

# BCCWJに含まれるウェブデータの特性について ——データ重複の諸相とBCCWJ使用上の注意点——

田野村忠温 (大阪大学大学院文学研究科)

## On Certain Properties of the Web-Based Subcorpora of BCCWJ

Tadaharu Tanomura (Osaka University)

### 1. はじめに

インターネット上に存在する日本語文書は膨大かつ多様で、言語研究資料としてきわめて大きな価値と魅力を有する。そのインターネット文書の短所と言えば、一般に予想されやすいのは特殊な言葉遣いの出現や書き誤りの多さといった点であろうが、実際に最も問題となるのはむしろ同一データの重複出現である(拙論(2010))。インターネット上にはさまざまな事情で同一の文章、段落、文、句が複製されて繰り返し現れる。

インターネット文書は「現代日本語書き言葉均衡コーパス(BCCWJ)」にも2つのサブコーパス「Yahoo!知恵袋」「Yahoo!ブログ」として収められ、両者を合わせてBCCWJ全体の約2割の分量を占めている。そしてこのたび、それらのサブコーパスにおいてもデータ重複の問題が思いのほか深刻であることが明らかとなった。

以下では、発表者が問題の認識に至った経緯と、2つのサブコーパスにおけるデータ重複の様相に関する調査・分析の結果について述べる。

### 2 BCCWJ N-gram分析ツール BNAalyzer

#### 2.1 ツールの概要

発表者は過日、BCCWJの検索結果を用いて語句のコロケーションを調べるための簡易な分析ツールBNAalyzer (Windows上で作動)を作成し、次のところで公開した。

<http://www.tanomura.com/research/BNAalyzer/>

BNAalyzerは、BCCWJ検索サイト「中納言」での検索結果をもとに、検索語の前後にどのような表現がよく現れるかを分析する。具体的には、検索語の前後文脈のN-gram (N個の短単位または長単位の連続)の一覧を作成し、エクセルで表示する。

#### 2.2 インストール方法および使用法

上記URLのページの説明に従ってBNAalyzerをインストールすると、デスクトップ上にアイコンが2つ作られる。それぞれを単純版、circumcollocate版と呼ぶ。

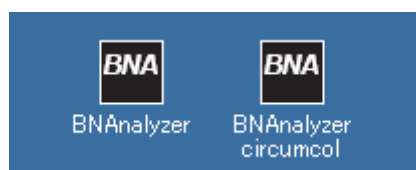


図1 BNAalyzerのデスクトップアイコン

BNAnalyzerを使ってBCCWJの検索結果からN-gramの一覧を得るには次のようにする。

- 1) 中納言で「検索結果のダウンロード」によって検索結果 (zipファイル) を取得
- 2) zipファイルをBNAnalyzerのアイコンの上にドラッグ&ドロップ

これによりエクセルの新しいブックが開かれ、単純版の場合は、検索語 (キー) の前文脈の末尾および後文脈の冒頭のN-gram (N=1~8) が頻度順に表示される。

