# Text Extraction from the Web via Text-to-Tag Ratio

## Tim Weninger and William H. Hsu

Department of Computing and Information Sciences

Kansas State University, Manhattan KS

DEXA 2008 Workshop on Text-based Information Retrieval

Turin, Italy

# Outline

- Introduction
  - › Motivation
  - › Related Work

- The Text-to-Tag Ratio
  - › Heuristic
  - › Worst Case

- Methodology
  - › Pre-processing
  - › Computing clusters

- Results
  - › Evaluation Metrics
  - › Results

- Conclusions and Future Work

# Introduction – Motivation [1]



Taken from The Hutchinson News on 8/14/2008

- Problem:
  - › Too much *junk* in a web page

- Goal:
  - › Extract only the content of a page

# Introduction – Motivation [2] – Example



Rendered HTML Document

Published online 8/13/2008

**A home away from school**
**Day care has after-school duties as some clients start academic year**
**By Kristen Roderick - The Hutchinson News - kroderick@hutchnews.com**

(Travis Morisse/The Hutchinson News) Mary Waln, 7, and Nija Morris, 6, read "The Magic Mat" together Wednesday at Hadley Day Care.

The doors at Hadley Day Care opened Wednesday afternoon, and children scurried in with tales of their first day of school.

Nija Morris, a 6-year-old attending Faris Elementary, smiled as she hung her pink-and-blue flowered backpack on a hook and talked to her classmates about her first day.

"I played and I did art and I played outside and I went to the gym, and I went inside and did centers," she said. "And then I went to meet the other classes and then we went home."

The school-aged children were a little more wound up on Wednesday, program director Christie Gardner said. The excitement is always higher the first day of school, and not everyone is in a routine.

Text content of the document

# Related Work [1]

- ## Naïve Approach
  - › Remove all HTML tags



Original, Rendered HTML Document

All Text of the Document

# Related Work [2]

- ## Tag Approach
  - › Use HTML tags as clues for content
  - › Problem: Style-sheets



Original, Rendered HTML Document

```
<div>
        <div>
        </div>
        <div>
                <div>
                        Eat at Joes
                </div>
        </div>
        <div>
                <div>
                        <div>
The doors at Hadley Day Care opened Wednesday
afternoon, and children scurried in with tales of their
first day of school.
                        </div>
                        <div>
Nija Morris, a 6-year-old attending Faris Elementary,
smiled as she hung her pink-and-blue flowered
backpack on a hook and talked to her classmates about
her first day.
                        </div>
                </div>
        </div>
</div>
```

