

## 7. ローカル環境での利用について

統語・意味解析コーパス (NPCMJ) チュートリアル  
@弘前大学

長崎郁

2019.5.11

1 / 10

## はじめに

- NPCMJ のデータ（ラベル付き括弧（Bracketed tree）形式）は、ダウンロードして利用することもできる。
- 例えば、ラベル付き括弧形式のファイルは次のようなツールを使って検索することができる。
  - Tgrep2 (<https://tedlab.mit.edu/~dr/Tgrep2/>)  
コマンドラインから利用，unix 環境が必要
  - Tregex (<https://nlp.stanford.edu/software/tregex.html>)  
java プログラム，windows や mac でも動く
- この時間は，tregex を使って NPCMJ のデータを検索する方法，および検索結果をエクスポートする方法について紹介する。

2 / 10

## NPCMJ データのダウンロード

- データのダウンロード（以下のいずれかから）

- トップページ

<http://npcmj.ninjal.ac.jp/>

[Bracketed Tree ファイルをダウンロードする](#) をクリック

- 概要ページ

<http://npcmj.ninjal.ac.jp/interfaces/cgi-bin/index.sh?db=npcmj&lang=jp>

[Download all bracketed trees](#) をクリック

- ダウンロードした zip ファイルを解凍

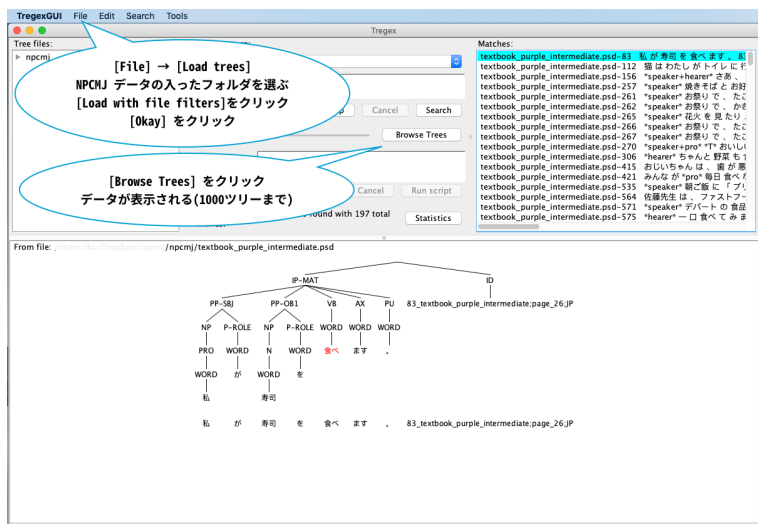
3 / 10

## tregex のインストールと起動

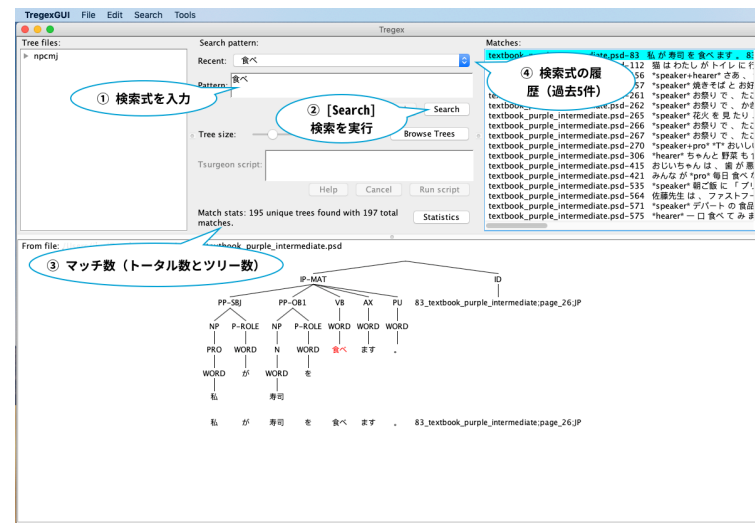
- <https://nlp.stanford.edu/software/tregex.html>
  - タイトル (Tregex, Tsurgeon and Serngrex) の下の  
| Download | をクリック
  - 最新版 (3.9.2) をダウンロード，適当な場所（例えば Desktop）に保存
  - zip ファイルを解凍し，フォルダの中の“stanford-tregex.jar”をクリック
  - tregex ブラウザーが立ち上がらなかったら，java をインストールする。
    - <https://www.java.com/en/download/>
    - [Free Java Download](#)
    - [Agree and Start Free Download](#)
    - ダウンロードしたファイルをクリックし，指示に従う。

