



Web Indexing on a Diet: Template Removal with the Sandwich Algorithm

Stephen Wan
stephen.wan@csiro.au

Paul Thomas
paul.thomas@csiro.au

Tom Rowlands
tom.rowlands@csiro.au

April 2009

Copyright and Disclaimer

© Copyright CSIRO 2009. To the extent permitted by law, all rights are reserved and no part of this publication covered by copyright may be reproduced or copied in any form or by any means except with the written permission of CSIRO.

Important Disclaimer

CSIRO advises that the information contained in this publication comprises general statements based on scientific research. The reader is advised and needs to be aware that such information may be incomplete or unable to be used in any specific situation. No reliance or actions must therefore be made on that information without seeking prior expert professional, scientific and technical advice. To the extent permitted by law, CSIRO (including its employees and consultants) excludes all liability to any person for any consequences, including but not limited to all losses, damages, costs, expenses and any other compensation, arising directly or indirectly from using this publication (in part or in whole) and any information or material contained in it.

Abstract

Web pages contain both unique text, which we should include in indexes, and template text such as navigation strips and copyright notices which we may want to discard. While algorithms exist for removing template text, most rely on first completing a crawl and then parsing each page. We present a cheap and efficient algorithm which does not parse HTML and which requires only a single pass of the document. We have used two web corpora to investigate the performance of a retrieval system using our algorithm and have found similar effectiveness with an index 9–54% smaller. Further experiments using a marked-up corpus have shown 97% of desired lines are returned.

1 Introduction

Retrieving information from within a web document is made more difficult by the presence of template text. Such templates include, for example, the header and footer information that sandwiches the real content of the document. These are typically inserted automatically by HTML authoring tools and scripts that dynamically generate HTML pages, in order to provide website with a consistent look-and-feel. Ideally, an information retrieval system would be able to discard such template material.

In this paper, we treat template detection and removal as a longest common subsequence (LCS) problem, giving an efficient solution. Our experiments with the WT10g corpus and an enterprise data set demonstrate gains in efficiency with low complexity.

Related work has been characterised as using either a *local* or a *global* approach [2]. A local approach examines a page in isolation to find the template material. In contrast, a global approach determines shared templates by examining two or more documents from a collection.

Most approaches handle templates with a two-pass algorithm: the first pass identifies the template and the second extracts it. Approaches to identifying the template have included structural comparisons, often using the document object model of the HTML document. Tree comparison methods have been used to examine similarities in HTML tag elements [5]. Similarly, Wang et al.[6] look for tables specifying layout.

In contrast to examining document structure, other approaches simply examine page text and are thus cheaper to run. Word-level features such as term frequency and word position statistics have been exploited to induce templates [4]. A similar approach using text fragment frequencies is explored by Gibson et al.[2]. Our work is similarly non-structural but does not require any statistical modelling.

2 The sandwich algorithm

We investigate template detection and removal from the viewpoint of improving the efficiency of a web search engine. As such, we start with the constraint that the solution must be able to operate as documents are crawled.

The algorithm is derived from the intuition that, given the prevalence of HTML authoring tools and website content management systems, documents in the same directory will likely share the same template. The template lines are detected by comparing the target file—line by line—with a sibling document in the directory, referred to here as a *peer*. The LCS is a non-contiguous set of common lines between a document and its peer. Our approach assumes this is a template and discards it. The remaining content is considered indexable material and kept. If there are no other pages in the directory, no template removal is attempted.

Our approach is global but reduces to a single pass. That is, identification of the

template is performed per document and at the same time, the template material is extracted. It can be implemented in a crawler before material is stored. If the crawl is breadth-first, in most cases an appropriate peer will simply be the last page crawled.

The LCS can be calculated in $O(mn)$ or less [3], where m and n are the number of lines in each document. No HTML parsing is required; the algorithm is entirely independent of the markup language. The algorithm can remove template text from ‘split’ content, where template material is injected in between portions of useful text. Our implementation is simpler than competing approaches, making template removal an option where engineering resources are limited.

3 Experiments

Our experiments consider two measures. First, we examine the effectiveness and efficiency of a retrieval system which employs the sandwich algorithm. Second, a corpus with templates explicitly marked allows us to investigate our algorithm’s accuracy.

To investigate the performance of a retrieval system which incorporates the sandwich algorithm, we used two corpora. The WT10g corpus, used by the TREC web track [1], includes about 1.7 million web pages from a variety of hosts. Peers were found for 92% of pages. We used three sets of associated queries (“topics”). Topics 451–500 (from TREC 2000) and 501–550 (from TREC 2001) are reverse engineered from search engine query logs. Topics EP1–145, also from TREC 2001, concentrate on finding home pages. Since by removing navigation blocks we will remove a number of links to each site’s entry page, performance on this latter set of queries seems likely to degrade.

The second corpus is in the media domain, and was collected from a large, national media organisation’s website. It comprises about 760,000 documents for which peers were found for 98%. 88 queries were used from a sample of the organisation’s query log, with judgements by an author who was familiar with the organisation. A subset of this corpus has templates explicitly marked.

The first question we ask is: how much more efficient can an index be if templates are removed? To our knowledge, template removal approaches have not been examined by this measure. Table 1 summarises the size of each corpus with and without processing; and the number of postings in an index of each.

Since a lot of templates are formatting or scripting instructions, which will not be indexed anyway, the savings in postings are less than the savings in corpus size—however even the smallest saving, 9% of postings for WT10g, seems worthwhile, and the figures for the media corpus represent a substantial savings. The figures for WT10g are smaller than the 40–50% suggested by Gibson et al., but the WT10g crawl is older than the one used there and the use of templates has been growing since [2]. By insisting on exact string matches we are also conservative in identifying possible templates.

Although a substantial fraction of the index has been removed, performance is

	As-is	Templates removed
WT10g	11,033 MB	9,169 (−17%)
	$1,367 \times 10^6$ postings	$1,238 \times 10^6$ (−9%)
media	12,592 MB	4,107 (−67%)
	$1,097 \times 10^6$ postings	504×10^6 (−54%)

Table 1: Corpus and index sizes for two corpora, before and after processing.

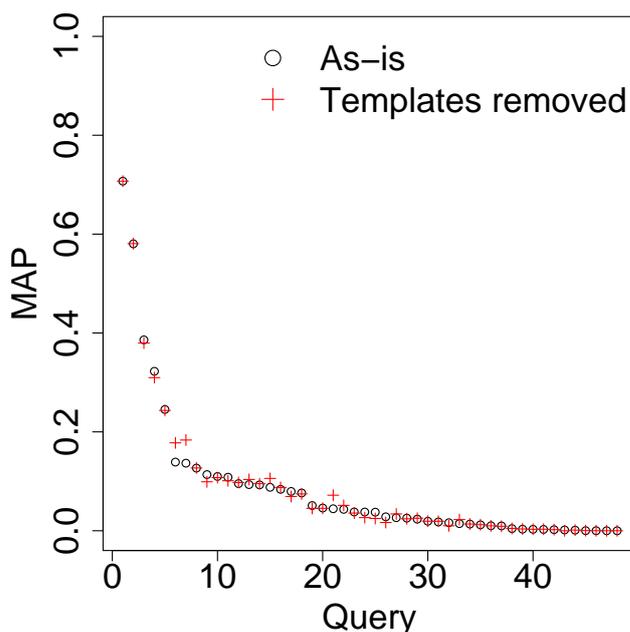


Figure 1: MAP scores for queries 451–500, processed with and without templates in the index.

unaffected. Figure 1 illustrates the MAP scores for each of topics 451–500: on most queries there is no discernable change and overall there is no significant difference (Wilcoxon $p > 0.99$). Topics 501–550 and EP1–145, and the media set, are similar ($p > 0.2$, $p > 0.5$, and $p > 0.4$ respectively).

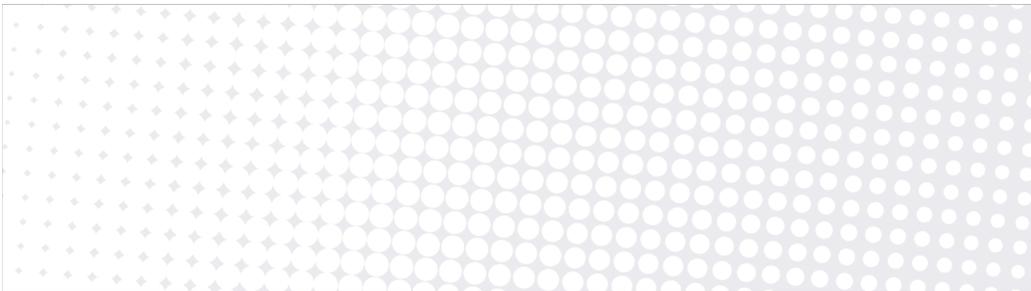
A further question is how accurate are we in removing templates? We compared our output, line by line, with a subset of the media corpus explicitly marked by the organisation. Blank lines and content-less HTML (e.g. a sole `<p>` on a line) were not considered in the comparison. The precision and recall of lines classified as non-template material (and hence kept) is 57% and 97% respectively, with an F1 score of 0.59. The algorithm is correctly keeping the great majority of interesting text, although our conservative approach means we are also keeping a portion of templates.

4 Conclusions and Future Work

Templates represent a substantial, though generally uninformative, portion of text on the web. Removing templates leads to a reduction in index size, without a drop in query performance. Line-based LCS comparison provides a cheap method for template detection and removal, allowing for easy integration within a web crawler. In future work, we intend to use the sandwich algorithm with question answering systems and automatic text summarisers, both of which can benefit greatly with the accurate removal of irrelevant template material.

References

- [1] P. Bailey, N. Craswell, and D. Hawking. Engineering a multi-purpose test collection for web retrieval experiments. *Info Proc & Management*, 39(6), 2003.
- [2] D. Gibson, K. Punera, and A. Tomkins. The volume and evolution of web page templates. In *Proc. WWW*, 2005.
- [3] J. W. Hunt and T. G. Szymanski. A fast algorithm for computing longest common subsequences. *CACM*, 20(5), 1977.
- [4] L. C. Liang, S. Ye, and X. Li. Template detection for large scale search engines. In *Proc. ACM Symposium on Applied Computing*, 2006.
- [5] K. Vieira, A. S. da Silva, N. Pinto, E. S. de Moura, J. ao M B Cavalcanti, and J. Freire. A fast and robust method for web page template detection and removal. In *Proc. CIKM*, 2006.
- [6] Y. Wang, B. Fang, X. Cheng, L. Guo, and H. Xu. Incremental web page template detection. In *Proc. WWW*, 2008.



Contact Us

Phone: 1300 363 400
+61 3 9545 2176

Email: enquiries@csiro.au

Web: www.csiro.au

Your CSIRO

Australia is founding its future on science and innovation. Its national science agency, CSIRO, is a powerhouse of ideas, technologies and skills for building prosperity, growth, health and sustainability. It serves governments, industries, business and communities across the nation.