

日本文の読みやすさの評価式

達石由佳、小野芳彦、山田尚勇
(東京大学理学部情報科学科)

日本文の表面の情報から、構文や意味によらないでその文章の読みやすさを評価する式を、読みやすさと関係のある表面情報のうちの4種類、すなわち(1)文の平均の長さ(文字数)、(2)各文字種(英字、ひらがな、漢字、カタカナ)について、その文字種の連(同一文字種の文字の一続き)の相対頻度、(3)文字種ごとの連の平均の長さ、(4)読点の数の句点の数に対する比、から線型式により求めた。主成分分析により、読みやすさに関係のある成分を見つけ、その計算式を評価式とした。この成分はサンプルとしてとった科学技術系の日本文におけるスコアの分布が、読みやすさについての経験的知識とよく一致した。また、このスコアを読みやすさの指標に使えることを、クローズ法と、それにかかる時間の計測とを用いた実験により確かめた。

Derivation of a Readability Formula of Japanese Texts

TATEISI Yuka, ONO Yoshihiko, YAMADA Hisao

Department of Information Science, Faculty of Science, University of Tokyo,
7-3-1, Hongo, Bunkyo-ku, Tokyo, 113, Japan

A readability formula is obtained that can be used by computer programs for style checking of Japanese texts that need not syntactic or semantic information. The formula is derived as a linear combination of the surface characteristics of the text that are related to its readability: (1) the average number of characters per sentence, (2) for each type of characters (Roman alphabets, kanjis, hiraganas, katakanas), relative frequencies of runs (maximal strings) that consists only of that type of characters, (3) the average number of characters per each type of runs, and (4) *tooten* (comma) to *kuten* (period) ratio.

To find the proper weighting, principal component analysis (PCA) was applied to these characteristics taken from 77 sample texts.

We have found a component which is related to the readability. Its scores match to the empirical knowledges of reading ease. We have also obtained experimental confirmation that the component is an adequate measure for stylistic ease of reading, by the cloze procedure and by the examination on the average time taken to fill out one blank of the cloze texts.

1. はじめに

この研究の目的は、計算機で日本語の文体をチェックするために利用できる、読みやすさの指標 (readability index) を計算するための式 (readability formula) を導出することである。読みやすさの指標とは、文章の表面情報から、構文、意味解析によらずにその文章の (内容や、レイアウトではなく) 文体による読みやすさ・読みにくさを推定する指標である。

英文の読みやすさの指標としては、Flesch (1949) の Reading Ease Score や、Smith と Kinkaid (1970) の Automated Readability Index (ARI) などが知られ、利用されている (Cherry 1981, Coke 1983)。これらの指標は、文章に含まれている文の1文あたりの平均の語数と、語の1語あたりのシラブル数 (Flesch) や文字数 (ARI) から、線型式を用いて計算される。Reading Ease Score は、文あたりの単語の数を w 、語あたりのシラブルの数を s として、

$$F = 206.835 - 1.015 \times w - 84.6 \times s$$

で表される。Flesch は、 F が80以上の文章は「難しい」、45以下の文章は「やさしい」としている。ARI は、文あたりの単語の数を w 、語あたりの字数を l として、

$$G = 4.71 \times w + 0.50 \times l - 21.43$$

で表される。 G は、その文章を読むのに何半年程度の読解力が必要を示す数値である。これらの式は、英語では長い単語は一般にはなじみのうすい語であることが多いことと、長い文は構文が複雑になりやすいことに基づいている。すなわち、語の長さは文中の語彙の難しさを、文の長さは構文の難しさを反映している。