次に示すのは「なかなか」の検索結果に基づく後文脈のN-gramの一覧である。

	A	B	C	D	E	F	
	1-gram	2-gram	3-gram	4-gram	5-gram	6-gram	
1	の(412)	できない(67)	うまくいかない(84)	出てこない(25)	うまくいきません。(19)	うまくいきません。_(6)	できるも
2	いい(161)	のもの(59)	出てこ(28)	うまくいきません(16)	出てこない。(10)	できることではない。(4)	できること
3	難しい(145)	うまくいか(46)	思うように(28)	うまくいかない。(12)	そうはいかない。(6)	できるものではありません(4)	うまくい
4	うまく(105)	出て(39)	できない。(19)	できません。(11)	のものだった。(6)	のものだった。_(4)	うまくいき
5	でき(101)	難しい。(38)	うまくいきませ(16)	そうはいかない(10)	うまくいかなかった。(5)	そうはいきません。(3)	そうはうま
6	、(98)	思うよう(28)	できません(15)	のものだ。(10)	お目にかかれない(5)	そうはうまくいきません(3)	できること
7	に(94)	の美人(27)	どうして(15)	のものだった(10)	のものである。(5)	思うようにはいかない(3)	できるも
8	むずかしい(68)	見つからない(25)	のもので(14)	帰ってこない(10)	出てこないの(5)	出てこない。_(3)	の美青年
9	出(57)	そうは(24)	見つからなかった(14)	思うようには(10)	できません。_(4)	出てこないの(3)	ふとんか
10	そう(54)	むずかしい。(21)	理解され(14)	できるものでは(9)	できることではない(4)	戻って来なかった。(3)	ピンとき
11	理解(52)	手に(21)	のものだ(18)	見つからなかった。(9)	できるものではありません(4)	うまくいかないことが多い(2)	気がつき
12	見つから(49)	いない(20)	帰ってこ(12)	寝つけなかった。(9)	できるものではない(4)	うまくいかないものです。(2)	気軽に作
13	手(41)	理解し(20)	手に入ら(12)	お目にかかれ(7)	の美人である。(4)	うまくいかなかった。_(2)	言うこと
14	大変(37)	興味深い(19)	そうはいか(11)	わかりません。(7)	帰ってこない。(4)	うまくいきませんでした(2)	思うよう
15	面白い(37)	いい。(17)	お目に(10)	手に入らない(7)	見えてこない。(4)	お目にかかることが(2)	出てこ
16	思う(35)	わからない(17)	できるもので(10)	消えなかった。(7)	見つからなかった。_(4)	お目にかかれない。(2)	消えな
17	い(34)	気が(17)	のものだ(10)	そうもいかない(6)	寝つかれなかった。(4)	できることではありません(2)	進みませ
18	ない(33)	進まない(17)	の美人で(10)	できることでは(8)	容易ではない。(4)	できるものではない。(2)	地元
19	おもしろい(31)	難しい(17)	わかりませ(10)	どうして。(6)	うまくいかないものです(3)	の見ものであった。(2)	容易では
20	よく(31)	うまくいき(16)	理解して(10)	のものである(6)	そうはいきません(3)	の美形である。_(2)	理解して
21	その(30)	大変な(16)	できないこと(9)	むずかしいものです。(6)	そうはうまくいきませ(3)	の美人でしたよ。(2)	恋愛など
22	気(30)	立派な(16)	見られない(9)	理解してもらえ(6)	のものであった(3)	の美青年であった(2)	
23	困難(30)	できませ(15)	寝つけな(9)	ありません。(5)	の見ものだった。(3)	ふとんから出る決心がつか(2)	
24	出(30)	出(15)	出(15)	出(15)	出(15)	出(15)	

図2 「なかなか」の後文脈のN-gram

この一覧は、「なかなか」の後には「の」「いい」「難しい」(1-gram)、「でき-ない」「の-もの」「うまく-いか」(2-gram)、「うまく-いか-ない」「出-て-こ」「思-う-よ-う-に」(3-gram)などの表現がよく現れることを示している。

circumcollocate版を使えば、例えば語彙素「惜しむ」の検索結果から、「寸暇を惜しんで」「努力を惜しまない」「骨身を惜しまず」のように「惜しむ」を前後からはさむように現れる2表現の慣習的な組合せがあることを知ることができる (circumcollocateの用語・概念については拙論(2010)を参照)。

以上のようなN-gramの頻度情報は、語句のコロケーションを分析する際の手がかりとなり得る。

### 3 BNAnalyzerによる不自然な分析結果とその原因

#### 3.1 不自然な分析結果

さて、BNAnalyzerを作ったあと、その動作確認のために随意に選んだ少数の検索語の検索結果を分析してみると、一見しておかしいと分かる結果が得られるケースが非常に多いことが分かった。次に2つの例を示す。

