

HTML 表データの構造認識と携帯端末表示*

塚本 修一[†], 増田 英孝[†], 中川 裕志[‡]

東京電機大学工学部[†] 東京大学情報基盤センター[‡]

1 はじめに

近年、携帯電話や PDA などの携帯端末から Web ページをブラウズしたいという要求が急増している。しかし、現状では、PC の高解像度大画面 (解像度が最低でも 640 × 480 以上) を前提として作られているページがほとんどである。携帯端末デバイスの画面解像度は年々高くなっているが、携帯端末の画面サイズは限られているために読める大きさで表示できる文字数には物理的限界がある。また、画面をスクロールさせるための操作量が増加する。さらに、Web ページ上の表を表示する際に、ブラウザによって <TABLE> タグの取り扱い方や、対応するタグの種類が異なるため、その表示に問題が発生する。そこで、本研究では高解像度大画面向けに作られた既存の Web ページを携帯端末でブラウズする際の表の表示に問題点をしぼり、まず、表の項目名、項目名に対応するデータ (以下、「項目データ」と呼ぶ) の境界を同定することにより、その構造を認識するアルゴリズムを提案し、評価した。次に、このアルゴリズムを用いて表を携帯端末に適した形に自動変換して表示するシステムを実装した。

2 携帯端末における表示の問題点

携帯電話、PDA などの携帯端末を用いて Web ページをブラウズする際には、小画面、低解像度のためにさまざまな問題が発生する。ここでは、具体的な問題点を挙げその解決方法の提案を行う。

2.1 携帯端末で表を表示する際の問題点

表は、本来情報を整理し分かりやすくするために作られている。しかし、小画面低解像度の携帯端末で表をブラウズすると、逆に可読性が低下し、読み誤りが生じることがある。また使用するブラウザによって表示が異なる場合がある。図 1 に PC で表を含むページを表示した例を示す。解像度が高く画面サイズが大きいため、表全体を見渡すことができる。図 2 に PalmsOS[1] 上の AvantGo[2] ブラウザ、図 3 に Palmscape[3] ブラウザで図 1 と同一の表を含むページを表示した例を示す。図 2 の AvantGo では、罫線が表示されないために表の行と列の関係を保持することが難しい。次に、図 3 の Xiino では罫線が表示されているので行と列の関係を認識できるが、小画面低解像度のために以下の問題が発生する。第 1 に、各セルの横幅が狭くなるためにセルデータの途中で折り返しが発生し、読み誤りを起こす可能性がある。第 2 に図 4 は、図 3 の画面をスクロールしたものであるが、表の項目名の部分が隠れてしまい、表の各セルが何を示すか見失ってしまう。その結果、スクロールしてページを戻さなくてはならない。表の行と列の数が大きくなればなるほどこれら 2 つの問題が顕著となる。また、表の <TD>, <TH> タグの colspan, rowspan オプションの値が増加すると、1 つのセルデータを 1 画面内に収めて表示できなくなり、さらに可読性が低下する。図 5 にその例が顕著に表れたものを示す。

3 表の構造認識システム

3.1 システムの概要

本システムは、[4] で述べた本質的な表のみを対象とする。これまでに、表の研究では言語的性質を点数化し表の表すドメインを認識する研究 [5, 6] があるが、複雑な表や、セルの複数に属性があるもの、また未知のドメインには対応できない。

*Table Transformation in HTMLdocument for Mobile Terminals

[†]Shuichi TSUKAMOTO, [†]Hidetaka MASUDA, [‡]Hiroshi NAKAGAWA

[†]Department of Electrical Engineering Tokyo Denki University, [‡]Information Technology Center The University of Tokyo

(人)

	総数	30～39歳	40～49歳	50～59歳	60～69歳	70歳以上
総数	8369	1480	1660	1995	1701	1533
男性	3854	682	777	928	832	635
女性	4515	798	883	1067	869	898

図 1: PC 画面での表の表示例

図 2: AvantGo での表示例

図 3: Xiino での表示例

図 4: スクロールした時の Xiino での表示例

図 5: rowspan オプションがあるページを表示した例

[7]の研究でタグの構造のみで項目名の認識を行っているが、正解率は60%である。これに対して、本研究ではセル間の類似度をベクトル空間法によって計算し、類似度の比を用いて、行と列の項目名と項目データを計算し区別する。表の i 行 j 列のセルを $Cell_{ij}$ として、各セルの N 個の言語的性質 $k = 1, \dots, N$ に対応して、その性質を持てば1、持たなければ0と値 w_k を定義する。 w_k を要素とするベクトルを式(1)のように定義し、表のセルデータをベクトル化し、計算する。

$$\overrightarrow{Cell_{ij}} = (w_1, w_2, \dots, w_N) \quad (1)$$

以下にベクトルの要素となる言語的性質を列挙する。今回の実験では N は合計で99であり、内容は以下に示す。

- 連続データ (1次元)
- 句読点 (2次元)
- 文字長 (3次元)
- 接頭辞 (14次元)
- 接尾辞 (43次元)
- 単位 (17次元)