日本語では、文あたりの文字数と漢字の文字全体に対する割合を測り、それぞれから独立に難しいかやさしいかの評価する方法が知られている。森岡 (1958) では、小・中学校の教科書と専門雑誌、総合雑誌、大衆雑誌、児童雑誌の調査から、文の長さが平均60字以上の文章を「難しい」、35字以上の文章を「やさしい」とし、また、漢字の割合が35%以上の文章を「難しい」、20%以下の文章を「やさしい」とする評価を与えている。また、安本 (1983) は、現代の文学作品の調査から、文の平均の長さが51字以上の文章は「文が長い」、35字以下の文章は「文が短い」とし、また、1000字あたりの漢字の数が405字以上の文章は「漢字が多い」、299字以下の文章は「漢字が少ない」として、文を短く、漢字を少なくすれば文章は読みやすくなる、としている。日本語でも、英文と同様、長い文は構文が複雑になりやすい。したがって、文の平均の長さは、文章の構文の複雑さを反映している。漢字の割合は、語彙の難しさを反映している。森岡 (1958) によれば、小・中学校の教科書では、学年が進むにつれて漢字が多くなる。また、安本 (1983) は、多くの漢字を含むことは多くの語を含むことと平行的であり、読解力を要するので、漢字の多い文章は、大人にとっても読みにくい、と述べる。

一般に文中の漢字をひらがなに直すと文字数は増える。文の (文字数で測った) 長さは文の表記のしかたに依存する。また、漢語は、対応する和語よりも字数が少ないことが多い。漢字の割合と文の長さを独立に評価すると、矛盾した結果を得る可能性がある。森岡・安本両者とも、2つの要素を1つの指標にまとめることはしていない。

阪本 (1967) は、子供向けの本が、対象となる子供の読解力に合うかどうかを評価する基準を提案した。これは、次の3つの独立した評価からなる (表1-1)。

- (1) 「教育基本語彙表」(阪本 1958) に基づく基本語彙の語全体に対する割合。
- (2) 10語以上からなる文の文全体に対する割合。
- (3) 漢字の文字全体に対する割合。

阪本の方法は、文の長さを語単位で測ることによって、文の長さが表記に依存するという問題を避けている。しかし、語を単位とすることは、日本語の場合は実用的でない。日本語は一般に分かち書きをしないので、文の語への切り分けが、特に機械的にはやさしくないからである。

既存の評価はまた、カタカナの存在を無視しているという点にも問題がある。カタカナは、外来語の表記に使われる。最近の、特に技術系の文章では、外来語が多数使われ、その結果、文中のカタカナの割合が増えている。渡辺ら (1983) によれば、「情報処理」誌第17巻に含まれる文字の約8分の1がカタカナである。また、佐竹 (1982) によれば、現代の各種の雑誌に使われる文字のうち4.44%から13.75%がカタカナである。既存の方法は、漢字しか考慮していないので、現代の文章を扱うには不十分である。

2. 読みやすさの要因

われわれは、日本語の読みやすさの要因として次の4つをとった。

- (1) 文字の種類ごとの文字の頻度
- (2) 連 (同一文字種の文字の一続き) の長さ
- (3) 文の長さ
- (4) 文あたりの読点の数

これらのうち、(1)と(2)は、語彙の難しさに関する要因であり、(3)と(4)は、構文の複雑さに関する要因である。文字の頻度は、少なくとも次のように文章の読みやすさに影響する。漢字は、前述のように文章を難しくすると考えられている。カタカナやアルファベット (ローマ字) は、外来語を表記するのに使われるので、これらの文字を多く含むことは、なじみのうすい語を多く含むことを意味している。ひらがなは、漢字と逆に、文章をやさしくすると考えられている。日本語には厳密な意味での正書法がなく、同じ語をひらがなで書くことも漢字で書くこともあるいはカタカナやローマ字を使って書くこともできる。しか

し、日本人の大人が普通に文を書くときに、どの語に対してどの文字を使うかはほぼ決まっている。漢字は名詞、あるいは、活用する語の語幹を表記するのに使う。ひらがなは、活用語尾やその他の機能語に使う。カタカナやアルファベットは、外来語（多くは名詞）に使う。したがって、ふつうの文章では文字の頻度は語彙の使い方を反映し、文章の読みやすさに影響する。

文字種の境は、語、あるいはそれより小さな文法的単位の境にはほぼ対応するので、同一文字種の、異種の文字で区切られた1統きはほぼ語と対応する。文字は種類によって外見が異なるので、普通の表記の文章では、文字種の境が読み手が文を語に切り分けるのを助ける。したがって、同一文字種が長く続けて現れると、語の境が見にくくなり文は読みにくくなる。同一文字種の、異種の文字で区切られた一統きを連 (run) と呼ぶ。

