

レイアウト情報に基づく複数のトピックを含む HTML ファイルの分割

3U-7

前田 英巳子 大野 潮満 黄瀬 浩一 松本 啓之亮

大阪府立大学 工学部 情報工学科

1 はじめに

現在, World Wide Web(WWW)には文書やデータベースなどの様々な情報が無数に存在しており, 利用者が必要な情報を捜し出すには利用者の代わりに情報を検索するシステムの利用が不可欠である. 我々は, 利用者が保存した Hyper Text Markup Language(HTML) ファイルから利用者の興味を抽出し, その結果を用いて検索を行なうシステムについて研究を進めている [1].

WWWのHTMLファイルの特徴として, 1つのファイルに複数のトピックを含むものが多いことがあげられる [2]. HTMLファイルから利用者の興味を正確に抽出するにはファイルをトピックごとに分割する必要がある.

本稿では, HTMLファイルのレイアウトの特徴を利用したファイルの分割手法を提案する. 本手法の特徴は, HTMLタグをレイアウトを示すものとしてとらえ, タグの並びを利用してファイルを分割すること, タグづけされていない部分にタグを挿入することである.

2 HTMLファイルの分割

HTMLファイルをトピックごとに分割するにはファイルの論理構造を知る必要がある. HTMLタグは文書型定義 (DTD: Document Type Definition) を与えることで論理構造を表すことができる. しかし, 実際にはレイアウトを示すタグとして扱われることが多く, この点がトピックを抽出する上で問題となる.

人間は様々な形式を用いられていてもレイアウトからある程度論理構造を認識できる [3]. とくに, 複数のトピックが並列的に書かれている場合, レイアウトの類似性からトピックの境界を発見していると考えられる. そこで本手法では, HTMLファイルからタグの並びをいくつか取り出し, その類似度を計算することによってファイルを分割する. 類似度の計算には, タグ列の伸縮を加味できる DP マッチングを用いる.

さて, HTMLタグには整形済みのデータを表す `<pre>`, `</pre>` がある. `<pre>` と `</pre>` で囲まれた部分には HTMLタグが含まれていないためにタグ列を利用した分割ができない. この問題を解決するため, 本手法では, `<pre>` と `</pre>` で囲まれた部分のレイアウトを解

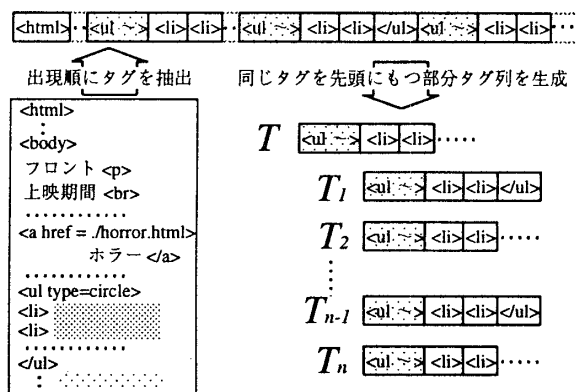


図1 タグ列の生成と部分タグ列の生成

析し, レイアウトを表現する HTMLタグを適切な箇所に挿入することによって他の部分と同等に扱う.

3 分割手法

具体的な手順を以下で示す.

3.1 タグの挿入

まず, `<pre>`, `</pre>` で囲まれたタグを用いていない部分のレイアウトをタグで表現する. レイアウトを表現するタグセットは特別に用意するのではなく, HTMLタグを利用する. これは HTMLタグでレイアウトの大部分を表現できると考えられるからである. 本手法では, センタリングを表現するために `<center>`, `</center>`, 空白行を表現するために `<p>`, 左寄せの行にインデントされた行が続く場合には `<dt>`, `<dd>` を用いる.

空白行は `<p>` に置き換える. `<center>`, `</center>` と `<dd>`, `<dt>` タグの挿入は行頭からの空白文字数によって判断する. ファイル f の行長の最頻値を l_f , i 行目の行長を l_i , i 行目の先頭からの空白文字数を s_i とすると, `<center>`, `</center>` の挿入は, 以下の式を満たす場合に行なう.

$$s_i > l_f/6$$

$$l_f - 3 \leq l_i + s_i \leq l_f + 3$$

また, $s_i = 0$ である左寄せされた行に

$$s_{i+1} = s_{i+2} = \dots = s_{i+m} = a (a > 0)$$

となる行が続けば, i 行目に `<dt>`, $i+k$ ($k = 1, 2, \dots, m$) 行目に `<dd>` を挿入する.

3.2 タグ列と部分タグ列の生成

図1のように HTMLファイルからタグを出現順に抽出してタグ列を生成する. 次に先頭のタグから順に1つずつタグを取り出す. 取り出したタグと同じタグがタグ

Segmentation of HTML Files into Topics Based on Their Layout

Emiko Maeda, Shiomi Ohno, Koichi Kise and Keinosuke Matsumoto

Department of Computer and Systems Sciences, College of Engineering, Osaka Prefecture University