# Text-to-Tag Ratio [1]

---

Algorithm 1: *Text-To-Tag Ratio* pseudocode

---

**input**

$h \leftarrow$ HTML source code

**begin**

Remove all `script, remark` tags and empty lines

**for** each line $k$ to *numLines*( $h$ ) **do**

$x \leftarrow$ number of non-tag ASCII characters in $h[k]$

$y \leftarrow$ number of tags in $h[k]$

**if** $y = 0$ **then**

*TTRArray*$[i] \leftarrow x$

**else**
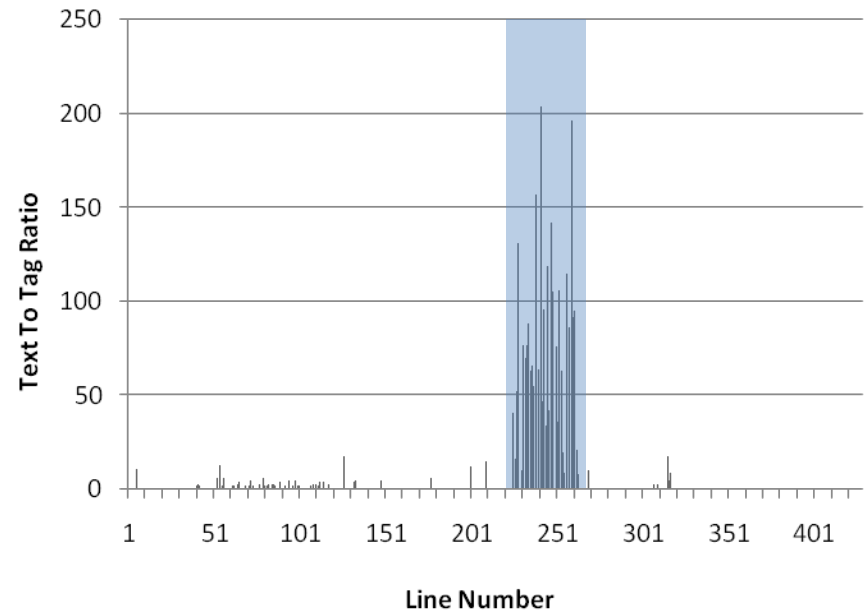
*TTRArray*$[i] \leftarrow x / y$

**end if**

**end for**

**return** *TTRArray*

**end**

---

# Text-to-Tag Ratio [2]

- ## Example

# Text-to-Tag Ratio [3]

- ## Worst Case [1]
  - › Non-HTML or all content pages



TIR'08 Paper

approximation

# Text-to-Tag Ratio [4]

- ## Worst Cases [2]
  - › American Declaration of Independence Web page



American Declaration of Independence TTR computed from digital copy at http://www.ushistory.org/declaration/document/index.htm

# Methodology [1]

- ## Preprocessing
  - › Content Blurring



$$e_k = \frac{\sum_{i=k-r}^{k+r} TTRArray_i}{2r+1}$$

# Methodology [2]

- ## Clustering [1]
  - › K-Means, Farthest First, Expectation







| Cluster | 1 cluster | 2 clusters | 3 clusters |
|---------|-----------|------------|------------|
| 1 | 6.85 | 0.56 | 10.12 |
| 2 | - | **53.40** | **70.42** |
| 3 | - | - | 0.59 |

K-Means clustering

# Methodology [3]

- ## Clustering [2]
  - › Threshold clustering based on standard deviation



Std. Dev. Is 20.3TTR for
Hutchinson News document

# Methodology [4]

- Clustering [3]
  - › Prediction clustering
    - Looks for jumps in the moving average of the TTRArray
    - Not formalized in this paper
    - Very good extension in ANNIE'08 paper.

# Methodology [5]

- ## Evaluation Metrics
  - ### Longest Common Subsequence (LCS)
    - Very Draconian
    - Treated as recall
  - ### Edit Distance Ratio (EDR)
    - Inverse Levenstein distance over  longest sequence
    - Treated as precision

$$EDR = \frac{1 - \text{edtDist}(o, m)}{\max(\text{len}(o), \text{len}(m))}$$

- ## Evaluation method
  - ### 176 Pages selected by querying Yahoo search for "the"
  - ### Gold standard for each page created by a CS undergraduate.
  - ### Metrics computed against gold standard and averaged

# Results [1]

- ## Threshold Only

# Results [2]

- ## Longest Common Subsequence

| | Threshold | EM | K-Means | Farthest First | Prediction |
|---|---|---|---|---|---|
| Mean (%) | **94.19** | 92.62 | 92.47 | 85.88 | 81.14 |
| Median (%) | 98.65 | **99.34** | 98.68 | 94.18 | 94.42 |
| Std Dev. | **14.03** | 17.60 | 16.57 | 21.32 | 24.85 |
| Matches | 34 | **43** | 35 | 25 | 22 |

- ## Edit Distance Ratio

| | Threshold | EM | K-Means | Farthest First | Prediction |
|---|---|---|---|---|---|
| Mean (%) | 56.21 | 48.77 | 57.44 | **62.53** | 52.40 |
| Median (%) | 61.63 | 48.98 | 61.17% | **77.03** | 55.30 |
| Std Dev. | 31.89 | 30.66 | 32.96 | 33.75 | **30.01** |

# Results [3]

- ## Space savings
  - › Mean file sizes

| | HTML | Extracted Text | GZip HTML | GZip Text |
|---|---|---|---|---|
| File Size (Kb) | 9,630.34 | 497.70 | 2,234.77 | 275.53 |

# Conclusions and Future Work

- Text-To-Tag Ratio Approach
  - › A valid content extraction technique
  - › But has Limitations

- Need for better evaluation metrics

- Prediction clustering
  - › Extended for ANNIE'08 in St. Louis, MO, USA
  - › General histogram clustering
    - Uses Gaussian Blurring
    - Analysis of the slope of the tangent line
    - Extracting dimensions and re-clustering
  - › Much better results exist, but were not available by the TIR deadline.

# Questions?