

日本文における付属語の統計的性質とその利用例

山内 佐敏 林 大川 舟部 敏行

(株) リコー 技術本部

1. はじめに

表1 事務文228例中の付属語出現度数

現在、一般のオフィスで、文書処理の対象となる現代日本文は論文、報告書等で用いられ、一方的に意志を伝える“である”調と、手紙文等で用いられ、相手の意向を伺う丁寧語、謙譲語を多く含む“です、ます”調とでは、文体が大きく異なる。この文体の違いを特徴づけている要素は、一つには“お”“ご”の接頭辞が使用されるか、されないかの違いがあるが、もう一つには付属語の用い方に大きな違いがある。

また、日本文での単語の概念は、欧文のそれに比べてあまりはっきりせず、定義も論者によってまちまちである。そのため単語単位で分かち書きをさせようとすると、人によってのバラツキが大きいだけでなく、一個人でも時によって変化してしまう。それに対して文節の概念は、形式名詞、補助用言の定義で若干のゆれはあるものの付属語と自立語との境界は比較的わかりやすいため、文節単位での分かち書きをさせた場合、人によるバラツキは小さい。近年急速に普及してきた日本語文書処理機に用いられている入力方式には上記のような要素も考慮されていることであろう、文節単位の仮名漢字変換方式を採用しているものが多い。

本稿では、日本文の文体を特徴づける付属語の用いられ方を頻度と付属語間の接続確率の面から調べ、その結果と、それを用いて仮名漢字変換ソフトの速度向上検討を行なった内容をあわせて報告する。

2. 調査対象文と付属語の定義

報告調文体の文書として技術論文¹⁾の要旨とまえがき部分45文例をとり上げ、丁寧語を多く含む文体の文書として事務文例²⁾228例をとり上げた。

NO.	頻度	種類	付属語
1	2571	格助詞	の
2	1924	丁寧助動詞	ます
3	1830	格助詞	に
4	1274	格助詞	を
5	1055	副助詞	は
6	1048	接続助詞	て
7	812	形式名詞	うえ(他17語)
8	784	形式動詞	する
9	707	格助詞	と
10	645	断定助動詞	だ
11	582	過去助動詞	た
12	424	格助詞	が
13	269	形容動詞	だ
14	267	副助詞	も
15	265	補助用言的	ついて、つき
16	226	接続助詞	が
17	219	形式動詞	おる
18	218	格助詞	から
19	212	形式動詞	ある
20	205	準体助詞	の
21	138	形式動詞	よる
22	133	副助詞	か
23	128	接続助詞	と
24	128	形式動詞	なる
25	127	打消助動詞	ない
26	122	受身助動詞	れる
27	114	副助詞	まで
28	108	格助詞	して
29	94	断定助動詞	です
30	88	希望助動詞	たい
31	74	形式動詞	いる
32	64	使役助動詞	せる
33	61	形式動詞	おる
34	58	打消助動詞	ん(ぬ)
35	51	接続助詞	ながら
36	49	副助詞	なぞ
37	45	接続助詞	ば
38	38	接続助詞	ところ
39	37	格助詞	へ
40	35	可能動詞語尾	る
41	29	受身助動詞	られる
42	26	接続助詞	から
43	23	格助詞	で
44	20	接続助詞	に
45	18	形式動詞	くる
46	18	形式動詞	いう
47	15	格助詞	より
48	14	副助詞	ばかり
49	13	推量助動詞	う
50	13	用言性接尾辞	はらう

また、調査量の影響がどの程度あらわれるかをみるために、事務文例228例中の約1/10にあたる21例をとり出し比較した。

またここで付属語として扱うのは、一般の国語文法で言う助詞、助動詞の他、形式名詞、形式動詞、用言性接辞等、仮名表記する単語（自立語以外）を含めており、その総数は212種である。

3. 付属語の出現確率分布

事務文228例中で使用された付属語は全部で107種であった。上位50位の内容を表1に示す（総頻度数17631）。またこの出現確率分布特性、およびその累積分布特性を図1に示す。

