

# A Study of Web Page Understandability for Consumer Health Search

Joao Palotti<sup>1</sup>, MSc; Guido Zuccon<sup>2</sup>; Allan Hanbury<sup>3</sup>

<sup>1</sup>Qatar Computing Research Institute, Doha, Qatar, Email: jpalotti@hbku.edu.qa

<sup>2</sup>Queensland University, Brisbane, Australia, Email: g.zuccon@qut.edu.au

<sup>3</sup>Vienna University of Technology, Vienna, Austria, Email: hanbury@ifs.tuwien.ac.at

## Abstract

**Background:**

**Objective:**

**Methods:**

**Results:** (make sure to include relevant statistics here, such as sample sizes, response rates, P-values or Confidence Intervals. Do not just say "there were differences between the groups").

**Conclusions:**

*OLD ABSTRACT: In this paper, we investigated methods to estimate the understandability of health Web pages and used these to improve the retrieval of information for people seeking health advice on the Web. Understandability plays a key role in ensuring that people accessing health information are capable of gaining insights that can assist them with their health concerns and choices. The access to unclear or misleading information has been shown to negatively impact on the health decisions of the general public. Our investigation considered methods to automatically estimate the understandability of health information in Web pages, and it provided a thorough evaluation of these methods using human assessments as well as an analysis of pre-processing factors affecting understandability estimations, and associated pitfalls. Furthermore, lessons learnt for estimating Web page understandability were applied to the construction of retrieval methods with specific attention to retrieving information understandable by the general public. We found that machine learning techniques were more suitable to estimate health Web page understandability than traditional readability formulas, which are often used as guidelines and benchmarking by health information providers on the Web. Learning to rank effectively exploited these estimates to provide the general public with more understandable search results. These results are important for specialised search services tailored to support the general public in seeking health advice on the Web.*

**KEYWORDS:** Consumer Health Search; Another Keyword

## Introduction

Search engines are concerned with retrieving relevant information to support a user's information seeking task. Commonly, signals about the topicality or aboutness of a piece of information with respect to a query are used to estimate relevance, with other relevance dimensions like understandability, trustworthiness, etc. [1] being relegated to a secondary position, or completely neglected. While this may be a minor problem for many information seeking tasks, there are some specific tasks in which dimensions other than topicality have an important

role in the information seeking and decision making process. The seeking of health information and advice on the Web by the general public is one such task.

A key problem when searching the Web for health information is that this can be too technical, unreliable, generally misleading, and can lead to unfounded escalations and poor decisions [2]. Where correct information exists, it can be hard to find and digest amongst the noise, spam, technicalities, and irrelevant information. In *high-stakes search tasks* such as this, access to poor information can lead to poor decisions which ultimately can have a significant impact on our health and well-

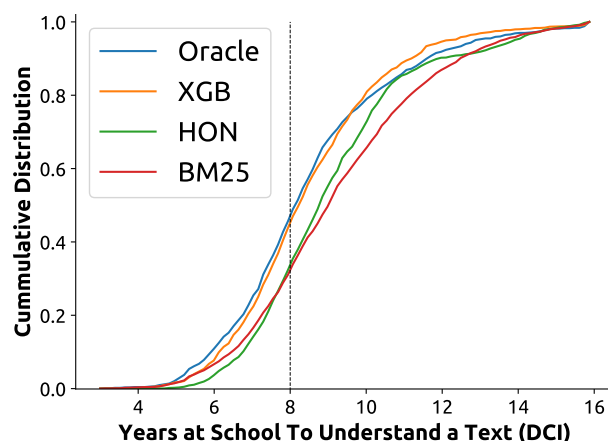
being [2, 3]. In this work we are specifically interested in the understandability of health information retrieved by search engines, and in improving search results to favour information understandable by the general public. We leave addressing reliability and trustworthiness of the retrieved information to future work; however this can be achieved by extending the framework we investigate here.

The use of general purpose Web search engines like Google, Bing and Baidu for seeking health advice has been largely analysed, questioned and criticised [4–10], despite the commendable efforts these services have put into providing increasingly better health information, e.g., the Google Health Cards [11].

Ad-hoc solutions to support the general public in searching and accessing health information on the Web have been implemented, typically supported by government initiatives or medical practitioner associations, e.g., HealthOnNet.org (HON) and HealthDirect.gov.au, among others. These solutions aim to provide *better* health information to the general public. For example, HON's mission statement is “to guide Internet users to reliable, understandable, accessible and trustworthy sources of medical and health information”. But, do the solutions these services currently employ actually provide this type of information to the health-seeking general public? As an illustrative example, we analysed the top 10 search results retrieved by HON<sup>1</sup> in answer to 300 search queries from CLEF 2016 eHealth (see Section ). Figure 1 reports the cumulative distribution of understandability scores for these search results (note, we did not assess their topical relevance here). Understandability scores were computed with the most effective readability formula and settings from Section 0.4 (Dale-Chall Index), and express how easy is to understand a Web page. Low scores correspond to easy to understand Web pages. We report also the scores for the “optimal” search results (Oracle), as found in CLEF 2016 (relevant results that have the highest understandability scores), along with the scores for the best retrieval method from Section 0.4. The results clearly indicate that, despite solutions like HON being explicitly aimed at supporting access to understandable health information, they often fail to do so.

In this article we aim to establish methods and best practice for developing search engines that retrieve *relevant and understandable* health advice from the Web. The overall contributions of this article can be summarized as:

1. We propose and investigate methods for the estimation of the understandability of health information in Web pages: a large number of medically-focused features are grouped in meaningful categories and their contribution to the understandability estimation task is carefully measured;



**Figure 1.** Distribution of Dale-Chall Index (DCI) of search results. DCI measures the years of schooling required to understand a document. The average US resident reads at or below an 8th grade level (dashed line)[14–17], which is the level suggested by NIH for health information on the Web [18]. The distribution for HON is similar to that of the baseline used in this article (BM25). Our best method (XGB) re-ranks documents to provide more understandable results; its distribution is similar to that of an “Oracle” system.

2. We further study the influence of HTML processing methods on these estimations and their pitfalls, extending our previous work that has shown how this often ignored aspect greatly impacts effectiveness [12];
3. We further investigate how understandability estimations can be integrated into retrieval methods to enhance the quality of the retrieved health information with particular attention to its understandability by the general public. New models are explored in this article, also extending our previous work [13];

This paper makes concrete contributions to practice, as it informs health search engines specifically tailored to the general public (for example the HON or HealthDirect services referred to above) about the best methods they should adopt, but they currently don't. These are novel and significant contributions, as no previous work has systematically analysed the influence of the components at play in this study and we show that these greatly influence retrieval effectiveness and thus delivery of relevant and understandable health advice.

## Related Work

Understandability refers to the ease of comprehension of the information presented to a user. Put in other words, health in-

<sup>1</sup>Results retrieved on 01/10/2017.

formation is understandable “when consumers of diverse backgrounds and varying levels of health literacy can process and explain key messages” [19]. Often the terms understandability and readability are used interchangeably: we use readability to refer to formulas that estimate how easy is to understand a text, usually based on its words and sentences. We use understandability to refer to the broader concept of ease of understanding: this is affected by text readability<sup>2</sup>, but may also be influenced by how legible a text is and its layout, including e.g., the use of images to explain difficult concepts.

There is a large body of literature that has examined the understandability of Web health content when the information seeker is a member of the general public. For example, Becker reported that the majority of health Web sites are not well designed for the elderly [21], while Stossel et al. found that health education material on the Web is not written at an adequate reading level [17]. A common finding of these studies is that, in general, health content available on Web pages is often hard to understand by the general public; this includes content that is retrieved in top ranked positions by current commercial search engines [4–9].

Previous Linguistics and Information Retrieval (IR) research has attempted to devise computational methods for the automatic estimation of text readability and understandability, and for the inclusion of these within search methods or their evaluation. Computational approaches to understandability estimations include (1) *readability formulas*, which generally exploit word surface characteristics of the text, (2) *machine learning* approaches, (3) matching with specialised *dictionaries or terminologies*, often compiled with information about understandability difficulty.

Measures such as Coleman-Liau Index (CLI) [22], Dale-Chall Index (DCI) [23] and Flesch Reading Easy (FRE) [24] belong to the first category. These measures generally rely on surface-level characteristics of text, such as characters, syllables and word counts [25]. While these measures have been widely used in studies investigating the understandability of health content retrieved by search engines (e.g., [4–9, 17, 21]), Palotti et al. found that they are heavily affected by the methods used to extract text from the HTML source [12]. They were able to identify specific settings of an HTML preprocessing pipeline that provided consistent estimates. We shall revisit this work in more details in Section 0.4, as we further investigate this problem by comparing the effect of HTML preprocessing on text understandability estimations in light of explicit human assessments.

The use of machine learning to estimate understandability forms an alternative approach. Earlier research explored

the use of statistical natural language processing and language modelling [26–28] as well as linguistic factors, such as syntactic features or lexical cohesion [29]. While we replicated here many of the features devised in these works, they focus on estimating readability of general English documents rather than medical ones. In the medical domain, Zeng et al. explored features such as word frequency in different medical corpora to estimate concept familiarity, which prompted the construction of the Consumer Health Vocabulary (CHV) [30–32].

The actual use of CHV or other terminologies such as the Medical Subject Headings (MeSH) belongs to the third category of approaches. The CHV is a prominent medical vocabulary dedicated to mapping layperson vocabulary to technical terms [31]. It attributes a score for each of its concepts with respect to their difficulty, with lower/higher scores for harder/easier concepts. Researchers have evaluated CHV in tasks such as document analysis [33] and medical expertise prediction [34]. The hierarchy of MeSH was previously used in the literature to identify hard concepts, assuming that a concept deep in the hierarchy is harder than a shallow one [35]. Other approaches combined vocabularies with word surface characteristics and syntactic features, like part of speech, into a unique readability measure [36].

