

# The Off-Topic Memento Toolkit

Shawn M. Jones  
Old Dominion University  
Norfolk, Virginia  
sjone@cs.odu.edu

Michele C. Weigle  
Old Dominion University  
Norfolk, Virginia  
mweigle@cs.odu.edu

Michael L. Nelson  
Old Dominion University  
Norfolk, Virginia  
mln@cs.odu.edu

## ABSTRACT

Web archive collections are created with a particular purpose in mind. A curator selects seeds, or original resources, which are then captured by an archiving system and stored as archived web pages, or mementos. The systems that build web archive collections are often configured to revisit the same original resource multiple times. This is incredibly useful for understanding an unfolding news story or the evolution of an organization. Unfortunately, over time, some of these original resources can go off-topic and no longer suit the purpose for which the collection was originally created. They can go off-topic due to web site redesigns, changes in domain ownership, financial issues, hacking, technical problems, or because their content has moved on from the original topic. Even though they are off-topic, the archiving system will still capture them, thus it becomes imperative to anyone performing research on these collections to identify these off-topic mementos. Hence, we present the Off-Topic Memento Toolkit, which allows users to detect off-topic mementos within web archive collections. The mementos identified by this toolkit can then be separately removed from a collection or merely excluded from downstream analysis. The following similarity measures are available: byte count, word count, cosine similarity, Jaccard distance, Sørensen-Dice distance, Simhash using raw text content, Simhash using term frequency, and Latent Semantic Indexing via the gensim library. We document the implementation of each of these similarity measures. We possess a gold standard dataset generated by manual analysis, which contains both off-topic and on-topic mementos. Using this gold standard dataset, we establish a default threshold corresponding to the best  $F_1$  score for each measure. We also provide an overview of potential future directions that the toolkit may take.

## CCS CONCEPTS

• Applied computing → Digital libraries and archives; • Information systems → World Wide Web;

## KEYWORDS

topic drift, Archive-It, web archive, similarity

## 1 INTRODUCTION

Public web archives are where web pages are preserved and accessible for curiosity, research, and evidentiary purposes [10, 25]. Some researchers go so far as to curate their own collections of archived web pages, or **mementos**. These curators will select **seeds**, or **original resources**, and create their own mementos from these seeds using variety of web archiving platforms, one of which is the Archive-It service [23] offered by the Internet Archive. Collections are created for some purpose and these seeds are chosen to support the collection’s topic. In order to understand the history of an event

or an organization, curators will often configure the platform to capture the same original resource multiple times, thus producing many mementos per seed. Many collections at Archive-It are like *Japan Earthquake*, with 81,014 seeds resulting in 486,227 mementos. Researchers examining these collections want to optimize the amount of time spent evaluating mementos, and the sheer quantity of mementos to evaluate makes it imperative that they not spend time on mementos with low information value.

Consider a collection about the Olympics. The top page of a sports site will be on-topic during the Olympics, but will cover other events once the Olympics have finished. Pages can go off topic for a variety of other reasons and web archives still capture these off topic mementos. Web sites have technical issues (Figure 1). Owners take sites down for maintenance (Figure 2). Hosting services suspend sites (Figure 3). New owners purchase existing domains and replace the site content (Figure 4). Political regimes change, resulting in news sites changing content [5]. Hackers deface pages (Figure 5). Owners restructure sites, resulting in broken links, which are off-topic. Detecting off-topic mementos is important for the development of automated collection summaries [3] or finding aids, where off-topic content needs to be detected and excluded, lest it alter the output. As many as 11% of the mementos in a collection can be off-topic [2]. To that end, we have developed the Off-Topic Memento Toolkit (OTMT)<sup>1</sup>. It is important to note that OTMT only identifies off-topic mementos, it does not remove them from the collection; indeed, when and how pages went off-topic may be of interest to researchers.

We discuss the different similarity measures available with OTMT version 1.0.0 alpha, how we arrived at the default threshold values for each measure, and mention some of the possible future features for the toolkit. We incorporated several different measures because curators may need to tailor the toolkit’s abilities to the content they are evaluating. Our contribution is the review of the effectiveness of different similarity measures for identifying off-topic mementos, especially additional similarity measures not covered in prior work [2], and the availability of a software package allowing users to discover off-topic mementos for themselves.

## 2 BACKGROUND AND RELATED WORK

Web archive collections, like those at Archive-It, are often created along a specific theme. Figure 6 displays a screenshot of the *Climate change and environmental policy* collection at Archive-It, collected by the Stanford University Social Sciences Research Group. The curators who built this collection intended for its mementos to contain information on this topic, and thus selected specific seeds to include in this collection. Each seed has a corresponding **TimeMap** which provides a record of all mementos for that seed [38]. Figure

<sup>1</sup><https://github.com/oduwsdl/off-topic-memento-toolkit>

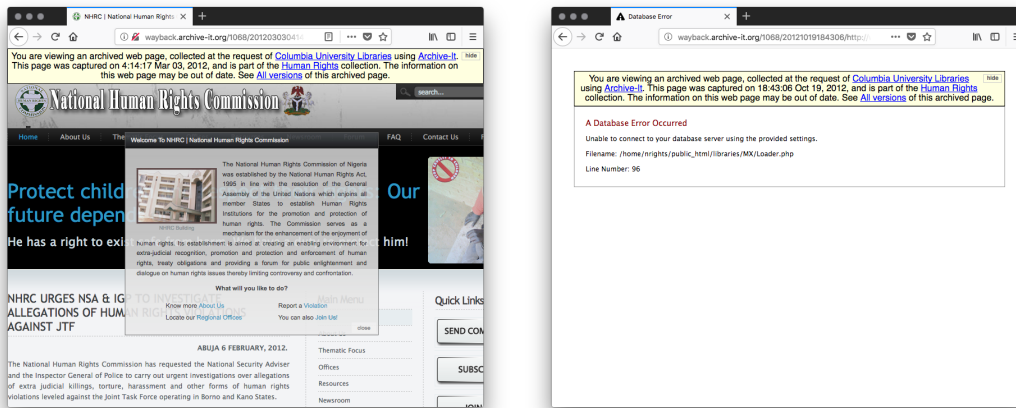


Figure 1: The seed <http://www.nigeriarights.gov.ng/> preserved in Archive-It's *Human Rights* collection (left) goes off topic due to technical issues (right).

