

Web の情報検索

- ディレクトリサーチサイトの分類法 -

清水由美子

Web 上の情報検索に際し、初心者がよく利用するといわれる階層構造を持った検索サイトを対象に、その分類法を検討する。3つのディレクトリサーチサイトを、従来の情報分類法の一つである日本十進分類法を基準に比較した。最上階層の類似度、サブカテゴリ数とその分布、同一カテゴリ名の重複使用、分類の観点の交叉などを調査した結果、いずれのサイトも基本的に従来の階層構造を踏襲していたが、中には意図的に新しい分類法(交叉分類、カテゴリの複数箇所への位置付け等)を取り入れているサイトもあった。

キーワード：Web ページの検索，ディレクトリサーチサイト，分類法

1 はじめに

ここ数年日本でもインターネットの使用者が激増している。その中でも Web 上の情報を有効に利用したいという要望は強い。しかし Web 上の情報はあまりにも膨大であり、自分の求める情報に的確に行き着くのはそう簡単ではない。また、得た情報が Web 上で得られる全てであったのかどうかという不安が常に付きまとう。情報利用に対する使用者の満足度は必ずしも高くないという調査結果もある(「Web 上の情報に関する利用者の評価は厳しく、全体の 80.5%が「情報を探すのに苦労している」と評価している[1])。初心者にとっての情報検索の難しさについて、木村は学生を対象とした実験結果から、画面には「メディアの文法」があり、それが利用経験の少ない学生にとって障壁となっていると述べる[2]。

情報検索の際、初心者が多く利用するといわれるディレクトリサーチ(Web ディレクトリを使った検索)をうまく利用するためには、「メディアの文法」を知ること、すなわち、それらディレクトリサーチサイトの構造を知ること必要であろう。

本稿では Web ページの検索方法について、その種類と特徴を概観し、続いて3つのディレクトリサーチサイトの分類法を検討する。

2 Web ページの検索方法

Web ページの検索方法には大きく分けて次の三種類がある。

2.1 Web ディレクトリ

Web ページの情報をその内容に従ってカテゴリ分けし、カテゴリを辿ることで検索対象に行き着くように工夫された検索方法。基本的に階層構造(木構造・tree 構造)

の分類システムが採用されている。カテゴリはいずれも人が作ったものであり、掲載サイトの選択・分類も人手で行われている。各カテゴリが階層の中に位置付けられているため、そのカテゴリ名がどのような文脈の中で使われているかが容易に分かる。使用者の直感に従った検索が可能であるため、初心者には使いやすい検索方法である(1994年にサービスを開始したYahoo! [3]が代表格)。

2.2 検索エンジン

検索エンジン(ロボット)と呼ばれる検索用ソフトが、Web 情報を蓄積しているサイトを自動的に順次訪問し、そのサイトの各 Web ページから自動収集したキーワード情報を自動的にデータベース化し、使用者はそうにして自動生成されたデータベースを対象に検索を行う方法。キーワード検索と全文検索がある。キーワード検索は、検索対象となる情報のそれぞれにキーワードを付け、そのキーワードに対して検索を行うものである。ただし、キーワード付けはhtml のタグから判断した特定部分について行われる。全文検索は、対象となる Web ページの文章すべてに対して検索をかけるものである。ただし、大規模な文書群に対しては、検索実行時に全テキストを走査するのは現実的ではなく、テキストに出現するすべての単語に対して転置インデックスを作成し、これを利用して全テキスト走査と同等のことを実現する方法が取られる[4]。

検索エンジンによる検索では、ほとんどのサーチサイトで検索語に関連が高いと推定される順に結果が並べられる。時に何千、何万単位のサイトを拾い出してしまうため、そのうちのどれが本当に自分の欲しい情報であるのかの見極めが難しい。的確な情報を得るためには論理演算子(and, or, not)を利用して、情報の範囲を絞り込む必要があるが、キーワードを入力すること自体がバリアとなったり、演算子の適切な使用法が分からなかったりで、検索エンジンによる検索は初心者には敬遠されがちであるという[5](日本のサイトでは goo[6]が代表格)。

2.3 メタ検索エンジン

使用者の検索要求に対し、複数の検索エンジンを利用し、その結果（どの検索エンジンを使えばよいか）を返すシステム。自分自身がインデックスを持っているわけではなく、メタな知識を持った検索システムである。

