

インターネットのサーチエンジンの評価尺度：
ESL(Expected Search Length)を使った検索実験

A Measure for Evaluating Search Engines on the World Wide Web:
Retrieval Test with ESL (Expected Search Length)

安形 輝, 野末道子, 服部紀彦, 上田修一
Teru Agata, Michiko Nozue, Norihiko Hattori, Shuichi Ueda

Résumé

Search engines are a kind of information retrieval systems on the Internet. However, they cannot be evaluated with the measures commonly applied in traditional retrieval experiments, i.e., recall and precision, since recall cannot be calculated in the essentially unrestricted resources of the Internet. The Expected Search Length (ESL) proposed by Cooper in 1968 is a measure for evaluating information retrieval systems and has the possibility of taking the place of recall and precision. The ESL calculates the cost paid by a user, i.e., number of retrieved sites the user must look through before (s)he gets sufficient number of relevant sites. This article proposes a version of the ESL in order to adapt it to evaluate search engines. An experiment was carried out to compare the validity of the ESL with that of precision and recall. Eight search engines of the domestic and overseas were evaluated by precision and ESL. The result of this retrieval experiment shows the adaptability of ESL for the evaluation of the search engines.

- I. 情報検索システムとしてのサーチエンジンの特色
- II. 情報検索システムの評価と評価尺度
 - A. 従来の評価研究
 - B. CooperによるESL

安形 輝：慶應義塾大学大学院文学研究科図書館・情報学専攻，東京都港区三田 2-15-45

Teru Agata: Graduate School of Library and Information Science, Keio University, 2-15-45, Mita, Minato-ku, Tokyo. itasan@slis.keio.ac.jp

野末道子：鉄道総合技術研究所輸送システム開発推進部，東京都国分寺市光町 2-8-38

Michiko Nozue: 2-8-38, Hikari-cho, Kokubunji-shi, Tokyo. michiko@rtri.or.jp

服部紀彦：慶應義塾大学文学部図書館・情報学科，東京都港区三田 2-15-45

Norihiko Hattori: School of Library and Information Science, Keio University, 2-15-45, Mita, Minato-ku, Tokyo.

上田修一：慶應義塾大学文学部図書館・情報学科教授，東京都港区三田 2-15-45

Shuichi Ueda: Professor, School of Library and Information Science, Keio University, 2-15-45, Mita, Minato-ku, Tokyo. ueda@slis.keio.ac.jp

受付日：1998年1月27日 改訂稿受付日：1998年7月27日 受理日：1998年9月8日

- C. ESL の問題点とその改善
- D. ESL と精度の関係
- III. サーチエンジン評価の問題点
- IV. ESL を用いた評価実験
 - A. 実験環境
 - B. 実験結果
- V. サーチエンジン評価における ESL の有効性
 - A. ESL の有効性に対する検討
 - B. まとめ

I. 情報検索システムとしての サーチエンジンの特色

現在、インターネット上の Web ページの探索には、一般的に「サーチエンジン」と呼ばれるシステムが利用されている。サーチエンジンは、発展途上にあるが、その基本は従来の情報検索システムの延長である。情報検索システムモデルでは、データベースを作成し、データベースを検索用に構造化し、特定の検索手法を用いて、個々の利用者の作成した個別の検索質問に応じて検索し、結果を表示する。

インターネット環境下では、同一の対象に対して異なるサーチエンジンが存在し、競争状態にある。そのため、サーチエンジンが今後発展するためには、これらの異なるサーチエンジンの評価が不可欠であり、既にいくつかの試みがなされている。情報検索システムの一つであるサーチエンジンの評価には、これまでの情報検索システムの評価で用いられてきた手法や尺度を導入することが可能である。

しかしながら、その評価においては、従来の情報検索システムでは考慮されなかった要素を含めなければならない。

第一は、インデックス作成の問題である。これまでの検索実験では、最初に特定のテスト集合を定め、それに対する検索性能を測定してきた。ところがサーチエンジンが扱うのは原則的にインターネット上のデータ全てであり、その意味でテスト集合はインターネット全体と考えられる。多くのサーチエンジンでは、インデックス作成もロ

ボットによって自動化している。従って、これまでは単にテスト集合全体からインデックスを作成するため、その作成手順、収録範囲等は大きな問題ではなかったが、制約の多いインターネット上ではいかに網羅的かつ効率的、さらには高い更新頻度でインデックスを作成できるかが重要である。つまり、サーチエンジンでは、サーバーから Web のページを収集するロボットの機能と性能を含めた評価が必要となる。