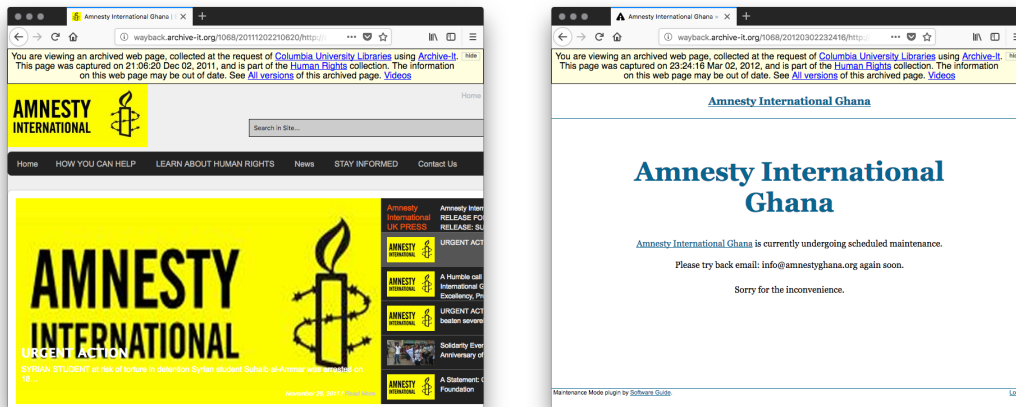


Figure 2: The seed <http://amnestyghana.org/> preserved in Archive-It's *Human Rights* collection (left) goes off topic due to site maintenance

7 displays a TimeMap for the seed URI <http://bloombergvillenow.org/>. Each entry with a relation of memento indicates that the URI on that line is a memento for this seed, and its corresponding datetime keyword indicates when this memento was captured, its **memento-datetime**. As noted above, seeds sometimes go off-topic, and anyone studying climate change will not want these off-topic mementos in the data that they review.

We will use the Memento terminology in the rest of this paper. TimeMap URIs are abbreviated as **URI-T** and memento URIs are abbreviated as **URI-M**.

AlSum [4] explored techniques for producing a thumbnail for each memento in a TimeMap. Each thumbnail is a resized image produced by taking a screenshot of a memento in a browser. Using AlSum's thumbnails, a curator can review all images visually and mark those that are off-topic. There are Archive-It collections, like the *Government of Canada Publications* with 314,032 mementos, that

make manual review of each memento a costly endeavor both in terms of time and personnel.

There are many methods of comparing the similarity of two documents. Seeking ways to improve web crawling, Manku [22] determined that Simhash [8] is effective at discovering documents that are similar. Adar [1] employed the Sørensen-Dice coefficient [12, 36] to understand the changes in content of the same resource over the course of a crawl. Sivakumar [35] tested the viability of the Jaccard Index [15] for improving search results by identifying duplicate advertisements and headers. Hajishirzi [13] applied cosine similarity [33] to the problem of identifying duplicate news articles. Zittrain [39] and Jones [19] have used these methods to analyze content drift in web archive collections. We adopt these **similarity measures** as a way to determine if a memento is off-topic.

Topic modeling techniques, such as Latent Dirichlet Allocation [7], typically break a corpus into clusters of documents. Each cluster

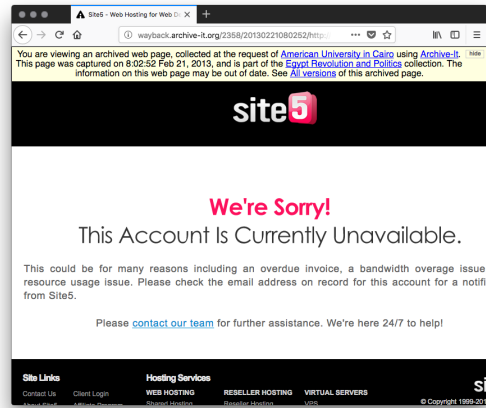
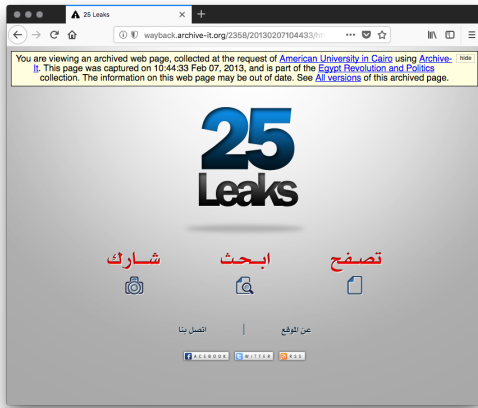


Figure 3: The seed <http://25leaks.com> preserved in Archive-It's *Egypt Revolution and Politics* collection (left) is later suspended due to non-payment.

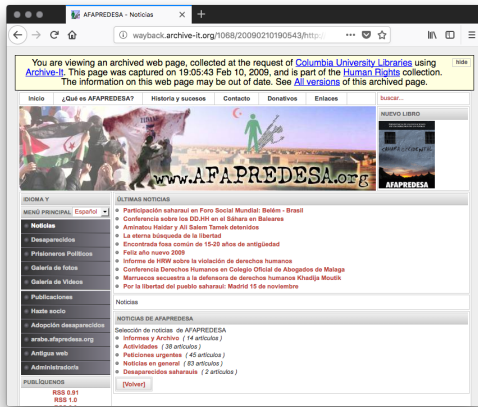


Figure 4: The seed <http://www.afapredesa.org> preserved in Archive-It's *Human Rights* collection (left) changes ownership to a different organization that publishes in Japanese (right).

contains documents that share some topic. We are looking for off-topic documents. Though it is conceivable that the smallest cluster may contain off-topic documents, we have not evaluated that topic modeling techniques behave this way and only consider cosine similarity informed by Latent Semantic Indexing [11], a process developed by Radim Řehůřek [30].

The mementos usually viewed by users of web archives have been augmented for usability and legal reasons, including banners to identify the containing archive. The extra content in these augmented mementos leaves them unsuitable for comparison. Fortunately, Archive-It provides access to **raw mementos** at special URIs [18, 37]. These raw mementos contain the original content that was observed by the web archive at the time of capture, without any rewriting. It is these raw mementos that the OTMT uses in its analysis.

Once the OTMT has the content of a raw memento, it still must preprocess it before comparing it using a given similarity measure. Most mementos consist of HTML, JavaScript, and CSS. This boilerplate provides no useful information for the decision as to whether a memento is on or off-topic, and must be removed [29]. The OTMT then tokenizes, stems, and removes stop words [9] from the resulting text. We refer to these steps as **preprocessing**.