In this work, we investigated approaches to estimate understandability from each of these categories. We further extended Palotti et al.’s work to understand the influence of HTML preprocessing on automatic understandability methods and establish best practices.

Some prior work has attempted to use understandability estimations for improving search results in consumer health search; as well as methods to evaluate retrieval systems that do account for understandability along with topical relevance. Palotti et al. [13] have used learning to rank with standard retrieval features along with features based on readability formulas and medical lexical aspects to determine understandability. Van Doorn et al. [37] have shown that learning a set of rankers that provide trade-offs across a number of relevance criteria, including readability/understandability, increases overall system effectiveness. Zuccon and Koopman [38], and later Zuccon [39], have proposed and investigated a family of measures based on the gain-discount framework, where the gain of a document is influenced by both its topical relevance and its understandability. They showed that, although generally correlated, topical-relevance evaluation alone provides differing system rankings compared to understandability-biased evaluation measures. In this work we further explored the development of retrieval methods that combine signals about topical relevance and understandability.

<sup>2</sup>Increasing readability tends to improve understanding [20].

## Data and Resources

### Data Collections

In this article we investigated methods to estimate Web page understandability, including the effect HTML preprocessing pipelines and heuristics have, and their search effectiveness when employed within retrieval methods. To obtain both topical relevance<sup>3</sup> and understandability assessments, we used the data from the CLEF 2015 and 2016 eHealth collections.

The CLEF 2015 collection contains 50 queries and 1,437 documents that have been assessed relevant by clinical experts and have an assessment for understandability [40]. Documents in this collection are a selected crawl of health Web sites, of which the majority are certified HON Web sites. The CLEF 2016 collection contains 300 queries and 3,298 relevant documents that also have been assessed with respect to understandability [41]. Documents in this collection belong to the ClueWeb12 B13 corpus, and thus are general English Web pages, not necessarily targeted to health topics, nor of a controlled quality (as are instead HON certified pages). Understandability assessments were provided on a 5-point Likert scale for CLEF 2015, and on a [0, 100] range for CLEF 2016 (0 indicates highest understandability).

To support the investigation on methods to automatically estimate the understandability of Web pages, we further considered correlations between multiple human assessors (inter-assessor agreement). For CLEF 2015, we used the publicly available additional assessments made by unpaid medical students and health consumers collected by Palotti et al. [42] in a study of how medical expertise affects assessments. For CLEF 2016 we collected understandability assessments for 100 documents. Three members of our research team, which did not author this paper, were recruited to provide the assessments. The Relevance tool [43] was used to assist with the assessments, mimicking the settings used in CLEF. **should we briefly show these results?**

### Evaluation Measures

In the experiments, we used Pearson, Kendall and Spearman correlations to compare the understandability assessments of human assessors with estimations obtained by the considered approaches, under all combinations of pipelines and heuristics. Pearson correlation is used to calculate the strength of the linear relation between two variables, while Kendall and Spearman measure the rank correlations between the variables. We opted to report all three correlation coefficients to allow for a thorough comparison to other work, as they are equally used in

the literature.

For the retrieval experiments in Section 0.4, we used evaluation measures that use both relevance and understandability. The uRBP measure [39] extends rank biased precision (RBP) to scenarios where multiple relevance dimensions are used. The measure is formulated as  $uRBP(\rho) = (1 - \rho) \sum_{k=1}^K \rho^{k-1} r(d@k)u(d@k)$ , where  $r(d@k)$  is the gain for retrieving a relevant document at rank  $k$  and  $u(d@k)$  is the gain for retrieving a document of a certain understandability at rank  $k$ ;  $\rho$  is the RBP persistence parameter. This measure was an official evaluation measure used in CLEF (we also set  $\rho = 0.8$ ).

A drawback of uRBP is that relevance and understandability are combined into a unique evaluation score, thus making it difficult to interpret whether improvements are due to more understandable or more topical documents being retrieved. To overcome this, we first separately calculated an RBP value for relevance and another for understandability, and then combined them into a unique effectiveness measure:

- $RBP_r@n(\rho)$ : uses the relevance assessments for the top  $n$  search results (i.e. this is the common RBP). We regarded a document as topically relevant if assessed as somewhat relevant or highly relevant.
- $RBP_u@n(\rho)$ : uses the understandability assessments for the top  $n$  search results. We regarded a document as understandable (1) for CLEF 2015 if assessed easy or somewhat easy to understand; (2) for CLEF 2016 if its assessed understandability score was smaller than a threshold  $U$  (we used  $U = 40$ <sup>4</sup>).
- $H_{RBP}@n(\rho) = 2 \times \frac{RBP_r@n \times RBP_u@n}{RBP_r@n + RBP_u@n}$ : combines the previous two RBP values into a unique measurement using the harmonic mean (in the same fashion that the  $F_1$  measure combines recall and precision).

For all measures we set  $n = 10$  because shallow pools were used in CLEF along with measures that focused on the top 10 search results (including  $RBP_r@10$ ).

Along with these measures of search effectiveness, we also reported the number of unassessed documents, the RBP residuals,  $RBP_r@10^*$ ,  $RBP_u@10^*$  and  $H_{RBP}^*$ , i.e. the corresponding measures calculated by ignoring unassessed documents. We did this to minimise pool bias since the pools built in CLEF were of limited size, and the investigated methods retrieved a substantial number of unassessed documents.

<sup>3</sup>We refer to this simply as relevance in the reminder of the paper, when this does not cause confusion.

<sup>4</sup>This choice for  $U$  was based on the distribution of understandability assessments. This distribution can be found in the online appendix.

### Preprocessing Pipelines and Heuristics

As part of our study, we investigated the influence the preprocessing of Web pages has on the estimation of understandability, when this is estimated using the methods in [Section 0.4](#). We did so by comparing the combination of a number of preprocessing pipelines, heuristics, and understandability estimation methods with human assessments of Web page understandability. Our experiments extended those by Palotti et al. [12] and provided a much thorough analysis, as they only evaluated surface level readability formulas and did not compare their results against human assessments.

To extract the content of a Web page from the HTML source we tested: BeautifulSoup<sup>5</sup> (*Naive*), which just naively removes HTML tags, Boilerpipe [44] (*Boi*) and Jstext [45] (*Jst*), which eliminates boilerplate text together with HTML tags. Palotti et al.'s data analysis highlighted that the text in HTML fields like titles, menus, tables and lists often missed a correct punctuation mark and thus the text extracted from them could be interpreted as many short sentences or few very long sentences, depending on whether a period was forced at the end of fields/sentences. We thus implemented the same two heuristics devised by Palotti et al. to deal with this: *ForcePeriod (FP)* and *DoNotForcePeriod (DNFP)*. The FP heuristic forces a period at the end of each extracted HTML field, while the DNFP does not.

### Additional Resources

Because of space limitations, in this article we only reported a subset of the results; the remaining results (which show similar trends to those reported here) are made available in an online appendix for completeness: <https://sites.google.com/view/www2018-sub341>. All data and code will be shared on GitHub upon acceptance. [Check what can we integrate from the google site and put here. Modify google site page to something that has no journal name.](#)

## Understandability Estimators

As reviewed in [Section](#), several methods have been used to estimate the understandability of health Web pages, with the most popular methods (at least in the biomedical literature) being readability formulas based on surface level characteristics of text. Next, we outline the categories of methods to estimate understandability used in this work; an overview is shown in [Table 2](#). Some of these methods further expand measures used in the literature.

**Traditional Readability Formulas (RF):** These include the

readability formulas mentioned in [Section](#), as well as other, less popular ones. A full list is provided in surveys by Collins-Thompson [46] and Dubay [25].

**Raw Components of Readability Formulas (CRF):** These are formed by the “building blocks” used in the traditional readability formulas; examples of such building blocks include the average number of characters per word and the average number of syllables in a sentence. Words are divided into syllables using the Python package Pyphen [47].

**General Medical Vocabularies (GMV):** These include methods that count the number of words with a medical prefix or suffix, i.e. beginning or ending with Latin or Greek particles (e.g., amni-, angi-, algia-, arteri-), and text strings included in lists of acronyms or in medical vocabularies such as the International Statistical Classification of Diseases and Related Health Problems (ICD), Drugbank and the OpenMedSpel dictionary<sup>6</sup>. An acronym list from the ADAM database [48] was used. Methods in this category were matched with documents using simple keywords matching. An acronym list from the ADAM database [48] was used and methods in this group were matched with documents using simple keywords matching.

**Consumer Medical Vocabulary (CMV):** The popular MetaMap [49] tool was used to map the text content of Web pages to entries in CHV [31]. We used the MetaMap semantic types to retain only concepts identified as symptoms or diseases. Similar approaches have been commonly used in the literature (e.g., [50–53]).

**Expert Medical Vocabulary (EMV):** Similarly to the CHV features, we used MetaMap to convert the content of Web pages into MeSH entities, studying symptom and disease concepts separately.