特に、漢字は連が長いと読みにくさの原因となる。漢語は、その構成要素を単につなげるだけで複合語を作ることができる。この場合、構成要素同士がどのような関連を持っているのかは明示されない。そのために、漢語は意味があいまいになることがある。

各文字種の連の連全体に対する頻度は、その文字種の文字の文字全体に対する頻度と相関が高い ($p \geq 0.6$) ので、連の頻度は文字の頻度に代わることができる。

日本語でも、文の長さは長いと読みにくくなることが多い。小・中学校の教科書では、学年が進むにつれて、文章の平均の長さが長くなる。阪本 (1963) は文中に含まれる語数の方が字数よりも学年の差を表してはいるが、両者は比例しており、文長を字数で測ってもよい、としている。

文あたりの句の数が多いと文は読みにくくなる。林 (1959) は中学生と高校生に対する実験で、修飾句のある文を書きかえて修飾句を独立な文にした文章もとの文章を比較して前者の方が読みやすいことを示している。また、小鶴 (1987) の調査によれば、教科書は学年が進むにつれて1文あたりの句の数が増える。読点は、句の終わりに打たれるので、1文あたりの読点の数は、1文あたりの句の数に対応する。

3. 解析の方針

文章から読みやすさの要因となる表面情報を抽出し、それらの数値の線型式として評価式を構成する。連や文の長さの単位には、計算を簡単にするために字数をとった。

文章によって表面情報の数値は異なるが、このばらつきはいくつかの要因を含む。たとえば、文章の内容によって語彙が変わるので、文章の内容あるいはその属する分野によって文字や連の頻度は異なる。また、文体によっても表面情報の数値は変わる。文体は、文章の目的によって変わりうる。たとえば、初心者向けの入門書は、専門家向けの論文などよりも、解かりやすくななければならないから、筆者は読みやすいように注意を払うであろう。また、翻訳文は、しばしば原語の構文そのままの文体で書か

れる。このような文体はしばしば日本語として不自然で読みにくくなる。このような、文章の目的による文体の特徴を抽出することができれば、それを基準として読みやすさを測ることもできるであろう。

われわれは、まず文体の特徴を抽出するために、とった表面情報の数値から主成分分析 (PCA) によって特徴となる成分を抽出した。次に、文章における各成分のスコアのばらつきを読みやすさに関する経験的知識と照らし、読みやすさを表す成分を選ぶ。主成分は線型式で計算するので、その計算式を評価式として使うことができる。

読みやすさの要因を表す変数として、次のもの (4種類、10変数) をとった。

(1) 各字種について、その字種の連の、連全体に対する頻度。前に述べたように、字種ごとの文字の頻度と連の頻度には相関が高いので、連の頻度で代表した。字種としては、アルファベット、ひらがな、漢字、カタカナの4つをとった。

(2) (1) であげた各字種について、文字種ごとの連の長さの平均。

(3) 文あたりの平均文字数。文の長さは、となりあった句点、疑問符、感嘆符の間の字数で測った。

(4) 読点の数の句点の数に対する比。これは、文あたりの平均の句の数に対応する。

データをとる文章 (これをサンプルという) は以前に述べたことから、大人が普通に書くような表記法のものでなければならない。不自然な表記の文章、たとえば小学校の教科書などは使えない。小学校の教科書は、学習者がまだ習っていない漢字をひらがなで書くために不自然な表記となっているからである。サンプルとして、われわれの研究室にある機械可読テキスト70編と、読みやすさと関わりのある軸を探るための補助となるテキスト7編の計77編をとった。70編の文章は、情報科学に関する論文、翻訳、雑誌記事、入門書である。7編のうち、5編は、作文技術についての本および日本語についてのエッセイからとった節で、これらは専門家向けの論文などよりは読みやすいと考えて加えた。残る2編はコンピュータプログラムの著作権に関する裁判の判決文と、コピーや磁気ディスクの文書性と偽造罪についての解説である。法律関係の文章は読みにくいと一般に言われる。