	E	F	G	H
1	5-gram	6-gram	7-gram	8-gram
2	見つけたいならYahoo!(25)	見つけたいならYahoo!縁結び	見つけたいならYahoo!縁結び	(見つけたいならYahoo!縁結び)(25)
3	わかった。_(「(10)	はわからなかった。_(4)	ダウンロードできます。◆メルマガ(	ダウンロードできます。◆メルマガやっ(4)
4	はわからなかった。(8)	ダウンロードできます。◆_(4)	駄目になるスキルでは(4)	駄目になるスキルではなく(4)
5	は答えなかった。(8)	ファーストメールを送りました(4)	売り切れてしまうそうなので(4)	売り切れてしまうそうなので、(4)
6	忘れてしまいます。(7)	食卓に出す。(4)	戻ってきた。_(「(4)	治療を受けたいときなど)医療(3)
7	戻ってきた。(7)	駄目になるスキルで(4)	いっぱいになってしまいます。(3)	食卓に出す。(4人分(3)

図3 「すぐに」の後文脈のN-gram

	D	E	F	G	H
1	4-gram	5-gram	6-gram	7-gram	8-gram
2	」といった(54)	に殉じて自分は(47)	気持ちに殉じて自分は(47)	の気持ちに殉じて自分は(47)	あなたの気持ちに殉じて自分は(47)
3	」などという(53)	のことを自分は(34)	すべてのことを自分は(29)	関するすべてのことを自分は(25)	に關するすべてのことを自分は(25)
4	次のような(53)	の気持ちに自分は(22)	あなたの気持ちに自分は(22)	あなたに自分はその成果を(11)	てあなたに自分はその成果を(11)
5	」という(51)	、次のような(20)	に自分はその成果を(11)	の気持ちに殉じて自分が(9)	あなたの気持ちに殉じて自分が(9)
6	ない」という(50)	ている」という(19)	気持ちに殉じて自分が(9)	【松下幸之助日々の(8)	【松下幸之助日々の(8)
7	殉じて自分は(47)	」という(17)	【松下幸之助日々の(8)	と言のように「その(8)	くればあなたの気持ちに自分は(8)

図4 「言葉」の前文脈のN-gram

いずれの例においても、とうてい一般性を持つとは考えられない「見つけたいならYahoo!縁結び」とか「あなたの気持ちに殉じて」といったN-gramがリストの上位を席卷している。

### 3.2 原因の調査

上の2例のうち、前者(図3)はYahoo!知恵袋かYahoo!ブログのいずれかのサブコーパスに原因がある可能性が高い。後者(図4)についてはこの分析結果だけからでは事情は分からない。

そこで、「BCCWJ-DVD版」——DVD媒体で頒布されているBCCWJ——とインターネットを用いて調べてみたところ、真相は次の通りであった。

まず、図3に見る「見つけたいならYahoo!縁結び」という高頻度N-gramは、ブログ記事の書き手が書いたものではなく、ヤフー株式会社の運営する出会い系サイトの宣伝用ブログ「そろそろ恋愛しませんか?」(<http://blogs.yahoo.co.jp/yjpartnerblog/>)の記事に自動的に張られるリンクのタイトル「[結婚相手をすぐに見つけたいならYahoo!縁結び]」の一部であった。リンクの実例は例えばインターネット上のYahoo!ブログの記事<http://blogs.yahoo.co.jp/yjpartnerblog/archive/2008/11/11>で見ることができる。この記事はBCCWJにサンプルID OY14\_29479として収録されている。

図4にある「殉じて」の奇異な用法を含む多数のN-gramはいずれも「世界で一番、誰よりも愛してる人へ」と題されたブログ(<http://blogs.yahoo.co.jp/geoburgher/>)に頻出するものであった。BCCWJにはこのブログから「殉じて」を含む記事が36件取られており、そこには「殉じて」が計243回も現れ、BCCWJ全体に含まれる「殉じて」計264例の実に92%を占めている。図4に見る、「自分」を含むほかのN-gramも同じブログに現れるものであった。

