

西塚 一樹* 西埜 覚* 苗村 憲司**

通信・放送機構 横浜コンテンツリサーチセンター* 慶應義塾大学環境情報学部**

1. はじめに

インターネットには、生徒に閲覧させることが授業目的の達成を妨げる情報(以下では有害な情報という)が多く公開されており、これらの閲覧を防ぐことが求められる。授業では利用目的が異なるので、公開される情報の判断基準を画一化することは困難である。画一的な基準を避け、目的にそって閲覧できるようにするために、情報のコンテンツ(ウェブページ)を複数のカテゴリに分けて、その各々に複数のレベルで有害な情報である可能性を提示する格付け方式(多段階レーティング方式)が考えられている。このような方式を用いれば、受信した側で有害な情報をブロックして、授業目的にあった情報のみを閲覧することが可能になる。WWWコンソーシアム標準のPICS(Platform for Internet Content Selection)は、この格付け情報を記述する構文である。PICSに基づく方式の一つであるRSACi (Recreational Software Advisory Council on the Internet)では、4カテゴリ・5レベルの格付け基準を定めており、国際的な基準作りも推進されている。

受信者が情報を選んで閲覧するためには、事前の格付けが必要である。格付けには情報提供者が自ら行なう「自己格付け」と、ISPや諸組織が独自に行なう「第三者格付け」がある。格付けの作業は人手に頼るため、多大な工数を必要としており、その省力化が求められている。

多段階レーティング方式を前提に、ウェブページをテキスト処理で有害な情報の数値化、ページのリンクの関連付けによって有害な情報の検出を行い、格付け作業の省力化を実現する技術の研究開発を行なっている[1],[2]。

2. テキスト処理による有害な情報の数値化

ページのテキスト処理を行ない、①キー単語、②単語の組み合わせ、③文節(短い文)、④自己URL中の語句、⑤リンク先ページURL中の語

句を抽出し、個々の重み付けと出現頻度を基に数値化を行い、4カテゴリ(アダルト、暴力、差別、悪い情報)で有害な情報の検出を試みた。約24,500ページから現在公開されている1,218ページについて、数値化と目視の結果を比較した(図1)。結果が一致したものが1,028ページ(84.4%)と大半を占めた。だが、目視の結果が有害な情報であるにもかかわらず甘い評価になったものが57ページ(2.2%)あった。

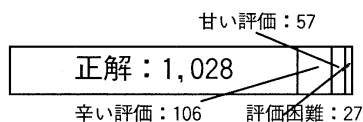


図1 数値化と目視の結果

3. 表示ページを組み込むリンク機能

2.の数値化では、表示される1ページを構成する個別のソースファイルを「ページ」として処理を行なった。図2を例にすると、ページAは、FRAMEタグによってソースファイルF1,F2を組み込み、画面にはAが表示される。この場合、ページAはテキストを含まないことも多い。METAタグのHTTP-EQUIV="refresh"、JAVA Script等でも同様に他のソースファイルを組み込んで表示するケースがある。

Page A (HTMLソース) Page A'(画面表示)

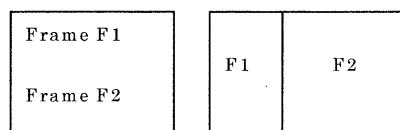


図2 Frameタグの例

2.で目視より甘い評価になったページでは、このようにソースファイルを組み込んで表示しているものが見られた。このような場合、有害な情報の検出には組み込まれるページ(図2ではページA'のF1,F2)との関連性と、その格付

¹A way of relating to the function linked other pages for detecting harmful contents on WWW, Kaduki NISHIDUKA, Satoru NISHINO, Telecommunications Advancement Organization of Japan, and Kenji NAEMURA, Keio University

け情報の利用（情報の先取りという）の可能性を明らかにする必要がある。

4. 表示と組み込まれるページの関連性の実験

FRAME タグ、META タグの HTTP-EQUIV="refresh"（以下では、META タグという）で、アダルトのカテゴリ（RSACi の SEX,NUDE に相当）を中心に関連性を実験した。

