

国立国語研究所学術情報リポジトリ

異なり語数の推定

メタデータ	言語: jpn 出版者: 公開日: 2020-06-29 キーワード (Ja): キーワード (En): 作成者: 山崎, 誠 メールアドレス: 所属:
URL	https://doi.org/10.15084/00002889

異なり語数の推定

山崎 誠

1. はじめに

2. 異なり語数と異なり表記形数

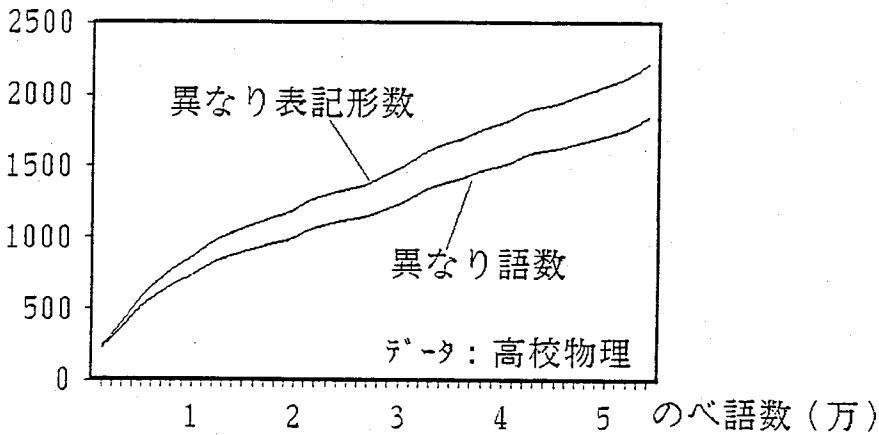


図1 異なり語数と異なり表記形数の推移

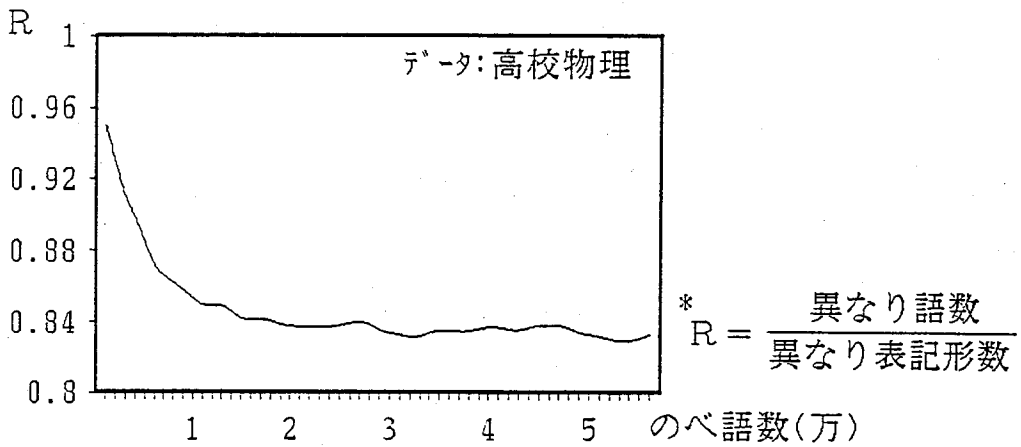
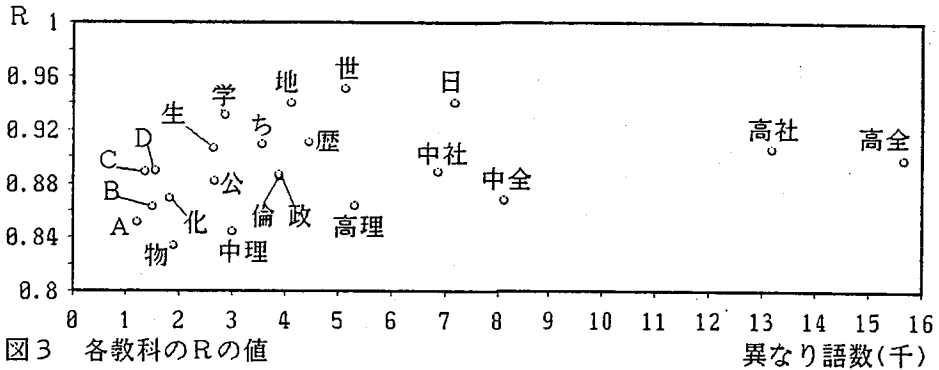