このように、不自然な分析結果の原因は、主にBCCWJのYahoo!ブログサブコーパスにおける同一データの重複出現にあることが判明した。

インターネットに高頻度で現れる一般性の低い表現は、人間がその都度書いたわけではなく、機械的に生成または複製されたものに過ぎない。言語研究上そのようなものを通常の言語データと同列に扱ってはならないことは明らかであるが、BCCWJ評価の観点から引き続き問われるべきは、コーパスにどのような重複がどれくらい含まれているのかという問題である。

以下においては、BCCWJのYahoo!ブログおよびYahoo!知恵袋の各サブコーパスの特性に



表 1 完全一致のサンプル一覧（部分）

サンプルID	字数
OY14_29957=OY14_34602	540
OY05_02386=OY05_02853	593
OY11_00268=OY14_10909	626
OY02_00210=OY02_00212	646
OY03_02330=OY03_03987	855
OY01_00055=OY11_00113	928
OY01_00506=OY01_00579=OY02_00125	947
OY01_00180=OY01_00214=OY01_00260=OY01_00552=OY01_00642= OY11_00260=OY11_00385=OY11_01547=OY11_01553	1038
OY04_01476=OY04_01665	1,172
OY01_01375=OY01_02646	1,365
OY14_40012=OY14_43821	2,467
OY14_49176=OY14_49589	2,467
OY11_03448=OY11_03449	2,863

なお、ここでの調査は、BCCWJ-DVD版のLUWディレクトリにあるタブ区切り形式データより復元したテキスト（文字コードはShift\_JISに変換）によった。他の方法で調査すれば統計（特に字数）に多少の違いが生じる可能性がある。

#### 4.2 サンプルの部分一致

サンプルが部分的に一致するものには、完全一致に近いものから小部分の一致にすぎないものまでさまざまな段階がある。また、一致・不一致のあり方も事例によって異なる。

部分的な一致を含むサンプルの数は完全一致のそれをはるかに上回るが、類似は程度問題であるので、部分一致を含むサンプルの範囲を明確にすることは原理的に不可能である。また、大量のテキストから部分一致を検出することには現実の処理上限界がある。

ここでは、BCCWJを使用するうえで特に深刻な問題を引き起こす可能性のある、データ一致の程度のはなはだしい事例を2つ見る。

第1の事例は、「犬幼稚園B u d d y D o g」のブログ (<http://blogs.yahoo.co.jp/lovedog111222>) から取られた以下のサンプル38件である（一部に完全一致のサンプルを含む）。

OY05\_00030, OY05\_00066, OY05\_00447, OY05\_00486, OY05\_00699, OY05\_01030, OY05\_01096,  
OY05\_01213, OY05\_01840, OY05\_02017, OY05\_02075, OY05\_02130, OY05\_02232, OY05\_02252,  
OY05\_02286, OY05\_02335, OY05\_02386, OY05\_02461, OY05\_02834, OY05\_02853, OY05\_03041,  
OY05\_03152, OY05\_03182, OY05\_03316, OY05\_03593, OY05\_03641, OY05\_03759, OY05\_04364,  
OY05\_04503, OY05\_04574, OY05\_04712, OY05\_04887, OY05\_04944, OY05\_05154, OY05\_05424,  
OY05\_06006, OY05\_06217, OY05\_06422

これらのサンプルにおいては、1文ないし数文のまとまりが、異なる組み合わせや順序において現れる。例えば、次はOY05\_02853=OY05\_02386のテキストであるが、