- 特殊文字 (11次元)
- テーブルタグの属性 (3次元)
- 文字種 (5次元)

3.2 認識アルゴリズム

m 行 n 列の表の行間、あるいは列間の類似度を計算するために、まず表の i 行 j 列のセルを $Cell_{ij}$ として表し、同じ列の $Cell_{kj} (k \neq i)$ との類似度の平均 $Sim_{row}(i, j)$ を次式で求める。

$$Sim_{row}(i, j) = \frac{1}{m-1} \sum \frac{\overrightarrow{cell_{ij}} \cdot \overrightarrow{cell_{kj}}}{|\overrightarrow{cell_{ij}}| |\overrightarrow{cell_{kj}}|} \quad (2)$$

ここで、 \sum の範囲は、 $k = 1, \dots, n$ 、但し、 $k = i$ は除く。 $\overrightarrow{cell_{ij}} \cdot \overrightarrow{cell_{kj}}$ は、 $\overrightarrow{cell_{ij}}$ と $\overrightarrow{cell_{kj}}$ の内積を表し、 $|\overrightarrow{cell_{ij}}|$ と $|\overrightarrow{cell_{kj}}|$ は、それぞれ $\overrightarrow{cell_{ij}}$ と $\overrightarrow{cell_{kj}}$ の絶対値を表す。したがって、 \sum の内側の式は、 $\overrightarrow{cell_{ij}}$ と $\overrightarrow{cell_{kj}}$ の cosine である。図6で $Cell(1,1)$ と第1列中のセルとの類似度の計算の様子を示した。次に、第 i 行のセル、即ち $Cell_{ij} (j = 1, \dots, n)$ のすべてについて $Sim_{row}(i, j)$ を計算し、その行と他の行との類似

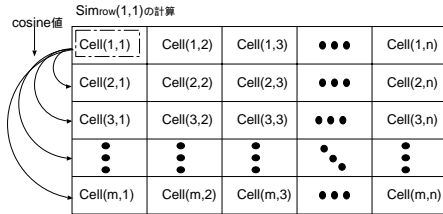


図 6: $Sim_{row}(1,1)$ の計算

度の平均 $Sim_{row}(i)$ を次式で求める。

$$Sim_{row}(i) = \frac{1}{n} \sum_{k=1}^n Sim_{row}(i,k) \quad (3)$$

式 (3) で計算した結果を図 7 に示す。 $Sim_{row}(i)$ の値は、第 i 行が、他の行と類似していれば大きく、類似していなければ小さくなる。項目名を表す行と項目データを表す行とは類似度が低い。一方、項目データを表す行同士は類似度が高い。また、項目名を表す行は Web ページでは上にくることが一般的である。 $i = 1$ が第 1 行である。例えば、1 行目と 2 行目の間が項目名と項目データの境界なら図 8 のようになる。

$Sim_{row}(1)$
$Sim_{row}(2)$
$Sim_{row}(3)$
\vdots
$Sim_{row}(m)$

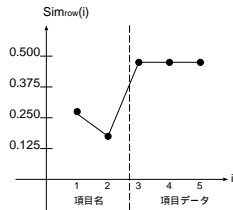


図 7: 式 (3) の計算結果 図 8: 項目名と項目データの境界における Sim_{row} の変化の様子

よって、 $Sim_{row}(i)$ と i 行以下の $Sim_{row}(i+1), \dots, Sim_{row}(m)$ の平均の比 $R(i)$ を、式 (4) のように定義し、

$$R(i) = \frac{Sim_{row}(i)}{\frac{1}{m-i} \sum_{k=i+1}^m Sim_{row}(k)} \quad (4)$$

項目名と項目データの行の境界 T は次のアルゴリズムで求まる。但し、 θ は、境界かどうかを判定する閾値である。

```

T = 0;
for(i=1;i<=m;i++){
    if(R(i)<θ){ T = i; }
    else{ break; }
}
if(T==0){ 縦方向に境界なし }
else{T 行までが項目名の行 }
    
```

以上は項目名の行と項目データの行の境界を求めるアルゴリズムだが、以上の導出において、縦横を交換すれば、 $Sim_{col}(j)$ を計算でき、そして項目名の列と項目データの列の境界を認識できる。以上のアルゴリズムによって、切り出された結果から表の型に当てはめる [4]。