なお、以下の記述において、ページが同一のドメインに存在するとき「内部サイト」、アダルトのカテゴリで数値化の結果がゼロでないとき「アダルトな情報」という。

(1)実験に用いたページは約 23,000 ページ。FRAME タグのページは 2,355 ページ(10.2%)、META タグは 417 ページ(1.81%)であった。

(2)FRAME タグのケース

①組み込まれるページの全てが内部サイトのみに構成されているものは 2,095 ページ(89.0%)。一部が外部サイトであるものは 260 ページ(11.0%)である。このことから FRAME タグは、ページ作成者がページの表示を工夫するために利用しているケースが多いと思われる。

内部サイトのみ 2095	外部サイトを含む 260
-----------------	-----------------

図3 組み込まれるページの存在するサイト

②組み込まれるページで有害な情報の数値化が行われたものは 468 ページ(20.0%)で、そのうちアダルトな情報のページが 134(28.6%)で、そうでないページが 334(71.4%)であった。(図4)

③134 ページのアダルトな情報を組み込む側のページでは、30 ページ(22.4%)がアダルトな情報を持っており、104 ページ(77.6%)がそうでないものである。(図5)

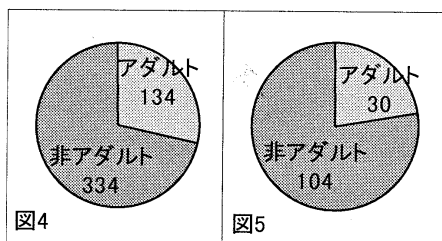


図4 組み込まれるページのアダルトな情報

図5 組み込む側のページのアダルトな情報

(3)META タグのケース

①組み込まれるページが外部サイトであるものは 284 ページ(68.1%)である(図6)。このことから META タグは、URL の変更の適用する

ために用いられるケースが多いと思われる。

内部サイト 133	外部サイト 284
-----------	-----------

図6 組み込まれるページの存在するサイト

②組み込まれるページで有害な情報の数値化が行われたものは 141 ページ(33.8%)。そのうちアダルトな情報を持つページは 62(44.0%)であった。(図7)

③62 ページのアダルトな情報を組み込む側のページでは、34 ページ(54.8%)がアダルトな情報を持っており、28 ページ(45.2%)がそうでないものである。(図8)

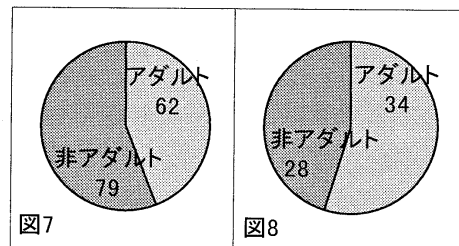


図7

図8

図7 組み込まれるページのアダルトな情報

図8 組み込む側のページのアダルトな情報

4. まとめ

FRAMEタグ、METAタグの評価に用いるデータの数を増やして同様の実験を継続して行なう必要があるが、アダルトな情報に関しては、組み込まれるページとの関連性が明らかになった。これらを用いたページにおいて、組み込まれるページの情報を先取りすることは、通常の別ページをリンク先にする場合と同様に、有害な情報を検出するために有効であると思われる。今後はJAVA Scriptを用いた場合についても同様の評価を行なう予定である。

テキスト処理による有害な情報の数値化、リンク先ページの情報の先取りによって、効率的に有害な情報を検出する仕組みのプロトタイプの開発に、本稿で述べた方式も適用したいと考えている。

参考文献

- [1] 西埜、他：“WWW 上の有害な情報を効率的に検出する一手順”，信学会 99 年総合大会,D-15-17(1999)
- [2] 西埜、他：“WWW の教育用格付けの効率化技術に関する検討”，情処学会 99 年後期全国大会,3X-04(1999)