サンプルID : OY05\_02853=OY05\_02386

犬幼稚園B u d d y D o gに愛犬を預ける飼い主さん。送り迎えの際、愛犬たちがじゃれあう広場でお茶をしつつ、その間に情報交換タイムが始まります。しつけや健康管理の話はもちろん、去勢手術や避妊手術について、詳しい説明や報告をしたり、不安な事を質問したり。フードやおやつ選び方・与え方、犬が喜ぶおもちゃや、留守中に便利なグッズについてなど、話題は多岐にわたります。犬の幼稚園「B u d d y D o g」は自由登校システムのため、毎回少しずつ違うメンバーが顔を合わせて、情報は増える一方。みんな子(犬)育て真っ最中で、お互いに相談もしやすいようです。仔犬は本来、

目を輝かせ好奇心旺盛・天真爛漫で元気すぎる程です！落ち着きがない・無関心、無反応・それは仔犬の本質ではありません。犬達は犬幼稚園 Buddy Dogで仲良くじゃれあったり、時にはおもちゃを取り合ってみたり・遊び疲れて寄り添って眠っていたり・愛くるしい表情をいっぱい見せてくれます。その姿は本当に純粋で愛しい程です。『犬の社会性』を身につけることが、将来に良い子になる秘訣。「三つ子の魂百までも」は、人間も犬も一緒なんです。犬幼稚園 Buddy Dogは、仔犬にとって世界を広める第一歩でもあるわけです。犬幼稚園 Buddy Dogは、きっとあなたと愛犬の間に新しい発見と更なる楽しみをもたらしてくれるはずです。お気軽にご相談ください。

同じブログから取られた他のサンプルOY05\_02286=OY05\_02232においては、このテキストの後ろに短いテキストが付け加えられている。また、サンプルOY05\_04364では冒頭に短いテキストが付け加えられ、かつ、「三つ子の魂百までも」以下のテキストが削除されている。

上で色分けして示したテキストの各部分は、このブログにおいてしばしばテキストの構成要素として——ときには、わずかに修正された形で——繰り返し現れる。それを“要素テキスト”と呼ぶことにすれば、次に示すサンプルでは、上記テキストの4つの要素テキストが使われ、かつ、波下線を施した別の要素テキストが付け加えられた形になっている。また、最初の要素テキストの冒頭には「仔」が加えられている。

#### サンプルID : OY05\_02461

仔犬達は犬幼稚園 Buddy Dogで仲良くじゃれあったり、時にはおもちゃを取り合ってみたり・遊び疲れて寄り添って眠っていたり・愛くるしい表情をいっぱい見せてくれます。その姿は本当に純粋で愛しい程です。『犬の社会性』を身につけることが、将来に良い子になる秘訣。「三つ子の魂百までも」は、人間も犬も一緒なんです。犬幼稚園 Buddy Dogは、仔犬にとって世界を広める第一歩でもあるわけです。犬幼稚園は犬をしつけるのではなく、犬とじゃれあうことにより社会性を形成する。家族以外の人と接することにより人への信頼・服従を確立する。飼い主は飼い主として必要な知識を学んでいただく所です。犬幼稚園 Buddy Dogはきっとあなたと愛犬の間に新しい発見と更なる楽しみをもたらしてくれるはずです。お気軽にご相談ください。

程度のはなはだしい部分一致の第2の事例は、先に見た「世界で一番、誰よりも愛する人へ」(<http://blogs.yahoo.co.jp/geoburgher/>)のブログの記事である。このブログからはBCCWJに少なくとも次の36件のサンプルが取られている。

OY14\_04922, OY14\_07061, OY14\_09969, OY14\_14614, OY14\_15863, OY14\_17848, OY14\_22996, OY14\_32704, OY14\_33189, OY14\_34427, OY14\_34837, OY14\_34962, OY14\_35316, OY14\_36252, OY14\_38580, OY14\_38725, OY14\_40012, OY14\_42324, OY14\_42502, OY14\_42913, OY14\_43821, OY14\_44162, OY14\_44932, OY14\_45460, OY14\_46975, OY14\_47226, OY14\_47461, OY14\_49176, OY14\_49589, OY14\_50563, OY14\_51147, OY14\_51273, OY14\_53034, OY14\_53047, OY14\_53827, OY14\_54064