文中の、図、表、参考文献、文と独立して示された数式・論理式は除いた。主成分分析には、東京大学大型計算機センタの Vax 8600 上の S (Becker 1984) を用いた。

4. 結果

最初の3成分 (累積寄与率70%) について検討する。各変数の負荷量を表4-1に示す。また、テキストの主成分スコアによるプロットを図4-1に示す。図において、i は入門書、m は論文でない解説記事など、p は論文、t と T は翻訳を示す。また、

D と E はそれぞれ「読みにくい」サンプルと「読みやすい」サンプルである。

第1成分は次のような特徴を持つ。まず、表4-1から、この成分にはアルファベットの頻度の正の寄与が大きく、ひらがな、漢字の負の寄与が大きい。また、図4-1では、同じ文字で表されているが、第1成分スコアの大きいものは英語の略語や、埋め込まれた式を多く含むものであった。従って、第1成分は、テキストの分野の差を反映していると考えられる。

第2成分は、漢字連の長さ、カタカナ連の長さ、文の長さの負の寄与が大きい。これらは、以前に述べたことから、文の読みやすさを下げる要因である。また、図4-1を見ると、「読みにくい」テキストと「読みやすい」テキスト、論文と入門書で値が異なり、「読みやすい」テキストと入門書で大きく、「読みにくい」テキストと論文で小さくなっている。この2つのことから、第2成分は読みやすさを反映していると考えられる。

第3成分は、漢字の多いテキストで高く、カタカナの多いテキストで低くなる成分である。

5. 第2成分と読みやすさ

第2成分が読みやすさを反映していることの傍証を2つ得ている。

まず、推蔽によって第2成分の値は大きくなる。サンプルのうち5編は、5人で分訳した論文を8人で推蔽した文章での、各人の担当分である。これらのスコアからそれぞれの推蔽前の原稿(77編の外)のスコアを引いた差を表5-1に示す。この表からわかるように、3成分とも、そのスコアは、推蔽の前後で一様変化している。第2成分の差は、他の2つの成分の差より平均して大きい。検定により第2成分の差の平均と第1成分の差の平均との差は5%で有意であった($p=0.044$)。第2成分と第3成分では、差は10%で有意であったが、5%で有意でなかった。 $(p=0.098)$ 。これらのことから第2成分は、推蔽により、他の成分より大きく変化するといえる。推蔽によって文章は読みやすくなるのが普通であるから、このことは、第2成分が読みやすさを反映することの1つの傍証となる。

第2に、第2成分のスコアは、その文章の文あたりの受身形の頻度と負の相関を持つ。表5-2に、各成分と受身形の頻度との相関係数を示す。

日本語の受身形は、可能な意味にも使われ、受身を多用すると意味があいまいになり、文が読みにくくなる。受身形は、牛島ら(1987)の方法に従って、パターンマッチングによりとり出すことができる。これを、文1000コあたりの頻度になおしたものと、第

2成分スコアとの相関を図5-1に示す。図中の直線は回帰直線である。図から「読みにくい」テキストでは受身の頻度が高く「読みやすい」テキストでは低いことがわかる。

6. 評価式

以上のことから、第2成分を読みやすさを反映する成分としてよい。まとめると、第2成分のスコアは、次の理由により読みやすさの指標として使うことができる。

- 1) 第2成分の値の大小は、読みやすい/読みにくいテキストに関する主観評価と一致する。
- 2) 入門のためにやさしく書かれたものは、専門家向けの論文よりスコアが高い。
- 3) スコアは推蔽により大きくなる。
- 4) スコアは受身の頻度と負の相関を持つ。

第1、第3成分では、これらのことが全て成り立つわけではない。第4成分以下は、テキストの表面情報のばらつきに関する寄与が小さい。従って、第2成分は、他の成分より読みやすさの指標として適当である。

第2成分の計算式を変換して、77編における平均が50、標準偏差が10になるようにしたもの(偏差値)が次の式である。

$$\begin{aligned}RS &= 0.05 \times pa + j0.25 \times ph - 0.19 \times pc - 0.61 \times p \\ &\quad - 1.34 \times ls \\ &\quad - 1.35 \times la + 7.52 \times lh - 22.1 \times lc - 5.3 \times lk \\ &\quad - 3.87 \times cp \\ &\quad - 109.1\end{aligned}\quad (1)$$