この出現確率分布特性は、指数分布

$$P(x) = a e^{-bx} \quad : a, b \text{ は定数} \quad (1)$$

の形で近似できる。

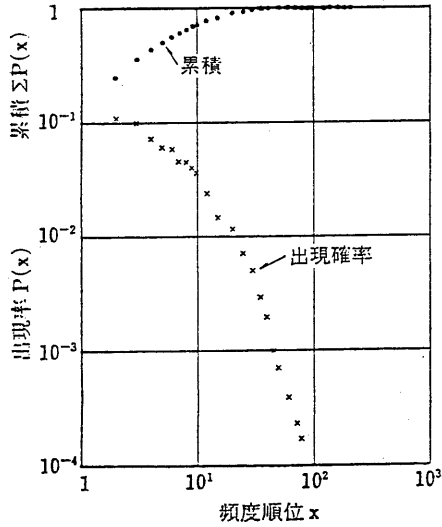


図1. 頻度順付属語出現確率 (事務文228例)

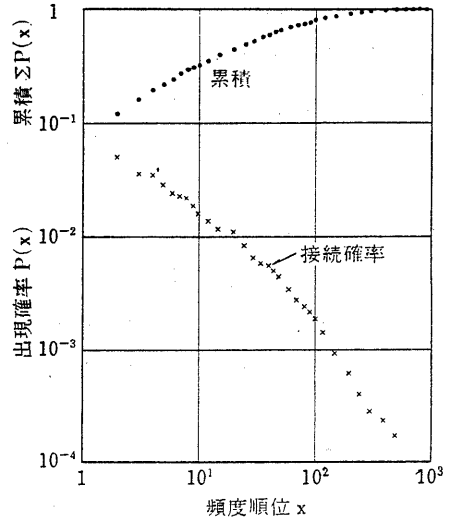


図3. 頻度順付属語接続確率 (事務文228例)