第二は、サーチエンジンの利用者がエンドユーザーであるという点である。従来の情報検索システムは、必ずしもエンドユーザーが利用するものではなかった。エンドユーザーのインターネット利用環境は、多くの場合、電話回線でインターネットプロバイダの提供するアクセスポイントに接続する形態をとっている。その場合に、電話回線使用料、インターネット接続料は原則的に接続時間の従量制で決まるため、接続時間を最小としようとする圧力が強く働くと考えられる。サーチエンジンで、検索結果が適切に順位付けがされていれば、検索結果すべてをみる必要がなくなる。結果として、接続時間が短くなるため、検索結果は順位付けがされていることが望ましい。現実には多数のサーチエンジンは、順位付けを行っている。しかし、情報検索における伝統的な評価尺度である精度と再現率は順位付け出力が考慮される以前の尺度であり、適切な評価を行えない恐れがある。

こうしたサーチエンジンの特色を考慮するならば、情報検索システムの伝統的な評価尺度である精度や再現率に代わる別の新しい評価尺度を考案

する必要があろう。以下では、これまでの評価手法と尺度を概観し、それらを検討した後、ESL と呼ばれる評価尺度に改良を加えた尺度を提案し、その有効性を確かめるために、8種のサーチエンジンの評価を試みる。

II. 情報検索システムの評価と評価尺度

A. 従来の評価研究

情報検索システムに対する評価は、伝統的に再現率(recall)と精度(precision)という尺度を使って行われている。これらの尺度は、Aslib によるクランフィールド(Cranfield II)プロジェクトにおける実験に端を発する。この実験では、どの程度のレバント文献が検索されたか、どの程度非レバント文献を検索しなかったかという点から評価され、前者は検索の網羅性を測る“再現率”として、後者は検索の正確さを測る“精度”として提案された。

再現率と精度による評価は、一般的にはクランフィールド型検索実験と呼ばれ、具体的には以下のような手順から構成される。

- 1) テスト集合と検索質問を用意する
- 2) 各検索質問に対して、どの文献がレバントであるか、レバンス判定を行う
- 3) 情報検索システムはテスト集合に対して検索を行う
- 4) 検索結果とレバンス判定の情報を使い、以下の式から再現率と精度を算出し評価を行う

$$\text{再現率} = \frac{\text{検索されたレバントな文献数}}{\text{全てのレバントな文献数}}$$

$$\text{精度} = \frac{\text{検索されたレバントな文献数}}{\text{検索された文献数}}$$

現在までに多くの研究で、クランフィールド型検索実験には多くの問題があることがたびたび指摘されてきた¹⁾²⁾。再現率と精度に焦点を絞り、問題点をまとめると以下ようになる。

- 1) 再現率算出には全レバント文献の情報が判明している必要がある。そのため、現実の情報検索では再現率の算出は非常に困難である

- 2) 再現率と精度が検索されたかどうかを問題にするために、順位付け出力を適切に評価できない
- 3) 1950年代のバッチ形式の実験環境における尺度であり、インタラクティブな情報検索システムへの応用可能性については疑問が残る
- 4) 評価には再現率、精度両方が必要であるが、二つの尺度は反比例の関係にあり、評価結果の解釈を困難にしている

サーチエンジンの評価においては、このような問題点のなかでも1)2)は致命的であり、再現率と精度両方による評価は困難であると言える。実際、後述のように再現率と精度両方からサーチエンジンを評価しているものはない。そこで、従来の情報検索評価研究において再現率と精度に代わるものとして提案はされたが、ほとんど使われてきていない尺度についても検討する。

情報検索評価研究において再現率と精度以外の評価尺度として現在までに様々なものが提案されてきた。それらの評価尺度の中で代表的なものとしては、A. R. Meetham³⁾によるCTI(Contingency Table Information), J. A. Swets⁴⁾によるArea, J. W. Wilbur⁵⁾によるレバンス情報(IdF), H. P. Frei と P. Shauble⁶⁾によるu と u*, W. S. Cooper⁷⁾によるESLがある。また、一部のオンライン検索システムに関する研究においては、利用者の満足度(satisfaction)や有用性(utility)といった基準から評価が行われることもある。これらの評価尺度が再現率・精度に置き換わらなかった理由としては、全ての尺度が何らかの問題を抱えており、少なくとも再現率・精度と同程度の妥当性しかなかったことがある⁸⁾。例えば、Wilburのレバンス情報は、その算出には再現率と同様に集合中の全レバント情報が必要であるという問題を抱えている。

全体的な傾向として、再現率と精度を含めほとんどの評価尺度は、検索結果から得られる情報あるいは便益を測定する点で、出力面を扱っているといえる。一方で、検索結果を吟味するような利用者の労力、時間といった利用者側のコス