ここで RS: 評価

- pa: アルファベット連の連全体に対する頻度 (%)
- ph: ひらがな連の連全体に対する頻度 (%)
- pc: 漢字連の連全体に対する頻度 (%)
- pk: カタカナ連の連全体に対する頻度 (%)
- ls: 文の平均長さ (文字)
- la: アルファベット連の平均長さ (文字)
- lh: ひらがな連の平均長さ (文字)
- lc: 漢字連の平均長さ (文字)
- lk: カタカナ連の平均長さ (文字)
- cp: 句点あたり読点の数

である。

RSの値のめやすを表6-1に示す。ここで「教科書」とは中学、高校の理科・社会の教科書から1単元分の文章をとったもの、各5編における値である。

7. RS の妥当性

文章の表面情報から、文章の読みやすさの指標の候補 RS を得たが、上にあげた傍証の他に、RS を指標としてよいという半断をうらづけるため、RS の大きいテキストと小さいテキストの間の差を調べるためにクローズ法 (Taylor53) を用いた実験を行った。

クローズ法 (cloze procedure) は、いくつかの文章の相対的な読みやすさを測る方法で、次のような手続きをとる。

- (1) ほぼ同じ長さの文章の一部の単語を、その文中でのほたらきや、意味を考えずに、一定の割合で抜き取る。抜き取る割合は、10%から20%の間とし、ランダムに抜き取るか、または、一定の間隔で抜き取る。
- (2) 単語を抜き取った場所には、一定の長さの空白を入れる。
- (3) このようにして作った文章を、何人かの被験者に見せ、残りの文脈から、空白に入るべき単語を書きこませる。
- (4) 各文章ごとに、もとの単語が正しく復元された数の、被験者全員についての合計を求める。これを、クローズスコアという。
- (5) また、クローズスコアを (空白の数×被験者の数) で割ったものを求める。これを平均的中率という。
- (6) 平均的中率が高い文章ほど、読みやすい文章である、とする。

日本語の文章についてクローズ法を適用する場合でも、語を抜き取りの単位とすればよいことがわかっている (芝 1958)。

実験は、次のようなことを行った。被験者は、東京大学理学部情報科学科の学生25人と大学院生3人の合計28人である。テキストは主成分の計算に用いた論文などの中から、RS の高いものと低いものを3つずつ (p1-p6) とった。このうち、p1 から p3 までは、RS が50をこえるもの、p4 から p6 までは、RS が50に満たないものである。各テキストは8語ごとに語を抜き取った。抜き取った場所には、5文字分のアンダースコアを置いた。各被験者には、このうちの、ランダムに与えられた3つについて復元を求めた。また、各テキストの復元を始めた時刻と、終了時刻を記録するように求めた。

さらに、臨時教育審議会の最終答申から、「科学技術の進展と教育」にかんする部分 (r1) および、それを、文を短くし、さらに漢語を和語におきかえるように書きかえたもの (r2) を加えた。r1 から r2 への書き換えの意図は、RS の値を大きくすることである。実際、r1 に対する RS は 27、r2 に対する RS は 47 であった。被験者は、このうち的一方についても、復元を求められた。

復元されたテキストのうち、途中でやめたもの、時間の記入のなかったものをのぞき、残ったものについて、平均的中率 (%)

と、各々の被験者が、そのテキストの復元にかかった時間を空白の数で割ったもの、すなわち、空白1か所を埋めるのににかかった時間 (秒) を計算した。

p1 から p6 について平均的中率と、空白1か所あたりの時間の中央値を表7-1に示す。RS と的中率との相関はみられなかった ($p = 0.295$)。

p1 から p6 を、RS の高い3つと RS の低い3つに分け、空白1か所あたりの時間に差があるかどうかの中央値検定を行った。表7-2に示すように、時間の差は有意であった ($\chi^2 = 6.722$, $df=1$, $p < 0.05$)。

r2 の的中率は59.6%で、r1 の的中率56.6%よりも大きかった。被験者ごとの的中率を中央値検定した結果、r1 と r2 の的中率の差は、有意ではなかった。また、空白1か所あたりの時間も、r1 では10.9秒、r2 では9.6秒で r2の方が短い。中央値検定の結果は有意でなかった。