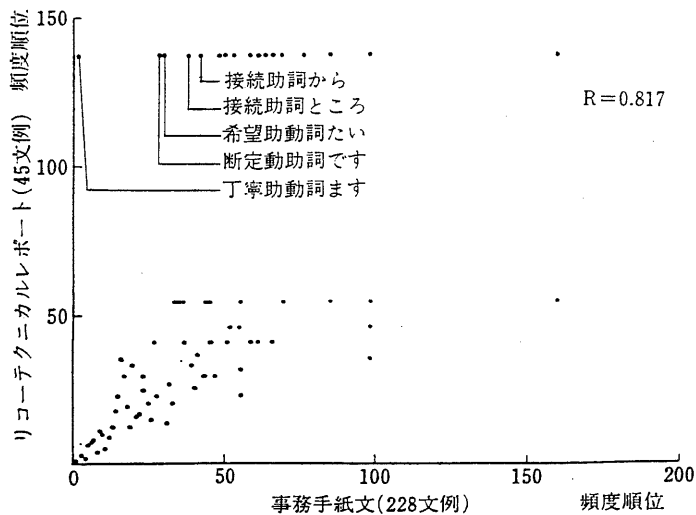
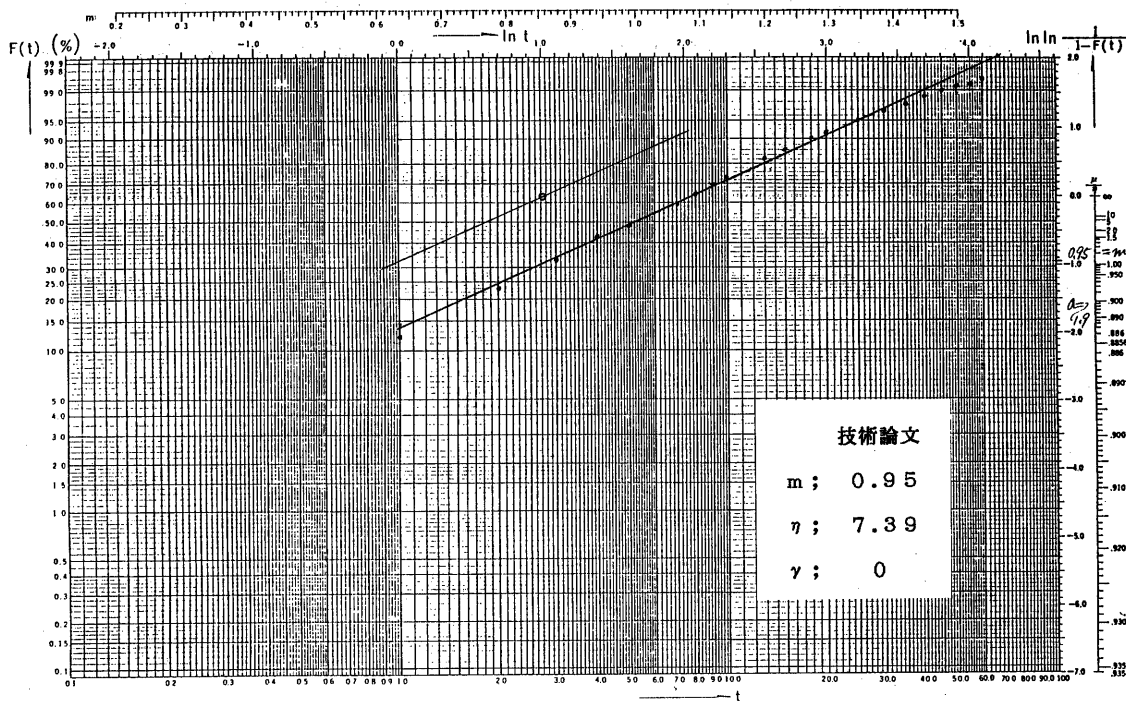
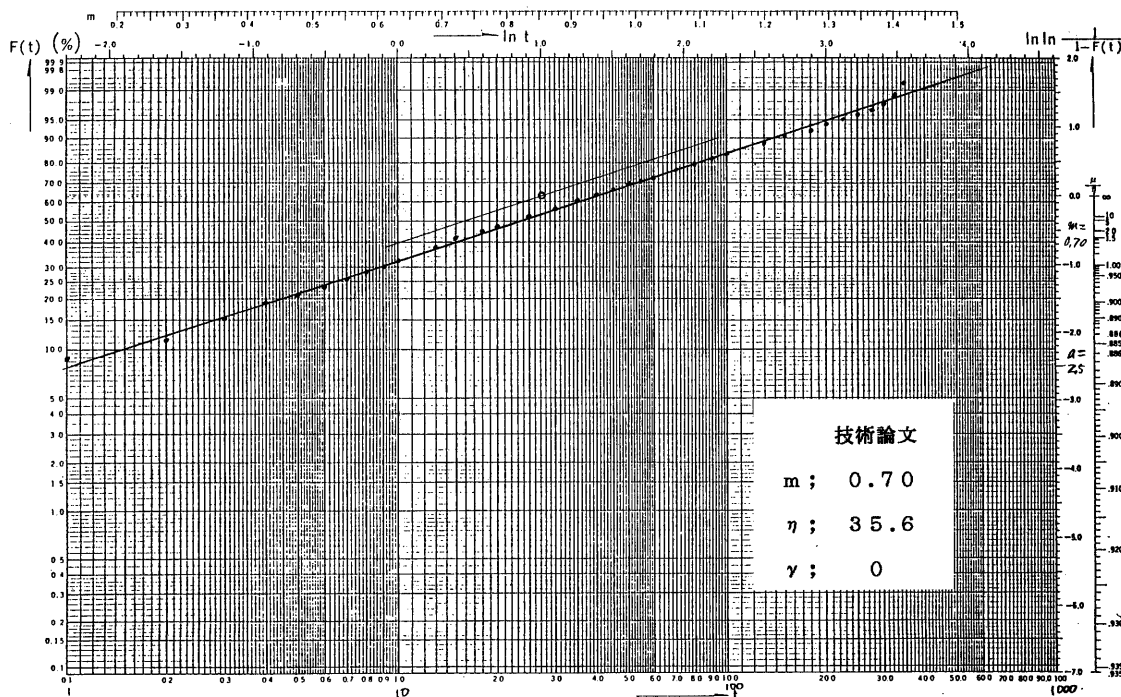