トは、固定としたり、あるいは無視している。唯一、ESLのみが、逆に利用者の満足する件数の形で得られる情報量を固定し、利用者側のコストを評価する尺度となっている。

B. Cooper による ESL⁷⁾

Cooper による ESL は、簡単に言えば、利用者の満足する文献数を得られるまで検討しなければならない非レバント文献数である。そのため、利用者の満足する件数が明らかであることを前提としている。ESL の概念を説明するために図 1 を使う。この図は検索システムにおける順位付け出力の結果を表しており、四角形は文献を、その横の数字は順位を示している。検索質問に対するレバント文献は灰色で示している。その検索質問に関して利用者が満足する文献数が 3 件だとするならば、その 3 件を利用者が得るまでに利用者が見なくてはならない非レバント文献は 2 件であるので、ESL は 2 となる。

実際の検索では複数の文献が同順位に出力され

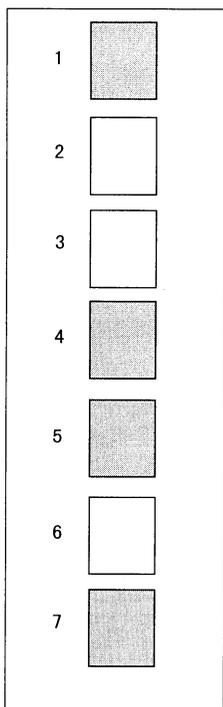


図 1 順位付け出力の例

る場合もあるため、そのことを考慮に入れた ESL は以下のような式 (1) から算出される。

$$esl(q) = \sum_{m=1}^M pr(l_m) \cdot l_m \quad (1)$$

この式は検索要求 q に対する ESL を表している。この式で、 M は同順位の文献が複数あるとき利用者が検討する順序の場合の数である。例えば、満足するレバント文献数が 4 であり、順位 5 までに得られたレバント文献数が 3 であり、順位 6 が 3 文献 A, B, C であったとき、その検討する順序には A-B-C, A-C-B, B-A-C, B-C-A, C-A-B, C-B-A の 6 つの状況が考えられるので、この場合 M は 6 である。 l_m はある状況 m における検討すべき非レバント文献数であり、 $pr(l_m)$ はその状況になる確率である。ここで、要求が満足される文献数に到達する順位までは単に非レバント文献の数を考慮すればいいことになるので、実際の計算は以下のような式 (2) で行われる。

$$esl(q) = j + \frac{i \cdot s}{r + 1} \quad (2)$$

ここで、 j は該当順位前までに出現した非レバント文献数であり、 i は利用者が満足する件数を得られる順位に存在する非レバント文献数、 r は該当順位に存在するレバント文献数である。この場合にレバント文献同士の間非レバント文献が何文献あるかを考えるとその期待値は $i/(r+1)$ となる。ここで s は該当順位の何件目のレバント文献で利用者が満足するかを示している。つまり、 $i \cdot s / (r + 1)$ は該当順位におけるレバント文献が満足する件数に達するまでに検討する必要のある非レバント文献となる。

実際の情報検索システム間の比較を行う場合には、多くの検索式に対する ESL を算出し、その平均を出すために以下の式を使うことになる。

$$\begin{aligned} \overline{esl} &= \frac{1}{N} \sum_{n=1}^N esl(q_n) \\ &= \frac{1}{N} \sum_{n=1}^N \left(j_n + \frac{i_n \cdot s_n}{r_n + 1} \right) \quad (3) \end{aligned}$$

この式で N は用意した検索質問数を表しており、その他の変数は式 (2) と同様である。

さらに、Cooper は、複数の実験環境において

テスト集合、検索式が異なる場合にも比較を行うことができるように ESL をさらに拡張した ESL 減少係数を提案している。しかし、この尺度は再現率と同様に集合中の全レバント文献数を必要とするため、インターネット上のサーチエンジンには応用することができない。実際にサーチエンジンの評価を考えた場合には、それぞれのエンジンについて複数のテスト集合があるわけではなく、インターネット全体は一つの大きな集合と考えられる。従って、ここでは、比較を行うためには ESL 減少係数でなく、ESL を使うことが妥当と考える。

ESL の特徴としては、1) 再現率・精度とは異なり単一の尺度である、2) 利用者の満足する件数を考慮に入れている、3) 利用者のコストに焦点を当てた尺度である、ことが挙げられる。

C. ESL の問題点とその改善

ESL には、利用者の満足する件数があらかじめ分かっている必要があるという欠点がある。現実の情報検索行動においては、事実検索ではこの満足する件数は容易に決定されるものであるが、主題検索のような場合には明確に決まらない場合も存在する。

ESL に関する Cooper の研究は多く引用されているが、この欠点のために、ESL 自体はほとんど使われてこなかった。ここでは、縦軸に ESL を取り、横軸に利用者の満足する件数を取るグラフを考えることで、その欠点を解消することを提案する。例えば、ある仮想的な環境において 2 つの情報検索システム A, B に対する ESL の平均と利用者の満足する件数が図 2 のようになる場合を考える。ESL の値は少なくなるほど性能が高いことになるが、この図ではシステム A, B の ESL は件数によって高低がある。この図からは、利用者の満足する件数が少ない場合にはシステム A の性能が高く、多い場合にはシステム B の性能が高いことが分かる。