3 ディレクトリサーチサイトの分類法

3.1 調査対象とそれぞれの特徴

日本語のディレクトリサーチサイトのうち Yahoo! JAPAN[7]、infoseek[8]、CSJ の厳選サイト集である CSJ What's Best![9]を対象とする。このほか、日本十進分類法[10]を分類法や分類カテゴリ検討の参考に用いる。

調査の対象とする階層と各カテゴリ数は次の通りである（なお、カテゴリは日々変化している。ここに挙げるカテゴリ数は2000年2月中の調査結果である）。

Yahoo! JAPAN：第1階層 14 カテゴリ、第2階層 431 カテゴリ、第3階層 4436 カテゴリ

infoseek：第1階層 18 カテゴリ、第2階層 271 カテゴリ、第3階層 1692 カテゴリ

CSJ What's Best!：第1階層 13 カテゴリ、第2階層 141 カテゴリ（第2階層から直接実際の Web ページにリンクされる）

日本十進分類法：第1次区分 10、第2次区分 100、第3次区分 941（1000 用意されている分類番号のうち、カテゴリ名が入らないものが59ある）

続いて、3つのディレクトリサーチサイトと日本十進分類法の特徴について概観する。

Yahoo! JAPAN：日本語での情報提供を目的とした世界中にいくつかある Yahoo! の一つ。第1階層はほぼ Yahoo! に対応するが、第2階層からは項目数などに違いが見られる。第2階層から直接実際の Web ページにリンクされているものもあれば（例：芸術と人文>外国のアート>実際のページ）、6階層以上の深い階層構造を持つもの（例：自然科学と技術>生物学>植物学>植物>樹木>林学、森林科学）もある。（「>」で階層構造を表すこととする。）

infoseek：トップに18のチャンネル（第1階層）を設け、そのそれぞれがチャンネルリンク集（第2階層）を持つ。さらにその一つ一つのカテゴリがリンク集を持つ（第3階層）という構造。深さは第3階層までのものが主になっている。

CSJ What's Best!：厳選サイト集。第2階層から実際の Web ページにリンクされている。約2000件のサイトを掲載している。

日本十進分類法：学校図書館、公共図書館等で広く採用されている階層構造を持った分類法。第1次区分であらゆる文献・資料の扱う情報・知識を9区分し、それらに含まれないものと、二つ以上の区分を含むものを「総記」

的な区分として設定する。第2次区分は、第1次区分のなかの一つの区分枝を、さらに9区分し、同じく9区分に含まれないものと、二つ以上の区分にわたるものを10番目の区分とする。以下、同様の考え方で、第3次区分、第4次区分、... 第n次区分と分けていく[11]。

日本十進分類法は次のような分類原則を持っている[10]。

- 1) 学術の体系に準拠
- 2) 分類原則の遵守
 - (1) 一貫性：一度に区別する原理は常に一つとし、無秩序な交叉分類を生じさせない、(2) 相互排他性：区分された部分集合（種）はお互いに領域を侵さない排他的なものとする、(3) 包括性：種の総和が被区分体（類）の全域を包含する、(4) 段階性：区分は段階的でなければならず、飛躍を認めない、の4原則に従う。

3) 分類項目名の明示

分類体系の各段階に列挙される分類項目名は、他の分類項目名と的確に区分される用語で明示されなければならない。単に用語が異なるだけでなく、その意味が明確で、紛らわしくないことも重要である。

4) 過去への対応力を持つ

学者が提示する学術の分類は、現在の時点で最新の知識の状態だけに着目しているが、図書の分類では、過去の文化遺産である図書も分類する責務があり、現在では廃れてしまった文化の記録も収まる体系でなければならない。

5) 将来の展望に備えて、新主題を収容できる余裕を持つ

図書の分類法は策定してから容易に変更できない固定性を持っており、それだけに将来の展望に備える余裕を確保すべきである。

3.2 最上階層の類似度

3つのディレクトリサーチサイト及び日本十進分類法の最上階層を対象に、四者の距離を見る。各サイト及び日本十進分類法の最上階層のカテゴリは次の通りである。Yahoo! JAPAN：芸術と人文/ビジネスと経済/コンピュータとインターネット/教育/エンターテインメント/政治/健康と医学/メディアとニュース/趣味とスポーツ/各種資料と情報源/地域情報/自然科学と技術/社会科学/生活と文化

infoseek：ニュース&天気/マネー/ビジネス/キャリア/学び/コンピュータ/インターネット/カルチャー&ホビー/政治・法律・社会/クルマ/トラベル/スポーツ/エンターテインメント/ヘルス&ビューティー/住まい/グルメ&クッキング/くらしの情報/ショッピング