In 2015, AlNoamany performed a study to detect off-topic mementos [2]. In that work, she reviewed the patterns that emerge when a collection goes off-topic. Pages can be always on topic or always off topic. Pages can start on topic and then drift off topic permanently at some point. This is usually the case when a page goes offline, falls under new ownership, or the content has drifted far from the original topic. Pages can also oscillate on or off-topic, often due to hacking or technical problems. She sampled mementos from three Archive-It collections and manually labeled these

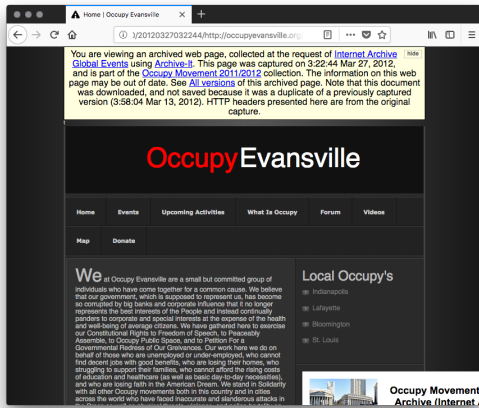


Figure 5: The seed <http://occupyevansville.org/> preserved in Archive-It's *Occupy Movement 2011/2012* collection (left) was later hacked by Anonymous (right)

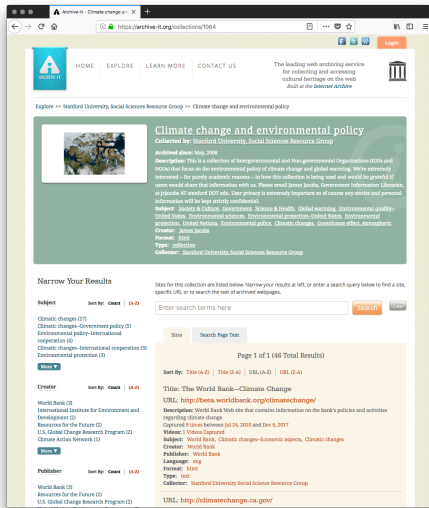


Figure 6: The *Climate change and environmental policy* Collection at Archive-It, collected by Stanford University

mementos as “on-topic” or “off-topic”. Using this dataset, she then evaluated the effectiveness of different similarity measures. This work was later used to remove off-topic mementos in consideration of generating summaries of Archive-It collections [3]. We build upon her work by evaluating additional similarity measures against this same dataset, which we refer to as the **gold standard dataset**, and have also developed the OTMT for curators to use in discovering off-topic mementos.

### 3 TIMEMAP MEASURES

The OTMT uses the TimeMap of each seed to group mementos for comparison. The OTMT supports different similarity measures

```
<http://bloombergvillenow.org/>; rel="original",
<http://wayback.archive-it.org/2950/timemap/link/http://
bloombergvillenow.org/>; rel="self"; type="application/link-
format"; from="Tue, 03 Jan 2012 01:43:26 GMT"; until="Thu, 31
May 2012 20:08:41 GMT",
<http://wayback.archive-it.org/2950/http://bloombergvillenow.org/>;
rel="timegate",
<http://wayback.archive-it.org/2950/20120103014326/http://
bloombergvillenow.org/>; rel="first memento"; datetime="Tue, 03
Jan 2012 01:43:26 GMT",
<http://wayback.archive-it.org/2950/20120109025617/http://
bloombergvillenow.org/>; rel="memento"; datetime="Mon, 09 Jan
2012 02:56:17 GMT",
<http://wayback.archive-it.org/2950/20120531200841/http://
bloombergvillenow.org/>; rel="last memento"; datetime="Thu, 31
May 2012 20:08:41 GMT"
```

Figure 7: An example TimeMap for the seed URI <http://bloombergvillenow.org/> in the Archive-It collection *Occupy Movement 2011/2012*

---

#### Algorithm 1 General algorithm used for all TimeMap measures

---

- 1: **for**  $timemap \in collection$  **do**  
 $f \leftarrow HTTPGET_{raw}(memento_{first})$   
 $f \leftarrow preprocess(f)$
  - 2: **for**  $memento \in timemap$  **do**  
 $m \leftarrow HTTPGET_{raw}(memento)$   
 $m \leftarrow preprocess(m)$   
 $score \leftarrow computeSimilarity(f, m, measure)$   
 $saveScore(score, memento, measure)$
  - 3: **end for**
  - 4: **end for**
- 

against the first (i.e., earliest) memento in a TimeMap. The assumption is that the first memento was on-topic when its seed was submitted to the archive and that automated crawling continued afterward at various intervals.

The general timemap measure algorithm used by the OTMT is shown in Algorithm 1. This algorithm iterates through all TimeMaps in the collection, dereferencing the raw version of the **first memento** (denoted by  $f$ ) in the TimeMap for its content. After preprocessing (if necessary), the algorithm iterates through the preprocessed version of every memento in the TimeMap, comparing the first memento to each additional **considered memento** (denoted by  $m$ ) with the selected similarity measure (denoted by *measure*).

The following sections provide more detail on these measures. As in Algorithm 1, the symbols  $f$  and  $m$  used in the following equations correspond to the first and considered memento.

### 3.1 Structural Measures

The OTMT provides two structural measures that execute much faster than the others. **Byte count** tallies the bytes within a memento’s content and compares the first memento’s bytes with the considered memento. Before calculating the score, only the content of each memento is dereferenced. No preprocessing is performed. Instead of bytes, **word count** tallies the number of words within a memento’s content and compares the number of words in the first memento with the considered memento. Before calculating the score, preprocessing is performed on the memento so that individual words can be counted. The score for each of these measures is based on the percentage difference between the size of the first memento and the considered memento, shown by Equation 1,

$$d_c(f, m) = \begin{cases} \frac{c(m)-c(f)}{c(f)} = \frac{c(m)}{c(f)} - 1 & \text{if } m < f \\ 0 & \text{if } m \geq f \end{cases} \quad (1)$$

where  $c(x)$  indicates the count (byte or word) of memento  $x$ .

The scores range from 0.0, meaning the two documents are the same size or larger, to  $-1.0$  meaning that the considered memento has reached a size of 0. We assume that adding content is common if a memento stays on topic. For this reason, we only consider scores that are negative because off-topic pages often only contain short sentences indicating a 404 message, that the web site has failed to pay its bills, or that there is a technical problem.

### 3.2 Set Operation Measures

Within the OTMT we provide set operations to evaluate the sets of words that make up each document. Because we are interested in the words of each document, the documents are preprocessed before using these measures.

The Jaccard Index (sometimes called Jaccard Coefficient) compares two sets [15]. To normalize this score, OTMT uses the **Jaccard distance** as defined by the Python distance library [24] to compare the two sets of words making up the documents. Jaccard distance calculates the percentage of overlap between the words in both documents, as shown in Equation 2,

$$d_J(f, m) = \frac{|t(f) \cup t(m)| - |t(f) \cap t(m)|}{|t(f) \cup t(m)|} \quad (2)$$

where  $t(x)$  indicates the tokens produced by the preprocessing of the content of memento  $x$ .