図2. 付属語の事務文228例と技術論文間の相関図



a) 出現頻度の累積分布



b) 接続確率の累積分布

図4. ワイブル確率紙による累積分布 (技術論文)

技術論文での出現確率分布特性および事務文21文例でのそれもほとんど同形であるが、各順位に対応する個々の語は異なる。

回帰した結果の各パラメータと実データと回帰式との相関係数 r は次のようになっている。

技術論文	事務文228例	事務文21例
$a; 8.7 \times 10^{-2}$	4.3×10^{-2}	6.1×10^{-2}
$b; 1.0 \times 10^{-1}$	7.1×10^{-2}	8.2×10^{-2}
$r; 0.98$	0.98	0.98

また出現確率の累積分布特性は (1) 式を積分しても良いがワイブル分布

$$F(x) = 1 - e^{-(x-\gamma)/\eta]^m} \quad (2)$$

ただし m : 形のパラメータ

η : 尺度のパラメータ

γ : 位置のパラメータ

で近似できる。日科技連ワイブル確率紙³⁾で回帰した結果の m と η の値を次に示す。

技術論文	事務文228例	事務文21例
$m; 0.95$	0.85	0.90
$\eta; 7.39$	8.82	7.56
$\gamma; 0$	0	0

図2に事務文228例と技術文との順位相関図を示す。特徴的な付属語としては、事務文228例で上位となった丁寧助動詞“ます”、断定助動詞“です”、希望助動詞“たい”等は、技術論文中には一度も使われていないが、その逆に技術論文中で上位となった付属語で事務文に使われていない語はない。

他の組合せの順位相関も調べたので、それぞれの相関係数 R を次に示す。

	R
事務文21例-事務文228例	0.880
事務文21例-技術論文	0.838
事務文228例-技術論文	0.817

4. 付属語の接続確率分布

図3に事務文228例の付属語の接続確率分布特性と累積分布特性を示す。

この接続確率分布特性も付属語の出現確率分布と同様指数分布式 (1) で近似できる。回帰した結果の各パラメータと回帰式との相関係数 r を同様の形で次に示す。

技術論文	事務文228例	事務文21例
$a; 7.7 \times 10^{-3}$	7.1×10^{-3}	5.0×10^{-3}
$b; 1.1 \times 10^{-2}$	4.8×10^{-3}	7.6×10^{-3}
$r; 0.92$	0.91	0.90

接続確率の累積分布特性もワイブル分布式 (2) で近似すると各パラメータは

技術論文	事務文228例	事務文21例
$m; 0.70$	0.59	0.63
$\eta; 35.6$	47.0	45.1
$\gamma; 0$	0	0

となる (このうち技術論文のワイブル分布図を、先の付属語の出現確率分布のそれとあわせて図4に示す)。

これも出現確率分布と同様、各分布は似たような分布となっているが、各順位に対応する個々の語は極端に異なる。各分布間の順位相関係数 R は

	R
事務文21例-事務文228例	0.630
事務文21例-技術論文	0.143
事務文228例-技術論文	0.325

となり、事務文と技術論文の間では無相関の方に近いことがわかる。

特徴的な接続の相違としては事務文で上位となった一般動詞連用形+丁寧助動詞“ます”、丁寧助動詞連用形“まし”+過去助動詞“た”、あるいは接続助詞“て”等の繋がりが技術論文中には一度も出