抜き取られた語を正確に復元することは、読者の背景知識と文章の内容の属する分野がかけ離れているときの方が、読者になじみのある文章の場合よりも難しい。平均的中率の差は、内容の難しさの差を含んでいる。

RS は、空白を埋めるのにかかる時間と相関がある。語を文字列として書きこむのにかかる時間はほぼ一定であるから、空白を埋めるための時間は、文脈を読んで、何が欠けているかを推測するのにかかる時間であると考えられる。したがって、RS の大小は、文脈のつかみにくさ、文章の字づらから内容をつかみとる難しさと関係すると考えることができる。このことから、RS は、平均的中率との相関は低いが、文章の表面的、文体的な読みやすさを反映しているといえる。

8. 式の簡略化

(1) 式を見ると、ls、la、lh、lc、lk に対する係数が、pa、ph、pc、pk に対する係数より大きいことがわかる。各項の値の、77編での平均の大きさ (絶対値) を表7-1に示す。この表から、連の長さに対する項は対応する文字種の頻度に対する項より大きくなる。このことは、pa~pk までの項を省略できる可能性を示唆している。

ls、la、lh、lc、lk、pc の6変数のみについて、10変数の場合と同様に主成分をとった。主成分同士の相関を表8-2に示す。

第1成分がもとの第2成分と相関が高い。また、スコアの分布も、第1成分は、10変数の第2成分と同様の傾向を示す。

6変数による第1成分をやはり77編の平均が50、標準偏差10となるように変換したものが、次の式である。

$$RS' = -0.12 \times ls$$

$$-1.37 \times la + 7.4 \times lh - 23.18 \times lc - 5.4 \times lk$$

$$-4.67 \times cp$$

$$+115.79$$

cpを除いた5変数に対する分析では、もとの主成分と相関の高い主成分は現われなかった。

10変数をとった根拠として、語いの差を文字の頻度と文字種ごとの連の長さが反映していることをあげたが、この結果は連の長さだけで十分であることを示している。

9. 結論と今後の課題

主成分分析の手法を用いて、読みやすさの指標を得た。これは、文体のもよその読みやすさを示すものである。英文の指標の場合もそうであるが、書きかえによって指標の値を高く(良く)しても、必ずして読みやすくなるとは限らない。例えば、文をすべてひらがなで書いた上、句ごとに独立した文にすれば、lhが増え、ls、lk、pc減るのでRSは大きくなる。RSは自然な表記で書かれたものにしか適用できない。ひらがなが不自然に多いための読みにくさを検出するためには、lh(あるいはph)についての2次項を含む式が必要かもしれない。普通の漢字かなまじり文では、かなの現れ方は文字ごとに一律ではない。漢字を単純にかなに直すと、普通の漢字かなまじり文では頻度の低い文字(たとえば、ゆ、よなど)が頻度が高くなる可能性もある。ひらがなの文字種だけではなく、特定の文字の頻度を調べることにより、不自然にひらがなの多い文章を検出することも可能であろう。

