

Report on the Lucene4IR workshop: Developing Information Retrieval Evaluation Resources using Lucene (L4IR2016)

Leif Azzopardi¹, Yashar Moshfeghi², Martin Halvey¹,
Rami S. Alkhawaldeh², Krisztian Balog³, Emanuele Di Buccio⁴,
Diego Ceccarelli⁵, Juan M. Fernández-Luna⁶,
Charlie Hull⁷, Jake Mannix⁸, Sauparna Palchowdhury⁹

¹ University of Strathclyde {*Leif.Azzopardi, Martin.Halvey*}@strath.ac.uk

² University of Glasgow {*Yashar.Moshfeghi, Rami.Alkhawaldeh*}@glasgow.ac.uk

³ University of Stavanger *krisztian.balog@uis.no*

⁴ University of Padova *dibuccio@dei.unipd.it*

⁵ Bloomberg *dceccarelli4@bloomberg.net*

⁶ University of Granada *jmfluna@decsai.ugr.es*

⁷ Flax *charlie@flax.co.uk*

⁸ LucidWorks *jake.mannix@lucidworks.com*

⁹ National Institute of Standards and Technology, USA *sauparna.palchowdhury@nist.gov*

October 28, 2016

Abstract

The workshop and hackathon on developing Information Retrieval Evaluation Resources using Lucene (L4IR) was held on the 8th and 9th of September, 2016 at the University of Strathclyde in Glasgow, UK and funded by the ESF Elias Network. The event featured three main elements: (i) a series of keynote and invited talks on industry, teaching and evaluation; (ii) planning, coding and hacking where a number of groups created modules and infrastructure to use Lucene to undertake TREC based evaluations; and (iii) a number of breakout groups discussing challenges, opportunities and problems in bridging the divide between academia and industry, and how we can use Lucene for teaching and learning Information Retrieval (IR). The event was composed of a mix and blend of academics, experts and students wanting to learn, share and create evaluation resources for the community. The hacking was intense and the discussions lively creating the basis of many useful tools but also raising numerous issues. It was clear that by adopting and contributing to most widely used and supported Open Source IR toolkit, there were many benefits for academics, students, researchers, developers and practitioners - providing a basis for stronger evaluation practices, increased reproducibility, more efficient knowledge transfer, greater collaboration between academia and industry, and shared teaching and training resources.

1 Introduction

Lucene and its expansions, Solr and ElasticSearch, represent the major open source Information Retrieval toolkits used in Industry. However, there is a lack of coherent and coordinated documentation that explains from an experimentalist's point of view how to use Lucene to undertake and perform Information Retrieval Research and Evaluation. In particular, how to undertake and perform TREC based evaluations using Lucene. Consequently, the objective of this event was to bring together researchers and developers to create a set of evaluation resources showing how to use Lucene to perform typical IR operations (i.e. indexing, retrieval, evaluation, analysis, etc.) as well as how to extend, modify and work with Lucene to extract typical statistics, implement typical retrieval models, etc. Over the course of the workshop participants shared their knowledge with each other creating a number of resources and guides along with a road map for future development.

2 Keynotes and Invited Talks

During the course of the workshops a series of talks on how Lucene is being used in Industry, Teaching and for Evaluation along with more technical talks on the inner workings of how Lucene's scoring algorithm works and how learning to rank is being included into Solr, were presented¹. A summary of each talk is below.

Introduction Talk: Why are we here?

Leif Azzopardi, University of Strathclyde: Leif explained how after attending the lively Reproducibility workshop [1] at ACM SIGIR 2015, he wondered where the Lucene team was, and why, if Lucene and the community is so big, why they don't come to IR conferences - he posited that perhaps we haven't been very inclusive or welcoming to such a large community of search practitioners. He further asserted that this has reduced our capacity to transfer our knowledge and experience into one of the largest Open Source toolkits available. He argued that if we as academics want to increase our impact then we need to improve how we transfer our knowledge to industry. One way is working with large search engines, but what about other industries and organisations that need search and use toolkits like Lucene? He argued that we need to start speaking the same language i.e. work with Lucene et al and look for opportunities on how we can contribute and develop resources for training and teaching IR and how to undertake evaluations and data science using widely used, supported and commonly accepted Open Source toolkits. He described how this workshop was a good starting point and opportunity to explore how academia and industry can better work together, where we can identify common goals, needs and resources that are needed to foster this relationship.

Keynote Talk: Apache Lucene in Industry

Charlie Hull, Flax: In his talk, Charlie first introduced Flax, and how it evolved over the years. Charlie explained that they have been building search applications using open search software since 2001. Their focus is on building, tuning and supporting fast, accurate

¹Slides are available from www.github.com/leifos/lucene4ir

and highly scalable search, analytics and Big Data applications. They are partners with Lucidworks, leading Lucene specialists and committers. When Lucene first came out clients were reluctant to adopt open source, but nowadays it is much more acceptable. Charlie notes that now you don't have to explain to clients what open source software is, and why it should be used. He described how Lucene-based search engines have risen in use - and that search and data analytics are available to those without six figure budgets. Charlie points out that Lucene is appealing because it is the most widely used open source search engine, which is hugely flexible, feature rich, scalable and performant. It is supported by a large and healthy community and backed by the Apache Software Foundation. Many of world's largest companies use Lucene including Sony, Siemens, Tesco, Cisco, LinkedIn, Wikipedia, WordPress and Hortonworks. Charlie notes that they typically don't use Lucene directly, instead they use the search servers, built on top of Lucene, i.e. Apache Solr (which is mature, stable, and crucially highly scalable), or Elasticsearch (easy to get started with, great analytics, scalable). He contrasts these products with some of the existing toolkits in IR [2, 7, 10, 14], and remarks on the latter, that "no one in industry has ever heard of them!". So even though they have the latest research encoded within them, it is not really viable for businesses to adopt them, especially as support for such toolkits is highly limited. He recommends that IR research needs to be within Lucene-based search services for it to be used and adopted.