CSJ What's Best：ニュース/音楽・芸術/旅行・地域/ス

スポーツ/趣味/娯楽・グルメ/生活・医療/車/ショッピング/
 コンピュータ/教育・学術/金融・投資/ビジネス
 日本十進分類法：総記/哲学/歴史/社会科学/自然科学/技
 術/産業/芸術/言語/文学

見た目の類似度を見るためにカテゴリ数，全文字数，
 1カテゴリ当たりの平均文字数，文字種（漢字・平仮名・
 片仮名・記号の延べ数と，全文字数に占める割合）を調
 べる（表1）。続いて，カテゴリ名の類似度を考察する（表
 2 もとになるサイトを縦に，対象とするサイトを横に
 とる）。カテゴリ名の類似は，他サイト等が全く同一のカ
 テゴリ名を持っていた場合を1，一部が同じまたは類似
 のカテゴリ名を持っていた場合を0.5としてカウントす
 る。なお，「エンターテインメント」と「エンタテインメ
 ント」，「クルマ」と「車」，「トラベル」と「旅」は同一
 カテゴリと扱う（例：Yahoo! JAPANの「社会科学」に対
 し，日本十進分類法「社会科学」は1。infoseekの「ニ
 ュース&天気」に対し，CSJ What's Bestの「ニュース」
 は0.5）。例えば，Yahoo! JAPANのカテゴリ名の自らのカ
 テゴリ名との類似度は14であり，Infoseekのカテゴリ名
 との類似度は6である。ここで，各サイトのカテゴリ数
 が違うため，比較のために，自らのカテゴリ名との類似
 度を10とした場合の数値を（ ）内に示した。これによ

ると，Yahoo! JAPANから見た日本十進分類法のカテゴリ
 名との類似度は1.8であるが，日本十進分類法から見た
 Yahoo! JAPANのカテゴリ名との類似度は2.5である。

見た目にはYahoo! JAPANとCSJ What's Bestとが近い
 が，カテゴリ名ではCSJ What's Bestに対するinfoseek
 が最も近い。

3.3 サブカテゴリ数とその分布

上位階層の各カテゴリに対して下位階層のいくつのカ
 テゴリが対応しているか（サブカテゴリ数がいくつ）か
 を見る。Yahoo! JAPANの第1階層1カテゴリに対するサ
 ブカテゴリ数は平均30.8（最少5カテゴリ，最多50カテ
 ゴリ）であり，第2階層1カテゴリに対するサブカテ
 ゴリ数は平均10.3（最少0カテゴリ，最多156カテゴリ）
 である。以下，各サイトの第1～第3階層までのカテ
 ゴリ数，上位階層1カテゴリに対するサブカテゴリ数，サ
 ブカテゴリ数の分布（最少から最多まで）を表3にまと
 める。なお，CSJ What's Best!は第2階層から実際のペ
 ージへリンクされているため，第2階層1カテゴリに対
 するサブカテゴリ数等は求められない。

表3に見るように，Yahoo! JAPANには第2階層から直
 接実際のWebページに飛ぶカテゴリ（サブカテゴリ数0）

表1 見た目の類似度（%は小数点第2位を四捨五入）

	カテゴリ数	全文字数	平均文字数	漢字 (%)	平仮名 (%)	片仮名 (%)	記号 (%)
Yahoo! JAPAN	14	86	6.1	41 (47.7)	9 (10.5)	36 (41.9)	0 (0)
Infoseek	18	94	5.2	10 (10.6)	7 (7.4)	71 (75.5)	6 (6.4)
CSJ What's Best	13	55	4.2	25 (45.5)	0 (0)	24 (43.6)	6 (10.9)
日本十進分類法	10	24	2.4	24 (100)	0 (0)	0 (0)	0 (0)

表2 カテゴリ名の類似度

	Yahoo! JAPAN	infoseek	CSJ What's Best	日本十進分類法
Yahoo! JAPAN	14(10)	6(4.3)	5(3.6)	2.5(1.8)
Infoseek	6(3.3)	18(10)	9(5)	0(0)
CSJ What's Best	5(3.8)	9(6.9)	13(10)	0(0)
日本十進分類法	2.5(2.5)	0(0)	0(0)	10(10)