図2 R^* の値の推移



- 高校
 物…物理 倫…倫理社会
 化…化学 政…政治経済
 生…生物 日…日本史
 学…地学 世…世界史
 地…地理
 高理…理科系
 高社…社会科系
 高全…全教科
- 中学
 A…理科1上 公…公民
 B…理科1下 ち…地理
 C…理科2上 歴…歴史
 D…理科2下
 中理…理科系
 中社…社会科系
 中全…全教科

3. 1語に対する表記形の数

表1 1語に対応する表記形数の多いもの(高校教科書)

度数	見出し語	類	表記形例
21	とる	2	取れる
18	つくる	2	造り
16	はかる	2	謀り
13	かわる	2	変わる
13	たつ	2	立つ
13	ひく	2	引く
13	ゆく	2	行く
12	あう	2	合う
12	おく	2	置く
12	おくる	2	贈る
12	かえる	2	返る
12	かく	2	書ける
11	あらわす	2	表わされる
11	おこなう	2	行なわれる
11	おこる	2	興る
11	きる	2	切れる
11	はなす	2	離れる
11	ひらく	2	開く
11	もつ	2	持つ

度数	見出し語	類	表記形例
10	したがう	2	従わ
10	する	2	支
10	たてる	2	立てる
10	つく	2	付くる
10	なる	2	成る
10	のぞく	2	除く
10	はたらく	2	働ける
10	はやい	3	速く

「類」は、「分類語彙表」での分類番号の1桁めの数字

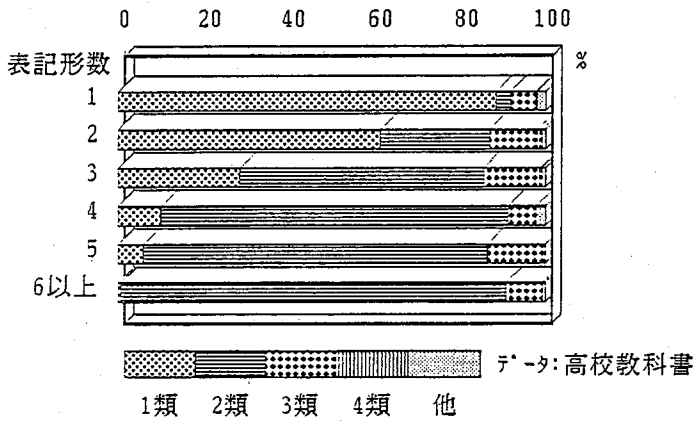


図4 1語あたりの異なり表記形数別の語類構成比

表2 Rと語類構成比との相関

	高校	中学
1類	0.913	0.680
2類	-0.554	-0.569
3類	-0.463	-0.279
4類	-0.669	-0.309

4. 1表記形に対する見出し語の数

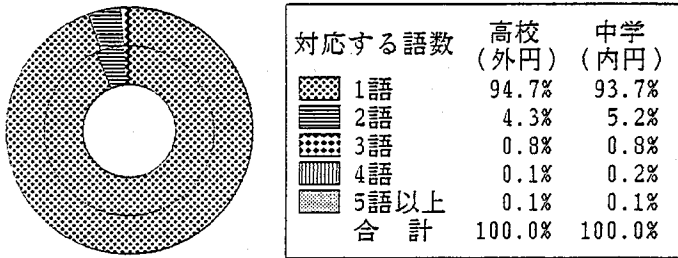


図5 1つの表記形がいくつの語に対応したか

表3 1表記形に対応する語の数が多いもの (高校教科書)

度数	表記形
7	より
6	かけ つい 下 元 日
5	かた から たち たて たら つき つけ で なり まき め 金 重 西 中 長 分
4	あい いき がけ さし さら し つぎ なし なれ み よ 家 間 京 月 原 源 出 上 人 生 土 文 米 本 名
3	あたり あり あわ いっ うち え おい おき か か かき かぎ かく かる かね か り が ら おき く ま く くり こめ かさ かる かつ け す す み ず み せ たえる たま え ます つみ づ ま っ な も と や り の り わ か オレソ 汗 ガ器 基 強 玉 近 史 軽 見 古 後 御 光 巻 き 行 上 高 常 神 細 世 成 功 星 青 石 赤 巢 大 御 出 光 広 床 値 文 治 平 歩 万 味 密 毛 木 夜 葉 両 封 じ 刀 風 林 和 公 商 端 物