参考文献

- (Becker 1984) Becker, R. A. and Chambers, J. M., "S: An Inter-active Environment for Data Analysis and Graphics", Wadsworth, Belmont, California, 1984.
- (Cherry 1982) Cherry, L. L., "Writing Tools", IEEE Transactions on Communications, Vol. 30, No. 1, pp. 100-104, 1982.
- (Coke 1983) Coke, E. U. and Koether, M. E., "A study for the Match Between Technical Documents and the Reading Skills of Technical Personnel", Bell System Technical Journal, Vol. 62, No. 6, pp. 217-226, 1983.
- (Flesch 1949) Flesch, R., "The Art of Readable Writing", Harper, 1949.
- (林 1959) 林四郎, 「読みの能力と読みやすさの要因と読まれた結果と」, 計量国語学, Vol. 11, pp. 20-33, 1959.
- (小鶴 1987) 小鶴康浩, 「日本語文の読みやすさの評価に関する基礎的研究」, 情報処理学会第34回全国大会, 1987.
- (森岡 1958) 森岡健二, 「リーダビリティ」, 遠藤嘉基他編「コトバの科学」第5巻「コトバの美学」, 中山書店, 1958.
- (阪本 1958) 阪本一郎, 「教育基本語彙表」, 学芸図書, 1958.
- (阪本 1962) 阪本一郎, 「国語教科書の文の長さとその測定法」, 読書科学, VII, 2, 1962.
- (阪本 1963) 阪本一郎, 「文の長さの比重の測定法」, 読書科学, VIII, 1, pp. 1-6, 1963.
- (阪本 1967) 阪本一郎, 「読みやすさの基準の一試案」, 読書科学, XIV, 1, 12, pp. 1-6, 1967.
- (佐竹 1982) 佐竹秀雄, 「各種文章の字種比率」, 国立国語研究所報告 71, p.p. 327-346, 1982.
- (芝 1957) 芝祐順, 「読み易さの測り方クローズ法の日本語への適用」, 心理学研究 Vol. 28 No. 2, pp. 67-73, 1957.
- (Smith 1970) Smith, E. A. and Kinkaid, P., "Derivation and Validation of the Automated Readability Index for Use with Technical Materials", Human Factors, Vol. 12, pp. 457-464, 1970.
- (Taylor 1953) Taylor, W. L., "Cloze Procedure: A New Tool for Measuring readability", Journalism Quarterly, Fall 1953.
- (牛島 1987) 牛島和夫他, 「日本語文章推敲ツールにおける受身形の抽出法」, 情報処理学会論文誌, Vol. 28, No. 8, 1987.
- (渡辺 1983) 渡辺定久, 大岸洋, 「情報処理における用語と用字」, 情報処理学会日本文入力方式研究会10-2, 1983年5月11日.
- (安本 1983) 安本美典, 「説得の文章技術」, 講談社現代新書 685, 講談社, 1983.

	No. 1	No. 2	No. 3	No. 4	No. 5	No. 6	No. 7	No. 8	No. 9	No. 10
Alpha. r. f.	0.87	0.03	0.03	-0.04	0.17	-0.39	0.10	-0.04	-0.22	-0.05
Hira. r. f.	-0.93	0.19	-0.13	0.03	-0.03	0.04	0.11	0.10	-0.22	0.11
Kanzi r. f.	-0.92	-0.14	0.24	0.0	-0.18	0.10	0.09	0.08	-0.08	-0.15
Kata. r. f.	0.01	-0.25	-0.85	-0.26	-0.11	-0.02	-0.33	0.12	-0.04	-0.03
Sent. length	-0.72	-0.34	-0.10	-0.05	-0.04	-0.55	0.16	0.07	0.13	0.02
Alpha. r. l.	0.34	-0.37	0.04	0.75	-0.39	-0.07	-0.12	0.06	-0.03	0.01
Hira. r. l.	-0.63	0.54	-0.22	0.25	0.02	-0.14	-0.19	-0.38	-0.01	-0.02
Kanzi r. l.	0.0	-0.78	0.25	-0.39	-0.30	0.02	-0.04	-0.28	-0.06	0.04
Kata. r. l.	-0.04	-0.63	-0.53	0.28	0.29	0.20	0.32	-0.14	0.0	-0.02
Tooten per Kuten	-0.43	-0.54	0.36	0.13	0.50	-0.03	-0.35	0.05	-0.03	0.01
Eigenvalue	3.66	1.95	1.34	0.95	0.65	0.53	0.45	0.29	0.13	0.04
Proportion (%)	36.60	19.50	13.40	9.50	6.50	5.30	4.50	2.90	1.30	0.40
Cumulative (%)	36.60	56.10	69.50	79.00	85.60	90.90	95.40	98.30	99.60	100.00