4 / 10

## tregex : ファイルの読み込み



## tregex : 検索



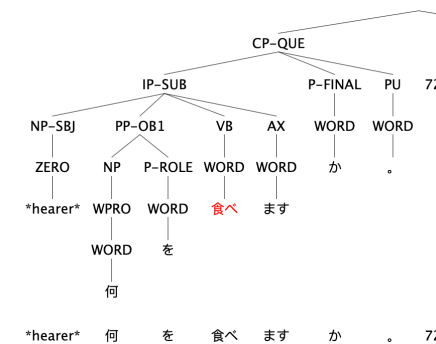
適当な語 (例 : 食べ) を検索してみましょう

## tregex : 検索式 (TGrep-lite との違い)

- TGrep-lite とは書き方が少し異なる。
  1. 終端ノードも非終端ノードも `"./.../"` や `"[...]"` で囲まずに書くことができる。その場合「完全一致」の意味になる。
  2. 終端ノードも非終端ノードも正規表現環境にするとときは `"./.../"` で囲む。単に文字列を囲めば「部分一致」の意味になる。
  3. 利用できる正規表現の種類が TGrep-lite よりも多い。
    - 例えば、選言の `"|"` は `./.../` で囲まなくても使える。
    - `"食べる|食べる|食べよ|食べ"` と `"/^(食べる|食べる|食べよ|食べ)$/"` は同じ (TGrep-lite は常に後者を使わなくてはならない)
    - 詳しくは [Help] を参照

## tregex : 検索式 (TGrep-lite との違い)

4. 終端ノード (空要素をのぞく) は `"WORD"` によって直接支配される (空要素は `"ZERO"` によって直接支配)。よって、終端ノードと品詞タグの支配関係を記述する場合は、`"A_<<:_B"`, `"B_>>:_A"` を使うと良い (`"_"` は半角スペース)。
  - `VB_<<:_/^食べ/` (VB (動詞) で「食べ」で始まるもの)
  - `か_>>:_P-FINAL` (「か」で P-FINAL (終助詞))
  - `か_!>>:_P-FINAL` (「か」で P-FINAL (終助詞) でないもの)



## tregex : 検索結果のエクスポート

- 例：助詞「しか」の用例を検索し、エクスポートしてみる
- 「しか」を検索 → 62 total matches
  - (「しか」はすべて P-OPTR (とりたて助詞) として扱われているか、確認したい場合は以下のようにして確認すると良い。  
“しか<\_!>>:\_P-OPTR”
- エクスポート
  - メニュー [File] → [Save matched trees] (タグあり)
  - メニュー [File] → [Save matched sentences] (タグなし)

9 / 10

## tregex: ツリーの一部を表示・エクスポートする

- データには非常に長いツリーも含まれる。そのため検索結果を確認するのが困難な場合は、必要な部分だけ表示させることもできる。
- 例：助詞「しか」を含む節 (IP)
- メニューから [TreeGUI] → [Preferences]  
[Show only matched portion of tree] をチェック
  - 検索式中のマスターノード (とそれに支配されるノード) を表示する。
- 助詞「しか」を支配する助詞句 (PP) をもつ節 (IP) を検索  
“/^IP/\_<\_ (/^PP/\_<\_(P-OPTR\_<<:\_しか))” (55 total matches)
- 上述と同じ方法でタグあり、タグなしのデータをエクスポートすることができる。

10 / 10