現せず、技術論文中で上位の、た行五段活用動詞連用形促音便“っ”+過去助動詞“た”、サ変動詞未然形“さ”+使役助動詞“せる”等の繋がり事務文中には一度も出現していない。

5. 文字種の頻度特性との比較

筆者等⁴⁾は以前に日本文務文書における文字種の調査を行なった。そのおり、文字種の頻度分布をワイブル分布で回帰している。それによると各パラメータは表6のようになる。

表6 日本文務文書中の各文字種のワイブル分布パラメータ

区分	m	η
記号	1. 1	4. 3 2
数字	1. 1	4. 5 3
英字	0. 8 4	1 5. 5
ひらがな	0. 9 7	1 2. 0
カタカナ	0. 8 3	1 6. 0
漢字	0. 7 0	1 5 0. 0
総合	0. 7 0	7 2. 2

注) 原典では $\alpha (= \eta^{-m})$ で表示されている。

この内容と付属語の内容を比較してみると、付属語の頻度分布特性では文字種の少ない文字セット、すなわち記号、数字、英字、ひらがな、カタカナに近く、付属語間の接続確率分布特性では文字種の多い文字セット、すなわち漢字、総合（日本文全体で使われる文字の総合）に近い特性であることがわかる。

6. 本調査の利用例

以上のような調査結果は、かな漢字変換装置等で行なう文法解析の効率向上に利用することができる。例えば、相沢ら⁵⁾が報告しているような方法、すなわち、自立語を最長一致法で仮確定し、接続行列表

の検索により順次接続可となっている付属語と後続の文字列を比較していき、文節末までの一致を確認する方法には有効である。

この方法では、接続行列表の検索において、接続の成功確率が大きいものを優先的に検索するほうが、検索時間を短縮させられる。そこで、最も適切な配列を決定する基準として、後続する付属語の並びを表わす列番号を、ランダムに配列した場合と使用される頻度の高い順に配列した場合とを比較し、どの程度検索時間が短縮されるかを検討した。

6-1. 接続行列表検索モデル

図5に検索モデルを示す。以下に検索モデルの動作を説明する。

まず、自立語が決定された段階で、その品詞をキーとして、Aを検索しCへ行く。Cにおいて、Dの行番号が決定されDへ行く。ただし、自立語が動詞の場合、活用形を決定してから、Dへ行く。Dの検索において、「1」の立っている列番号のところ、それをキーとして、Bを検索しCへ行く。Cにおいて付属語を検索する。まだ付属語候補の文字列が残っている場合はDへ行き、それ以外は、検索を終了する。