r. f. = run frequency, r. l. = run length

Fig. 4-1. Principal Component Scores

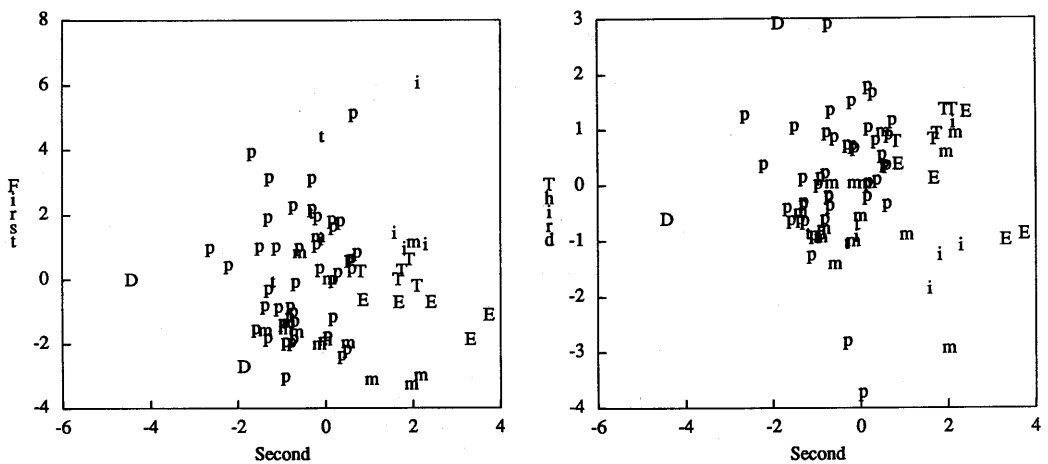


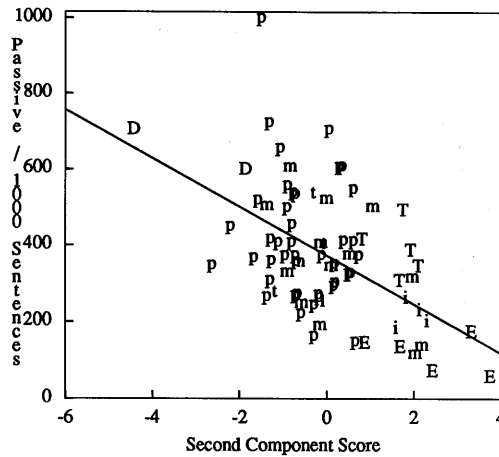
Table 5-1. Score Change by Improvement

	No. 1	No. 2	No. 3
text 1	0.47	1.18	0.34
text 2	0.11	0.60	0.87
text 3	0.41	0.38	0.23
text 4	0.0	0.42	0.20
text 5	0.50	1.33	0.43
mean	0.30	0.78	0.41
sdev	0.10	0.20	0.12

Table 5-2. Correlations to the Frequency of Passive Forms

	No. 1	No. 2	No. 3	No. 4	No. 5	No. 6	No. 7	No. 8	No. 9	No. 10
r	-0.25	-0.53	0.02	-0.15	0.09	-0.06	-0.20	-0.01	-0.09	0.03

Fig. 5-1. Frequencies of Passives



	p1	p2	p3	p4	p5	p6
RS	61.44	54.54	51.86	43.72	37.87	35.64
cloze %	66.58	63.81	56.28	56.25	64.69	60.46
time/blank (sec.)	6.89	6.19	8.18	8.75	9.06	8.05

	RS > 50	RS < 50
long	12	24
short	24	12

$\chi^2 = 6.72, p < 0.05$

Text Type		Max	Mean	Min
PCA Samples	Easy Indicators	76.8	67.2	56.2
	Difficult Indicators	36.7	27.5	18.3
	Technical Documents	66.4	49.4	31.2
Textbooks	Junior High School	59.9	55.2	48.5
	Senior High School	58.0	49.2	39.5

	variable	multiplier	term
Alpha r. f.	2.91	0.06	0.17
Hira r. f.	36.89	0.25	9.22
Kan r. f.	30.97	-0.19	5.88
Kata r. f.	3.83	0.61	2.34
Sent. len	51.19	-1.34	68.60
Alpha r. l.	3.31	-1.35	4.47
Hira r. l.	2.87	7.52	21.58
Kan r. l.	1.96	-22.10	43.39
Kata r. l.	4.11	-5.30	21.78
Tooten per Kuten	1.87	-3.87	7.24