表3 各階層のカテゴリ数，及び上位階層1カテゴリに対するサブカテゴリ数とその分布

	第1階層 カテゴリ数	第1対第2 サブカテゴリ数 (平均)	サブカテゴリ数 分布	第2階層 カテゴリ数	第2対第3 サブカテゴリ数 (平均)	サブカテゴリ数 分布	第3階層 カテゴリ数
Yahoo! JAPAN	14	30.8	5~50	431	10.3	0~156	4436
Infoseek	18	15.1	5~25	271	6.2	0~47	1692
CSJ What's Best	13	10.8	5~20	141			
日本十進分類法	10	10	10	100	9.41	1~10	941

もあれば、第2階層1カテゴリ当たり156というサブカテゴリを持つものもあり、サブカテゴリ数のばらつきが大きい。100以上の大きなサブカテゴリ数は次のような分類法によるものである。「地域情報(第1階層)>世界の国と地域(第2階層)」の下位階層には「アイスランド」「アイルランド」「アフガニスタン」といった国の名前が並び、これらの国名を合わせると140カテゴリとなる。infoseekにも「トラベル(第1階層)>世界の地域情報(第2階層)」のような類似カテゴリがあるが、このサブカテゴリは11である。こちらの分類は「ヨーロッパ」「北欧」「中東」のように、かなり大きなレベルで行われている。また、Yahoo! JAPANの「健康と医学(第1階層)>病気、症状(第2階層)」は129のサブカテゴリを持つ。ここには「インフルエンザ」「肝炎」「咽頭癌」といった多くの病名が並べられると共に「肝臓病」「癌」のような、その上位レベルのカテゴリ名も同列に置かれている。

3.4 同一カテゴリ名の重複使用

一つのカテゴリ名が同一階層内で、あるいは別階層で複数回使用されることがある。CSJ What's Best!の第1階層カテゴリ名「コンピュータ」は、同じく第1階層「ショッピング」のサブカテゴリ名としても用いられる。また、infoseekでも「キャリア>資格取得」「学び>資格取得」のように同一カテゴリ名が別の上位階層のサブカテゴリ名として使われている。日本十進分類法の第1区分カテゴリ名も第2区分以下で用いられる(例 第1区分「1哲学」、第2区分「10哲学」、第3区分「100哲学」)。ただし、これらはいずれもカテゴリ名が同じであるからといって同一のサブカテゴリを持つわけではない。What's Best!の第1階層カテゴリ「コンピュータ」のサブカテゴリは「全般/マック/PCWindows/UNIX/ハードウェア/ソフトウェア/コンピュータ雑誌/周辺機器」であり、「ショッピング」のサブカテゴリ「コンピュータ」は実際のWebページにリンクされている。

カテゴリ名の重複使用の度合いをそれぞれのサイトで見た。(表4 例: Yahoo! JAPANの第1階層カテゴリ名は、同一階層では重複使用されないが、第2階層では延べ18回出現する。Yahoo! JAPAN第2階層カテゴリ名は、第1階層で延べ18回、同一階層内では延べ561回の重複使用がある。)

表4 同一カテゴリ名(出現階層と延べ数)

	Yahoo! JAPAN		infoseek		CSJ What's Best!		日本十進分類法	
	1	2	1	2	1	2	1	2
1	0	18	0	5	0	2	0	8
2	18	561	5	8	2	10	8	0

ここで特徴的なのが、Yahoo! JAPANにおける同一カテゴリ名重複使用の多さである。これはYahoo! JAPANが次のような階層構造のルールを持っているためである[12]。

あるカテゴリやサイトを配置する際、一つのカテゴリに置くだけでは不十分な場合があります。「スポーツニュース」のあるべき場所は「趣味とスポーツ:スポーツ:メディアとニュース」でしょうか、それとも「メディアとニュース:スポーツ」でしょうか。どちらも可能性があり、どちらか片方だけでは不十分です。このような場合、Yahoo! JAPANでは「カテゴリリンク」という方法を用いて、一つのカテゴリを複数のカテゴリに表示できるようにしています。例えば、「エンターテインメント>ゲーム」のサブカテゴリと、「趣味とスポーツ>ゲーム」のサブカテゴリは、共に「アーケード/囲碁/イベント/インターネットゲーム/...」である。これは、一つのカテゴリにアプローチするのに複数の方法がある、または一つのカテゴリを階層構造の複数箇所に置くということであり、他のサイト及び日本十進分類法には見られない分類法である。Yahoo! JAPANでは、積極的にこのような分類法を用いているため、同一カテゴリ名の重複使用が多くなっている。

