

3X-05

ページ内リンクリスト情報を利用した
Web ページアクセス制御用 URL リスト拡張手法

牧野 俊朗 杉崎 正之 稲垣 博人

日本電信電話株式会社
NTT サイバーソリューション研究所

1 はじめに

インターネット上には、多種多様な情報があるが、子供がインターネットを使用する場合、親が子供に見せることが望ましくないと思うコンテンツに子供がアクセスしてしまうことがあり、このようなアクセスを規制したいという要望がある。また、企業によっては、社員の生産性の低下を防ぐために、業務に無関係なサイトへのアクセスを規制する場合がある。このようなアクセス規制を実現するために、様々なフィルタリングソフトが提供されている。フィルタリングソフトは、基本的に人手で作成した URL リストを元に、アクセスの可否を決定している。インターネット上では、新サイトの出現および既存サイトの消滅が頻繁の行われているため、この URL リストの作成およびメンテナンスは多大な労力を必要とする。一方、インターネット上には、様々なリンク集が存在し、ページ作成者がメンテナンスをしている。本稿では、このリンク集を利用することによって、Web ページへのアクセスを制御するための URL リストを拡張する方法について述べる。

2 フィルタリングソフトの現状

インターネットへのアクセスを制限するためのフィルタリングソフトは、ページ製作者等がページに書き込んだジャンルや表現の程度に関する識別情報のラベルとページ閲覧者等が設定した情報とを比較してアクセスの制限をしたり、フィルタリングソフト会社が作成した URL のリストを用いてアクセスの制限をしたりする。また、単語やフレーズの辞書を持ち、辞書内の語とページ内に出現する語を比較することにより、制限するか否かの判断をするものもある。ページ内のラベルを利用する方法は、適切にラベルが貼られていれば確かであるが、現状では、適切なラベルが貼られているページはまだ少ない。また、辞書を使う方法は、適切な

辞書を作成することが困難であり、また、単語やフレーズによる機械的な判断は間違えることも多い。また、画像のみのページでテキストのないページも存在し、そのようなページには対処できない。このため、フィルタリングソフトは、URL リストを必要とするが、新サイトの出現や既存サイトの移転などが頻繁におこるため、そのメンテナンスには多大な労力が必要とされている。

3 ページ内リンクリスト情報を利用した URL リスト拡張手法

インターネット上には、様々なリンク集的なページが存在し、ページ製作者によってメンテナンスが行われている。これを利用して、URL リストを拡張する方法を提案する。

3.1 手法の概要

次のような手順で、URL リストに追加する URL を決定する。

- (1) ページ内リンクリストの作成
まず、ロボットで収集した Web ページの HTML ファイルから、HTML タグを利用して、リンク先の URL の一覧を作成する。これをページ内リンクリストと呼ぶことにする。
- (2) リンク集ページの推定
次に、規制対象ジャンルのリンク集を見つける必要がある。ここでは、あらかじめ規制対象ジャンルの URL リストをある程度人手で作成しておき (これを初期 URL リストと呼ぶ)、ページ内リンクリストと比較し、初期 URL リスト中の URL をある数以上含むページをリンク集と見なすことにする。
- (3) 各リンク先 URL の評価値の算出
リンク集と判断したページのページ内リンクリス

URL database expansion method for web filtering by using the link lists in web pages

Toshiro MAKINO, Masayuki SUGIZAKI,
and Hirohito INAGAKI
NTT Cyber Solutions Laboratories, NTT Corporation

ト中のURLのうち、初期URLリストに存在しないURLに関して、評価値をつける。

(4) 追加するURLの決定

ロボットで収集したすべてのページで上記の(1)～(3)を行い、評価値がある閾値以上のURLを追加するURLとする。

3.2 評価値の算出方法

評価値の算出に関しては、以下の点を考慮した。

- 1ページ内に複数のジャンルのリンクを掲載しているリンク集が存在すること。
- リンク集ページ内にも、リンク集とは無関係なリンクが存在すること。

1点目に関しては、まず、ページ内リンクリスト中のURLのうち、初期URLリストに存在するURLと存在しないURLの割合を考え、存在するURLの割合が高い場合に、存在しないURLの評価値が高くなるようにした。これにより、複数ジャンルのリンク集が存在するページでは各URLの評価値が小さくなることが期待できる。また、複数ジャンルのリンク集が存在する場合、目的のジャンル以外のジャンルはページによって異なることが期待できるので、ページ毎の評価値の和を最終的な評価値とすることにより、さらに他ジャンルのURLの評価値を相対的に小さくすることができる。これは2点目の無関係なリンク先の評価値を下げるためにも有効である。2点目に関しては、さらに各サイトの被リンク数をカウントし、リンク集からのリンク数との割合を考え、目的のジャンルのリンク集以外からの被リンク数の割合が高いサイトに関しては、評価値が低くなるようにした。上記の考え方に従い、以下の式により評価値 v を求めた。

$$v = \frac{\sum_{i=1}^M \frac{b_i}{n_i}}{N}$$

ただし、 b_i は、あるページ内リンクリスト i 中の初期URLリストに含まれるURLの数、 n_i は、ページ内リンクリスト i 中のURLの数、 M はそのURLを含むリンク集ページの数、 N は、そのURLの被リンク数である。

4 実験と考察

ロボットで収集したWebページからランダムサンプリングした70万強のページを対象に、フィルタリ

ングソフトで必ず規制対象とされるアダルトジャンルに関して、本手法によりURLリストの拡張の実験を行った。初期URLリスト中のURLの数は197、初期URLリスト中のURLを2つ以上含むページをリンク集ページと見なした。結果は、リンク集として判定されたページ数が325、閾値を0とした場合の追加URL候補数は3779であった。図1に閾値を変化させた場合の追加URL候補の数および、そのうち実際にアダルトジャンルのURLであった数、接続ができなかったURLの数を示す。なお、評価値の絶対値は初期URLリストによって大きく変化するので、図では最大値を1として割合で表示してある。

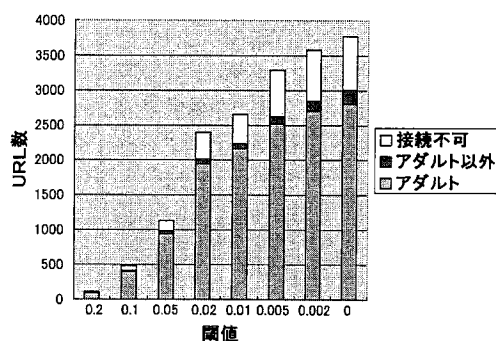


図1 閾値によるURL数の変化

閾値に関わらず接続不可のURLを除いて求めた適合率は0.93以上あり、本手法で有効にURLの拡張ができることが確認できた。また、閾値が0.2,0.1の場合に、一時的に適合率が低下するが、それ以降は閾値を下げるるとともに適合率が低下しており、本評価値の設定方法が、無関係のURLを取り除く効果を持つことが確認できた。一時的な適合率の低下は、同一作者が作成したと思われるほぼ同じ内容のアダルトリンク集ページからリンクされているアダルト以外のページで他のページからあまりリンクされていないページが存在し、そのページの評価値が高くなってしまっているのが原因であった。今後、処理ページ数を増やして、同様の問題が生じるか否かの検証を行うとともに、よりよい評価値の算出方法を検討したい。

5 おわりに

ロボットで収集したWebページ内のリンク情報を利用して、アクセス制御用のURLリストを拡張する方法を提案した。本手法は再帰的に適用できるため初期URLリストが小さくても大量のURLリストを自動的に作成できる。