

tf-idf法の改良による Web情報検索結果の表示法に関する研究

能登研究室

阿部剛仁 (16002)

1 はじめに

World Wide Web の情報空間の内容の充実と拡大に伴って、情報を探手段としてのサーチエンジンの重要性はますます高まっている。しかし検索するページ数が膨大であった場合、ユーザが求めている内容のページがページのトップに出現するとは限らないなどの問題点もある。

本研究ではユーザの求めているページを検索結果から排除しユーザにとって見易い検索結果を表示させることを目的とする。

2 サーチエンジン

サーチエンジンには、大きく分けてディレクトリ検索とロボット検索の2種類が存在する。前者は人の手によってホームページの内容を登録し、カテゴリからページを絞り込んでいく検索で、Yahoo!などに代表される。後者はWWWサーバ上をロボットが常に巡回してページのソースを自動収集し、ユーザが入力したキーワードが入っているページを出力するものであり、gooなどに代表される。

サーチエンジンはWWWの情報を検索するのに有効であるが、ユーザが求めている情報がページのトップに出てこないことが多いようである。これはHP作成者がアクセス数を増やすためにページに手を加えているというのが原因の1つとして考えられる。

本研究では検索結果が膨大な量になってしまうロボット検索について研究を行う。

3 問題点の改善

ユーザが目的のページを探しにくいのは検索結果が膨大な数になってしまうからである。よって先にあげた問題点を解決することで見易い検索結果の表示をすることができると考え、それを改善する方法を以下に提案する。

3.1 tf-idf法

ページ内で繰り返し使われている語句はそのページにおいて重要と考えられる。よって出現回数 (term frequency: tf) をランク付けに使用する。しかし多くのページで使用されている語句は索引語としての有用性は低い。そこで出現ページ数の逆数 (inverse document frequency: idf) も使用する。具体的にはtfとidfを掛け合わせることで重みを定めることとなる。この手法をtf-idf法と言う。

3.2 問題点の改善方法

あるキーワードをサーチエンジンで検索したとき、キーワードと全く関係のないページが検索されたことはないだろうか？

もし背景と同じ色で文字が書かれていた場合、ユーザの閲覧するページには何の支障もなくキーワードをそのページに埋め込むことができる。サイトへのアクセス数が一定数を越えることでプロバイダから宣伝料が支払われるなどの金銭的な問題もあり、アクセス数を増やすためにHP作成者がこのような行為をすることは少なくない。これが検索結果にキーワードと関係のないページが出力される原因の1つである。これを改善するためには背景と同じ色で書かれた文字を出現文字として換算しないようなtf-idf法を開発すれば良い。

通常、検索を行う場合、ページのソースをデータベースに蓄積する。そしてHTML言語などの命令文

字を排除した状態で検索を行う。だが、命令文字を排除した時点で背景と同色の文字も排除し検索を行えば、このような問題点を改善できると考えられる。

具体的にはHTML言語において背景を設定するBGCOLOR="XXXX"とフォントを設定するFONT COLOR="XXXX"の両者のXXXXの部分が一致したとき、そのフォントで書かれた文字をtfとして換算しないというプログラムになる。

しかし色の設定には"black"や"white"などのように英語で書かれたものと"#000000"や"#ffffff"などのように光の三原色を16進数で表したものがある。この"black"と"#000000"、"white"と"#ffffff"はそれぞれ同じ色を表示することなどでどちらにも対応するようなプログラムを構築する必要がある。

4 実験

今回は既存の「goo」を使って実験を行った。キーワードとして使われることが多い「MP3」というキーワードを使用し、またページに手を加えているのはアダルトサイトに多いと言うことで、「MP3」と「アダルト」のAND検索をする。その結果、検索結果数は1126件が検索され、実際にトップページに表示された10ページ(1ページは存在せず)を開いたところ、背景と同じ色で書かれたページが6ページ存在した。開いたのはトップページの10ページだけだが残りの1116件のページにトップページと同じ割合でこのようなページが入っていると考えると背景と同じ色で書かれたページ数Xは

$$X = A \times \frac{M}{10 - N}$$

ただし

A: 検索結果数

M: トップページ内で背景と同じ色で書かれた文字が存在したページ数

N: トップページ内で存在しなかったページ数

となり全検索結果の中でこのようなページ数は750件存在すると予想される。半数以上のページが該当することからこれによってキーワードと関係ないページのランクを下げることができ、ユーザの見やすい検索結果が表示できると考えられる。図1はgooの検索結果ページの種類の分類である。

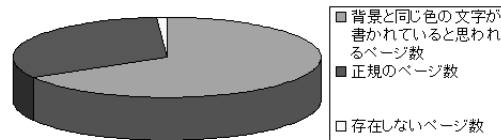


図1: 実験結果

5 おわりに

本研究では、サーチエンジン「goo」を用いて背景と同じ色で書かれた文字を出現文字として換算しないシステムを提案し、ユーザにとって有用な検索結果を表示することがわかった。今後は、ホームページのURLからソースを取得し、そこに本稿で述べたようなシステムを適用してランク付けを行えるサーチエンジンのプロトタイプシステムの構築が必要である。また、既存のサーチエンジンと比較し、このシステムが有効であることを示す。