The Sørensen-Dice Coefficient is another method of comparing two sets [12, 36]. The OTMT uses the **Sørensen-Dice distance** as defined by the Python distance library to compare the two sets of

words making up the documents. Sørensen-Dice is different in that it takes twice the number of words in common and divides them by the number of total words, shown in Equation 3.

$$d_S(f, m) = 1 - \frac{2|t(f) \cap t(m)|}{|t(f)| + |t(m)|} \quad (3)$$

Both distance measures have scores ranging from 0.0, meaning that the documents are the same, to 1.0, meaning that the documents are completely dissimilar. Both are different from word count because the individual words in each document are considered.

### 3.3 Simhash Measures

The OTMT provides Simhash as implemented by the Python Simhash library [34]. That library’s functions allow for multiple types of input: term frequencies or raw content.

For the **Simhash of term frequencies** the term frequencies (TF) are provided as input to a corresponding cryptographic hash function, thus preprocessing is needed. The hash of each term frequency then makes up part of a larger hash representing the document.

If the **Simhash of raw content** is desired, then the raw memento content is supplied to the function, and no preprocessing is needed. The function converts the raw document into 4-grams, strings of 4 characters long. A hash is computed on each 4-gram and these hashes make up the resulting Simhash. This form of Simhash is influenced by the position of each 4-gram in the document.

Differences between these smaller hashes are reflected in the resulting Simhash, allowing one to compare two documents by comparing their Simhash values. Simhash scores are calculated based on the number of bits different between the two Simhashes. A score of 0 indicates that the two strings of bits are the same. A score of 64 indicates that the two strings are completely different. Because of this, it is highly unlikely that all 64 bits of a Simhash will be different.

### 3.4 Cosine Similarity Measures

Cosine similarity compares two documents as the distance of two vectors. These vectors can be constructed different ways, but the source mementos always require preprocessing.

The first vector construction method supported by the OTMT is **TF-IDF**. With this method each document is converted into a vector that represents each word and its term frequency. The values of the vector are further weighted with the inverse document frequency (IDF) [17] of every term of the other mementos in the TimeMap. This way, each document vector is informed not by just the first and considered memento, but also all other mementos in the TimeMap. The OTMT uses the TF-IDF functionality of the scikit-learn library [28].

The second vector construction method uses the Latent Semantic Indexing (**LSI**) [11] capability of the gensim library [32]. The OTMT implementation for computing these vectors follows the gensim Similarity Queries tutorial [31]. In this case, the vector of each document is informed by LSI.

The cosine of the resulting angle of these vectors (either via TF-IDF or LSI) is then used to generate a distance score. Equation 4 shows how these vectors produce a score,

**Table 1: Similarity measures supported by the OTMT**

| Measure                             | Fully Equivalent Score | Fully Dissimilar Score | Preprocessing Performed | OTMT -tm keyword |
|-------------------------------------|------------------------|------------------------|-------------------------|------------------|
| Byte Count                          | 0.0                    | -1.0                   | No                      | bytecount        |
| Word Count                          | 0.0                    | -1.0                   | Yes                     | wordcount        |
| Jaccard Distance                    | 0.0                    | 1.0                    | Yes                     | jaccard          |
| Sørensen-Dice Distance              | 0.0                    | 1.0                    | Yes                     | sorensen         |
| Simhash of Term Frequencies         | 0                      | 64                     | Yes                     | simhash-tf       |
| Simhash of Raw Memento Content      | 0                      | 64                     | No                      | simhash-raw      |
| Cosine Similarity of TF-IDF Vectors | 1.0                    | 0.0                    | Yes                     | cosine           |
| Cosine Similarity of LSI Vectors    | 1.0                    | 0.0                    | Yes                     | gensim_lsi       |

$$d_c(f, m) = \frac{v(f) \cdot v(m)}{\|v(f)\| \|v(m)\|} \quad (4)$$

where  $v(x)$  indicates the vector produced by the content of memento  $x$ . Cosine similarity ranges from 1, most similar, to 0 indicating that the vectors are completely different.

#### 4 TOOLKIT USAGE

The OTMT allows a user to select an input type, one or more similarity measures, and an output file. These arguments indicate the input-measure-output architecture of the OTMT which attempts to separate the concerns of acquiring mementos for comparison (input), measuring those mementos (measure), and then producing results (output). This architecture facilitates the addition of future input types, measures, and outputs.

The toolkit is run from the command line. For example, to evaluate Archive-It collection 7877 using both measures Jaccard distance and bytecount and then save the output to outputfile.json, one would run the command:

```
detect_off_topic -i archiveit=7877 -o outputfile.json -tm
jaccard=0.80,bytecount=-0.50
```

where  $-i$  indicates the input type followed by its arguments, and the  $-o$  indicates the name of the output file, and  $-tm$  (for *TimeMap Measure*) indicates that the next argument is a list of measures and thresholds. Table 1 shows the available similarity measures, their score ranges, whether preprocessing is performed, and the OTMT keyword used to specify the measure on the command line.

The toolkit supports the following forms of input:

- One or more Memento TimeMaps (input type `timemap` followed by a = and then the URIs of TimeMaps separated by commas)
- One or more WARC [14] files (input type `warc` followed by a = and then the filenames of the WARCs separated by commas)

```
"http://wayback.archive-it.org/1068/timemap/link/http://www.badil.org/": {
  "http://wayback.archive-it.org/1068/20130307084848/http://www.badil.org/": {
    "timemap measures": {
      "cosine": {
        "stemmed": true,
        "tokenized": true,
        "removed boilerplate": true,
        "comparison score": 0.10969941307631487,
        "topic status": "off-topic"
      },
      "bytecount": {
        "stemmed": false,
        "tokenized": false,
        "removed boilerplate": false,
        "comparison score": 0.15971409055425445,
        "topic status": "on-topic"
      }
    },
    "overall topic status": "off-topic"
  },
  ...
}
```

**Figure 8: Example JSON output record from the OTMT for the memento at URI-M <http://wayback.archive-it.org/1068/20130307084848/http://www.badil.org/> from the TimeMap at URI-T <http://wayback.archive-it.org/1068/timemap/link/http://www.badil.org/>**

- An Archive-It collection ID (input type `archiveit` followed by a = and then the Archive-It collection ID)

If an Archive-It collection ID is supplied, then the OTMT extracts all seeds from that collection’s Archive-It pages and constructs URIs to discover all mementos for those seeds. If a WARC is supplied, then its contents are extracted and TimeMaps are generated for each original resource.