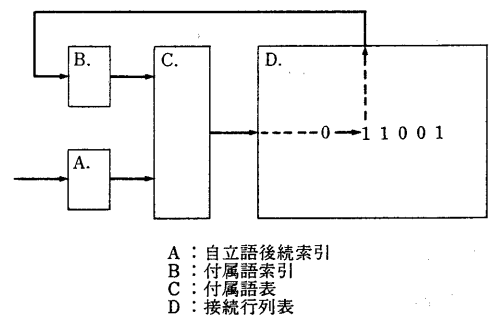


図5 接続行列表検索モデル

6-2. 検出時間の算出

検索モデルにおいて、一文節に対する平均検索時間を求める。平均検索時間 \tilde{T} は、次のように定義できる。

$$\tilde{T} = \tilde{T}_1 + \tilde{T}_2 + \tilde{T}_3 + \tilde{T}_4 \quad (3)$$

\tilde{T}_1 ; 接続行列表 (D) 内の平均検索時間

\tilde{T}_2 ; D→B→C→Dの平均検索時間

\tilde{T}_3 ; A→C→Dの平均検索時間

\tilde{T}_4 ; 自立語の切り出しに失敗した時の平均検索時間

\tilde{T}_1 は

$$\tilde{T}_1 = \hat{n} E(x) t_1 \quad (4)$$

\hat{n} ; Dの最初の列番号から検索する平均回数

$E(x)$; 接続が成功する列番号の期待値

t_1 ; x列からスタートして、x+1列を検索する時間

である。

\tilde{T}_2 は

$$\tilde{T}_2 = \hat{n} (3 \hat{f} t_2 + \hat{h} t_3) \quad (5)$$

\hat{f} ; 接続が成功する列 (その列も含む) までの1の平均個数

\hat{h} ; C内部で実際に検索する平均付属語候補数

t_2 ; 表間、索引・表間の移動の検索時間

t_3 ; 付属語候補1個の検索時間

である。

\tilde{T}_3 は

$$\tilde{T}_3 = 2 t_2 + \hat{h} t_3 \quad (6)$$

である。

\tilde{T}_4 は

$$\tilde{T}_4 = (\tilde{T}_1 + \tilde{T}_2 + \tilde{T}_3) \hat{m} \quad (7)$$

\hat{m} ; 自立語の切り出しに失敗した平均回数

よって式(3)は、

$$\tilde{T} = (1 + \hat{m}) \{ \hat{n} E(x) t_1 + 3 \hat{f} t_2 + \hat{h} t_3 \} + 2 t_2 + \hat{h} t_3 \quad (8)$$

となる。ここで $t_1 = t$ 、 $t_2 = a t$ 、 $t_3 = b t$ とすると (a、bは定数)

$$\tilde{T} = (1 + \hat{m}) \{ \hat{n} E(x) + 3 a \hat{f} + b \hat{h} \} t + 2 a + b \hat{h} \quad (9)$$

となる。

6-3. 接続行列表の構成

前述の検索時間のうちで接続行列表のみが配置順によって、大きく変動をうける。

1) ランダム構成

接続行列表をランダムに配列した時の接続成功列番号の期待値 ($E_r(x)$)、標準偏差 ($D_r(x)$) を求める。

期待値 ($E(x)$) は次式で算出する。

$$E(x) = \sum_{i=0}^n x_i \cdot p_i \quad (10)$$

x_i : 列番号

p_i : 各列番号で接続が成功する確率

標準偏差 ($D(x)$) は次式で求める。

$$D(x) = \sqrt{E(x^2) - \{E(x)\}^2} \quad (11)$$

ランダム構成において、 p_i は各列とも等しいので、式(10)は次のように変形できる。

$$E_r(x) = \frac{1}{I} \sum X_i \quad (12)$$

I : 列番号数

ここで、 $I = 212$ として計算すると、

$$E_r(x) = 107$$

となる。列を検索した場合、平均107番目の列で接続が成功することが期待できる。

標準偏差は、

$$D_r(x) = 6.1$$

となる。

2) 頻度順構成

接続行列表を頻度順に配列した時の接続成功列番号の期待値 ($E_o(x)$)、標準偏差 ($D_o(x)$) を同様に算出する。

式(10)は 次のように変形できる。

$$E_o(x) = \sum_{i=0}^n x_i \cdot P(x_i) \quad (13)$$

$P(x_i)$ は、図1の確率曲線で示される確率である。

式(1)の回帰式は、事務文228例を用いると、

$$P(x_i) = 4.3 \times 10^{-2} \times e^{-7.1 \times 10^{-2} x_i} \quad (14)$$

よって

$$E_o(x) = 1.5, D_o(x) = 1.4 \text{ となる。}$$

列を検索した場合、平均1.5番目の列で接続が成功することが期待できる。

3) 実用構成

上記のように接続行列表を頻度順に配列することにより、検索時間を短縮できるので頻度順構成にするのが良い。ただし、接続行列表のメンテナンス等を考慮すると、多少頻度順を犠牲にせざるをえない。

