

テキストデータ群の重ね合わせによるヒストリと特徴のインタラクティブ表示(3)

松原 伸人、土屋 正人

matubara@sra.co.jp, m-tsuchi@sra.co.jp

◆はじめに

大量のデータをインタラクティブに操作する、Web ブラウザ上で動作するアプリケーションのプロトタイプ開発を行っています。データとしてGSLetterNeoを使って、

- ドキュメントタイトル、著者等を表示するレイヤー
- 特徴語群を表示するレイヤー
- 希少語群を表示するレイヤー

の3層を、画面奥に向かって重ねて配置することで、ドキュメントの時間関係と特徴がひとまとまりに見えるようにしたプロトタイプのインタフェースと実装を、Vol.95とVol.96で紹介しました。

今回はデータを解析するために必要な辞書を中心に紹介します。

◆形態素解析¹と辞書

各単語のドキュメント内での使用頻度と、単語が使われているドキュメント数から、tf-idfが計算できます。

tf-idfはドキュメント内でたくさん使われていて、他のドキュメントではそれほど使われていない単語の度合いを表します。

tf-idf - [Wikipedia]

<https://ja.wikipedia.org/wiki/Tf-idf>

希少語群の検出では、tf-idfの計算式の分子tfを逆にして単語の出現頻度の低い単語を計算しています。

$$tfidf = tf / idf$$
$$minority = (1 - tf) / idf$$

実際に形態素解析を行ってみると分かりますが、解析対象に合う辞書の作成が必要になります。

今回はIPA辞書をベースにして利用して、解析結果を見ながら必要に応じて辞書に単語を追加したり、不明な語の定義の変更を行ったりしています。

辞書に1個単語を追加すると、解析結果が単語の計測結果にも影響します。

◆解析結果と計測結果の可視化

このような辞書作成して単語計測結果を見るためのツールが図1～図7です。

このツールはHTMLとJavaScriptとCSSでプログラムしてあります。

GSLetterNeoのテキストの解析結果を見ていく過程で作成し、使っていました。

画面左側にあるドキュメントリストで対象ドキュメントをクリックすると、画面中央のテキストエリアにテキスト内容を表示、画面右側に解析結果を表示します。

テキストエリア内で、太字で大きな字になっている単語はtfidf値が高い単語で、値が大きいほど大きな字で表示されるようになっています。

¹ 文法的な情報の注記の無い自然言語のテキストデータ(文)から、対象言語の文法や、辞書と呼ばれる単語の品詞等の情報にもとづき、形態素(Morpheme、おおまかにいえば、言語で意味を持つ最小単位)の列に分割し、それぞれの形態素の品詞等を判別する作業(Wikipediaより引用)。

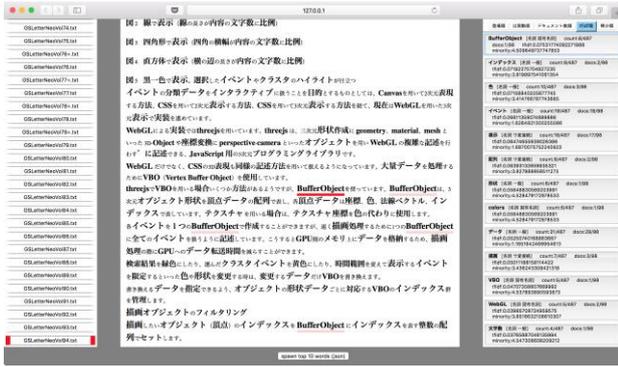


図 1 選んだ単語にアンダーラインをつけて示す

右側のリストの単語をクリックすると、テキストエリア上の単語にアンダーラインをつけて示します(図1)。

クリックするたびに登場する順に単語を移動して見ることができます。

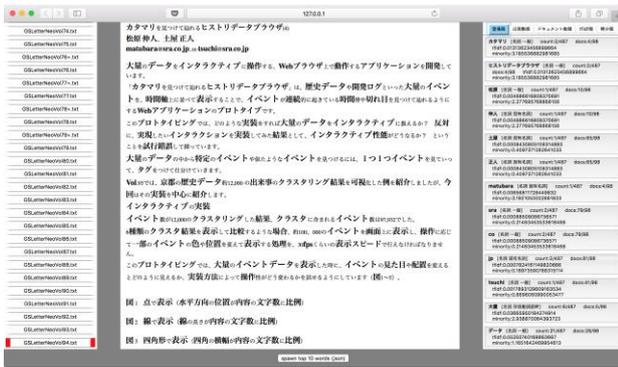


図 2 単語計測結果を見るツールの画面 (単語の登場順)

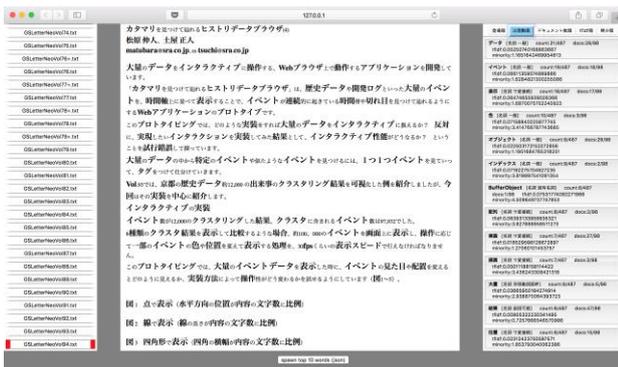


図 3 出現教順

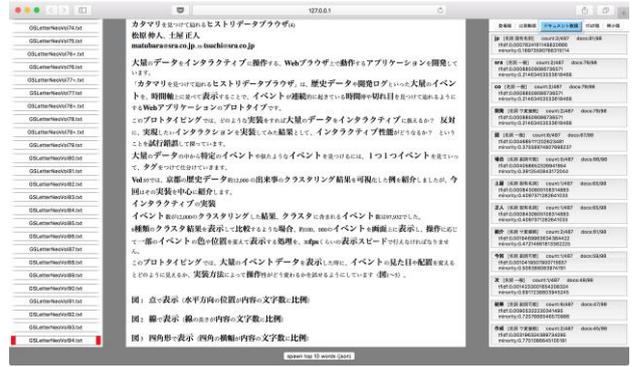


図 4 ドキュメント数順

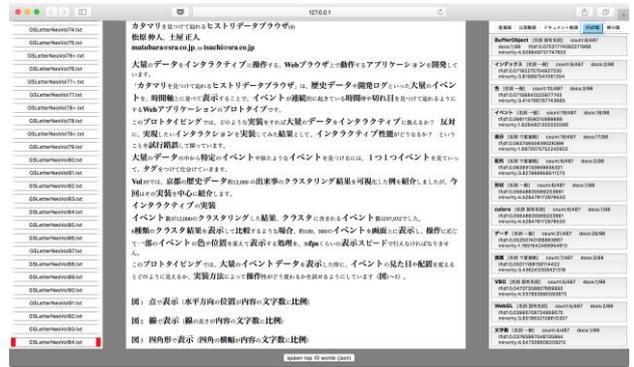


図 5 tdfdf 順

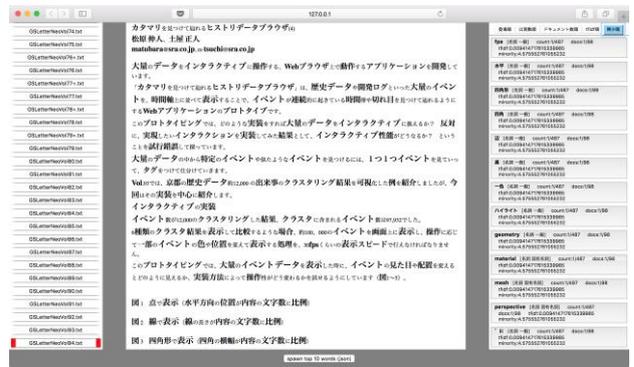


図 6 希少順

画面下部にある「spawn top 10 words」ボタンを押すと、登場順、出現数順、ドキュメント数順、tfidf順、希少順にそれぞれスコアの高い10語を抽出し、JSON形式で解析結果を出力します(図7)。

出力したJSONデータを用いて、特徴語レイヤーと希少語レイヤーを作成しています。

```
{
  "GSLetterNeoVol1": {
    "name": "GSLetterNeoVol1",
    "filename": "GSLetterNeo_text_exports/GSLetterNeoVol1.txt",
    "appearance": [
      "アジャイル",
      "計画",
      "コンサルタントファシリテーター",
      "野島",
      "勇",
      "NOJIMA",
      "Isamu",
      "nojima",
      "sra",
      "co"
    ],
    "termfrequency": [
      "計画",
      "野島",
      "勇",
      "手塚",
      "アジャイルソフトウェア",
      "プロジェクト",
      "アジャイル",
      "変化",
      "連絡",
      "従来",
      "実業"
    ],
    "documentfrequency": [
      "SRA",
      "sra",
      "co",
      "野島",
      "必要",
      "勇",
      "連絡",
      "紹介",
      "今回",
      "プロジェクト"
    ]
  }
}
```

図7 解析結果をJSON形式で出力した画面

◆終わりに

解析結果のリストを見るだけでは、なぜこういう分解がなされたのかわからない場合も、元のテキストと合わせて見られると理由がなんとなくわかり、試しに辞書に単語を追加してみると、正しそうな結果になることが多かったです。

多くの場合は、もともと辞書には載っていないような単語が現れた時に、分解がうまくいかないようです。

<p>GSLetterNeo Vol. 97 2016年8月20日発行 発行者●株式会社SRA 先端技術研究所 編集者●土屋正人</p> <p>バックナンバーを公開しています●http://www.sra.co.jp/gsletter ご感想・お問い合わせはこちらへお願いします●gsneo@sra.co.jp</p>	<p>夢を。</p> 
--	--

株式会社SRA

〒171-8513 東京都豊島区南池袋2-32-8

夢を。Yawaraka Innovation
やわらかいのバージョン