5. 表記の種類とRの値

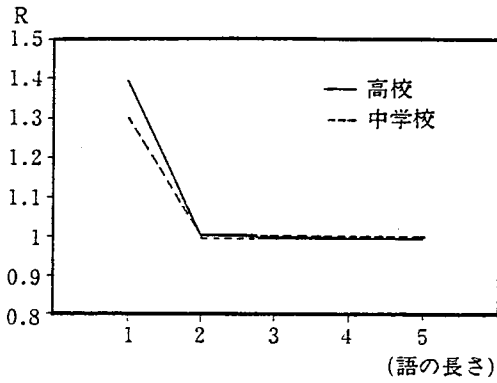


図6 Rの値と語の長さ（漢字表記形）

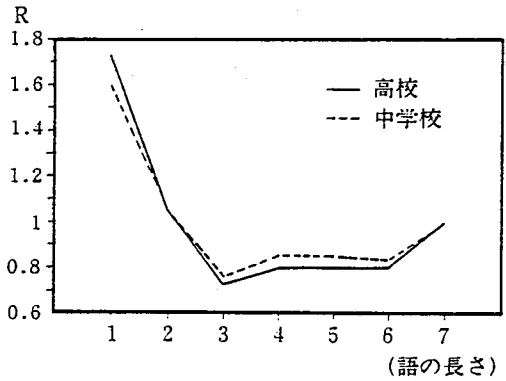


図7 Rの値と語の長さ（平仮名表記形）

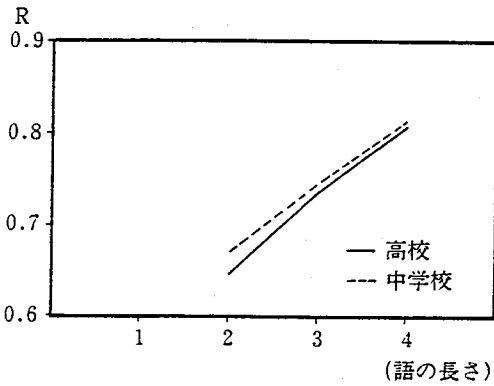


図8 Rの値と語の長さ（混ぜ書き表記形）

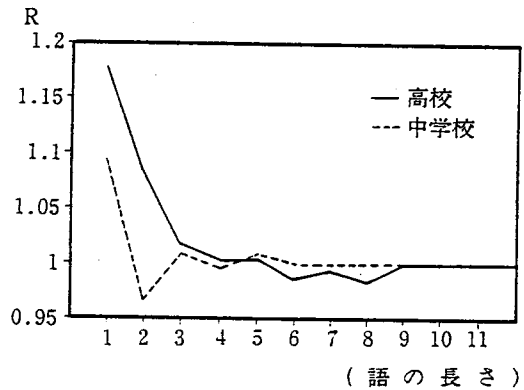


図9 Rの値と語の長さ（片仮名表記形）

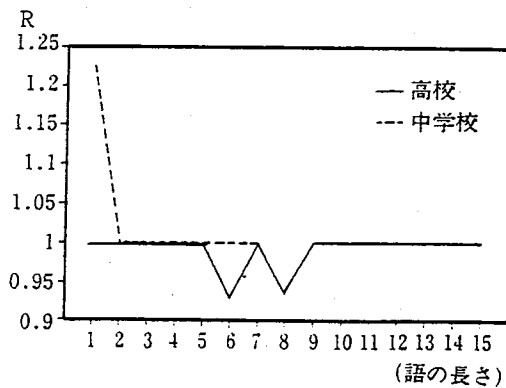


図10 Rの値と語の長さ（その他の表記形）

推定値 (単位：千)

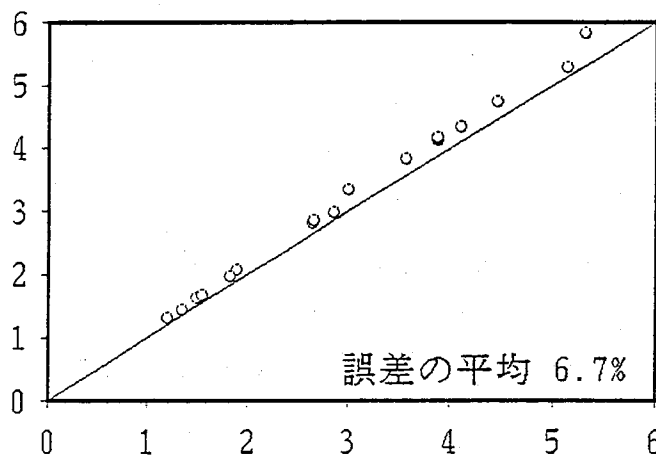
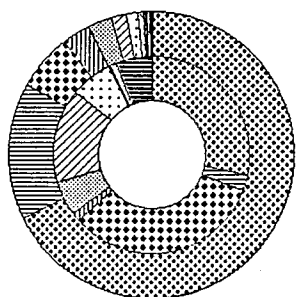