The default output is in JSON format, as shown in Figure 8. Each URI-T key contains a dictionary of URI-M keys. Each URI-M key contains a dictionary of similarity measures run against that URI-M. For each similarity measure the output indicates which preprocessing was performed on the memento. Each measure record indicates the topic status (on or off-topic) based on the threshold supplied. For each URI-M an overall topic status is listed, which is determined based on whether or not one of the measures determined that the memento was off-topic. The toolkit also supports a CSV version of this same data.

If a threshold value is not provided on the command line for a given measure, then the OTMT uses a reasonable default. In the next section, we detail how we arrived at reasonable default values for each measure.

#### 5 EVALUATION OF REASONABLE DEFAULT THRESHOLDS

To acquire reasonable default thresholds for each measure, we used the dataset from AlNoamany’s work [2]. Information about the



**Table 2: Definitions of the conditions used to calculate the  $F_1$  score for threshold determinations**

| Gold Standard Data | OTMT      | Condition      |
|--------------------|-----------|----------------|
| On-Topic           | On-Topic  | True Negative  |
| Off-Topic          | On-Topic  | False Negative |
| On-Topic           | Off-Topic | False Positive |
| Off-Topic          | Off-Topic | True Positive  |

**Table 3: Distribution of the Gold Standard Data Set**

| Collection ID | Collection Name               | # seeds in sample | # mementos in sample | # off-topic |
|---------------|-------------------------------|-------------------|----------------------|-------------|
| 1068          | Human Rights                  | 199               | 2302                 | 95 (4%)     |
| 2358          | Egypt Revolution and Politics | 136               | 6886                 | 384 (6%)    |
| 2950          | Occupy Movement 2011/2012     | 255               | 6569                 | 458 (7%)    |

gold standard dataset used in this study is shown in Table 3 and the version used in this study is available on GitHub<sup>2</sup>.

For evaluating the OTMT, we focused on the  $F_1$  measure, seen in Equation 5,

$$F_1 = \frac{2TP}{2TP + FP + FN} \quad (5)$$

where  $TP$  indicates the number of true positives,  $FP$  indicates the number of false positives, and  $FN$  indicates the number of false negatives. Table 2 shows the conditions we used to calculate these values for each measure and threshold combination.

For comparison with AlNoamany’s results, we provide values for a second metric of **accuracy**, shown in Equation 6 using the same symbols as Equation 5, where  $TN$  indicates the number of true negatives and the other symbols have the same meaning as Equation 5.

$$Accuracy = \frac{(TP + TN)}{(TP + FP + FN + TN)} \quad (6)$$

AlNoamany tested each similarity measure with 21 thresholds [2]. To get more precise threshold values for each measure we ran the OTMT and then saved the resulting scores from the comparison of each first memento and considered memento. We then iterated through each memento in the output. Starting with a lower limit as the threshold, we tested each score against that limit and declared that memento as on or off-topic depending on the value of the threshold, the score, and the direction of the comparison operator (e.g.,  $>$  or  $<$ ) for that measure. We then incremented the threshold and compared again, saving the off-topic determination again. This process was repeated until we reached a designated upper limit. Visualizations of the results for each measure are shown in Figure 9 with each threshold on the x-axis and the resulting  $F_1$  score on the y-axis.

<sup>2</sup><https://github.com/oduwsdl/offtopic-goldstandard-data/tree/5139aca762e1ddac76da628436dbc48ae38807f2>

For example, with Byte Count, we ran the OTMT and saved the scores for each memento. We declared a memento off-topic if its score was less than -0.99. We then took each memento’s byte count score and declared it off topic if its score was less than -0.98. We then repeated for each memento with a threshold set at -0.97. This process was repeated, in increments of 0.01 until we reached 0.

We then compared the off-topic determinations per threshold with the gold standard data. From there, we were able to generate a corresponding  $F_1$  score for each threshold value.

We had assumed that testing the thresholds in this way would help us discover threshold values close to those found by AlNoamany, and we did get close in some cases. The  $F_1$  scores, however, are often worse. The OTMT uses the justext library [29] for boilerplate removal, whereas AlNoamany used boilerpipe [20]. OTMT uses nltk [6] for tokenization and stemming whereas AlNoamany used scikit-learn [28]. These subtle differences in libraries combined with gold standard data set updates, download errors, changes in how Archive-It handles mementos now compared with 2015, and differences in preprocessing techniques, are likely the reason for these differences.

Byte count threshold scores between -1 and 0 were tried, in increments of 0.01. A threshold score of -0.39 produces the best  $F_1$  score. This means that the off-topic memento is 39% smaller than the first memento in the TimeMap. AlNoamany’s findings suggested a threshold value of -0.65, making the OTMT results more strict. This is one case where our  $F_1$  score of 0.756 is higher than AlNoamany’s finding of 0.584.

We used the same range and increments to test word count. A threshold score of -0.70 produces the best  $F_1$  score. This means that the off-topic memento has 70% fewer words than the first memento in the TimeMap. AlNoamany’s findings suggested a threshold value of -0.85, again making the OTMT results of -0.70 more strict.

For Jaccard and Sørensen-Dice we used the same range of threshold values from 0 to 1, in increments of 0.01. A threshold score of 0.94 has the best  $F_1$  value for Jaccard. AlNoamany discovered that a value of 0.05 was best, but her work used the pure Jaccard Index rather than the Jaccard distance. Seeing as the Jaccard distance is  $1 - \text{Jaccard index}$ , this threshold is consistent with her findings. The  $F_1$  score of 0.651, however, is greater than her  $F_1$  result of 0.538 for Jaccard.

Sørensen-Dice was not attempted by AlNoamany. Its best threshold value is close to that of Jaccard at 0.88. In both cases, the document must be quite dissimilar to be marked as off-topic.

Simhash scores range from 0 to 64, based on the number of bits different between simhashes. AlNoamany did not test Simhash. For Simhash based on term frequencies, the best  $F_1$  score is achieved if one uses a threshold value of 28 bits. If just the raw content is fed into the function, the best  $F_1$  score is achieved at 25 bits.

As noted above, cosine similarity scores range from 1 to 0. Very close to AlNoamany’s threshold result of 0.15 for the cosine of TF-IDF vectors, the highest  $F_1$  score is achieved by the OTMT at a threshold of 0.12. This score is very close to the score of complete dissimilarity between documents. AlNoamany achieved an  $F_1$  score of 0.881 compared to our 0.766.