**Natural Language Features (NLF):** These included commonly used natural language heuristics such as the ratio of part-of-speech (POS) classes, the height of the POS parser tree, the number of entities in the text, the sentiment polarity and the ratio of words found in English vocabularies. The Python package NLTK<sup>7</sup> was employed for sentiment analysis, POS tagging and entity recognition. The GNU Aspell<sup>8</sup> dictionary was used as a standard English vocabulary and a stop word list was built by merging those of Indri [54] and Terrier [55]. Discourse features, such as the distribution of POS classes and density of entity in a text, were previously studied in the task of understandability prediction [56] yielding being superior to complex features such as entity co-reference and entity grid ([57]). To the best of our knowledge, sentiment polarity were never investigated. Our intuition is that laypeople produced content (patient

<sup>6</sup><http://extensions.openoffice.org/en/project/openmedspel-en-us>

<sup>7</sup><http://www.nltk.org/>

<sup>8</sup><http://www.aspell.net/>

<sup>5</sup><https://www.crummy.com/software/BeautifulSoup/>



forums or blogs) might content a larger number of emotional content, whereas scientific publications might not

**HTML Features (HF):** These included the identification of a large number of HTML tags, which were extracted with the Python library BeautifulSoup<sup>9</sup>. The intuition for these features is that Web pages with many images and tables may explain and summarise health content better, thus providing more understandable content to the general public.

**Word Frequency Features (WFF):** Generally speaking, common and known words are usually frequent words, while unknown and obscure words are generally rare. This idea is implemented in readability formulas such as the DCI, which uses a list of common words and counts the number of words that fall outside this list (complex words) [23] and has shown success in other recent approaches [58, 59]. We extended these observations by studying corpus-wide word frequencies. We modelled word frequencies in a corpus in a straightforward manner: we sorted the word frequencies and normalized word rankings such that values close to 100 are attributed to common words and values close to 0 to rare words. Three corpora were analysed to extract word frequencies:

- **Medical Reddit:** Reddit<sup>10</sup> is a Web forum with a sizeable user community which is responsible for generating and moderating its content. Any user can start a discussion or reply to a discussion. This forum is intensively used for health purposes: for example in the Reddit community AskDocs<sup>11</sup>, licensed nurses and doctors (subject to user identity verification) advise help seekers free of charge. We selected six of such communities (medical, AskDocs, AskDoctorSmeeee, Health, WomensHealth, Mens\_Health) and downloaded all user interactions available until September 1st 2017 using the Python library PRAW<sup>12</sup>. In total 43,019 discussions were collected.
- **Medical English Wikipedia:** after obtaining a recent Wikipedia dump<sup>13</sup> (May 1st 2017), we filtered articles to only those containing an Infobox<sup>14</sup> in which at least one of the following words appeared as a property: ICD10, ICD9, DiseasesDB, MeSH, MeSHID, MeshName, MeshNumber, GeneReviewsName, Orphanet, eMedicine, MedlinePlus, drug\_name, Drugs.com, DailyMedID, LOINC. In doing so, we followed the method by Soldaini et al. [60], which favours precision over recall when identifying a health-related article. This resulted in a collection of 11,868

articles.

- **PubMed Central:** PubMed Central<sup>15</sup> is an online database of biomedical literature. We used the collection distributed for the TREC 2014 and 2015 Clinical Decision Support Track [61, 62], consisting of 733,191 articles.

A summary of the statistics of these three collections is reported in Table 1. Unless explicitly stated otherwise, we ignored out of vocabulary words in our feature calculations.

**Machine Learning on Text - Regressors (MLR) and Classifiers (MLC):** These include machine learning methods for estimating Web page understandability. While Collins-Thompson highlighted the promise of estimating understandability using machine learning methods, a challenge is identifying the background corpus to be used for training [46]. To this aim, we used the three corpora detailed above, and assumed understandability labels according to the expected difficulty of documents in these collections:

- **Medical Reddit (label 1):** Documents in this collection are expected to be written in a colloquial style, and thus the easiest to understand. All the conversations are in fact explicitly directed to assist inexperienced health consumers;
- **Medical English Wikipedia (label 2):** Documents in this collection are expected to be less formal than scientific articles, but more formal than a Web forum like Reddit, thus somewhat more difficult to understand;
- **PubMed Central (label 3):** Documents from this collection are expected to be written in a highly formal style, as the target audience are physicians and biomedical researchers.

Models were learnt using all documents from these collections after features were extracted using Latent Semantic Analysis (LSA) with 10 dimensions (empirically set based on document word counts in the three collections). We modelled a classification task as well as a regression task using these collections. Thus, after applying the same LSA transformation to test documents from CLEF, a continuous score was assigned to each document by a regressor<sup>16</sup>, while each classifier assigned the documents to one of the three classes.

## Evaluation of understandability estimators

Using the CLEF eHealth 2015 and 2016 collections, we studied the correlations of methods to estimate Web page under-

<sup>9</sup><https://www.crummy.com/software/BeautifulSoup/>

<sup>10</sup><https://www.reddit.com>

<sup>11</sup><https://www.reddit.com/r/AskDocs/>

<sup>12</sup><https://praw.readthedocs.io/>

<sup>13</sup><https://dumps.wikimedia.org/enwiki/>

<sup>14</sup>A Wikipedia infobox is a structured template that appears on the right of Wikipedia pages summarizing key aspects of articles.

<sup>15</sup><https://www.ncbi.nlm.nih.gov/pmc/>

<sup>16</sup>In principle, regressors should output a continuous value between 1 and 3, but no restrictions are set and potentially any value can be assigned to a document.

**Table 1.** Statistics for the collections used as background models for understandability estimations.

Statistic	Medical Wikipedia	Medical Reddit	PubMed Central
Number of Docs.	11,868	43,019	733,191
Number of Words	10,655,572	11,978,447	144,024,976
Number of Unique Words	467,650	317,106	2,933,167
Avg. Words per Doc.	898.90 $\pm$ 1351.76	278.45 $\pm$ 359.70	227.22 $\pm$ 270.44
Avg. Char per Doc.	5107.81 $\pm$ 7618.57	1258.44 $\pm$ 1659.96	1309.11 $\pm$ 1447.31
Avg. Char per Word	5.68 $\pm$ 3.75	4.52 $\pm$ 3.52	5.76 $\pm$ 3.51

standability (Table 2), compared with human assessments. For each category of understandability estimation, Table 3 reports the methods with highest Pearson, Spearman or Kendall correlations. For each method, we used the best preprocessing settings; a study of the impact of preprocessing is reported in Section 0.4.

Overall, Spearman and Kendall correlations obtained similar results (in terms of which methods exhibited the highest correlations): this was expected as, unlike Pearson, they are both rank-based correlations.

For traditional readability measures, SMOG had the highest correlations for CLEF 2015 and DCI for CLEF 2016, regardless of correlation measure. These results resonated with those obtained for the category of raw components of readability formulas. In fact, the polysyllable words measure, which is the main feature used in SMOG, had the highest correlation for CLEF 2015 among these methods. Similarly, the number of difficult words, which is the main feature used in DCI, had the highest correlation for CLEF 2016 among these methods.

When examining the expert vocabulary category, we found that the number of MeSH concepts obtained the highest correlations with human assessments; however its correlations were significantly lower than those achieved by the best method from the consumer medical vocabularies category, i.e. the scores of CHV concepts. For the natural language category, we found that the number of pronouns, the number of stop words and the number of out of vocabulary words had the highest correlations – and these were even higher than those obtained with MeSH and CHV based methods. In turn, the methods that obtained the highest correlations among the HTML category (counts of P tags and list tags) exhibited overall the lowest correlations compared to methods in the other categories. P tags are used to create paragraphs in a Web page, being thus a rough proxy for text length. Among methods in the word frequency category, the use of Medical Reddit (but also of PubMed) showed the highest correlations, and these were comparable to those obtained by the readability formulas.

Finally, regressors and classifiers exhibited the highest correlations amongst all methods: in this category, the eXtreme Gradient Boosting (XGB) regressor and the multinomial Naive Bayes best correlated with human assessments.

## Evaluation of Preprocessing Pipelines and Heuristics

Results from experiments with different preprocessing pipelines and heuristics are shown in Figure 2 (top: CLEF 2015; bottom: CLEF 2016). For each category of methods and combination of preprocessing and heuristics, we report their variability in terms of Spearman rank correlation with the human assessments<sup>17</sup>. We further report the summary results across all understandability assessment methods and sentence ending heuristics for each of the preprocessing pipelines. Finally, we also report the inter-assessor correlation (last box) when multiple assessors provided judgements about the understandability of Web pages (details about this data in Section ). This provides an indication of the range of variability and subjectiveness when assessing understandability, along with the highest correlation we measured between human assessors.

We first examined the correlations between human assessments and readability formulas. We found that the *Naive* preprocessing resulted in the lowest correlations, regardless of readability formula and heuristics (although *DoNotForcePeriod* performed better than *ForcePeriod*). Using Justext or Boilerplate resulted in higher correlations with human understandability assessments, and the *ForcePeriod* heuristic was shown to be better than *DoNotForcePeriod*. These results confirm the speculations of Palotti et al. [12]: they found these settings to produce lower variances in understandability estimations and thus hypothesised that they were better suited to the task.

Overall, among readability formulas, the best results (highest correlations) were obtained by SMOG and DCI (see also Table 3). Although no single setting outperformed the others in both collections, we found that the use of CLI and FRE with Justext provided the most stable results across the collections, with correlations as high as the best ones in both collections. These results confirmed the advice put forward by Palotti et al. [12], i.e. in general, if using readability measures, then CLI should be preferred, along with an appropriate HTML extrac-

<sup>17</sup>Results for Pearson and Kendall correlations are reported in the online appendix, but showed similar trends.

**Table 2.** Methods used to estimate understandability. \*: raw values were used.  $\diamond$ : values normalised by number of words in a document were used.  $\dagger$ : values normalised by number of sentences in a document were used.