図 1 1 異なり語数の単純推定値 異なり語数 (単位：千)

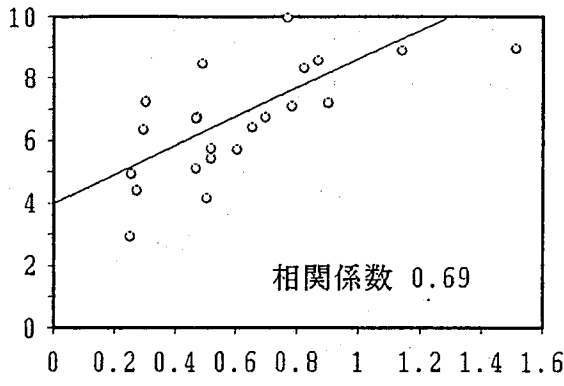


	異なり語数 (外円)	のべ語数 (内円)
漢	67.6%	28.2%
片	15.5%	2.6%
平	6.3%	32.5%
混	3.2%	1.4%
平混	2.6%	5.6%
他	1.8%	15.5%
漢平	1.3%	6.7%
漢片	0.6%	0.2%
漢混	0.4%	0.3%
漢平混	0.3%	0.9%
平片	0.2%	0.1%
漢平片	0.1%	6.1%
その他	0.1%	0.0%
合計	100.0%	100.0%

図 1 2 1 語の持つ表記パターン

データ：高校教科書

推定誤差*



A…漢字，片仮名，その他の
異なり表記形数の和
B…平仮名，混ぜ書きの異
なり表記形数の和

$\frac{B}{A}$

図 1 3 推定誤差と表記種の割合との相関

$$* \text{推定誤差} = \frac{\text{異なり語数推定値} - \text{異なり語数}}{\text{異なり語数推定値}} \times 100$$

推定値（単位：千）

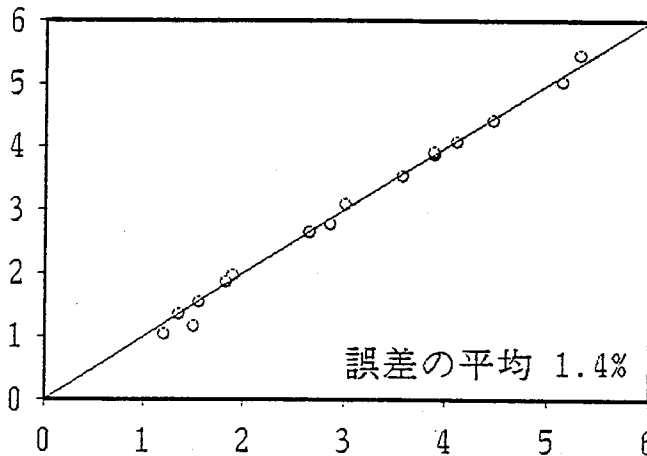


図 1 4 異なり語数の推定値

異なり語数
(単位：千)

$$K' = E - \frac{E}{100} \left[4.647 \frac{B}{A} + 3.978 \right]$$

$$(E = 1.279C_1 + C_2 + 1.513H_1 + 0.992H_2 + 0.731H_3 + \\ 0.648M_2 + 0.735M_3 + 0.845M_4 + 1.01K_1 + 1.02K_2 + \\ K_3 + O)$$

但し、 K' ：異なり語数推定値、 $C_1 \sim B$ は以下のそれぞれの表記語に対応する異なり表記形数。

C_1 ：漢字1字 C_2 ：漢字2字以上 H_1 ：ひらがな1字 H_2 ：ひらがな2字

H_3 ：ひらがな3字以上 M_2 ：混ぜ書き2字 M_3 ：混ぜ書き3字

M_4 ：混ぜ書き4字以上 K_1 ：カタカナ1字 K_2 ：カタカナ2字

K_3 ：カタカナ3字以上 O ：その他

A ：漢字+カタカナ+その他 B ：ひらがな+混ぜ書き

E ：異なり語数の単純推定値

6. まとめと課題