このようなグラフを用いて評価を行うことで、ESL の欠点を解消するとともに、利用者の満足する件数の変化に伴った詳細な分析を行うことがで

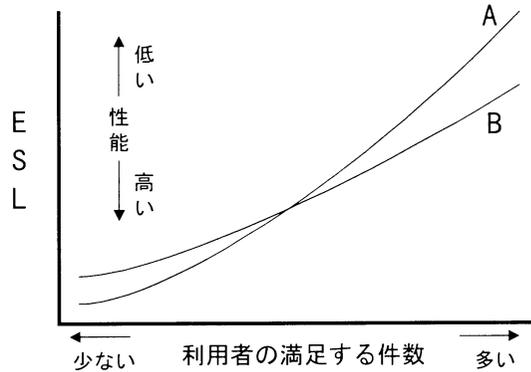


図 2 ESL と利用者の満足する件数

きる。

本研究と同様に Cooper の ESL を取り上げて 1997 年の M. D. Dunlop による研究⁹⁾でも、ESL のグラフ化手法が提案されている。これは ESL に対してシステムごとの検索時間の算出を行い評価することを提案しているが、各システムに応じた複雑な操作を行っているが故に、ESL の利点である汎用性、単純さを失っていると指摘できる。しかし、この研究はコストに焦点を当てた ESL という尺度が近年見直されてきたことを示すものと考えられる。

D. ESL と精度の関係

ESL と従来使われてきた精度は、ESL が検索された文献中にある非レバント文献数、精度が検索された文献中にあるレバント文献の割合を表すという点で、同じものを表現しているとも考えられる。しかし、以下の点で二つの尺度の概念は異なるものである。

- 1) ESL は利用者のコストを表すのに対し、精度は正確さを表している
- 2) ESL は単一の尺度であるのに対して、精度は再現率との関係で用いられる

ESL と精度の違いをより明確に示すために、ある仮想的な情報検索システム A, B による極端な検索結果の評価を考える。システム A, B 共に検索結果 100 件中のレバント文献数が 50 件であるが、システム A は順位 1~50 にレバント文献が集中し、システム B は順位 51~100 にレレ

バント文献が集中しているものとする。精度は100件出力された時点で評価することとし、ESLでは利用者の満足する件数を50件としている。この場合、明らかにシステムAの方が優秀なシステムであるが、精度で見るとシステムA、B共に50%と評価されてしまう。一方、ESLにおいては、システムAは0件、システムBは50件となり、ESLではシステムAの性能が高いと評価される。これはESLが順位付け出力を考慮に入れた尺度であるためである。

III. サーチエンジン評価の問題点

近年、インターネット上の情報資源、なかでもWWW (World Wide Web) による情報は激増しており (<http://www.mit.edu/people/mkgray/net/internet-growth-summary.html>), 将来的にもさらに増加の一途を辿ると考えられる。情報資源の増加は、有用な情報だけでなく、不要な情報の増加も意味している。そのため利用者がインターネット上の情報の中から必要な情報を選択することは次第に困難になりつつある。そこで、必要な情報を重要と思われる順に提示するサーチエンジンの重要性が高まっている。それとともに、サーチエンジンの評価に関する研究もここ数年数多く行われている。Louise T. Su はインターネット上でしか公開されていないような評価研究も含めレビューを行っている¹⁰⁾。

従来のサーチエンジンに対する研究では、多くの場合検索性能ではなく、ある検索語によって検索された件数、インターフェースの使いやすさ、検索機能、応答時間、収録数、更新頻度といった点から評価が行われてきた。例えば、国内のサーチエンジン評価研究において著名な浅井は「検索力」という評価尺度を提案している。これはある決められた検索語で何件検索されるかという非常に単純なものである¹¹⁾。

しかし、サーチエンジンは情報検索システムであり、情報検索システムは利用者にとって必要な情報を取捨選択することを目的とするものである。従って、当然、何らかの形でその検索性能は評価されなければならない。

従来のクランフィールド型検索実験で検索性能は再現率と精度から測定されてきたが、サーチエンジンの場合には、精度は算出できるが再現率は算出できない。インターネット上の資源は膨大であり、再現率を算出するためにすべての情報を吟味することは不可能なためである。実際に、サーチエンジンの評価を行っている研究では、再現率は放棄されているか省略されている^{12) 13)}。その結果として、サーチエンジンに関する評価研究においては精度のみが使われているが、元来再現率と精度は両方の関係から評価を行うものであるため、精度だけによる評価では十分とは言えない。また、サーチエンジンのほとんどは順位付け出力を行うが、前述のように再現率、精度は順位付け出力には対応できない尺度であり、あるシステムが高順位にレバントなWebページを多く出力したとしても、それが適切に評価に反映できないという問題がある。