このうちの1件のサンプルOY14\_53034の冒頭の2文——第2の句点までの内容——だけを示せば次の通りである。

#### サンプルID : OY14\_53034

(3からつづく)今夜から明日にかけて一部の地域で雨に注意する必要があり今日は昨日よりさむいのでさむさにも注意する必要があり花粉症も明日は少ないとなっているけれどまだ終息しないようで黄砂についても気をつけなければならないようでインフルエンザもまだかなりはやっていて個人的な経験で申し訳がないのだけどこの時期も風邪を引きやすく油断がならないからどうかさむさによる悪影響や花粉症による悪影響や雨や急な強い雨による悪影響やインフルエンザや風邪による悪影響が絶対に絶対に絶対に絶対に絶対に絶対に絶対に防がれ絶対にあなたが暖かい健康でゆとりのある毎日を過ごしてほしいとほんとうにほんとうにほんとうにほんとうにほんとうにほんとうにほんとうに強くおもう。永久に絶対にどんな場合も現実としても可能性としても当為としてもあらゆるすべてのことについてあなたの気持ちに殉じてあなたのためになるように自分は言葉だけじゃなくて実証されているようにどんな犠牲を払っても自分の命を犠牲にしても自分のみがすべての責



また、少納言・中納言による検索における“文境界無視”とでも呼ぶべき問題がある。これは特にウェブデータの場合に限ったことではないが、ブログ記事には句点を使わないものが多く、そのような場合、文の境界がないかのように扱われる。

サンプルID : OY11\_06057

今日も晴れたので散歩に松本城のお堀に二羽の白鳥が仲良く泳いでいたお城もなんか寒そうお城の堀越しに北アルプスの常念岳が白く映えていた緑町まで歩き上土の「花月」ホテルでカレーライスを食べた帰りは縄手通りを通過して帰宅 少し欲張って歩いたので左脚の脹脛が痛くなった午後からは曇りになり寒くなった

この記事は実際のブログでは次のように表示される (<http://blogs.yahoo.co.jp/sinnshuu28/archive/2008/12/8>)。

今日も晴れたので散歩に [見出し]

松本城のお堀に二羽の白鳥が仲良く泳いでいた

[白鳥の写真]

お城もなんか寒そう

[松本城の写真]

お城の堀越しに北アルプスの常念岳が白く映えていた

[常念岳の写真]

緑町まで歩き上土の「花月」ホテルでカレーライスを食べた

[カレーの写真]

帰りは縄手通りを通過して帰宅 少し欲張って歩いたので左脚の脹脛が痛くなった

午後からは曇りになり寒くなった

この例の場合、中納言の短単位検索ではサンプル全体が長い1文として扱われ、「泳いでいたお城」や「食べた帰りと」といった短単位の連続の“用例”があるものとして処理される。中納言の長単位検索や少納言による検索の場合も同様である。

これはBCCWJにおけるテキスト構造に関わる情報の付与の仕方に関係している。一般論として、行が句点以外の文字で終わる場合、文がそのまま次の行に続く可能性と、文が句点なしでそこで終わる可能性とがある。発表者の断片的な確認によればBCCWJの情報付与では両者の可能性が区別されておらず、その結果として“文境界無視”の現象が生じるものと見られる。

ほかにも細かい問題はあるが、日本語研究上特に問題となり得ると思われるYahoo!ブログサブコーパスの特性としてこれまでに気付いたのは以上である。

## 5 Yahoo!知恵袋の特性

### 5.1 サンプルの完全一致

Yahoo!知恵袋サブコーパスにはYahoo!ブログサブコーパスに見られるほどのデータの重複はないが、サンプルの完全一致が1例だけあった。

サンプルID : OC08\_00706=OC08\_05640

国勢調査って何のためにするの? 国勢調査の人口は、議員定数や地方交付税算定の基準など、法定人口として利用されます。また、男女・年齢別人口・産業別帯・高齢者のいる世帯などの統計は、国や市町村の社会福祉雇用政策、環境設備政策、号際対策などの行政資料として利用されます。