Based on Charlie's experience he provided us with a number of home truths:

- Open source does not mean cheap
- Most search engines are the same (in terms of underlying features and capabilities)
- Complex features are seldom used - and often confusing
- Search testing is rarely comprehensive
- Good search developers are hard to find

Charlie reflected on these points considering how we can do better. First, learn what works in industry and how industry are using search - there are lots of research challenges which they rarely get to solve and address but solutions to such problems would have real practical value. Second, improve Lucene et al with ideas from academia - faster - for example, it took years before BM25 replaced TFIDF as the standard ranking algorithm, where as toolkits like Terrier [11] already have infrastructure for Learning to Rank, while this is only just being developed in Lucene. Third, he pointed out that testing and evaluation of Lucene based search engines is very limited, and that thorough evaluations by search developers is poor. He argued that this could be greatly improved, if academics and researchers, contributed to the development of evaluation infrastructure, and transferred their knowledge to practitioners on how to evaluate. Lastly, he pointed that the lack of skilled and knowledgeable search developers was problematic - having experience with Lucene, Solr and Elasticsearch are highly marketable skills, especially, when there is a growing need to process larger and larger volumes of data - big data requires data scientists! So there is the pressing need to create educational resources and training material for both students and developers.

Using Lucene for Teaching and Learning IR: The University of Granada case of study

Prof. Juan M. Fernández-Luna(University of Granada) : In his talk, Juanma ex-

plained the Bologna Process, the establishment of the European Higher Education Area, and how the University of Granada (UGR) has adopted its study programmes, changing the existing undergraduate and master degrees and introducing a few new ones. Currently, Computer Science studies at UGR are composed of an undergraduate degree of four years and a master degree of one, with three different options in this last case: one professional master in Computer Science and two research masters (Data Science and Computer Engineering, and Software Development). With a lot more attention focused on Information Retrieval, a number of courses have been introduced within their undergraduate and masters courses:

- *Information Retrieval* at the undergraduate Computer Science degree (6 ECTS, 4th year). The objective of this subject is that the students learn the foundation of IR (document preprocessing, indexing, retrieval models, evaluation and text classification and clustering).
- *Information Management in the Web* at Master in Computer Science (4 ECTS). This subject is composed of three parts: social network analysis, IR and recommender systems. The basic aim is to show the students different ways of managing and accessing the information in the Web. The part related to IR is focused just in briefly explaining the IR foundations.
- *Information Retrieval and Recommender Systems* at Master in Data Science and Computer Engineering (3 ECTS). In the context of this master, IR foundations are shown to the students, but with a perspective closer to Data Science (preprocessing of large document collections, clustering and classification, etc.).
- *Undergraduate and Master Thesis* (12 ECTS). Each degree contains a final thesis where the students have to show the skills they have acquired during their studies by means of the development of a project. These are proposed by the lecturers and some of them are related to Information Retrieval.

Juanma explained that before introducing programming details, the lecturers at UGR thought that it would be better if students could understand the core IR process itself. To facilitate the teaching and learning process, *SulaIR* [4] was designed. This is a desktop tool that covers the different IR stages: web crawling, document pre-processing, indexing, retrieval and relevance feedback. The tool lets students interact with all these processes and secure the concepts from a practical point of view.

In any of these IR-related subjects, the same question came about when the lecturers were planning lab work and exercises, is it better to: (1) create IR projects from the scratch, programming even the more basic classes and focusing on the implementation details, or (2) use an existing IR library (Lucene, Terrier, Lemur, MG, etc.) and focusing on the process?

At the very beginning, they opted for the first alternative as they thought that this could help to understand the details of the search engines, but this created two learning challenges. The students had to understand the IR process, in a first abstraction step, and then, to transfer it to code, in a second step. Often the second step interfered with their understanding of the first step as they faced many programming challenges. So the learning process was not very effective.

Therefore, Juanma and the lecturers made the decision of using an Open Source library, where the problem is reduced to learning about an API and identifying the classes and methods that need to be used. In this case, students do not need to care about the implementation details, per se, instead they could focus on the process (or at least this was their theory). In

addition, they realised that the typical professional developer with real needs regarding IR will use an API/Toolkit and there will not be programming IR-related modules from scratch. From all of available APIs, Lucene was, without any doubt, the first choice. This was because it is the most popular IR toolkit openly available and most widely used in industry.