つまり、従来のサーチエンジンの評価に関する研究では、検索性能以外の尺度から評価を行っている場合がほとんどであり、さらに、検索性能に関する評価が行われている場合でも、精度しか使われていないため、十分な評価が行われていない。

また、インターネット上の検索では実際に利用者が電話などの低速回線を使って接続していることが多い。このため、利用者は今までの情報検索のように検索の正確さや網羅性というよりは、経済的・時間的コストをいかに減らせるかにより重きをおくと考えられる。

本研究ではESLという情報検索領域でほとんど注目されてこなかった尺度を取り上げ、その特徴を生かしながら欠点を改善する評価手法を提案している。この評価手法は、他の評価手法と比べコスト面に焦点をあてている点の特徴であり、よりサーチエンジンの特性にあった評価を行うことができると考えられる。

IV. ESLを用いた評価実験

CooperのESLを使って実際のサーチエンジンを検索し、各サーチエンジンの性能について評

価を行った。その実験の結果から ESL の有効性について検討した。

A. 実験環境

実験対象としたサーチエンジンは、国内 4 種、海外 4 種の合計 8 種である (表 1)。これらは順位付け出力のできること、1997 年秋の時点で代表的なサーチエンジンであること、できるだけ海外版と国内版があることを条件として選択した。さらに、順位付け出力は行えないが、代表的なディレクトリサービスである YAHOO!, YAHOO! JAPAN の二つを参考サービスとして加え、精度を使つての比較評価を行った。実験は 1997 年 10 月 15 日から 11 月 24 日の時期に行った。

検索質問は先行研究¹⁴⁾ で使われた質問に基づき国内外それぞれ 6 件ずつを用意した。ただし、質問の 1 件は国内のサーチエンジンでの検索結果が少なすぎ除いたため、国内 5 件、海外 6 件の合わせて 11 件を使って評価を行った。サーチエンジン検索に使われた 6 件の検索質問は以下のようになっている。

- ・ギリシャ哲学 (Greek philosophy)
- ・遺伝的アルゴリズム (genetic algorithm)
- ・ボランティアの募集 (volunteer)
- ・自己破産 (voluntary bankruptcy)
- ・長野オリンピック (Nagano Olympic)
- ・世界の七不思議 (Seven Wonders of the World) [海外のみ]

各検索質問に対し実際にサーチエンジンを検索する際には、各サーチエンジンで用意されている機能、文法に適した形で検索式を構築した。調査者はインターネットの利用経験が 3 年以上あり、各サーチエンジンに習熟していることを条件として選択した 3 名とした。レlevance判定は調査者

が行った。各質問に対する判定者は 1 名であり、これは判定方針を一貫させるためである。また、現実のレlevance判定では、情報検索過程の進行につれ、判定基準の変化することが知られている。変化しては一貫した判定を行うことができなくなるため、最初の検索式構築のさいに、その検索質問に対する判定基準を記述しておくことで、判定基準が変化することを防いでいる。

レlevance判定は、各サーチエンジンの検索結果に対して、リンクを辿った最初のページの情報から行った。1 ページ単位で判定することとした理由は、サーチエンジンにおける情報単位のとり方と、時間的な制約の二つである。第一に、現在、サーチエンジンにおける情報の単位は物理的な 1 ページである。そのため、システムの意図された動作を評価するならば、評価も 1 ページ単位で行うことになる。第二に、WWW 上の情報、およびサーチエンジン内のインデックスは、ある程度の期間で更新されてしまう。そのため、各質問について複数のサーチエンジンの検索結果を比較する際には、できる限り短期間で行うことが望ましい。ここでは、各質問について 1 サーチエンジンすべての検索結果を 1 日でレlevance判定することとした。各質問で検討した件数は多いもので 200 件以上にも及んだ。回線状況にもよるが、すべての検索結果についてページ中のリンクも辿り、複数ページを見ることは現実的に難しいため、ここでは検索結果の判定を最初の 1 ページのみで行うこととした。

実際の Web ページには、1 ページで内容が完結しているページもあれば、他のページと組み合わせられてある内容を表現しているページも存在する。今後は 1 ページという物理的な単位でなく、内容的な単位から扱うサーチエンジンが登場すると思われる。そのようなサーチエンジンを評価するさいには、Web ページにおける情報の単位をどのように規定するべきかを検討する必要があるだろう。

判定のレベルは、レrelevant、非レrelevantの 2 段階とした。そのページに存在するリンクを辿れば、レrelevantな情報が入手できると予想され

表 1 実験対象サーチエンジン

	対象サーチエンジン	参考サービス
日本	goo, Infonavigator, Infoseek Japan, Excite Japan	YAHOO! JAPAN
海外	Alta Vista, HotBot, Infoseek, Excite	YAHOO!