Cat.	Method	Cat.	Method	Cat.	Method
RF	Automated Readability Index (ARI) [63]	WFF	25th percentil English Wikipedia	HF	# of Abbr tags
	Coleman-Liau Index (CLI) [22]		50th percentil English Wikipedia		# of A tags
	Dale Chall Index (DCI) [23]		75th percentil English Wikipedia		# of Blockquote tags
	Flesch-Kincaid Grade Level (FKGL) [24]		Mean Rank English Wikip.		# of Bold tags
	Flesch Reading Ease (FRE) [24]		Mean Rank English Wikip. - Includes OV		# of Cite tags
	Gunning Fog Index (GFI) [64]		25th percentil Medical Reddit		# of Div tags
	Lasbarhetsindex (LIX) [65]		50th percentil Medical Reddit		# of Forms tags
	Simple Measure of Gobbledygook (SMOG) [66]		75th percentil Medical Reddit		# of H1 tags
CRF	# of Characters $\ast\diamond\dagger$	WFF	Mean Rank Medical Reddit	HF	# of H2 tags
	# of Words $\ast\dagger$		Mean Rank Medical Reddit - Includes OV		# of H3 tags
	# of Sentences $\ast\diamond$		25th percentil Pubmed		# of H4 tags
	# of Difficult Words (Dale Chall list [23]) $\ast\diamond\dagger$		50th percentil Pubmed		# of H5 tags
	# of Words Longer than 4 chars $\ast\diamond\dagger$		75th percentil Pubmed		# of H6 tags
	# of Words Longer than 6 chars $\ast\diamond\dagger$		Mean Rank Pubmed		# of Hs (any H above)
	# of Words Longer than 10 chars $\ast\diamond\dagger$		Mean Rank Pubmed - Includes OV		# of Img tags
	# of Words Longer than 13 chars $\ast\diamond\dagger$		25th p. Wikipedia+Reddit+Pubmed		# of Input tags
NL	# of Number of Syllables $\ast\diamond\dagger$	GMV	50th p. Wikipedia+Reddit+Pubmed	MLR	# of Link tags
	# of Polysyllable Words (>3 Syllables) $\ast\diamond\dagger$		75th p. Wikipedia+Reddit+Pubmed		# of DL tags
	# of Entities $\ast\diamond\dagger$		Mean R. Wiki.+Reddit+Pubmed		# of UL tags
	# of verbs $\ast\diamond\dagger$		Mean R. Wiki.+Reddit+Pubmed - w. OV		# of OL tags
	# of nouns $\ast\diamond\dagger$		# of Words with Medical Prefix $\ast\diamond\dagger$		# of List (DL + UL + OL)
	# of pronouns $\ast\diamond\dagger$		# of Words with Medical Suffix $\ast\diamond\dagger$		# of Q tags
	# of adjectives $\ast\diamond\dagger$		# of Acronyms $\ast\diamond\dagger$		# of Scripts tags
	# of adverbs $\ast\diamond\dagger$		# of ICD Concepts $\ast\diamond\dagger$		# of Spans tags
NL	# of adpositions $\ast\diamond\dagger$	CMV	# of Drugbank $\ast\diamond\dagger$	MLC	# of Table tags
	# of conjunctions $\ast\diamond\dagger$		# of Words in medical dict. (OpenMedSpel) $\ast\diamond\dagger$		# of P tags
	# of determiners $\ast\diamond\dagger$		CHV Mean Score for all Concepts $\ast\diamond\dagger$		Linear Regressor
	# of cardinal numbers $\ast\diamond\dagger$		# of CHV Concepts $\ast\diamond\dagger$		Multi-layer Perceptron Regressor
	# of particles or other function words $\ast\diamond\dagger$		CHV Mean Score for Symptom Concepts $\ast\diamond\dagger$		Random Forest Regressor
	# of other POS (foreign words, typos) $\ast\diamond\dagger$		# of CHV Symptom Concepts $\ast\diamond\dagger$		Support Vector Machine Regressor
	# of punctuation $\ast\diamond\dagger$		CHV Mean Score for Disease Concepts $\ast\diamond\dagger$		Gradient Boosting Regressor
	Height of part-of-speech parser tree $\ast\diamond\dagger$		# of CHV Disease Concepts $\ast\diamond\dagger$		Logistic Regression
NL	# of Stopwords $\ast\diamond\dagger$	EMV	# of MeSH Concepts $\ast\diamond\dagger$	MLC	Multi-layer Perceptron Classifier
	# of words not found in Aspell Eng. dict. $\ast\diamond\dagger$		Average Tree of MeSH Concepts $\ast\diamond\dagger$		Random Forest Classifier
	Positive Words $\ast\diamond\dagger$		# of MeSH Symptom Concepts $\ast\diamond\dagger$		Support Vector Machine Classifier
	Negative Words $\ast\diamond\dagger$		Average Tree of MeSH Symptom Concepts $\ast\diamond\dagger$		Multinomial Naive Bayes
	Neutral Words $\ast\diamond\dagger$		# of MeSH Disease Concepts $\ast\diamond\dagger$		Gradient Boosting Classifier
			Average Tree of MeSH Disease Concepts $\ast\diamond\dagger$		

tion pipeline, regardless of the heuristic for sentence ending<sup>18</sup>.

When considering methods beyond those based on readability formulas, we found that the highest correlations were achieved by the regressors (MLR) and classifiers (MLC), independently of the preprocessing method used. There is little difference in terms of effectiveness of methods in these categories, with the exception of regressors on CLEF 2015 that exhibited not negligible variances: while for the Neural Network Regressor the Pearson correlation was 0.44, for the Support Vector Regressor it was only 0.30.

A common trend when comparing preprocessing pipelines is that the Naive pipeline provided the weakest correlations with human assessments for CLEF 2016, regardless of estimation methods and heuristics. This result however was not confirmed for CLEF 2015, where the Naive preprocessing negatively influenced correlations for the readability formula category (RF), but not for other categories, although it was generally associated with larger variances for the correlation coefficients.

<sup>18</sup>We provide detailed plots to compare our results with Palotti's in the online appendix.

## Integrating Understandability into Retrieval

We then investigated how understandability estimations can be integrated into retrieval methods to increase the quality of search results. Specifically, we considered three retrieval methods of differing quality for the initial retrieval. These included the best two runs submitted to each CLEF task, and a plain BM25 baseline (default Terrier parameters:  $b = 0.75$  and  $k_1 = 1.2$ ). As understandability estimators we used the eXtreme Gradient Boosting (XGB) regressor<sup>19</sup>[67], as well as SMOG for CLEF 2015 and DCI for CLEF 2016. These were the best performing approaches from Section 0.4.

<sup>19</sup>For assessed documents, we used 10-fold cross validation, training XGB on 90% of the data, and used its predictions for the remaining 10%. For unassessed documents, we trained XGB on all assessed data, and applied this model to generate predictions. Different machine learning methods and feature selection schemes were experimented with; results are available in the online appendix. XGB was selected because its results were the best, although other methods followed similar trends.



**Table 3.** Methods with the highest correlation per category.

Cat.	CLEF 2015					CLEF 2016				
	Method	Preproc.	Pears.	Spear.	Kend.	Method	Preproc.	Pears.	Spear.	Kend.
<b>RF</b>	SMOG Index	Jst NFP	<b>0.438</b>	<b>0.388</b>	<b>0.286</b>	Dale Chall Index	Jst FP	<b>0.439</b>	0.381	0.264
		Boi FP	0.437	<b>0.382</b>	<b>0.264</b>		Boi FP	<b>0.431</b>	<b>0.379</b>	<b>0.262</b>
<b>CRF</b>	Avg. Num. of Polysyl. Words per Word	Jst FP	<b>0.429</b>	0.364	0.268	Avg. Difficult Words Per Word	Boi FP	<b>0.431</b>	<b>0.379</b>	<b>0.262</b>
	Avg. N. of Polysyl. Words per Sentence	Jst NFP	0.192	<b>0.388</b>	<b>0.286</b>					
<b>GMV</b>	Avg. N. Medical Prefixes per Word	Naive FP	<b>0.314</b>	0.312	0.229	Avg. Prefixes per Sentence	Jst FP	<b>0.263</b>	0.242	0.164
	Number of Medical Prefixes		0.131	<b>0.368</b>	<b>0.272</b>	ICD Concepts Per Sentence	Jst NFP	0.014	<b>0.253</b>	<b>0.172</b>
<b>CMV</b>	CHV Mean Score for all Concepts	Naive FP	<b>0.371</b>	<b>0.314</b>	<b>0.228</b>	CHV Mean Score for all Concepts	Jst FP	<b>0.329</b>	0.313	0.216
						CHV Mean Score for all Concepts	Boi FP	0.329	<b>0.325</b>	<b>0.224</b>
<b>EMV</b>	Number of MeSH Concepts	Naive FP	<b>0.227</b>	<b>0.249</b>	<b>0.178</b>	Number of MeSH Concepts	Boi NFP	<b>0.201</b>	0.166	0.113
						Number of MeSH Disease Concepts		0.179	<b>0.192</b>	<b>0.132</b>
<b>NLF</b>	N. of words not found in Aspell Dict.	Jst NFP	<b>0.351</b>	0.276	0.203	Avg. Stopword Per Word	Boi FP	<b>0.344</b>	0.312	0.213
	Number of Pronouns per Word	Naive FP	0.271	<b>0.441</b>	<b>0.325</b>			0.341	<b>0.364</b>	<b>0.252</b>
<b>HF</b>	Number of P Tags	None	<b>0.219</b>	<b>0.196</b>	<b>0.142</b>	Number of Lists	None	<b>0.114</b>	0.021	0.015
						Number of P Tags		0.110	<b>0.123</b>	<b>0.084</b>
<b>WFF</b>	Mean Rank Medical Reddit - Includes OV	Jst NFP	<b>0.435</b>	0.277	0.197	Mean Rank Medical Reddit	Boi NFP	<b>0.387</b>	0.312	0.214
	25th percentil Pubmed	Jst NFP	0.330	<b>0.347</b>	<b>0.256</b>	50th percentil Medical Reddit	Jst NFP	0.351	<b>0.315</b>	<b>0.216</b>
<b>MLR</b>	eXtreme Gradient Boosting (XGB) Regressor	Boi NFP	<b>0.602</b>	0.394	0.287	eXtreme Gradient Boosting (XGB) Regressor	Jst NFP	<b>0.454</b>	<b>0.373</b>	0.258
	eXtreme Gradient Boosting (XGB) Regressor	Jst FP	0.565	<b>0.438</b>	<b>0.324</b>	Random Forest Regressor	Boi NFP	0.389	0.355	<b>0.264</b>
<b>MLC</b>	Multinomial Naive Bayes	Naive FP	<b>0.573</b>	<b>0.477</b>	<b>0.416</b>	Multinomial Naive Bayes	Jst FP	<b>0.461</b>	<b>0.391</b>	<b>0.318</b>