The cosine similarity of LSI vectors has an  $F_1$  score of 0.711. The LSI algorithm requires that one specify the number of topics into which one should break up the corpus. We tried values of 2, 3, 5, 7,

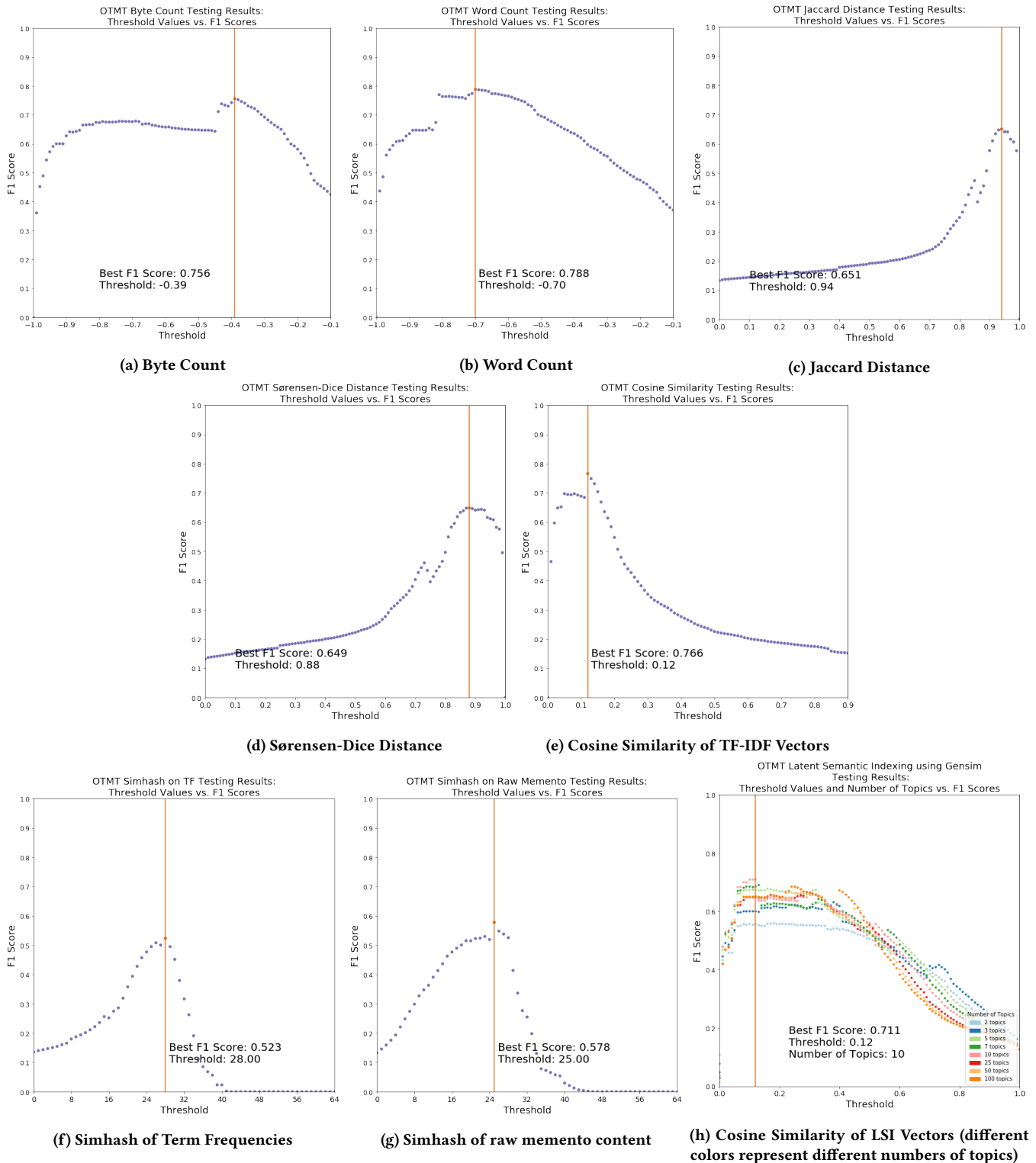


Figure 9: Scatter plots of threshold  $F_1$  testing results for different similarity measures



10, 25, 50, and 100 topics. The value of 10 worked best for testing with our gold standard data. Unfortunately, there is an element of randomness in the results produced by LSI in gensim [27]. Five different runs with LSI produced  $F_1$  scores near or at 0.711, but their corresponding thresholds ranged from 0.08 to 0.12. We took the mean of the scores and set the default threshold at 0.10.

The best test scores for each measure are shown in Table 4. AlNoamany’s results are shown for comparison. Word Count has the best  $F_1$  score, followed by Cosine Similarity of TF-IDF Vectors. Byte Count and Cosine Similarity of LSI Vectors score at third and fourth place, respectively.

Can we do better by using different measures together? Structural measures require less time to execute than more semantic measures like cosine similarity. If we can short-circuit the process with a structural measure then we can eliminate those off-topic mementos prior to review with a more time-intensive measure. AlNoamany did so and found that Word Count and Cosine Similarity worked best. The OTMT accepts multiple measures and evaluates their results as a *logical or*. If at least one of the measures scores a memento as off-topic, then that memento is marked as off topic. We tried different combinations of thresholds just as before and recorded the corresponding  $F_1$  and accuracy scores. Table 5 displays these combinations. The cosine of LSI vectors scores a slightly higher  $F_1$  and the same accuracy as word count. The word count score appears to exert more influence than its partner measures in all cases where it is present, making both cosine measures require stricter thresholds to be successful. It does not appear that combining these measures improves the  $F_1$  score in the OTMT.

## 6 FUTURE WORK

We intend to improve the OTMT over time. For example, we intend to explore making LSI scores reproducible by establishing a specific random number generator seed [27]. We have also considered additional TimeMap measures like Spamsum [21], which works like Simhash and was used to demonstrate memento content drift by Jackson [16]. The OTMT also includes an experimental implementation of cosine of Latent Dirichlet Allocation (LDA) [7] vectors from the Gensim library [32]. This implementation generates errors for some TimeMaps possibly due to a mismatch between the number of topics and the number of features generated by gensim. We have not yet addressed this issue and do not recommend the use of this measure at this time.

Where the previous list of measures were run against all mementos in a TimeMap, it is also conceivable that one can compare each memento in a collection against the collection as a whole. The OTMT does not yet support any measures with this concept. We do envision that such measures could be easily introduced to the OTMT using its existing input-measure-output architecture.

We anticipate the need for curators to have finer grained control over removing boilerplate, stemming, stop word removal, and tokenization. It would also be useful for curators to be able to select different boilerplate removal libraries. We selected justext based on a survey of boilerplate removal methods [26], but curators may find that other boilerplate removal libraries work better for their use cases.

## 7 CONCLUSION

For researchers, identifying off-topic mementos is an important first step to analyzing any web archive collection corpus. We have implemented the Off-Topic Memento Toolkit (OTMT) version 1.0.0 alpha as a way for researchers to identify mementos that are off-topic so that these mementos can be excluded from downstream processing. Because different collections have different needs, we have provided the following similarity measures for use in detecting off-topic mementos:

- byte count
- word count
- Jaccard distance
- Sørensen-Dice distance
- Simhash of term frequencies
- Simhash of raw memento content
- cosine similarity of TF-IDF vectors
- cosine similarity of Latent Semantic Indexing (LSI) vectors

Using a gold standard dataset from a prior study, we have evaluated each measure in terms of effectiveness at determining whether a memento is off-topic. We iterated through many threshold values for each measure, and recorded the  $F_1$  scores at each threshold. We discovered that word count has the best  $F_1$  score, followed by cosine of TF-IDF vectors. Combining the measures did not improve the result, as suggested by prior work.

We present the Off-Topic Memento Toolkit to the world in the hopes that it will help web archive collection curators save time and resources by identifying off-topic mementos.

## ACKNOWLEDGMENTS

This work has been supported in part by the Institute of Museum and Library Services (LG-71-15-0077-15) and the Andrew Mellon Foundation through the Columbia University Libraries Web Archiving Incentive program.

## REFERENCES

- [1] Eytan Adar, Jaime Teevan, Susan T. Dumais, and Jonathan L. Elsas. 2009. The Web Changes Everything: Understanding the Dynamics of Web Content. In *Proc. of the 2nd ACM Int. Conf. on Web Search and Data Mining (WSDM '09)*. ACM, Barcelona, Spain, 282–291. <https://doi.org/10.1145/1498759.1498837>
- [2] Yasmin AlNoamany, Michele C. Weigle, and Michael L. Nelson. 2016. Detecting off-topic pages within TimeMaps in Web archives. *International Journal on Digital Libraries* 17, 3 (2016), 203–221. <https://doi.org/10.1007/s00799-016-0183-5>
- [3] Yasmin AlNoamany, Michele C. Weigle, and Michael L. Nelson. 2017. Generating Stories From Archived Collections. In *Proc. 2017 ACM on Web Science Conference (WebSci '17)*. ACM, Troy, New York, USA, 309–318. <https://doi.org/10.1145/3091478.3091508>
- [4] Ahmed AlSum and Michael L. Nelson. 2014. Thumbnail Summarization Techniques for Web Archives. In *Proceedings of the 36th European Conference on IR Research on Advances in Information Retrieval*. 299–310. [https://doi.org/10.1007/978-3-319-06028-6\\_25](https://doi.org/10.1007/978-3-319-06028-6_25)
- [5] Mohamed Aturban. 2016. 2016-11-05: Pro-Gaddafi Digital Newspapers Disappeared from the Live Web! <http://ws-dl.blogspot.com/2016/11/2016-11-05-pro-gaddafi-digital.html>.
- [6] Steven Bird, Ewan Klein, and Edward Loper. 2009. *Natural Language Processing with Python*. O’Reilly Media. ISBN: 978-0596516499.
- [7] David M. Blei, Andrew Y. Ng, and Michael I. Jordan. 2003. Latent Dirichlet Allocation. *J. Mach. Learn. Res.* 3 (March 2003), 993–1022. <http://dl.acm.org/citation.cfm?id=944919.944937>
- [8] Moses S. Charikar. 2002. Similarity Estimation Techniques from Rounding Algorithms. In *Proc. 34th ACM Symposium on Theory of Computing (STOC '02)*. ACM, Montreal, Quebec, Canada, 380–388. <https://doi.org/10.1145/509907.509965>
- [9] W. Bruce Croft, Donald Metzler, and Trevor Strohman. 2015. *Information Retrieval in Practice*. Pearson Education. ISBN: 978-0136072249.

**Table 4: OTMT TimeMap measures sorted by best  $F_1$  score and compared with AlNoamany’s results**

| Similarity Measure                  | AlNoamany’s results |                        |                         | Results of this study |                        |                         |
|-------------------------------------|---------------------|------------------------|-------------------------|-----------------------|------------------------|-------------------------|
|                                     | Best $F_1$ Score    | Corresponding Accuracy | Corresponding Threshold | Best $F_1$ Score      | Corresponding Accuracy | Corresponding Threshold |
| Word Count                          | 0.806               | 0.982                  | -0.85                   | 0.788                 | 0.971                  | -0.70                   |
| Cosine Similarity of TF-IDF Vectors | 0.881               | 0.983                  | 0.15                    | 0.766                 | 0.965                  | 0.12                    |
| Byte Count                          | 0.584               | 0.962                  | -0.65                   | 0.756                 | 0.965                  | -0.39                   |
| Cosine Similarity of LSI Vectors    | Not tested          |                        |                         | 0.711                 | 0.965                  | 0.10 with 10 topics †   |
| Jaccard Distance                    | 0.538               | 0.962                  | 0.95*                   | 0.651                 | 0.953                  | 0.94                    |
| Sørensen-Dice Distance              | Not tested          |                        |                         | 0.649                 | 0.953                  | 0.88                    |
| Simhash on raw memento content      | Not tested          |                        |                         | 0.578                 | 0.934                  | 25                      |
| Simhash on TF                       | Not tested          |                        |                         | 0.523                 | 0.942                  | 28                      |

\* A derived value is shown for easier comparison. AlNoamany’s threshold was 0.05, but she used Jaccard Index rather than Distance.

† LSI is non-deterministic. This threshold value is a mean of several runs.

**Table 5: The top 4 scoring measures combined in groups of 2**

| Measure                         | Best $F_1$ Score | Corresponding Thresholds | Corresponding Accuracy |
|---------------------------------|------------------|--------------------------|------------------------|
| Cosine of LSI, Word Count       | 0.789            | (0.01, -0.70)            | 0.971                  |
| Cosine of TF-IDF, Word Count    | 0.788            | (0, -0.70)               | 0.971                  |
| Word Count, Byte Count          | 0.788            | (-0.70, -0.94)           | 0.971                  |
| Cosine of TF-IDF, Byte Count    | 0.766            | (0.12, -0.95)            | 0.965                  |
| Cosine of LSI, Cosine of TF-IDF | 0.766            | (0.12, 0.12)             | 0.965                  |
| Cosine of LSI, Byte Count       | 0.759            | (0.01, -0.39)            | 0.965                  |

[10] Renata Gonçalves Curty and Ping Zhang. 2011. Social commerce: Looking back and forward. *Proc. of the Am. Soc. for Inform. Sci. and Tech.* 48, 1 (2011), 1–10. <https://doi.org/10.1002/meet.2011.14504801096>

[11] Scott Deerwester, Susan T. Dumais, George W. Furnas, Thomas K. Landauer, and Richard Harshman. 1990. Indexing by latent semantic analysis. *J. of the Am. Soc. for Info. Science* 41, 6 (1990), 391–407. [https://doi.org/10.1002/\(SICI\)1097-4571\(199009\)41:6<391::AID-AS11>3.0.CO;2-9](https://doi.org/10.1002/(SICI)1097-4571(199009)41:6<391::AID-AS11>3.0.CO;2-9)

[12] Lee R. Dice. 1945. Measures of the Amount of Ecologic Association Between Species. *Ecology* 26, 3 (1945), 297–302. <https://doi.org/10.2307/1932409>

[13] Hannaneh Hajishirzi, Wen-tau Yih, and Aleksander Kolcz. 2010. Adaptive Near-duplicate Detection via Similarity Learning. In *Proc. of the 33rd International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR '10)*. ACM, Geneva, Switzerland, 419–426. <https://doi.org/10.1145/1835449.1835520>

[14] ISO 28500:2017 2017. *Information and documentation – WARC file format*. Standard. International Organization for Standardization, Geneva, Switzerland.

[15] Paul Jaccard. 1912. The Distribution Of The Flora In The Alpine Zone. *New Phytologist* 11, 2 (1912), 37–50. <https://doi.org/10.1111/j.1469-8137.1912.tb05611.x>

[16] Andrew Jackson. 2015. Ten years of the UK web archive: what have we saved? <http://anjackson.net/2015/04/27/what-have-we-saved-ipc-ga-2015/>.

[17] Karen Sparck Jones. 1972. A Statistical Interpretation Of Term Specificity And Its Application In Retrieval. *J. of Documentation* 28, 1 (1972), 11–21. <https://doi.org/10.1108/eb026526>

[18] Shawn M. Jones, Herbert Van de Sompel, and Michael L. Nelson. 2016. 2016-04-27: Mementos in the Raw. <http://ws-dl.blogspot.com/2016/04/2016-04-27-mementos-in-raw.html>.

[19] Shawn M. Jones, Herbert Van de Sompel, Harihar Shankar, Martin Klein, Richard Tobin, and Claire Grover. 2016. Scholarly Context Adrift: Three out of Four URI References Lead to Changed Content. *PLOS ONE* 11, 12 (2016), 1–32. <https://doi.org/10.1371/journal.pone.0167475>

[20] Christian Kohlschütter, Peter Fankhauser, and Wolfgang Nejdl. 2010. Boilerplate Detection Using Shallow Text Features. In *Proc. 3rd ACM Int. Conf. on Web Search and Data Mining (WSDM '10)*. ACM, New York, New York, USA, 441–450. <https://doi.org/10.1145/1718487.1718542>

[21] Jesse Kornblum. 2006. Identifying almost identical files using context triggered piecewise hashing. *Digital Investigation* 3, 91 – 97. <https://doi.org/10.1016/j.diin.2006.06.015>

[22] Gurmeet Singh Manku, Arvind Jain, and Anish Das Sarma. 2007. Detecting Near-duplicates for Web Crawling. In *Proc. of the 16th Int. Conf. on World Wide Web (WWW '07)*. ACM, Banff, Alberta, Canada, 141–150. <https://doi.org/10.1145/1242572.1242592>

[23] Marji McClure. 2006. Archive-It 2: Internet Archive Strives to Ensure Preservation and Accessibility. <http://www.econtentmag.com/Articles/News/News-Feature/Archive-It-2-Internet-Archive-Strives-to-Ensure-Preservation-and-Accessibility-18132.htm>. *EContent* (Oct. 2006).

[24] Michaël Meyer. 2013. Distance 0.1.3. <https://pypi.python.org/pypi/Distance/>.

[25] Ian Milligan. 2012. Mining the ‘Internet Graveyard’: Rethinking the Historians’ Toolkit. *J. of the Can. Hist. Assoc.* 23, 2 (Mar. 2012), 21–64. <https://doi.org/10.7202/1015788ar>

[26] Alexander Nwala. 2017. 2017-03-20: A survey of 5 boilerplate removal methods. <http://ws-dl.blogspot.com/2017/03/2017-03-20-survey-of-5-boilerplate.html>.

[27] Tom O’Hara and Radim Rehůřek. 2015. making LSI reproducible across machines. [https://groups.google.com/forum/#!topic/gensim/upiK51Hs\\_Pc](https://groups.google.com/forum/#!topic/gensim/upiK51Hs_Pc).

[28] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. 2011. Scikit-learn: Machine Learning in Python. *J. of Mach. Learn. Res.* 12 (2011), 2825–2830.

[29] Jan Pomikálek. 2011. *Removing Boilerplate and Duplicate Content from Web Corpora*. Ph.D. Dissertation. Masaryk University.

[30] Radim Rehůřek. 2011. *Scalability of Semantic Analysis in Natural Language Processing*. Ph.D. Dissertation. Masaryk University.

[31] Radim Rehůřek. 2018. Similarity Queries. <https://radimrehurek.com/gensim/tut3.html>.

[32] Radim Rehůřek and Petr Sojka. 2010. Software Framework for Topic Modelling with Large Corpora. In *Proc. LREC 2010 Workshop on New Challenges for NLP Frameworks*. ELRA, Valletta, Malta, 45–50. <http://is.muni.cz/publication/884893/en>.

[33] C. J. Van Rijsbergen. 1979. *Information Retrieval*. Butterworth-Heinemann. ISBN:0408709294.

[34] Leon Sim. 2013. A Python Implementation of Simhash Algorithm. <https://leonsim/posts/a-python-implementation-of-simhash-algorithm/>.

[35] P. Sivakumar. 2015. Effectual Web Content Mining using Noise Removal from Web Pages. *Wireless Personal Communications* 84, 1 (01 Sep 2015), 99–121. <https://doi.org/10.1007/s11277-015-2596-7>

[36] Thorvald Julius Sørensen. 1948. *A method of establishing groups of equal amplitude in plant sociology based on similarity of species content and its application to analyses of the vegetation on Danish commons*. København, I kommission hos E. Munksgaard.

[37] Herbert Van de Sompel, Michael L. Nelson, Lyudmila Balakireva, Martin Klein, Shawn M. Jones, and Harihar Shankar. 2016. 2016-08-15: Mementos in the Raw, Take Two. <http://ws-dl.blogspot.com/2016/08/2016-08-15-mementos-in-raw-take-two.html>.

[38] Herber Van de Sompel, Michael L. Nelson, and Robert Sanderson. 2013. RFC 7089: HTTP Framework for Time-Based Access to Resource States – Memento. <https://tools.ietf.org/html/rfc7089>.

[39] Jonathan Zittrain, Kendra Albert, and Lawrence Lessig. 2014. Perma: Scoping and Addressing the Problem of Link and Reference Rot in Legal Citations. *Legal Inform. Mgmt.* 14, 2 (2014), 88–99. <https://doi.org/10.1017/S1472669614000255>