検索質問：Genetic Algorithm (遺伝的アルゴリズム)

検索式例：AltaVista -- “genetic algorithm”

判定基準：遺伝的アルゴリズムの概要、定義、使い方、実現方法のいずれかについて書いてあるページをレバントとする。遺伝的アルゴリズムという語が単に載っているページは除く。また、説明なしに計算・図があるものは除く。

実験結果：

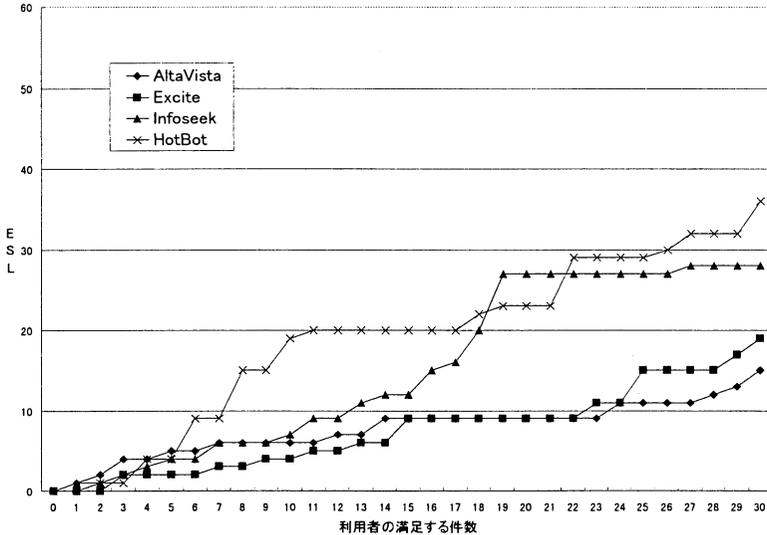


図3 検索質問「遺伝的アルゴリズム」

る場合には、レバントとした。検索結果の出力に対するレバンス判定は、各質問、各システムについて利用者の満足する件数が30件に到達し、さらに精度の算出点である検索件数が50件になるまで行った。検索結果に重複があれば、利用者にとってはその分検討する必要のあるページが増えると考え、上位に既に出力されているページはすべて非レバントとした。各エンジンの検索結果からページまでたどり着けない場合には、インターネット上のトラフィックが混雑しているなどの理由でタイムアウトしてしまう場合と、検索結果中のリンク情報が古くリンク先が存在しない場合とがある。タイムアウトの場合は、しばらく時間をおいて再びアクセスを試みた。リンク先が存在しない場合には非レバントとした。

以上のような形で実験を行った検索質問「遺伝的アルゴリズム」の例を図3に示す。

B. 実験結果

各検索質問に関する実験結果の平均をグラフで表現したものが、国内のサーチエンジンに関する図4と海外のサーチエンジンに関する図5になる。この図では下にあるほど性能が高いことになる。

図4と図5において明らかなのは、海外のサーチエンジンの方が国内のサーチエンジンと比較して、全体的に性能が高いことである。この理由としては、1) 形態素解析、構文解析等において日本語は処理が難しいこと、2) 国内のWebページ全体の質が海外のWebページと比べて低いこと、が考えられる。前者の問題は、特に形態素解析の面で顕著である。これは英語では単語が空白によって区切られているが、日本語は何らかの手法で単語を分ける必要があるために、起こる問題である。あるサーチエンジンでは、検索式を処理するさいに細かすぎる単位まで区切ったことが原因