Juanma points out that Teaching and Learning IR with Lucene in these subjects is not without its own problems and challenges. After using it with their courses, the lecturers came to a number of realisations:

- From *Information Retrieval*:
 - Lucene is a “monster”, with lots of classes and methods. This is because it is a large-scale production system, and so students often are frightened by it, unsure of how to work with it.
 - Considering Java as the programming language to work with Lucene, our students usually have to learn this language first, as they are used to work with C++ in the degree. This language could be a possibility but the Lucene C++ API is poorly documented in comparison to its Java version, so it was discarded.
 - There is a large amount of documentation about how basic tasks are carried out with Lucene, but sometimes the students select sources from different versions of Lucene. And this causes many headaches when understanding and debugging the problems faced when programming with Lucene.
 - The students’ learning curve is very steep when working with the basic process. However, when they progress to more advanced topics such as fielded and complex queries, they find it extremely difficult to progress and really struggle.
- From *Information Management in the Web*:
 - Due to the length of the course, the ABCs of IR with Lucene are shown to the students. They consider this a good approach because in case of needing to build a search engine in their professional career, they have the basic knowledge to use a toolkit to configure a search engine.
- From *Information Retrieval and Recommender Systems*:
 - In the context of a Data Science Master, students are not so interested in the retrieval process itself, but in the pre-processing and indexing stages as bases of further tasks related to machine learning. Other toolkits, as Tika, are shown.
 - Students have different programming backgrounds (Computer Science, Statistics, Mathematics, Information Science, Electronics, etc.), so it is a real problem to use Lucene for developing lab projects. Considering that R is the programming language on which most of the subjects from the Data Science master are based, an initiative such as RLucene² could be very interesting for the students.
- From *Undergraduate and Master Thesis*:
 - It is difficult to find Lucene advanced documentation for more specific topics, so they usually spend too much time trying to work out how such functionality works (usually through trial and error). In addition, it is very difficult for them to develop new retrieval models or ranking functions.

²<https://github.com/s-u/RLucene>

-
- As Lucene has not got native support for documents, they have to make great efforts to build the piece of software required to extract the text of the real documents that they find in their projects.
 - There is no support for IR evaluation in Lucene, so students have to write their own tools to evaluate the performance.

As a conclusion, Juanma and the other lecturers at UGR recognised that Lucene is a great tool for teaching and learning IR, but there was scope for improvement:

- Lucene documentation or tutorials from the point of view of teaching and learning were available as well as material describing advanced tasks.
- Visual and/or Command-Line Tools for teaching and learning the IR process based on Lucene, a kind of SulaIR-L, would really be very useful.
- Visual and/or Command-Line Tools for IR evaluation (TREC-based) were available, or at least, classes for these purposes were included in the API.

Black Boxes are Harmful

Sauparna Palchowdhury (National Institute of Standards and Technology, Gaithersburg, Maryland, USA): Having seen students and practitioners in the IR community grapple with abstruse documentation accompanying search systems and their use as a black box, Sauparna, in his talk, argued why Lucene is a useful alternative and how and why we must ensure it does not become another black box. In establishing his views, he described the pitfalls in an IR experiment and the ways of mitigation. The suggestions he put forth, as a set of best practices, highlighted the importance of evaluation in IR to render an experiment reproducible and repeatable and the need for a well-documented system with correct implementations of search algorithms that are traceable to a source in IR literature. In the absence of such constraints on experimentation students are misled and learn little from the results of their experiments and it becomes hard to reproduce the experiments. As an example, the talk cited a wrong implementation of the *Okapi BM25* term-weighting equation in a popular research retrieval system (Table 1). Following this was a brief how-to on implementing *BM25* (or any $TF \times IDF$ weighting scheme) in Lucene (Table 2). This also explained Lucene’s way of computing the similarity between two text documents (usually referred to as *Lucene’s scoring formula*³).

Some of the points of failure mentioned in the talk were misplaced test-collection pieces (document-query-qrel triplet), counterintuitive configuration interfaces of systems, poor documentation that make systems look enigmatic and lead to the creation of heuristics passed around by word-of-mouth, naming confusion (a myriad of $TF \times IDF$ model names), blatant bugs and a obscure parser. As mitigation, Sauparna listed some of the things he did as an experimenter. He wrote a script (TRECBOX⁴) to abstract parts of the IR experiment pipeline and map them to configuration end-points of the three systems; Indri [13], Terrier [11], and Lucene [5]. This would enable documenting and sharing an experiment’s design in plain text files. He constructed a survey of term-weighting equations titled *TF×IDF Repository*⁵ meant to be a single point of reference to help disambiguate the variants in the wild. All

³<https://goo.gl/ZOMVYe>

⁴<https://github.com/sauparna/TRECBOX>

⁵<http://kak.tx0.org/IR/TFxIDF>

equations mentioned in this repository are traceable to a source in IR literature. He also showed how to visually juxtapose evaluation results obtained using a permutation of a set of systems, retrieval models and test-collections on a chart that would act as a sanity check for the system’s integrity. As a part of these investigations he modified Lucene for use with TREC collections (the mod was named LTR⁶) which is available for others to use. The “mod” is also accompanied by notes to augment Lucene’s documentation. The gamut of Sauparna’s work is collected on a website⁷.

Lucene’s documentation does not use a well-defined notation to represent its way of computing the similarity score between a query Q and document D . Equation (1) denotes Lucene’s scoring formula as described in Lucene’s documentation. In the equation, T denotes a term. The functions, in order from left to right, on the right-hand-side of the equation is the *coordination factor*, *query normalization factor*, *term-frequency transformation*, *document-frequency transformation*, *query boost* and *document-length normalization factor*. A well-defined, generalized, notation for Lucene’s scoring, in step with the definition from Lucene’s documentation, is Equation (2) (function names were shortened appropriately).

$$score(Q, D) = coord(Q, D) \cdot qnorm(Q) \cdot \sum_{T \in Q} (tf(T \in D) \cdot idf(T)^2 \cdot boost(T) \cdot norm(T, D)) \quad (1)$$

$$score(Q, D) = f_c(Q, D) \cdot f_q(Q) \cdot \sum_{T \in Q \cap D} (tf(T) \cdot df(T) \cdot f_b(T) \cdot f_n(T, D)) \quad (2)$$

To explain Lucene’s scoring, Sauparna picked two popular $TF \times IDF$ variants, broke them down into meaningful components (a term-frequency transformation, a document-frequency transformation and a length-normalization coefficient) and plugged these components into Lucene’s equation. The components in Lucene’s equation that were left unused were replaced by the integer 1, meaning, the functions returned 1; which would have no effect on the chain of multiplications. Table 1 lists the variants and components and Table 2 shows where the components were transplanted to.

Making a reference to the SIGIR 2012 tutorial on *Experimental Methods for Information Retrieval* [9], Sauparna stated that we need to take a more rigorous approach to the IR experimental methodology. A list of best practices were recommended that would add more structure to IR experiments and prevent the use of systems as black boxes. These were:

1. Record test-collection statistics.
2. Provide design documentation for systems.
3. Use a consistent naming scheme and a well-defined notation.
4. Use an evaluation table as a sanity check.
5. Isolate shareable experimental artifacts.
6. Ensure that implementations are traceable to a source in IR literature.

In conclusion, Sauparna suggested that if we, the IR research community, were to build and work with Lucene, it would be helpful to consider these points when introducing new features into Lucene.

⁶<https://github.com/sauparna/LTR>

⁷<http://kak.tx0.org/IR>

TF×IDF Variants: What’s correct and what’s not.

| Name | w_{ik} | w_{jk} |
|----------------------|--|------------------------------------|
| <i>BM25</i> (A) | $\frac{f_{ik}}{k_1((1-b)+b\frac{dl_i}{avdl})+f_{ik}} \times \log(\frac{N-n_k+0.5}{n_k+0.5})$ | $\frac{(k_3+1)f_{jk}}{k_3+f_{jk}}$ |
| <i>BM25</i> (B) | $\frac{(k_1+1)f_{ik}}{k_1((1-b)+b\frac{dl_i}{avdl})+2f_{ik}} \times \log(\frac{N-n_k+0.5}{n_k+0.5})$ | $\frac{(k_3+1)f_{jk}}{k_3+f_{jk}}$ |
| <i>Okapi BM25</i> | $\frac{(k_1+1)f_{ik}}{k_1((1-b)+b\frac{dl_i}{avdl})+f_{ik}} \times \log(\frac{N-n_k+0.5}{n_k+0.5})$ | $\frac{(k_3+1)f_{jk}}{k_3+f_{jk}}$ |
| components | $TF \times DF$ | QTF |
| <i>SMART dtb.nnn</i> | $\frac{(1+\log(1+\log(f_{ik}))) \times \log(\frac{N+1}{n_k})}{1-s+s \cdot \frac{b_i}{avgb}}$ | f_{jk} |
| components | $TF \times DF \div LN$ | QTF |

Table 1: The similarity score; $score(D_i, D_j) = \sum_{k=1}^t (w_{ik} \cdot w_{jk})$ $\forall i \neq j$, combines the weight of a term k over the t terms which occur in document D_i and D_j . Since a query can also be thought of as a document in the same vector space, the symbol D_j denotes a query. *BM25*(A) and *BM25*(B) are the two incorrect implementations found in a popular retrieval system. Comparing them to *Okapi BM25* on the third row shows that A has the $k_1 + 1$ factor missing in the numerator, and B uses twice the term-frequency, $2f_{ik}$, in the denominator. Neither can they be traced to any source in IR literature, nor does the system’s documentation say anything about them. The *Okapi BM25* and the *SMART dtb.nnn* variants are known to be effective formulations developed by trial and error over eight years of experimentation at TREC 1 through 8. Their forms have been abstracted using the abbreviations TF , DF , LN and QTF (term-frequency, document-frequency, length-normalization and query-term-frequency) to show how these components fit in Lucene’s term-weight expression.

| Implementing TF×IDF variants in Lucene | | | | | | | |
|--|-------------|---|----------|---|---------------------------|---------|--------------------------------------|
| Lucene | $f_c(Q, D)$ | · | $f_q(Q)$ | · | $\sum_{T \in Q \cap D} ($ | $tf(T)$ | · $df(T)$ · $f_b(T)$ · $f_n(T, D)$) |
| BM25 | 1 | · | 1 | · | $\sum_{T \in Q \cap D} ($ | TF | · DF · QTF · 1) |
| dtb.nnn | 1 | · | 1 | · | $\sum_{T \in Q \cap D} ($ | TF | · DF · QTF · LN) |

Table 2: Plugging components of the TF×IDF equation into Lucene’s scoring equation; the first row is the generalized form and the following two rows show the components of two popular TF×IDF equations from Table 1 transplanted to Lucene’s equation.

Deep Dive into the Lucene Query/Weight/Scorer Java Classes

Jake Mannix, Lucidworks: In this more technical talk, Jake explained how Lucene scores a query, and what classes are instantiated to support the scoring. Jake described, first, at a high level how to do scoring modification to Lucene-based systems, including some “Google”-like questions on how to score efficiently. Then, he went into more details about the BooleanQuery class and its cousins, showing where the Lucene API allows for modifications of scoring with pluggable Similarity metrics and even deep inner-loop, where ML-trained ranking models could be instantiated - *if you’re willing to do a little work*.

Learning to Rank with Solr

Diego Ceccarelli, Bloomberg On day two of the workshop, Diego started his talk by explaining that tuning the relevance of a search system is often performed by “experts” who hand tune and craft the weightings used for the different retrieval features. However, this approach is manual, expensive to maintain, and based on intuitive or deep domain knowledge, rather than data. His working goal behind this project was to automate the process. He motivated the use of Learning To Rank, a technique that enables the automatic tuning of an information retrieval system. He pointed out that sophisticated learned models can make more nuanced ranking decisions than a traditional ranking function when tuned in such a manner. This is the reason why, at Bloomberg, they have integrated a Learning to Rank component directly into Solr and contributed the code back⁸ enabling others to easily build their own Learning To Rank systems. During his talk, Diego presented the key concepts of Learning to Rank, how to evaluate the quality of the search in a production service, and finally described how the Solr Learning to Rank component works.

3 Hackathon

As part of the workshop, a day and half was dedicated to the hackathon, where numerous attendees contributed to the Lucene4IR GitHub Repository - <http://github.com/leifos/lucene4ir/>. In the repository, three main applications were developed and worked on:

- **IndexerApp** - enables the indexing of several different TREC collections, e.g. TREC123 News Collections, Aquaint Collection, etc.

⁸<https://issues.apache.org/jira/browse/SOLR-8542>

-
- **RetrievalApp** - a batch retrieval application when numerous retrieval algorithms can be configured, e.g. BM25, PL2, etc
 - **ExampleStatsApp** - an application that shows how you can access various statistics about terms, documents and the collection. e.g. how to access the term posting list, how to access term positions in a document, etc.

The repository also contained a sample test collection (documents, queries and relevance judgements) was provided (CACM), so that participants could try out the different applications.

During the workshop, a number of different teams undertook various projects:

- **Customisation of the tokenisation, stemming and stopping during the indexing process:** this enabled the IndexerApp to be configured so that the collections can be indexed in different ways i.e. they could change the stemming algorithm, include a stop list, enable positions to be recorded, etc. The idea being that students would be able to vary the indexing parameters and then see the effect on the collection and performance .
 - **Implementation of other retrieval models:** inheriting from Lucene's BM25Similarity Class, BM25 for Long documents was implemented BM25L [6], OKAPI BM25's was also implemented to facilitate the comparison between how it is currently implemented in Lucene versus an implementation of the original BM25 weighting function.
 - **Alternative Scoring Mechanism:** Rather than scoring through Lucene's mechanics (Query \rightarrow Weight \rightarrow Scorer), others attempted to implement BM25 by directly accessing the inverted index through the Lucene's index API. The objective was twofold. First, to provide a "template" to implement a retrieval model where document matching is performed through a Document At A Time (DAAT) strategy and access to term vocabulary (via Terms and TermsEnum) and to posting lists (via PostingsEnum) is made more explicit; in some scenarios (e.g. for teaching activities) it could be useful to provide sample code that relies only on general concepts (term vocabulary, posting lists, etc.) and it is not tailored to specifics of the library — e.g. the Lucene scoring model. The second objective was to provide an "easier" way to control the variables in the experimental settings or to investigate if choices made for efficiency purposes (e.g. document length approximation) significantly affect effectiveness.
 - **Query Expansion:** A QueryExpansionRetrievalApp was developed that expanded queries using word synonyms - a parameter was introduced to mix together the original query with the expanded query.
 - **Hacking the Inner-Loop:** The break-out group focused on inner loop scoring wanted to try something that was simultaneously simple, practical, and yet required some inner loop scoring magic. Based on the interests of the group members, we decided on "cross-field phrase queries": an extension of the idea of a sloppy phrase query where the "slop" allowed for a pair of terms occurring in *different* fields to be part of a phrase (but with a parametrizably lower score than terms in the same field). We worked out the design (delegating most of the work to Query / Weight / Scorer classes already in Lucene, but then combining them together across fields), and stepped through much of the iteration implementation. While we got most of the plumbing done, we only had enough time for our "score()" method to be implemented as naively as imaginable, and did not get
-

it fully working in the time of the workshop. Some participants expressed interest in working on it further, to see how efficient it was, and what effect on scoring it would have (if a QueryParser was configured to explicitly spit out queries of this form sometimes).

- **Working with Lucene’s Index and Reader:** Additional Examples on how to access and work with Lucene’s index were also added to the ExampleStatsApp. These code snippets showed how it was possible to iterate through the term postings list, how to iterated through documents, and access various document, term and collection statistics. The purpose of the app was to demonstrate how to work with the Lucene index in order to perform various operations, which are often difficult to work out from the code base or existing documentation.

Along with the code some documentation was also produced to explain various aspects and how to set up and run the different applications. However, now, post workshop, there is the need to bring together all the elements developed and collate all the documentation so that these resources can be used for teaching and learning IR.

4 Discussion

During the course of the workshop, two breakout groups were formed to discuss how we can use Lucene when teaching and learning, and what were the main challenges in bridging the industry/academia along with what opportunities it could bring about. Finally, we asked participants to provide some feedback on the event.

4.1 Teaching and Learning

To seed the discussion for this working group, various members explained how the Information Retrieval course was taught at their institute. As Juanma had already discussed how they teach at the University of Granada (see above), others described their courses and experiences.

Emanuele Di Buccio, University of Padova: Emanuele described one of the courses taught as part of the Master Degree in Statistical Science at the University of Padua, called *Information Systems (Advanced)*⁹. The course covered both basic IR topics: indexing and retrieval methods, retrieval models, and evaluation, along with more advanced topics such as Web Search or Machine Learning for IR. A detailed description of the course contents can be found in [8], which is an IR book developed from the experiences teaching the course. The course was designed so that, for most of the topics, lessons at a theoretical level on a specific topic were followed by a laboratory assignment on that topic. The topics covered in laboratory assignments were: creation of a test collection, indexing, retrieval, relevance feedback, link analysis, learning to rank, and optimization of ranking functions with parameters.

Emanuele explained that students were asked to propose their own methodology to carry out the laboratory activities. For instance, when considering the topic of “relevance feedback”, each student could propose their own methodology to perform feedback, e.g. through a query expansion method or term re-weighting. Each assignment, then, involved the experimental evaluation on a shared test collection. Indeed, the objective of the assignments was three-fold:

⁹The description refers to the course editions in the Academic Years 2011/12-2014/15. The professor in charge was Massimo Melucci.

-
1. to better understand the topic;
 2. to become familiar in the design and the implementation of experimental methodologies to evaluate methods and/or components and,
 3. more generally, to test research hypotheses.

Students were allowed to use a manual approach (when possible), a software library or build their own software modules to achieve the assignment objective; a list of software libraries were provided before the first laboratory assignment to make the students aware of possible options. However, the adoption of a manual approach for some of the laboratory activities was mandatory. For instance, in the case of the assignment on indexing, the use of a manual approach aimed at a better understanding of the conceptual mechanisms to identify the most effective descriptors to retrieve relevant documents. When the students proposed their own methodology for indexing, they were asked to present their approach as a set of steps that can be automated.

The availability of software libraries or resources to easily use the basic operations is crucial to allow the students to test their methodology with little/less effort. In an edition of the course, a lesson was dedicated to a general introduction to Apache Lucene, where sample code was provided. Along with Apache Lucene, and introduction to Elasticsearch [3] was also presented, particularly how to index documents, perform retrieval, and how to customize the scoring mechanism via scripting¹⁰ The main reason for the introduction to Elasticsearch was that the students could index, retrieve, and customize the retrieval algorithm – and therefore test some of their methodologies – without writing actual code but only through the use of REST requests.

Another aspect Emanuele commented on was the heterogeneous background of the students, and how they came from various disciplines. This was one of the reasons, why they did not restrict the laboratory activities to a single software library and allowed students to select the tool they felt most comfortable with. While students in Computer Science and Computer Engineering were familiar with Java, students in Statistical Science tended to prefer the R language because it was used in many courses within their course degree. Therefore, one of the resources that could be useful for teaching is a R wrapper for Apache Lucene — wrappers in other programming languages exist, e.g. PyLucene [12] for Python. Software libraries such as Elasticsearch, could be useful tools to support teaching: for instance, they provide functionalities – in the event of Elasticsearch a REST request – to display how a specific fragment of text is processed given a pipeline, e.g. a specific tokenizer and a set of filters (lowercase, porter stemming, etc).

Prof. Krisztian Balog, University of Stavanger: The Web Search and Data Mining course is part of the Computer Science master’s programme at the University of Stavanger, but it is also offered to (advanced) bachelor students. The data mining part of the course includes data processing, classification and clustering methods. The IR part consists of indexing, retrieval models, evaluation, link analysis, query modeling, and entity linking and retrieval. The course is 6 hours per week, which is divided to 2 lectures (2x2 hours) and a practical session (2 hours).

Krisztian explained that during the lectures, after presenting the theory, students get a small paper exercise sheet where they need to apply the said theory on toy-sized input data. Examples of such exercises include constructing an index from some input text, computing

¹⁰ElasticSearch allows to evaluate a custom score via scripts — see the *Scripting* module. Apache Solr provides similar functionalities via *Function Queries*.

term weights and scoring a small number of documents, calculating PageRank scores, etc. They can use a calculator and/or a spreadsheet program, but the input is simple enough for pen-and-paper. The reference solutions to these exercises are made available after the class. Student feedback has been very positive; they are really appreciative of this element in the lectures. While completing these exercises, interesting questions (both practical and theoretical) can often spring up and be discussed.

The practical sessions involve implementing methods from the lectures and applying them on small (but real) datasets. Typically it happens on two levels. First, students need to implement one or two of the simpler methods for a given problem, e.g., decision trees or Naive Bayes for classification. Second, they get to use a ready-made third-party implementation; for the previous example, it would be SVM and Random Forests (from the scikit-learn Python package). For retrieval, ElasticSearch is used; the RESTful API is well documented and can easily be used from Python (or from any programming language for that matter). Evaluation is a core element of these exercises, so they need to measure and compare the performance of different approaches according to some metric. In order to allow students to focus on the more interesting parts of the problem as opposed to more “mechanical” tasks (e.g., reading in data from a file), they get the skeleton of the code along with explanations as an iPython notebook, and they only need to complete the missing parts.

Finally, students have a handful of larger (obligatory) assignments throughout the semester that they need to complete in teams. These assignments involve larger-scale datasets (note that this scale is still in accordance with academic standards, not with industrial ones). The assignments are set up as competitions on Kaggle,¹¹ with a (hard) deadline and a minimum performance threshold (e.g., a certain MAP score for a ranking task). There are no restrictions on the choice of the programming language or libraries used. The members of the best performing team for each assignment are rewarded with some bonus points that they can “cash in” during the final exam.

Martin Halvey, University of Strathclyde: The Department of Computer and Information Science at the University of Strathclyde has two modules relevant to the discussion. The first module is Information Access & Mining (ISA) is delivered to final year undergraduates and cover a range of techniques for extracting information from textual and non-textual resources, modelling the information content of resources, detecting patterns within information resources and making use of these patterns. The second is Information Retrieval and Access (IRA) which is delivered to Masters students. This module is a required module for students on Strathclyde’s Information & Library Studies and Information Management Masters Programmes, as well as being an optional module for other Masters students. To offer a contrast to other modules Martin described IRA in detail, as the cohort is different to others described. Typically, with some exceptions, the students do not have experience in programming or mathematics in their undergraduate degree. This presents a number of challenges when teaching some of the core concepts, where the syllabus includes: information seeking and behaviour, indexing, term weighting, retrieval models, IR evaluation, multimedia retrieval, user interfaces and interaction, web retrieval.

Martin explained that in laboratory and tutorial sessions that students were given problems to solve on paper. The intention is that students understand how different concepts, models, evaluation measures, etc. work. For some problems students are provided with

¹¹Kaggle in Class (<https://inclass.kaggle.com/>) is provided free of charge for academics.

spreadsheets that automatically calculate some of the equations discussed in lectures so that students can see the relationship between different inputs and outputs.

Martin outlined how developing some demonstrators using Lucene could replicate what he currently does with spreadsheets, with the benefit being that these demonstrators would be based on a real toolkit and also be more adaptable to use a wider range of retrieval models, evaluation measures etc. There is also the possibility in future years that students will be introduced to tools like Apache Lucene and ElasticSearch in a different module to IRA. Here, it was pointed out by Ian Ruthven, that often these students won't need to modify such toolkits, but they will need to know how they work, how to configure them, and how to evaluate their configuration choices.

Discussion: From the various perspectives, it was clear that there was a number of key concepts that were felt to be fundamental to teaching Information Retrieval. From the discussion it was also clear that the lecturers wanted to give students hands-on experience so that they could see the impact and effect of the different components i.e. what does tokenization and stemming do to the size of vocabulary, the size of the index, and the influence on precision and recall. Also, from the descriptions there was a consensus towards teaching IR in an inquiry led manner - focusing mainly on the science (rather than the engineering). In this way the lectures and course work would be focused on presenting experimental contexts in which the students could go off and conduct experiments to gain insights and understanding into the effect and influence of different factors (rather than just told what would happen). In this way, students would become more scientific in their approach, and know how to conduct an experiment (aim, method, results, conclusion). This was seen to be an important skill to learn, both from a research point of view, but also from a very practical point of view, i.e. many students will become data scientists, and so they need to be trained to be methodical in their approach.

With this in mind, we drew up a table of different facets of the courses, roughly split into Indexing, Retrieval, Analysis and Evaluation (see Table 3). We considered the two different levels, high (i.e. using/configuring the apps) and low (i.e. programming and coding algorithms). The reason for considering the different levels, was that different cohorts of students have different skill sets or the focus of the course maybe more oriented towards technical under the hood skills versus high level understanding and usage. At the different levels, we considered the different apps being built during the hackathon along with other apps that would also help support teaching and learning in IR. Table 3 summarises the different apps and some of the things that we would like students to be able to do with them. For example, consider a lecture on stemming and the following questions: what is the influence of stemming and what are the differences between no stemming, porter stemming and krovetz stemming, in terms of vocabulary and index size, and performance? Having apps that let students easily configure and run different stemmers, then be able to conduct retrieval experiments, measure the performance, and analyze/inspect the index, would then enable them to form their own insights into the effects of stemming. On the other hand, during the retrieval algorithms lectures, different comparisons between models can be made, and if required new algorithms developed. In terms of further analysis, an app (ResultAnalyzerApp) could be used to help analyze the influence of document length normalization - does changing the b parameter in BM25 actually effect the length of documents that are retrieved? While these are simple examples - it was felt that this on-boarding stage helps students to contextualize and understand the taught material - and enables them to go on to conduct more advanced

| Apps | High Level | Low Level |
|-------------------------|--|---|
| IndexerApp | Modify how the indexer is performed i.e. different tokenizers, parsers, etc | Can modify parsers, tokenizers, etc |
| IndexAnalyzerApp | Inspect the influence of indexer | |
| RetrievalApp | Try out different retrieval algorithms Change retrieval parameters | Implement new retrieval algorithms |
| ResultAnalyzerApp | Inspect and analyze the results returned | Customise the analysis, put out other statistics of interest |
| ExampleApp | Examples of how to work with the Lucene | index, to make modifications |
| Batch Retrieval Scripts | Configure to run a series of standard batch experiments | Customize to run specific retrieval experiments |
| RetrievalShellApp | n/a | Implement retrieval algorithms not using Lucene's scorer assuming term independence |

Table 3: A summary of different apps and what could be varied at different levels.

evaluations.

4.2 Challenges and Opportunities

During the workshop the challenges and opportunities between academia and industry were discussed, focusing on research, teaching and learning, and graduate attributes. Below is a summary of the main points stemming from the discussion.

The first point discussed were research opportunities and challenges that can arise between industry and academia in information retrieval, text mining and big data domain. One of the first problems discussed was regarding the access to data and research problems. First, not having access to data often precludes academic investigations, but even when available, the data needs to be of high enough quality and representative of the research problems faced by the company, for useful solutions to be developed. It was noted that in industry they often hit upon really interesting problems, but often, do not have the time, to investigate - it was suggested that this is where academia could potentially help - big problems are hard to find - however, finding funding to help companies (especially small companies) work with academia was seen as difficult to acquire and very time consuming with high overheads - and so more efficient knowledge transfer mechanisms were needed (at funding body level). Another aspect regarding the data, is that often it difficult to disclosure data because Non-Disclosure Agreements (NDA) are in place with clients or the data is in-house. While one solution would be if the client would agree to its release (unlikely), or that a research project between the industrial partner, the client and the academic be formed (long lead time). Another alternative suggested would be to abstract the problem away from the data and client so that the problem can be exposed - without disclosing the data or compromising the client. While this means that the data is still an issue - but at least then the problem can be further examined - and funding sought to solve it.

Another point discussed was the need for a common Open Source platform that can be used for teaching and learning information retrieval and big data. For example, both academic and industry participants agreed upon the need for a better and general documentation for Lucene. The industry participants felt that such teaching and learning materials

could be developed by academics with the help of the Lucene community. It was unclear whether there were funding schemes available to help facilitate this, though.

The final point discussed were the graduate attributes and the skills required by Computer Scientist and Software Engineering Graduates working in information retrieval, text mining and big data. The obvious and core skills required in terms of being able to : (i) develop high quality, robust code, (ii) understand and think about complex systems/problems, (iii) communicate, discuss and resolve issues and (iv) have a good understanding of software engineering principles and practices, were seen as mandatory. More specifically to the field of search and big data, industry participants felt that there was a lack of skill graduates that knew about search technologies, how to process large scale data sets, and how to work with big (text) data. They expected candidates to understand more about the processes and core concepts of search and big data - to be able to describe it at a conceptual level, at least - but with ideally some practical skills (i.e. Lucene, Solr, Spark, Pig, etc). It was noted that the learning curve can be very steep when picking such technologies - but such skills are seen to be increasingly valuable by employers. For example, Lucene is a very large and complex project, that has evolved over years, so it can feel very opaque and daunting to begin with, however, given that it one of the largest OS toolkits for information retrieval and data mining it is a skill worth learning. It was felt that there was a strong need for developing more training and resources for beginners to learn how to use such toolkits - and this was seen to be an area where more investment from funding agencies and universities could be directed.

4.3 Feedback

At the end of the workshop, we asked participants to provide feedback on the following questions: “What would you suggest us to Stop/Start/Continue to do in our next workshop?”. Below is a summary of the points contributed by participants.

First of all, participants were pleased by the workshop, specially they enjoyed the hackathon part as well as the talks from the industry. They also encouraged us to organise a follow up workshop on this topic. Participants liked the idea of the workshop were it brought together academics and industry and encouraged us to continue inviting people from industry and in particular Lucene developers. They also suggested to invite more undergraduate students to the hackathon.

Some participants suggested defining the goals and objectives by asking participants before the hackathon. In addition to this point, other participants suggested that they would prefer to have a more well-defined structured hackathon, which means, setting more well-designed goals, having tools and data tested and ready, as well as making sure all participants use the same data. One proposal on this front was to create a standardised mini-competition between workshop participants, e.g. providing a template code with a challenge to increase MAP with the expectation that each team formally presents its results.

Since, a number of the participants were not very familiar with using Lucene, it was suggested that in subsequent hackathons a mini-tutorial is included to help participants get up to speed. For example, some demos of various Lucene modules and how they work. Others suggested that some tutorials for the IR community would be a great way to encourage people to start working with Lucene - and in particular - help students to learn the basics and how to work with the toolkit, instead of against it. Some more detailed technical talks on the various modules were also encouraged.

Finally, some participants were requesting to include Solr into the program and promote the event in the Lucene and Solr Communities. And that the workshop could have been recorded and/or stream to increase its reach and dissemination.

We are very encouraged by the suggestions and feedback from participants as they provide a number of ways in which this initiative can be further developed and improved to support research and industry in this area.

5 Summary

This report provides an overview of the workshop on developing teaching and training resources for Information Retrieval evaluation. While, the event was informative and enjoyable, it was clear that there is lot of scope for developing greater links between industry and academia - and that tools like Lucene, Solr and ElasticSearch - provide an obvious way to foster collaboration. Furthermore, such toolkits, being supported by a large community, and widely used, mean that researchers working with such tools can effect more efficient knowledge transfer and their research can have greater impact. Furthermore, students moving into industry can benefit from learning such technologies - which are not only relevant - but because they provide a realistic introduction to using large scale production systems that are used for processing, searching and mining big data. With more focus, and more events, of this nature, we can develop common tools that support research through first developing learning and teaching resources for working with such toolkits. A logical next step would be to bring relevant parties together to bring together the work done during the hackathon, and to integrate the applications within an IR course.

6 Acknowledgments

We thank the European Science Foundation / ELIAS Network for funding the workshop (Grant No. SM 5916). We would also like to thank our speakers as well as Bloomberg, Flax, LucidWorks and the University of Strathclyde. Finally, we would like to thank all the participants, Pablo Arteaga, Martynas Buivys, Matteo Catena, Aidan O'Grady, Florence Kolberg, Florian Meier, Mohammad Alian Nejadi, Aigars Reimers, Wim Vanderbauwhede, Colin Wilkie, Charlotte Wilson, and Alan Woodward, for their contributions to the workshops and hackathon. Also, thanks to Henry Enfield, Jimmy Lin, Casper Petersen, Ian Soboroff, Manisha Verma, Guido Zucco, for their offline contributions and discussions.

References

- [1] ARGUELLO, J., CRANE, M., DIAZ, F., LIN, J., AND TROTMAN, A. Report on the sigir 2015 workshop on reproducibility, inexplicability, and generalizability of results (rigor). *SIGIR Forum* 49, 2 (Jan. 2016), 107–116.
 - [2] DOWIE, D., AND AZZOPARDI, L. *Re-leashed! The PuppyIR Framework for Developing Information Services for Children, Adults and Dogs*. 2013, pp. 824–827.
 - [3] ELASTICSEARCH. <https://www.elastic.co/products/elasticsearch>.
-

-
- [4] FERNÁNDEZ-LUNA, J. M., HUETE, J. F., RODRÍGUEZ-CANO, J. C., AND RODRