

R

relativeFrequency 相対頻度について

サンプルデータ 100万語

<http://www.thegrammarlab.com/?nor-portfolio=1000000-word-sample-corpora>

<http://micusp.elicorpora.info/>

MICUSP Sample

<http://www.thegrammarlab.com/?nor-portfolio=1000000-word-sample-corpora#>

全体の高頻度語

1, 10, 100, 1000位の語と頻度

topfeatures()の結果の中身

- ・ names 属性のついた数値ベクトル
 - ・ 要素番号のように、要素を指定すると、その数値が表示される。
 - ・ 要素番号の指定もできる。(この場合、要素番号が頻度順位となる)

チャンク内の検索と頻度

"the"の頻度分布

"the"の累積頻度分布

100語ずつ100万語まで1万回サンプリングしてみるとほぼ直線的に増えていることがわかる。

最初の100語での出現頻度は、7回。これを1万倍したら70000。実際は、68659回。1.02の精度。

1000位のsitesについて

sitesの頻度分布

siteの累積頻度分布

分散をしてみる

theはどこにおいても均一に出現している

sitesの出現は、偏っている。

ゆえに、sitesについては、その一部に基づき、相対頻度を出すことは、結果がゆがむ恐れが大きい。

thereforeについて

累積頻度