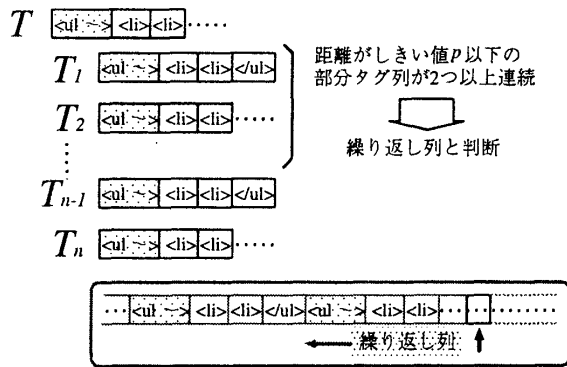


図2 繰り返し列の判定

列にあれば、部分タグ列を生成する。図1の場合はまず、 が繰り返し使われていることが分かるのでで始まる部分タグ列を生成することになる。1つめのから次のが出るまでのタグを抽出し、部分タグ列Tを作る。同様に、 $i+1 (i \geq 1)$ 回目のから $i+2$ 回目のまでのタグを抽出して部分タグ列 T_i を作る。から始まる部分タグ列を全て作り終えたらDPマッチングを用いてTと T_i の距離を計算する。

3.3 DP マッチングを用いた距離計算

部分タグ列Tと類似したタグ列が繰り返されているかどうかを調べるためにTと T_i との距離を計算する。距離計算にはタグ列の伸縮を加味できるDPマッチングを用いる。DPマッチングにはタグ列の伸縮を2つまでとする傾斜制限を設ける。また、タグは1対1対応であるとして、タグの脱落を許した対応づけを行なう。脱落したタグを記号*で表すとき、タグtと t' との距離 $d(t, t')$ を次のように定める。

$$d(t, t') = \begin{cases} 0 & (t = t' \text{ のとき}) \\ d_1 & (t \neq t', t \neq *, t' \neq * \text{ のとき}) \\ d_2 & (t = * \text{ または } t' = * \text{ のとき}) \end{cases}$$

ただし、 $d_1 > d_2$ である。 $d(t, t')$ に基づいたDPマッチングで求められるTと T_i との距離を $DP(T, T_i)$ 、Tの長さを $l(T)$ とすると、Tと T_i との距離 $d(T, T_i)$ は次式で計算される。

$$d(T, T_i) = \frac{DP(T, T_i)}{d_1 \min\{l(T), l(T_i)\} + d_2 |l(T_i) - l(T)|}$$

$DP(T, T_i)$ の値は部分列の長さに依存するため、上式では $DP(T, T_i)$ がとり得る最大値(分母)によって正規化している。 $d(T, T_i)$ がしきい値p以下であれば、 T_i をTの繰り返し列の候補とする。

3.4 繰り返し列の判定

図2のように距離 $d(T, T_i)$ がしきい値p以下である T_i が連続して $n (n \geq 2)$ 個続けば、 $T_1 \sim T_n$ をTの繰り返し列と判断する。

繰り返し列を見つけた場合は、タグ列から T_n の次のタグを取り出して3.2、3.3で説明した手順を繰り返す。

一方、全ての T_i に対して $d(T, T_i) > p$ である、あるいは、先頭のタグ(図2の場合は)から始まる部分列がない場合には、先頭の次のタグを取り出して3.2、3.3で説明した手順を繰り返す。

最終的には以上の処理で抽出された $T, T_1 \sim T_n$ を個々のトピックに対応する分割部とする。分割部はHTMLファイルから各部分タグ列に当たる部分を取り出して生成する。なお、タグ列に繰り返し列が見つからなかった場合は、ファイル全体を1つの分割部とする。

4 実験および考察

情報処理学会、人工知能学会、電子情報通信学会、電気学会のサイトから得たHTMLファイル156個に対して、本手法と<pre>と</pre>で囲まれた部分にタグを挿入しない手法との比較実験を行なった。各パラメータは、 $d_1 = 2, d_2 = 1, p = 0.21$ に設定した。評価は人手で作成した分割部3551個に対する抽出率を用いた。

本手法で生成した分割部は4516個、抽出率は63.9%であった。一方、タグを挿入しない手法で生成した分割部は4110個、抽出率は57.5%であった。この結果から、以下のことが分かる。

- 複数のトピックが並列的に書かれたファイルを対象とすれば、本手法のような単純な処理でも60%程度のトピックを抽出できる。
- <pre>と</pre>で囲まれた部分へのタグの挿入により、分割の精度を向上できる。
- タグに基づくファイル分割は、過度に分割する傾向がある。

本手法では、部分タグ列間の距離計算にDPマッチングを用いているので、部分タグ列間の多少の違いを吸収できる。しかし、改行を示す
や空行を表示する<p>が連続して用いられると、1つのトピックが複数行にわたっていても各行を分割部としてしまう。この問題に対処するためには、タグ以外の特徴である文字列中の文字種の規則性などを考慮した分割が必要となるだろう。

5 まとめ

本稿では、HTMLタグをレイアウトを示すタグをしてとらえ、タグの並びを利用してファイルをトピックごとに分割する手法を提案した。今後の課題としては、レイアウトだけでなく、文字種の違いも利用してファイル分割の精度を向上させることがあげられる。

参考文献

- [1] 大野潮満 他: “HTMLファイル分割に基づくユーザの興味把握とWWW検索”, 情報処理学会第57回全国大会, 5V-3, pp.3-279-3-280(1998).
- [2] 野本豊裕 他: “文脈情報を利用したブラウジング支援—Context-Sensitive Filtering—”, 人工知能学会研究会資料, SIG-FAI-9702-16(9/30), pp.91-96 (1997).
- [3] 大久保雅且 他: “話題が混在するテキストからの話題セグメントの抽出方式”, 情報処理学会第57回全国大会, 1V-5, pp.3-209-3-210(1998).