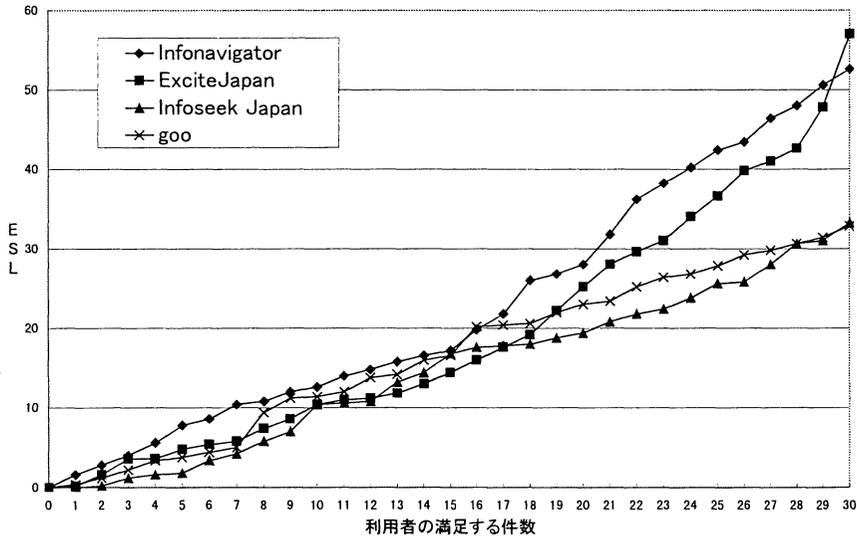


図4 ESL: 国内のサーチエンジン

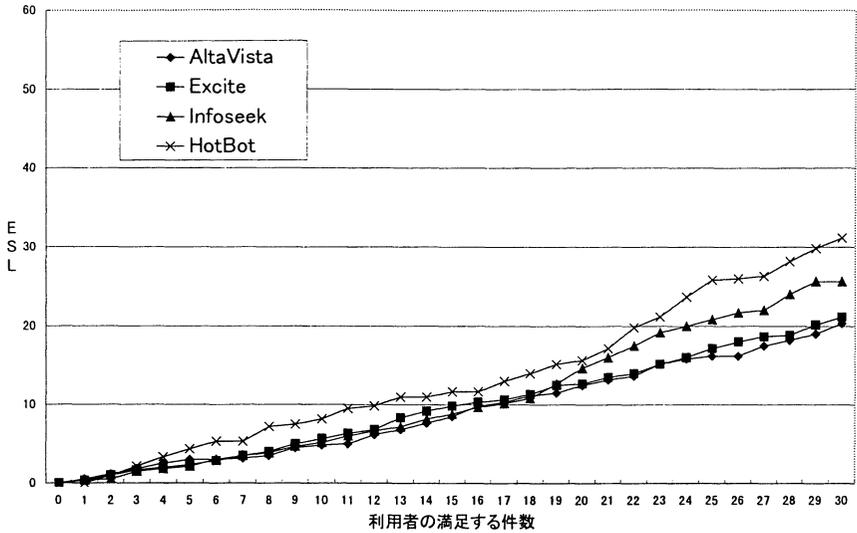


図5 ESL: 海外のサーチエンジン

と考えられるノイズが含まれていた。

今回の実験ではレlevance判定のデータしか収集してないため、後者の問題はあくまで印象である。しかし、日本語ということ国内の利用しか想定していないためか、情報を公開する意識の低いページが国内には多く存在していた。例えば、ページ内の画像へのリンクが切れているといった単純に技術的な問題のあるページや、「ホーム

ページを作ってみました」という程度の内容の乏しいページなどがあつた。

図4から国内のサーチエンジンについて検討すると、必要なレlevance文献数が20件程度までの検索においては各エンジンとも大差ないが、20件以上必要な検索においては Infoseek Japanあるいはgooを利用した方が検討する件数が少なくてすむことがわかる。次に図5から海外の

インターネットのサーチエンジンの評価尺度

表2 50件の時点での精度：国内サーチエンジン

	Infonavigator	Excite Japan	Infoseek Japan	goo	YAHOO! JAPAN
精度 (50 件)	38.8%	45.2%	50.8%	50.0%	44.2%

表3 50件の時点での精度：海外サーチエンジン

	AltaVista	Excite	Infoseek	HotBot	YAHOO!
精度 (50 件)	60.0%	59.0%	55.3%	50.7%	55.6%

サーチエンジンについて検討する。どの時点においても HotBot を使うと他の三つのエンジンと比較して検討すべき件数が多い結果となっている。つまり、海外のサーチエンジンに関しては全体的に HotBot の性能が悪いと考えられる。

各サーチエンジンについて 50 件の時点での精度を示したのが表 2, 3 である。この表で YAHOO! と YAHOO! JAPAN については、検索件数が全体的に少なかったため、検索結果が 50 件以上のものはその時点での精度、50 件に満たない場合には全検索結果に対する精度となっている。

V. サーチエンジン評価における ESL の有効性

A. ESL の有効性に対する検討

今回の実験によって、図 4 のように必要な件数が少ない場合にはどのエンジンを使ってもほぼ同じコストしかかからないが、必要な件数が多い場合には Infoseek Japan と goo は他の二つに比べると性能が高くなるというように、ESL によってサーチエンジンごとの性能の差異を明らかにできた。つまり、ESL によってサーチエンジンの評価を行うことができると考えられる。

ここでは、ESL の有効性をさらに検討する上で、従来の代表的な評価尺度である精度との比較、さらには再現率を含めての比較も考える。

精度に対して、ESL は必要な件数に応じた評価を行うことができる点が優れていると言える。精度は検索件数を横軸に取ることで ESL と同様にグラフを作成することも可能であるが、これには二つの点で問題がある。まず、精度は再現率との