以下に、実用を考慮した接続行列表の期待値

($E_m(x)$)、標準偏差 ($D_m(x)$) を算出する。

$$E_m(x) = 2.6, D_m(x) = 2.1$$

平均2.6番目の列で接続が成功することが期待できる。

また、参考に、当初作成した(以下「初期構成」と言う)接続行列表を用いた場合の期待値 ($E_n(x)$)、標準偏差 ($D_n(x)$) を同様の手法で算出する。

$$E_n(x) = 5.0, D_n(x) = 3.7$$

平均5.0番目の列で接続が成功することが期待できる。

4) 各構成の比較

接続行列表内の平均検索時間 \bar{T}_1 について上記4つを比較する。表3にランダム構成の \bar{T}_1 を基準とした各 \bar{T}_1 の比率を示す。初期構成を実用構成と比較すると、

$$\text{初期構成 : 実用構成} = 1 : 0.51$$

で実用構成の平均検索時間 \bar{T}_1 は約半分になることが期待できる。

表3 平均検索時間の比較

構成	K
ランダム構成	1.0
頻度順構成	0.14
初期構成	0.47
実用構成	0.24

K…接続行列表内の平均検索時間の比率

6-4. 実文による検索処理速度向上の確認

事務文例より10文節をランダムに抽出し(表4にしめす)、初期構成と実用構成の接続行列表を用いて、検索回数を調べ、比較し、どの程度検索時間が短縮するかを検討した。

表4 抽出した文節

No.	文節	No.	文節
1	皆さまに	6	部分にも
2	上りましたが	7	郵便に
3	願えれば	8	検討の
4	拝見いたし	9	実績を
5	申し上げます	10	対しまして

表5. 接続行列表の検索時間 (単位; 検索回数)

構成	NO.	1	2	3	4	5	6	7	8	9	10	計	回数比
初期構成	46	83	82	234	33	79	46	45	48	86	782	1	
実用構成	25	53	56	118	11	37	25	24	26	52	429	0.55	

注: 2個以上付属語を持つものは最後の付属語の検索時間を取る。

表6. 全検索時間 (単位; 検索回数)

構成	NO.	1	2	3	4	5	6	7	8	9	10	計	回数比
初期構成	67	538	335	351	46	198	67	60	73	1184	2919	1	
実用構成	37	365	196	225	14	101	37	33	39	678	1725	0.59	

1) 接続行列表内の検索時間

表5は、表4の10文節について、初期構成と、実用構成の各接続行列表を用い、接続行列表内の実際の検索回数を調べた結果である。表5から、実用構成は初期構成と比較すると、検索時間は0.55倍に減少していることがわかる。これは6-3で述べた期待値0.51倍に近い。

2) 全検索時間

表6は、表4の10文節について、初期構成と、実用構成の各接続行列表を用い、実際の全検索回数を調べた結果である。表6から、実用構成は、初期構成と比較すると、全検索時間は、平均0.59倍に減少していることがわかる。

以上のように、配列に頻度順を考慮した接続行列表を採用することで、当初のものと比較し実用的にも半分近くに検索時間を短縮している。

7. おわりに

以上述べたように、付属語の使われ方が日本文中の文字の使われ方に似ていることがわかった。また、このような調査を行なうことにより、具体的な装置の改善を机上で行なえることを示した。今後は自立語や接辞の特性も調査し、日本文中の文字や単語の

ふるまいについて考察していきたい。

最後に本調査について御協力いただいた関係各位に、特にデータの整理に献身的につくしてくれた島田美佐子嬢に深謝いたします。

参考文献

- 1) リコー: "RICOH TECHNICAL REPORT" NO.1 ~NO.6, (1979 ~ 1980)
- 2) 山城, 安田: "企業経営文例全書" 1~4 ぎょうせい (昭53)
- 3) 日科技連編: 信頼性データの解析 (日科技連ワイブル確率紙の使い方), 日科技連ライブラリー 12
- 4) 村山: 日本文事務文書における字種の解析と応用, RICOH TECHNICAL REPORT NO.6 (昭和56年11月25日)
- 5) 相沢, 江原: 計算機によるカナ漢字, NHK技術研究, Vol 25, NO. 5, pp. 261~298 (1973)