3.5 分類の観点

Yahoo! JAPANには、観点の異なる分類が同一ページに掲載されている場合がある。例えば「芸術と人文>人文>文学>販売」のサブカテゴリには「イベント/書店/タイトル/取次/本の検索サービス」等が並ぶ。また同一ページにライン一本で分けられて、「SF,ファンタジー,ホラー/アダルト/エンターテインメント/音楽/技術/教育/ミステリ/料理」等、本のジャンル別カテゴリが並んでいる。このように、一つの事柄についての異なった観点からの分類を同時に掲載しているページがいくつか見られる。

4 おわりに

4.1 まとめ

今回調査した3つのディレクトリサーチサイトは、階層構造を持つという点で従来の分類法の一つ、日本十進分類法を踏襲している。infoseek, CSJ What's Best!の

第1階層はそのカテゴリ名からかなり斬新な印象を受けるものの、基本的に分類法においては日本十進分類法を大きく逸脱していない。

しかし、日本十進分類法のカテゴリ名に最も近い第1階層カテゴリ名を持つYahoo! JAPANは、「1度に区別する原理を常に一つとする」[10]といったその分類原則に反して交叉分類を採用したり、一つのカテゴリを複数箇所に位置付けるなど、カテゴリ分類法で意図的に新しい方法を取り入れている。

4.2 今後の課題

今回の調査では、それぞれのサイトがカバーするサイト数の規模が異なるため、必ずしも正確な比較ができとは言えない。今後、同程度の規模のディレクトリサーチサイトを対象とした調査を行い、より正確な対比を試みたい。

また、各サイトについて、上位階層と下位階層との意味関係の分類を行う等、更に詳しい分析によりそれぞれのサイトの分類法の特徴をいっそう明確にしていきたい。

Yahoo! JAPANでは、従来になかった新しい分類法を取り入れている。こうした新しい分類法が、本当に初心者にとって検索しやすいものであるかどうかの検証も今後の課題である。例えば、いずれのディレクトリサーチサイトも、自分が今居る位置を様々な仕方で表示している。Yahoo! JAPANではページ上部左に「ホーム>健康と医学>女性の健康」のような形で現在位置を示しているが、一つのカテゴリを複数箇所に位置付けた場合、次のような事が起こる。「ホーム>エンターテインメント>芸能人、タレント>アーティスト」と辿ると、次の段階でのページの位置表示が「ホーム>エンターテインメント>音楽>アーティスト」となる。自分が辿った通りの道筋が表示されないことから、使用者が戸惑いを感じる可能性はないか等、実際の検索行動の観察なども併せて行いたい。

参考文献

- [1] 白澤基紀・新垣紀子・野島久雄・石崎雅人：“WWW 検索行動における「戻る」行動と検索方針の変化との関係”，情報処理学会ヒューマンインタフェース研究会発表資料，99-HI-83，pp.61-66，1999
- [2] 木村忠正：“「次の10件」に気付かない学生たち - メディアの文法とネットワークへの受動的意識構造 - ”，日本語学九月臨時増刊号，第十七巻第11号，明治書院，pp.188-203，1998
- [3] Yahoo! ，<http://www.yahoo.com/>
- [4] 長尾真・黒橋禎夫・佐藤理史・池原悟・中野洋：岩波講座言語の科学9 言語情報処理，岩波書店，1998
- [5] 新垣紀子・野島久雄：“電子メディア社会における人の情報検索プロセス”，CMCC 研究会第1回シンポジウム論文集，pp.3-12，1999
- [6] goo ，<http://www.goo.ne.jp/>
- [7] Yahoo! JAPAN ，<http://www.yahoo.co.jp/>
- [8] infoseek ，<http://japan.infoseek.com/>
- [9] CSJ What's Best! ，
<http://www.csj.co.jp/whatsbest/>
- [11] 柴田正美：資料組織概説，日本図書館協会，1998
- [10] もり・きよし原編 社団法人日本図書館協会：日本十進分類法新訂9 版本表編，社団法人日本図書館協会，1995
- [12] Yahoo! HOW-TO? ，
<http://www.yahoo.co.jp/docs/howto/chapters/8/6.htm>