実装せよ。

なお、問題37, 38, 39は[matplotlib](http://matplotlib.org/) (<http://matplotlib.org/>) もしくは [Gnuplot](http://www.gnuplot.info/) (<http://www.gnuplot.info/>) を用いるとよい。

30. 形態素解析結果の読み込み

形態素解析結果 (`neko.txt.mecab`) を読み込むプログラムを実装せよ。ただし、各形態素は表層形 (`surface`)、基本形 (`base`)、品詞 (`pos`)、品詞細分類1 (`pos1`) をキーとするマッピング型に格納し、1文を形態素 (マッピング型) のリストとして表現せよ。第4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。

31. 動詞

動詞の表層形をすべて抽出せよ。

32. 動詞の基本形

動詞の基本形をすべて抽出せよ。

33. 「AのB」

2つの名詞が「の」で連結されている名詞句を抽出せよ。

34. 名詞の接続

名詞の接続 (連続して出現する名詞) を最長一致で抽出せよ。

35. 単語の出現頻度

文章中に出現する単語とその出現頻度を求め、出現頻度の高い順に並べよ。

36. 頻度上位10語

出現頻度が高い10語とその出現頻度をグラフ（例えば棒グラフなど）で表示せよ.

37. 「猫」と共起頻度の高い上位10語

「猫」とよく共起する（共起頻度が高い）10語とその出現頻度をグラフ（例えば棒グラフなど）で表示せよ.

38. ヒストグラム

単語の出現頻度のヒストグラムを描け. ただし, 横軸は出現頻度を表し, 1から単語の出現頻度の最大値までの線形目盛とする. 縦軸はx軸で示される出現頻度となった単語の異なり数（種類数）である.

39. Zipfの法則

単語の出現頻度順位を横軸, その出現頻度を縦軸として, 両対数グラフをプロットせよ.

🌱 Updated: May 20, 2020