

テキストデータの統計科学入門

- 第4章 -

テキストにおける集計モデルと集計ツール

茨城大学情報工学科
菊池 裕紀

テキストにおける集計モデル

何らかの処理や統計を行う際...

単位を決めることが必要になる！

記号列を集計するときの基本

→それぞれの記号がテキスト中に現れる度数(頻度)

- ・「き」という記号

「きしやのきしやがきしやできしやする。」

→ 文章中に「き」が4つ検出される

この記号の数を増やす事で集計するモデルを拡張できる！

n-gram : n個の記号の度数を集計する方法

(nは集計のため切り取った隣接する記号列の長さ)

1つの記号 : 「き」 (unigram)

2つの記号 : 「きし」 (bigram)

3つの記号 : 「きしや」 (trigram)

4つの記号 : 「きしやが」 (four-gram)

データ集計のツール

MLTP (multi lingual text processor)

文字単位のn-gramを抽出するソフト
複数のテキストを同時に処理できる。

RMeCab (Rパッケージ)

R言語上で日本語テキストを統計的に分析する

茶器(ChaKi)

KH-coder

etc...

MLTP

・タグ

File List

Summary

n-gram

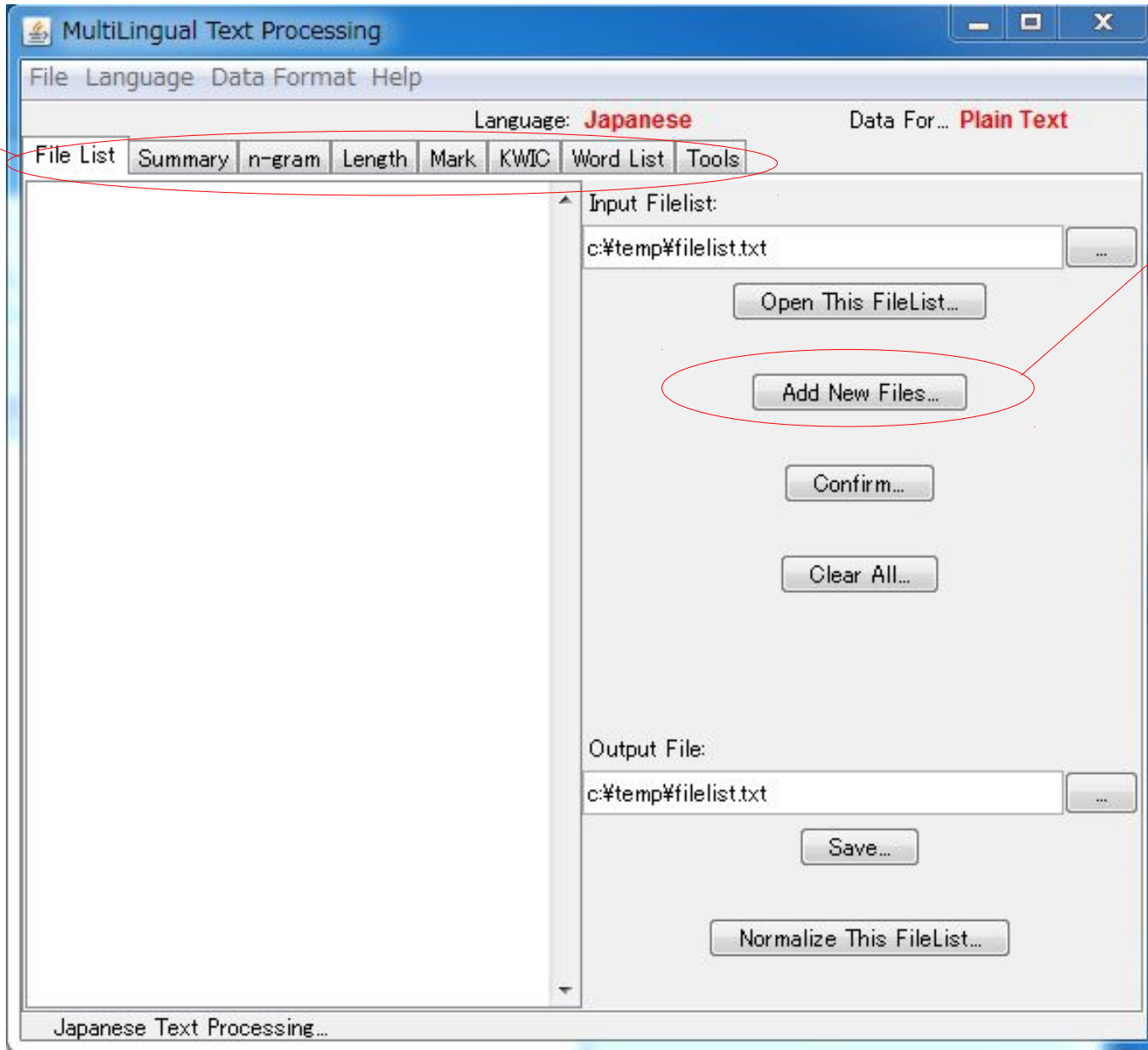
Length

Mark

KWIC

Word List

Tools



ファイルの追加