Yahoo!知恵袋は、利用者どうしの知識の交換を目的とする問答形式の掲示板である。中納言では質問と回答の境界が考慮されず、問答が上記のような一続きのテキストとして扱われるが、上のサンプルの場合、冒頭の「国勢調査って何のためにするの?」が質問で、残



りの部分が回答である。また、Yahoo!知恵袋のサイトで1つの質問に対して複数の回答があった場合は、BCCWJにはそのうち「ベストアンサー」とされたものだけが回答として収録されている(丸山・柏野・田中(2011))。

それにしても、問答の完全一致がなぜ生じたのか。上記のようにそれなりの長さを持つ問答が一字一句変わらず2度繰り返されることは常識的に考えがたい。

発表者の推測をあえて通俗風に表現すれば、これはYahoo!知恵袋を運営するヤフー株式会社が“自作自演”の問答を半ば不手際で2度掲載してしまったものである。

そのことを理解するには、BCCWJのYahoo!知恵袋サブコーパスのデータが何物であるかを知る必要がある。BCCWJのマニュアルには正確な記載がないが、BCCWJに収録されているのは、実は“Yahoo!知恵袋”のデータではなく、その準備段階の“Yahoo!知恵袋ベータ版”のデータである。

Yahoo!知恵袋は2005年11月7日に正式運用を開始したが、その約1年半前の2004年4月7日にはその試験段階としてYahoo!知恵袋ベータ版の運用が始められた。<sup>2</sup> Yahoo!知恵袋ベータ版では必ずしも実際の利用者が問答のやり取りをしたわけではなく、一部の質問や回答はヤフー株式会社によって用意されたのであった。<sup>3</sup>

現行のYahoo!知恵袋のサイトにはYahoo!知恵袋ベータ版の時期の問答も併せて掲載されており、当該の2サンプルは今も次の異なる問答として参照することができる。

[http://detail.chiebukuro.yahoo.co.jp/qa/question\\_detail/q116243661](http://detail.chiebukuro.yahoo.co.jp/qa/question_detail/q116243661)

(質問日時：2005/9/23 10:27:12、解決日時：2005/9/23 10:35:34、回答数：1)

[http://detail.chiebukuro.yahoo.co.jp/qa/question\\_detail/q136197063](http://detail.chiebukuro.yahoo.co.jp/qa/question_detail/q136197063)

(質問日時：2005/9/26 13:42:59、解決日時：2005/9/27 09:38:14、回答数：7)

Yahoo!知恵袋の正式運用開始を間近に控えた2005年9月にあつて、1度目の掲載は問答の模範例を示すために行われ、その3日後の2度目の掲載は複数回答からのベストアンサーの選び出しの実験あるいは例示の目的で行われた——ただし、同一の問答を使ったために重複が生じてしまった——といったところかと推測される。

以上、本発表の目的の中心を外れるが、“Yahoo!知恵袋”サブコーパスに収められたデータが正確には“Yahoo!知恵袋ベータ版”のデータであることを明らかにする目的も兼ねて、サンプルの完全一致の生じた背景に関する推測を述べた。

## 5.2 サンプルの部分一致

サンプルの部分一致についても、Yahoo!知恵袋サブコーパスにはYahoo!ブログサブコーパスにおけるような極端なデータの重複はない。ただし、相手の発言をコピーして引用するインターネット掲示板の慣習がデータの重複を多数生じている。