3.3 認識アルゴリズムの評価実験

本アルゴリズムの評価には、 θ の最適化を含め 10 fold 交差検定によって 300 個の表を評価し、結果、最適な θ の値 (行の閾値 $\theta = 0.90$, 列の閾値 $\theta = 0.70$) を適用し 80% の正解率となった [4]。また、評価を行った表の大きさの平均は 9 行 4 列であり、それぞれの型の個数とその内訳を表 3 に示す。

この表 3 の結果の内、時間割型となるものは 67 個

表 1: 行方向の結果

データの種類	正解率
トレーニングデータ	83.23%
テストデータ	82.11%

表 2: 列行方向の結果

データの種類	正解率
トレーニングデータ	79.11%
テストデータ	78.11%

表 3: 交差検定によるテストデータとして評価をした 300 表の内訳

		切れ目の行 (or 列)				合計
		0	1	2	3	
型	縦	70	202	25	3	300
	横	183	115	2	0	300

あり、切れ目の内訳は 1 行目 1 列目が 44 個、2 行目 1 列目が 22 個、2 行目 2 列目が 1 個である。表 1、表 2 の結果から、システムはおよそ 80% の正解率で表の項目名を認識することができる。残りの 20% の表は項目名の部分にもかかわらず、言語的類似度がす

べて高く認識できない表 (40%)、逆に項目データの部分にもかかわらず、言語的類似度が低い表 (60%) の2つに大別できる。

4 表示変換

3.3で認識した項目名と項目データを携帯端末で理解しやすい形に表示するための変換の方針としては、常に項目名と項目データをペアで表示することにした。これは、2における考察の結果、項目名と項目データが乖離して読み難くなっていることが分かったので、それを回避するための方策である。これによって、スクロールしても表が表示内容を見失うことがなくなる。表示領域の制限が緩和され、単語途中の折り返しにより可読性が低下することを避けることができる。システムが求めた結果を使って図1、図5の表を変換した例を図9、図10に示す。図9では、はじめに列の項目名の“男性”を表示し、次にそれに付随する行の項目名と、そのペアの値を表示してあり、図2、図3よりは理解しやすい。図10では、スクロールによって、見えなくなってしまった、項目名の“種類”、“内容”、“重量”、“料金”、とそれらペアの値を表示することにより、表の最上部にページ戻すことなく値が何を示すのか理解できる。

平成12年第5次循環器病...	
【男性】	
【総数】	3854
【30～39歳】	682
【40～49歳】	777
【50～59歳】	928
【60～69歳】	832
【70歳以上】	635
【女性】	
【総数】	4515
【30～39歳】	798
【40～49歳】	883

図9: 図1の表をシステムで変換した例

5 まとめ

本稿では表形式データの変換のために表の項目名と項目データを切り出すシステムについて述べた。提案したアルゴリズムを適用したシステムはおおよそ80%の正解率で項目名と項目データを認識することができる。今後の課題として、現段階では各々のベクトル要素の値は1か0としているが、実際に強く働

郵便料金表 通常郵便物	
【種類】	第三種郵便物
	(認可を受けた定期刊行物・開封)
【内容】	心身障害者団体の発行する定期刊行物を内容とし、発行人から差し出されるもの
【内容】	毎月3回以上発行する新聞紙
【重量】	50gまで
【料金】	18円

図10: 図5の表をシステムで変換した例

いているものを調査し、ベクトル値を最適化し、項目名の認識率を向上させる。

参考文献

- [1] パームコンピューティング株式会社.
<http://www.palm-japan.com/>.
- [2] AvantGo, Inc.: AvantGo4.2.
<http://avantgo.com/>.
- [3] 株式会社イリンクス: Xiino2.1/SJ.
<http://www.ilinx.co.jp/>.
- [4] 塚本修一, 増田英孝, 中川裕志. HTMLの表形式データの変換と携帯端末表示への応用. 第151回自然言語処理研究会, Vol. 142, pp. 35-42, 9 2002.
- [5] Matthew HURST and Shona DUGLAS. Layout and language: Preliminary Experiments in assigning logical structure to table cells. *Proceedings of the Fifth Conference on Applied Natural Language Processing*, pp. 217-220, 1997.
- [6] 伊藤史朗, 大谷紀子, 上田隆也, 池田祐治. 属性オントロジーの抽出と統合を用いた実空間と情報空間のナビゲーションシステム. 人工知能学会, Vol. 14, No. 6, pp. 69-77, 1999.
- [7] Hidetaka MASUDA, Daisuke YASUTOMI, and Hiroshi NAKAGAWA. How to transform tables in html for displaying on mobile terminals. *6th NLPRS2001 Workshop of Automatic Paraphrasing: Theories and Applications*, pp. 29-36, 2001.