To integrate understandability estimators into the retrieval process, we first investigated *re-ranking* search results retrieved by the initial runs purely based on the understandability estimations. If all the search results from a run were to be considered, then such a re-ranking method may place at early ranks Web pages highly likely to be understandable, but possibly less likely to be topically relevant. To balance relevance and understandability, we only re-ranked the first  $k$  documents. We explored rank cut-offs  $k = 15, 20, 50$ . Because evaluation was performed with respect to the first  $n = 10$  rank positions, the setting  $k = 15$  provided a conservative re-ranking of search results, while,  $k = 50$  provided a less conservative re-ranking approach. Results are presented in Section 0.5.

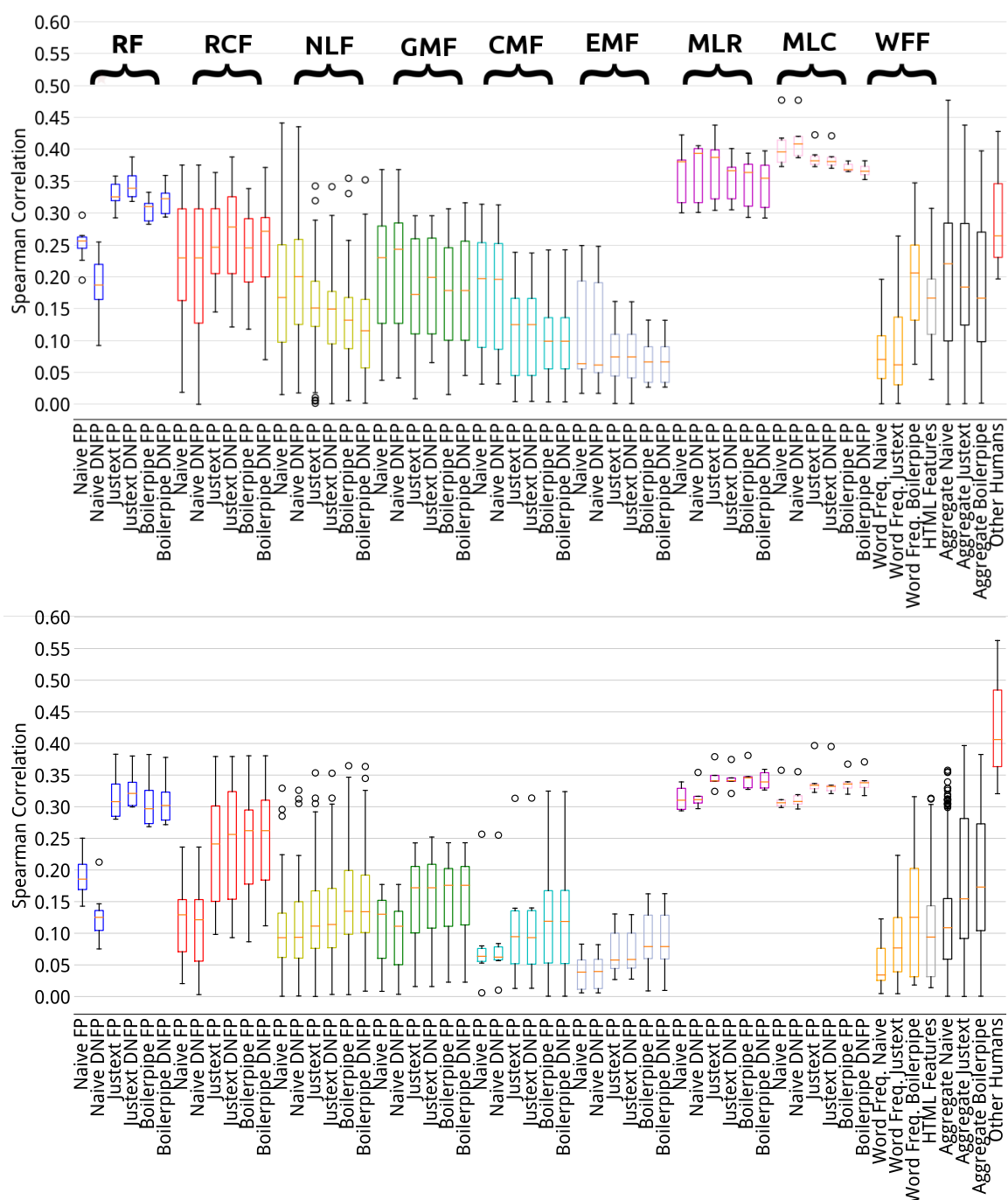
As an alternative to the previous two-step ranking strategy for combining topical relevance and understandability, we explored the *fusion* of two search results lists separately obtained for relevance and understandability. For this, we used the Reciprocal Rank Fusion (RRF) method [68], which was shown effective for combining two lists of search results based on their documents *ranks*, rather than scores. This approach was selected above score-based fusion methods because of the different scoring strategies and distributions employed when scoring for relevance compared to for understandability. For relevance, we used, separately, the three methods used for re-ranking (ECNU [69] and KISTI [70] for CLEF2015, GUIR [71] and ECNU [72] for CLEF 2016, and BM25 for both collections). For understandability, we used, separately, the estimations from SMOG/DCI and XGB. Also for this approach, we studied limiting the ranking of results to be considered by the methods across the cut-offs  $k = 15, 20, 50$ . Results are presented in

**Table 4.** Learning to rank settings.

Name	Feature Set	Labeling Function
Combo 1	IR features	$F(R,U) = R$
Combo 2	IR + Unders. features	$F(R,U) = R$
Combo 3	IR + Unders. features	$F(R,U) = R \times (100 - U)/100$

## Section 0.6.

Finally, we considered a third alternative to combine relevance and understandability: using *learning to rank* with features derived from retrieval methods (IR features) and understandability estimators. With the CLEF 2016 collection, we explored five combinations of label attribution and feature sets, maintaining the same pairwise learning to rank algorithm based on tree boosting (XGB). These combinations are listed in Table 4, with  $R$  being the relevance of documents and  $U$  their understandability estimation. While the definitions of Combo 1 and 2 are straightforward, the other methods deserve some further explanation. In Combo 3, a penalty was proportionally assigned to documents according to how far their understandability score was from a target score  $U$  (we used  $U = 40$ ). For example, a document with understandability 100 received no penalty, as 100 was the easiest level of understanding, while another with understandability 50 received a 50% penalty, meaning that its relevance score was halved. Combo 4 and 5 were based on a fixed threshold applied to the understandability score: if the score was higher than the threshold  $U = 40$ , then the original relevance score (for Combo 4) or a boosted value (for Combo 5) was assigned to the corresponding document.



**Figure 2.** Correlations between understandability estimators and human assessments for CLEF 2015 (top) and 2016 (bottom). For example, the first boxplot in top part of this figure is the result of correlating with Spearman method each one of the features in the category Readability Features (Table 2), obtained with the Naive *ForcePeriod* preprocessing, with human assessments collected during CLEF 2015. Each box extends from the lower to the upper quartile values, with the red marker representing the median value for that category. Whiskers show the range of the data in each category and circles represent values considered outliers for the category (e.g. Spearman correlation for SMOG index was 0.296 and for ARI was 0.194, which were considered outliers for the example category which had a distribution of values highly close to the median value of 0.25).

## Evaluation of Understandability Retrieval

Results for the considered retrieval methods are reported in Table 5. We report only the results for CLEF 2016 for brevity;

those for CLEF 2015 exhibited similar trends and are included in the online appendix. The effectiveness of the top two submissions to CLEF 2016 and the BM25 baseline are reported at indices 1-3 of Table 5. Statistically significant differences compared to the best CLEF 2016 run, GUR, are indicated with  $\diamond$ ;

**Table 5.** Results obtained by integrating understandability estimations within retrieval methods on CLEF 2016. Baseline runs are reported at table indices 1-3. Re-ranking experiments are reported at indices 4-21. Fusion experiments are reported at indices 22-30. Learning to rank experiments are reported at indices 31-35. All measures were calculated up to rank  $n = 10$ .

Index	Rerank	Run	Official CLEF 2016 Metrics				New Metrics to Evaluate Underst. in Retrieval - Sec.									
			$RBP$	Res.	$uRBP$	Res.	$RBP_u$	Res.	$H_{RBP}$	Res.	Unj	$RBP_r^*$	$RBP_u^*$	$H_{RBP}^*$		
1	No Rerank	GUIR [71] (Best Run)	<b>28.11</b>	7.65	<b>18.12</b>	7.19	<b>45.69</b>	8.86	<b>25.61</b>	6.50	0.01	<b>28.29</b>	<b>46.03</b>	<b>25.79</b>		
2		ECNU [72] (Runner Up)	27.70	7.37	17.55	<b>7.34</b>	43.89°	8.66	25.35	6.26	0.01	27.77	44.18°	25.48		
3		Plain BM25 Baseline	25.28°	<b>8.24</b>	16.05°	6.94	42.08°	<b>10.97</b>	22.97°	<b>7.19</b>	<b>0.06</b>	26.01°	43.89°	23.93°		
4	Dale-Chall Top 15	Based on GUIR	24.70†°	8.70	16.83†°	7.27	49.10†°	10.62	24.94	7.50	0.03	25.24†°	50.33†°	25.54		
5		Based on ECNU	24.78†°	7.83	16.64†°	7.16	48.88†°	9.71	24.80	6.50	0.02	25.12†°	49.64†°	25.21		
6		Based on BM25	23.22†°	8.78	15.85°	6.94	47.09†°	11.83	24.01	7.42	0.07	24.04†°	48.60†°	24.82		
7	Dale-Chall Top 20	Based on GUIR	22.19†°	9.37	15.36†°	6.98	48.71†°	12.30	23.21†°	8.12	0.06	23.26†°	51.39†°	24.45†°		
8		Based on ECNU	23.01†°	8.93	15.70†°	6.91	48.99†°	11.69	23.73†°	7.80	0.05	23.84†°	51.00†°	24.66		
9		Based on BM25	21.58†°	9.51	14.83†°	7.02	46.99†°	13.00	22.89°	8.06	0.09	22.93†°	49.55†°	24.26		
10	Dale-Chall Top 50	Based on GUIR	16.18†°	15.24	11.56†°	6.80	41.79†°	24.49	18.10†°	14.42	0.22	20.90†°	53.28†°	23.27†°		
11		Based on ECNU	16.88†°	17.37	11.78†°	<b>7.30</b>	40.76†°	23.77	18.30†°	<b>15.57</b>	<b>0.24</b>	21.34†°	52.07†°	23.33†°		
12		Based on BM25	15.06†°	15.35†°	10.55	6.62	40.03°	23.88	16.55†°	13.83	<b>0.24</b>	19.42†°	51.69†°	21.59†°		
13	XGB Top 15	Based on GUIR	<b>25.16</b> †°	8.09	<b>17.27</b> †°	7.12	<b>50.96</b> †°	10.11	<b>25.16</b>	6.89	0.02	<b>25.61</b> †°	52.00†°	<b>25.68</b>		
14		Based on ECNU	24.18†°	7.69	16.54°	7.09	50.00†°	9.91	24.56	6.65	0.02	24.56†°	50.74†°	25.01		
15		Based on BM25	22.33†°	8.14	15.46	6.76	47.90†°	12.13	22.89°	7.25	0.07	23.11†°	49.43†°	23.69°		
16	XGB Top 20	Based on GUIR	22.38†°	9.49	15.61†°	7.05	50.45†°	12.08	23.30†°	8.16	0.05	23.62†°	52.98†°	24.68		
17		Based on ECNU	22.95†°	8.82	15.95†°	7.02	50.42†°	11.70	23.97°	7.56	0.04	23.68†°	52.15†°	24.73		
18		Based on BM25	20.65†°	9.42	14.46†°	6.84	47.74†°	13.56	21.93°	8.34	0.09	21.98†°	50.28†°	23.27°		
19	XGB Top 50	Based on GUIR	16.65†°	15.73	12.39†°	6.84	43.49†°	23.63	18.70†°	13.74	0.22	21.13†°	<b>55.07</b> †°	23.58†°		
20		Based on ECNU	16.19†°	<b>17.01</b>	11.82†°	7.27	43.05°	<b>24.75</b>	18.27†°	14.41	<b>0.24</b>	20.16†°	54.70†°	22.96†°		
21		Based on BM25	15.43†°	15.37	11.33†°	6.48	41.93°	23.65	17.43†°	13.40	0.26	19.58†°	54.04†°	22.17†°		
22	RRF (XGB & Orig.) Top 15	Based on GUIR	<b>27.23</b> †°	7.76	<b>18.31</b>	<b>7.23</b>	49.69†°	9.18	26.49†°	6.62	0.01	<b>27.46</b> †°	50.07†°	<b>26.69</b> †°		
23		Based on ECNU	26.60†°	7.41	17.81	7.19	48.67†°	8.80	26.02	6.09	0.01	26.76†°	49.10†°	26.27†		
24		Based on BM25	24.57°	8.15	16.51°	6.91	46.76†	11.23	24.16†	7.20	0.06	25.32°	48.52†°	25.08†		
25	RRF (XGB & Orig.) Top 20	Based on GUIR	26.21†°	7.96	17.73	7.19	50.29†°	9.58	25.89	6.73	0.03	26.53†°	50.98†°	26.25		
26		Based on ECNU	26.15†°	7.64	17.69	7.09	49.70†°	9.28	<b>26.07</b>	6.39	0.02	26.38†°	50.32†°	26.35		
27		Based on BM25	24.04†°	8.24	16.32°	6.87	47.69†°	11.40	24.08†°	7.35	0.06	24.82†°	49.52†°	25.01†		
28	RRF (XGB & Orig.) Top 50	Based on GUIR	24.09†°	<b>9.44</b>	16.85†°	7.02	50.55†°	11.76	24.76	<b>8.01</b>	0.07	25.08†°	<b>52.84</b> †°	25.84		
29		Based on ECNU	24.17†°	8.67	16.75°	7.12	<b>50.63</b> †°	11.66	25.00	7.61	0.07	24.90†°	52.50†°	25.84		
30		Based on BM25	22.28†°	8.87	15.50	6.76	48.79†°	<b>12.90</b>	23.13†°	7.82	<b>0.10</b>	23.46†°	51.89†°	24.57		
31	XGB LeToR	Combo 1 on BM25	20.42†°	17.61	13.00†°	7.41	32.17†°	24.61	18.39†°	14.41	0.28	25.25°	43.19°	23.83°		
32		Combo 2 on BM25	24.98†°	19.83	15.30†°	8.09	35.09†°	25.14	22.26°	17.50	0.24	30.41	46.09	28.28†°		
33		Combo 3 on BM25	26.35†	<b>20.48</b>	15.88†°	8.16	34.73†°	24.69	21.81†	17.41	0.22	32.25°	45.44	28.22†°		
34		Combo 4 on BM25	16.16†°	19.48	10.76†°	7.27	<b>36.75</b> †°	<b>28.51</b>	16.77†°	<b>17.80</b>	<b>0.29</b>	22.20†°	<b>50.06</b> †°	23.32°		
35		Combo 5 on BM25	<b>26.76</b> °	<b>20.48</b>	<b>16.19</b> °	<b>8.34</b>	35.26†°	24.13	<b>22.96</b>	17.59	0.22	<b>32.60</b> †	45.87	<b>29.20</b> †°		

differences between an original run (indices 1-3) and its modifications are indicated with † (paired, two-tail t-test,  $p < 0.05$ ). Note that the baseline BM25 is significantly worse than GUIR across all measures.

### Re-ranking

Indices 4-12 of Table 5 report the results of re-ranking methods applied to the runs listed at indices 1-3. Re-ranking was applied based on the DCI score of each document calculated using the preprocessing combination of Boilerpipe and ForcePeriod (best according to Pearson correlation, from Table 3). We found that the relevance of the re-ranked runs (as measured by  $RBP_r$  and  $RBP_r^*$ ) significantly decreased, compared to the original runs: e.g., re-ranking the top 15 search results using DCI made  $RBP_r$  decreasing from 25.28 to 21.58. However, these re-ranked results were significantly more understandable: for the previous example,  $RBP_u$  passed from 42.08 to 47.09.

In the experiments, we also studied the influence of the num-

bers of documents considered for re-ranking (cut-off). Indices 4-6 refer to re-ranking only the top  $k = 15$  documents from the original runs; 7-9 refer to the first  $k = 20$ ; and 10-12 to the first  $k = 50$ . The results show that the more documents are considered for re-ranking, the more degradation in  $RBP_r$  effectiveness. Considering understandability-only in the evaluation shows mixed results. Similar trends were observed for evaluation measures that consider understandability ( $RBP$  and  $RBP_u$ ), however with some exceptions. For example, an increase in  $uRBP$  was observed when re-ranking ECNU using the top 50 results.

Note that with the increase of the number of documents considered for re-ranking, there is an increase in number of unassessed documents being considered by the evaluation measures. Both the RBP residuals and the column  $Unj$  quantify the effect unassessed documents have on evaluation. Nevertheless, we note that if unassessed documents are excluded from the evaluation, similar trends are observed, e.g., compare findings

with those for  $uRBP^*$ ,  $RBP_r^*$ ,  $RBP_u^*$  and  $H_{RBP}^*$ .

Indices 13-21 refer to using the XGB regressor trained with all features listed in Table 2 to estimate understandability. Similarly to when using DCI, as the cut-off increased, e.g., from  $k = 15$  to  $k = 50$ , the documents returned were more understandable but less relevant. For the same cut-off value, e.g.,  $k = 15$ , the machine learning method used for estimating understandability consistently yielded more understandable results than DCI (higher  $RBP_u$  and  $RBP_u^*$ ).

Overall, statistical significant improvements over the baselines were observed for most configurations and measures.

### Rank Fusion

Next, we report the results of automatically combining topical relevance and understandability through rank fusion (indices 22 to 30). We used the XGB method for estimating understandability, as it was the one yielding highest effectiveness for the re-ranking method. Runs were thus produced by fusing the re-ranking with XGB and the original run. (Results for DCI are reported in the online appendix and confirm the superiority of XGB.)

Like as for re-ranking, also for the rank fusion approaches we found that, in general, higher cut-offs were associated to higher effectiveness in terms of understandability measures on one hand, but higher losses in terms of relevance-oriented measures on the other. Overall, results obtained with rank fusion were superior to those obtained with re-ranking only, though most differences were not statistically significant. Statistical significant improvements over the baselines were instead observed for most configurations and measures.

### Learning to Rank

Last, we analyse the results obtained by the learning to rank methods (indices 31-35). Unlike with the previous methods, we did not impose a rank cut-off on learning to rank. Learning to rank was only applied to the BM25 baseline, as we had no access to the IR features for the runs submitted at CLEF (i.e. GUIR and ECNU for CLEF 2016).

When considering  $RBP_r$  and  $uRBP$ , learning to rank exhibited effectiveness that was significantly inferior to that of the GUIR and ECNU baseline runs, though higher than those for the BM25 baseline (for some configurations). The examination of the RBP residuals (and the number of unassessed documents) revealed that this may have been because measures were affected by the large number of unassessed documents retrieved in the top 10 ranks. For example, the  $RBP_r$  residual for learning to rank methods was about double than that of the baselines or other approaches. In fact, among the documents retrieved in the top 10 results by learning to rank, there were

20% that were unassessed, compared to an average of 3% for the other methods. (Excluding XGB with cut-off 50, which also exhibited high residuals).

We thus should carefully account for unassessed documents through considering the residuals of RBP measures, as well as the measures that ignore unassessed documents. When this was done, we observed that learning to rank methods overall provided substantial gains over the original runs and other methods (when considering  $RBP_r^*$ ,  $RBP_u^*$  and  $H_{RBP}^*$ ), or large potential gains over these methods (when considering the residuals). Next, we analyse these results in more detail.

No improvements over the baselines were found for Combo 1 (index 31), and the high residuals for  $RBP_r$  were not matched by other residuals or by considering only assessed documents. Combo 1 was a simple method that used only IR features<sup>20</sup> and was trained only on topical relevance. Although simple, this is a typical learning to rank setting.

Compared to Combo 1, Combo 2 (index 32) included the understandability features listed in Table 2. This inclusion was as beneficial to the understandability measures as to the relevance measures, with  $RBP_r^*$ ,  $RBP_u^*$  and  $H_{RBP}^*$  all showing gains over the baselines. Combo 3 obtained similar  $H_{RBP}^*$  values, though with higher effectiveness for relevance measures ( $RBP_r^*$ ) than for understandability ( $RBP_u^*$ ).

Combos 4 and 5 were devised based on a set understandability threshold  $U = 40$ . While Combo 4 took into consideration only documents that are easy-to-read (understandability label  $\leq U$ ), Combo 5 considered all documents, but boosted the relevance score. Combo 4 reached the highest understandability score for the learning-to-rank approaches ( $RBP_u^* = 50.06$ ), but it failed to retrieve a substantial number of relevant documents ( $RBP_r^* = 22.20$ ). In turn, Combo 5 reached the highest understandability-relevance trade off ( $H_{RBP}^* = 29.20$ ). Compared to the BM25 baseline (on which it was based), Combo 5 largely increased both relevance ( $RBP_r^*$  from 26.01 to 32.60 – a 25% increase) and understandability ( $RBP_u^*$  from 43.89 to 45.87 – a 4% increase). Note that Combo 5 was also significantly better than the best run submitted to CLEF 2016 for both  $RBP_r^*$  (15% increase) and  $H_{RBP}^*$  (13% increase).

## Conclusion

In this paper we have examined approaches to estimate the understandability of health Web pages, including the impact of HTML preprocessing techniques, and how to integrate these within retrieval methods to provide more understandable

<sup>20</sup>We devised 24 IR features using the Terrier framework. The score of various retrieval models was extracted from a multi-field index composed of title, body and whole document.

search results for people seeking health information.

The empirical experiments suggested that:

1. machine learning methods based on regression are best suited to estimate the understandability of health Web pages;
2. preprocessing does affect effectiveness (both for understandability prediction and document retrieval), although, compared to other methods, ML-based methods for understandability estimation are less subject to variability caused by poor preprocessing;
3. learning to rank methods can be specifically trained to promote more understandable search results, while still providing an effective trade off with topical relevance.

This paper makes a clear contribution to improving search engines tailored to consumer health search because it thoroughly investigates promises and pitfalls of understandability estimations and their integration into retrieval methods. The paper further highlights which methods and settings do provide better search results to health information seekers. As shown in Figure 1, these methods would clearly improve current health-focused search engines.

## References

1. Zhang Y, Zhang J, Lease M, Gwizdka J. Multidimensional relevance modeling via psychometrics and crowdsourcing. In: Proceedings of the 37th international ACM SIGIR conference on Research & development in information retrieval. ACM; 2014. p. 435–444.
2. White RW, Horvitz E. Cyberchondria: Studies of the Escalation of Medical Concerns in Web Search. *ACM Transactions on Information Systems*. 2009 Nov;27(4):23:1–23:37.
3. White R. Beliefs and biases in web search. In: Proceedings of the 36th international ACM SIGIR conference on Research and development in information retrieval. SIGIR '13. New York, NY, USA: ACM; 2013. p. 3–12.
4. Graber MA, Roller CM, Kaebler B. Readability levels of patient education material on the World Wide Web. *Journal of Family Practice*. 1999;48(1):58–59.
5. Fitzsimmons P, Michael B, Hulley J, Scott G. A readability assessment of online Parkinson's disease information. *The journal of the Royal College of Physicians of Edinburgh*. 2010;40(4):292–296.
6. Wiener RC, Wiener-Pla R. Literacy, pregnancy and potential oral health changes: The internet and readability levels. *Maternal and child health journal*. 2014;18(3):657–662.
7. Patel CR, Cherla DV, Sanghvi S, Baredes S, Eloy JA. Readability assessment of online thyroid surgery patient education materials. *Head & neck*. 2013;35(10):1421–1425.
8. Atcherson SR, DeLaune AE, Hadden K, Zraick RI, Kelly-Campbell RJ, Minaya CP. A Computer-Based Readability Analysis of Consumer Materials on the American Speech-Language-Hearing Association Website. *Contemporary Issues in Communication Science & Disorders*. 2014;41.
9. Meillier A, Patel S. Readability of Healthcare Literature for Gastroparesis and Evaluation of Medical Terminology in Reading Difficulty. *Gastroenterology Research*. 2017;10(1):1–5.
10. Ellimoottil C, Polcari A, Kadlec A, Gupta G. Readability of websites containing information about prostate cancer treatment options. *The Journal of urology*. 2012;188(6):2171–2176.
11. Gabrilovich E. Cura Te Ipsum: answering symptom queries with question intent. In: Second WebQA workshop, SIGIR 2016 (invited talk); 2016. .
12. Palotti Ja, Zuccon G, Hanbury A. The Influence of Preprocessing on the Estimation of Readability of Web Documents. In: Proceedings of the 24th ACM International on Conference on Information and Knowledge Management. CIKM '15. New York, NY, USA: ACM; 2015. p. 1763–1766.
13. Palotti J, Goeuriot L, Zuccon G, Hanbury A. Ranking health web pages with relevance and understandability. In: Proceedings of the 39th international ACM SIGIR conference on Research and development in information retrieval. ACM; 2016. p. 965–968.
14. Cowan CF. Teaching patients with low literacy skills. *Fusard's Innovative Teaching Strategies in Nursing*. 2004;p. 278.
15. Wallace LS, Lennon ES. American Academy of Family Physicians patient education materials: can patients read them? *Family medicine*. 2004;36(8):571–574.
16. Davis TC, Wolf MS. Health literacy: implications for family medicine. *Family Medicine*. 2004;36(8):595–598.
17. Stossel LM, Segar N, Gliatto P, Fallar R, Karani R. Readability of patient education materials available at



- the point of care. *Journal of general internal medicine*. 2012;27(9):1165–1170.
18. ) NCIUS. Clear & Simple: Developing Effective Print Materials for Low-literate Readers. National Institutes of Health, National Cancer Institute; Accessed: 2017-09. <https://www.nih.gov/institutes-nih/nih-office-director/office-communications-public-liaison/clear-communication/clear-simple>.
  19. Shoemaker SJ, Wolf MS, Brach C. Development of the Patient Education Materials Assessment Tool (PEMAT): a new measure of understandability and actionability for print and audiovisual patient information. *Patient education and counseling*. 2014;96(3):395–403.
  20. Ley P, Florio T. The use of readability formulas in health care. *Psychology, Health & Medicine*. 1996;1(1):7–28.
  21. Becker SA. A study of web usability for older adults seeking online health resources. *ACM Transactions on Computer-Human Interaction (TOCHI)*. 2004;11(4):387–406.
  22. Coleman M, Liao TL. A Computer Readability Formula Designed for Machine Scoring. *Journal of Applied Psychology*. 1975;.
  23. Dale E, Chall JS. A Formula for Predicting Readability: Instructions. *Educational Research Bulletin*. 1948;27(2):37–54.
  24. Kincaid J, Fishburne R, Rogers R, Chissom B. Derivation of New Readability Formulas for Navy Enlisted Personnel. National Technical Information Service; 1975.
  25. Dubay WH. *The Principles of Readability*. Costa Mesa, CA: Impact Information. 2004;.
  26. Liu X, Croft WB, Oh P, Hart D. Automatic Recognition of Reading Levels from User Queries. In: *Proceedings of the 27th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*. SIGIR '04. ACM; 2004. p. 548–549.
  27. Collins-Thompson K, Callan J. Predicting reading difficulty with statistical language models. *Journal of the Association for Information Science and Technology*. 2005;56(13):1448–1462.
  28. Heilman M, Collins-Thompson K, Callan J, Eskenazi M. Combining lexical and grammatical features to improve readability measures for first and second language texts. In: *Human Language Technologies 2007: The Conference of the North American Chapter of the Association for Computational Linguistics; Proceedings of the Main Conference*; 2007. p. 460–467.
  29. Pitler E, Nenkova A. Revisiting readability: A unified framework for predicting text quality. In: *Proceedings of the conference on empirical methods in natural language processing*. Association for Computational Linguistics; 2008. p. 186–195.
  30. Zeng Q, Kim E, Crowell J, Tse T. A text corpora-based estimation of the familiarity of health terminology. *Biological and Medical Data Analysis*. 2005;p. 184–192.
  31. Zeng QT, Tse T. Exploring and developing consumer health vocabularies. *Journal of the American Medical Informatics Association*. 2006;13(1):24–29.
  32. Zeng-Treitler Q, Goryachev S, Tse T, Keselman A, Boxwala A. Estimating consumer familiarity with health terminology: a context-based approach. *Journal of the American Medical Informatics Association*. 2008;15(3):349–356.
  33. Leroy G, Helmreich S, Cowie JR, Miller T, Zheng W. Evaluating online health information: Beyond readability formulas. In: *AMIA Annual Symposium Proceedings*. vol. 2008. American Medical Informatics Association; 2008. p. 394.
  34. Palotti J, Hanbury A, Muller H. Exploiting Health Related Features to Infer User Expertise in the Medical Domain. In: *Proceedings of WSCD Workshop on Web Search and Data Mining*. John Wiley & Sons, Inc.; 2014. .
  35. Yan X, Lau RYK, Song D, Li X, Ma J. Toward a semantic granularity model for domain-specific information retrieval. *ACM Transactions on Information Systems*. 2011 Jul;29(3):15:1–15:46.
  36. Kim H, Goryachev S, Roseblat G, Browne A, Keselman A, Zeng-Treitler Q. Beyond surface characteristics: a new health text-specific readability measurement. In: *AMIA Annual Symposium Proceedings*. vol. 2007. American Medical Informatics Association; 2007. p. 418.
  37. van Doorn J, Odijk D, Roijers DM, de Rijke M. Balancing relevance criteria through multi-objective optimization. In: *Proceedings of the 39th International ACM SIGIR conference on Research and Development in Information Retrieval*. ACM; 2016. p. 769–772.
  38. Zuccon G, Koopman B. Integrating Understandability in the Evaluation of Consumer Health Search Engines. In: *MedIR*; 2014. .

39. Zuccon G. Understandability biased evaluation for information retrieval. In: *European Conference on Information Retrieval*. Springer; 2016. p. 280–292.
40. Palotti J, Zuccon G, Goeuriot L, Kelly L, Hanbury A, Jones GJF, et al. ShARe/CLEF eHealth Evaluation Lab 2015, Task 2: User-centred Health Information Retrieval. In: *Working Notes for CLEF 2015 Conference*, Toulouse, France, September 8–11, 2015.; 2015. .
41. Zuccon G, Palotti J, Goeuriot L, Kelly L, Lupu M, Pecina P, et al. The IR Task at the CLEF eHealth evaluation lab 2016: user-centred health information retrieval. In: *CLEF 2016-Conference and Labs of the Evaluation Forum*. vol. 1609; 2016. p. 15–27.
42. Palotti J, Zuccon G, Bernhardt J, Hanbury A, Goeuriot L. Assessors Agreement: A Case Study across Assessor Type, Payment Levels, Query Variations and Relevance Dimensions. In: *Experimental IR Meets Multilinguality, Multimodality, and Interaction: 7th International Conference of the CLEF Association*, CLEF'16 Proceedings. Springer International Publishing; 2016. .
43. Koopman B, Zuccon G. Relevation!: An open source system for information retrieval relevance assessment. In: *Proceedings of the 37th international ACM SIGIR conference on Research & development in information retrieval*. ACM; 2014. p. 1243–1244.
44. Kohlschütter C, Fankhauser P, Nejdl W. Boilerplate detection using shallow text features. In: *Proceedings of the third ACM international conference on Web search and data mining*. ACM; 2010. p. 441–450.
45. Pomikálek J. Removing Boilerplate and Duplicate Content from Web Corpora; 2011.
46. Collins-Thompson K. Computational assessment of text readability: A survey of current and future research. *ITL-International Journal of Applied Linguistics*. 2014;165(2):97–135.
47. PyPhen. Python module to hyphenate text; 2017. [Online: accessed 21-October-2017]. <http://www.pyphen.org/>.
48. Zhou W, Torvik V, Smalheiser N. ADAM: Another Database of Abbreviations in MEDLINE. *Bioinformatics*. 2006;22(22):2813–2818.
49. Aronson AR, Lang F. An overview of MetaMap: historical perspective and recent advances. *JAMIA*. 2010;17(3):229–236.
50. Pang CI. Understanding Exploratory Search in Seeking Health Information; 2016.
51. Agraftotes C, Arampatzis A. Augmenting Medical Queries with UMLS Concepts via MetaMap. In: *Proceedings of The Twenty-Fifth Text REtrieval Conference, TREC 2016*, Gaithersburg, Maryland, USA, November 15–18, 2016; 2016. .
52. Palotti J, Hanbury A, Müller H, Kahn CE. How users search and what they search for in the medical domain. *Information Retrieval Journal*. 2016 Apr;19(1):189–224.
53. Yates A, Goharian N. ADRTrace: detecting expected and unexpected adverse drug reactions from user reviews on social media sites. In: *European Conference on Information Retrieval*. Springer; 2013. p. 816–819.
54. Strohman T, Metzler D, Turtle H, Croft WB. Indri: A language model-based search engine for complex queries. In: *Proceedings of the International Conference on Intelligent Analysis*. vol. 2. Amherst, MA, USA; 2005. p. 2–6.
55. Ounis I, Amati G, V P, He B, Macdonald C, Johnson. Terrier Information Retrieval Platform. In: *Proceedings of the 27th European Conference on IR Research (ECIR 2005)*. vol. 3408 of Lecture Notes in Computer Science. Springer; 2005. p. 517–519.
56. Feng L, Jansche M, Huenerfauth M, Elhadad N. A comparison of features for automatic readability assessment. In: *Proceedings of the 23rd International Conference on Computational Linguistics: Posters*. Association for Computational Linguistics; 2010. p. 276–284.
57. Barzilay R, Lapata M. Modeling Local Coherence: An Entity-based Approach. *Comput Linguist*. 2008 Mar;34(1):1–34. Available from: <http://dx.doi.org/10.1162/coli.2008.34.1.1>.
58. Elhadad N. Comprehending technical texts: Predicting and defining unfamiliar terms. In: *AMIA annual symposium proceedings*. vol. 2006. American Medical Informatics Association; 2006. p. 239.
59. Wu DT, Hanauer DA, Mei Q, Clark PM, An LC, Proulx J, et al. Assessing the readability of ClinicalTrials. gov. *Journal of the American Medical Informatics Association*. 2015;23(2):269–275.
60. Soldaini L, Cohan A, Yates A, Goharian N, Frieder O. In: *Retrieving Medical Literature for Clinical Decision Support*. Springer International Publishing; 2015. p. 538–549.

61. Roberts K, Simpson M, Demner-Fushman D, Voorhees E, Hersch W. State-of-the-art in biomedical literature retrieval for clinical cases: a survey of the TREC 2014 CDS track. *Information Retrieval Journal*. 2016;19(1):113–148.
62. Roberts K, Simpson MS, Voorhees EM, Hersch WR. Overview of the TREC 2015 Clinical Decision Support Track. In: *Proceedings of The Twenty-Fourth Text REtrieval Conference, TREC 2015, Gaithersburg, Maryland, USA, November 17-20, 2015*; 2015. .
63. Smith EA, Senter RJ. Automated Readability Index. AMRL-TR-66-220. Aerospace Medical Research Laboratories; 1967.
64. Gunning R. *The Technique of Clear Writing*. McGraw-Hill; 1952.
65. Björnsson CH. Readability of Newspapers in 11 Languages. *Reading Research Quarterly*. 1983;18(4):480–497.
66. McLaughlin GH. SMOG Grading - a New Readability Formula. *Journal of Reading*. 1969;.
67. Chen T, Guestrin C. XGBoost: A Scalable Tree Boosting System. In: *Proceedings of the 22Nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining. KDD '16*. New York, NY, USA: ACM; 2016. p. 785–794.
68. Cormack GV, Clarke CLA, Buettcher S. Reciprocal Rank Fusion Outperforms Condorcet and Individual Rank Learning Methods. In: *Proceedings of the 32Nd International ACM SIGIR Conference on Research and Development in Information Retrieval. SIGIR '09*. New York, NY, USA: ACM; 2009. p. 758–759.
69. Song Y, He Y, Hu Q, He L, Haacke EM. ECNU at 2015 eHealth Task 2: User-centred Health Information Retrieval. In: *Working Notes of CLEF 2015 - Conference and Labs of the Evaluation forum, Toulouse, France, September 8-11, 2015*; 2015. .
70. Oh H, Jung Y, Kim K. KISTI at CLEF eHealth 2015 Task 2. In: *Working Notes of CLEF 2015 - Conference and Labs of the Evaluation forum, Toulouse, France, September 8-11, 2015*; 2015. .
71. Soldaini L, Edman W, Goharian N. Team GU-IRLAB at CLEF eHealth 2016: Task 3. In: *Working Notes of CLEF 2016 - Conference and Labs of the Evaluation forum, Évora, Portugal, 5-8 September, 2016*; 2016. p. 143–146.
72. Song Y, He Y, Liu H, Wang Y, Hu Q, He L. ECNU at 2016 eHealth Task 3: Patient-centred Information Retrieval. In: *Working Notes of CLEF 2016 - Conference and Labs of the Evaluation forum, Évora, Portugal, 5-8 September, 2016*; 2016. p. 157–161.