サンプルID：OC02\_07077

DVDをプレーヤーに入れたところ、全く動きません。以前正常に作動したプレーヤーもどれも動かなくなりました。原因は何でしょうか。ちなみに6月中旬にNortonを更新しています。それ以降に使えなくなりました。>正常に作動したプレーヤーもどれもどれもって、何台あるのですか？ウィルスではないですか。スキャンして下さい。

<sup>2</sup> それぞれの日付は <http://chiebukuro.yahoo.co.jp/docs/whats2004.html>、<http://chiebukuro.yahoo.co.jp/docs/whats2005.html> による。

<sup>3</sup> より詳しくは [http://detail.chiebukuro.yahoo.co.jp/qa/question\\_detail/q1011740930](http://detail.chiebukuro.yahoo.co.jp/qa/question_detail/q1011740930) などを参照。

また、この問答はYahoo!知恵袋のサイトでは次のように表示されるものであるが ([http://detail.chiebukuro.yahoo.co.jp/qa/question\\_detail/q135001663](http://detail.chiebukuro.yahoo.co.jp/qa/question_detail/q135001663))、

質問：

DVDをプレーヤーに入れたところ、全く動きません。  
以前正常に作動したプレーヤーもどれも動かなくなりました。  
原因は何でしょうか。  
ちなみに6月中旬にNortonを更新しています。  
それ以降に使えなくなりました。

回答：

>正常に作動したプレーヤーもどれも  
どれもって、何台あるのですか？  
ウイルスではないですか。スキャンして下さい。

4.3で見た“文境界無視”の事例と共通の理由により、中納言による検索では「どれ-も-どれ」という短単位連続の“用例”があるものとして扱われる。

## 6 おわりに——BCCWJ使用上の教訓

BCCWJの検索結果に基づくN-gram分析ツールの作成が契機となって明らかになった、BCCWJのウェブデータに含まれるデータ重複の問題ほかについて粗い調査・分析を行ってみた。

Yahoo!ブログサブコーパスにおけるデータの重複は、用例の頻度に着目する研究にしばしば破壊的な影響を与える。Yahoo!ブログサブコーパスのサンプル52,680件に完全一致の相手を持つサンプルが410件含まれるというのは一見小さな比率のようでもあるが、上で見た「殉じて」に関わる事例などが示す通り、全体の用例数が少ないときにデータ重複による“用例”が多数あれば、分析は致命的にゆがんだものになってしまう。そして、サンプルの部分一致は完全一致よりはるかに多いので、問題は410件のサンプルにとどまらない。Yahoo!知恵袋サブコーパスも相対的に軽度ながら同様の問題をはらむ。機械生成や単なる複製などによる表現の“用例”をそれと知らずに通常の用例と同列に扱ってしまうことのないよう十分な注意が必要である。

今回の調査・分析から得られたBCCWJ使用上の教訓を最後に一般的な形にまとめておく。均衡性の考慮から明確な基準に基づいて収集された出版物（特に書籍）のデータに、それとは採録の手順が異なるだけでなく言語的に異質でデータの重複や信頼性の問題も大きいYahoo!ブログ、Yahoo!知恵袋のデータを単純に加えて使うのは、“均衡”コーパスのまっとうな用法ではない。中納言では検索対象——少納言ではメディア/ジャンル——を目的に応じて適切に指定したうえで検索しなければならぬ。そして付言すれば、話しことばの書き起こしであり、やはり異質性の高い国会会議録のデータも、“書き言葉”コーパスの他の部分との不用意な併用は避けるべきであろう。

## 文献

- 田野村忠温(2010)「日本語コーパスとコロケーション——辞書記述への応用の可能性——」『言語研究』138, pp.1-23.
- 丸山岳彦・柏野和佳子・田中牧郎(2011)「第3章 サンプルング」『「現代日本語書き言葉均衡コーパス」利用の手引』第1.0版 (BCCWJ-DVD版所収のPDF文書), pp.21-38. 国立国語研究所コーパス開発センター.