関係で考えられた尺度であり、グラフ化する場合には再現率を横軸に取るべきであるという点である。二つ目としては、横軸に検索件数を取り精度をグラフ化したとしても、それが利用者の役に立つのかということがある。つまり、利用者が情報検索システムの評価をシステム選択の参考とする場合に、検索を行う前には、どの程度の分量の情報が欲しいかはある程度予想できるが、何件検索するかは予想できないため、ESL のグラフと比較すると解釈がより困難であるという点である。

また、再現率をも含めて ESL との比較を考えると、評価の基準が、従来の再現率と精度では検索の網羅性と正確さであるのに対して、ESL では検索にかかるコストである。主に研究のために用いられてきた従来の文献検索では検索の網羅性に対する要求は高かったと考えられる一方で、用途が研究に限られないサーチエンジンについてはインターネット利用においては、網羅性や正確さというよりはむしろコスト面に対する要求が高いと考えられる。この点でも、サーチエンジン評価においては ESL の方がより妥当な尺度として考えることができる。

さらに、再現率と精度では相反する関係にある二つの尺度をペアにして評価を行うため、結果の解釈が時に難しい場合があるのに対して、ESL は単一の尺度であり、その内容もコストという単純なものであるため、解釈が容易である。

B. まとめ

本研究では、Cooper により 1960 年代に提案された尺度である ESL をグラフ化することでサーチエンジン評価に応用し、評価実験を行っ

た。その結果、サーチエンジン評価における ESL の有効性を明らかにすることができた。

ESL の特徴としては、単一の尺度である、単純であり解釈が容易である、コストに焦点を当てた尺度である、順位付け出力に対応している、ことがあげられる。従来の情報検索システムの代表的な評価尺度は再現率と精度であるが、データベースの大規模化、インタラクティブな情報検索システムの登場によって、それらに代わる尺度の模索が情報検索領域の重要な課題となってきた。今回 ESL の有効性を示すことができたことから、ESL の元来の目的であった情報検索システム評価一般に対する応用について、ESL のグラフ化も含め再検討していく必要があると考えられる。

引用文献

- 1) Ellis, David. "The Dilemma of Measurement in Information Retrieval Research". *Journal of the American Society for Information Science*. Vol. 47, No. 1, p. 23-36 (1996)
- 2) Su, Louise T. "Evaluation Measures for Interactive Information Retrieval". *Information Processing and Management*. Vol. 28, No. 4, p. 503-516 (1992)
- 3) Meetham, A. R. "Communication Theory and the Evaluation of Information Retrieval Systems". *Information Storage and Retrieval*. Vol. 5, p. 129-134 (1969)
- 4) Swets, J. A. "Measuring the Accuracy of Diagnostic Systems". *Science*. No. 240, p. 1285-1293 (1988)
- 5) Wilbur, John W. "An Information Measure of Retrieval Performance". *Information Systems*. Vol. 17, No. 4, p. 283-298 (1992)
- 6) Frei, H. P.; Schauble, P. "Determining the Effectiveness of Retrieval Algorithms". *Information Processing and Management*. Vol. 27, No. 2/3, p. 153-164 (1995)
- 7) Cooper, William S. "Expected Search Length: a Single Measure of Retrieval Effectiveness Based on the Weak Ordering Action of Retrieval Systems". *American Documentation*. Vol. 19, No. 1, p. 30-41 (1968)
- 8) Tague-Sutcliffe, J. "The Pragmatics of Information Retrieval Experimentation, Revisited". *Information Processing and Management*. Vol. 28, p. 467-490 (1992)
- 9) Dunlop, M. D. "Time, Relevance and Interaction Modeling for Information Retrieval". *SIGIR 1997 Proceedings*, p. 206-213 (1997)
- 10) Su, Louise T. "Search Engines on the World Wide Web and Information Retrieval from the Internet: a Review and Evaluation". *Online and CDROM Review*. Vol. 21, No. 2 (1997)
- 11) 浅井勇夫. "サーチエンジンからみた Web の世界". *情報の科学と技術*, Vol. 47, No. 9, p. 453-458 (1997) [実際に「検索力」を使ったサーチエンジン評価は検索デスクで公開している。http://www.bekkoame.ne.jp/~asaisan/]
- 12) Leighton, H. V. "Performance of Four World Wide Web Index Services: Infoseek, Lycos, Web Crawler and WWW Worm". <http://www.winona.msus.edu/is-f/library-f/webind.htm>, 1997/9/4
- 13) Chu, H. T.; Rosenthal, M. "Search Engines for the World Wide Web: a Comparative Study and Evaluation Methodology". *Proceedings of the 59th Annual Meeting of the American Society for Information Science*, Vol. 33, Baltimore, 1996-10, ASIS, Baltimore, 1996, 127-135
- 14) Ding, W.; Marchionini, G. "A Comparative Study of Web Search Service Performance". *Proceedings of the 59th Annual Meeting of the American Society for Information Science*, Vol. 33, Baltimore, 1996-10, ASIS, Baltimore, 1996, 126-140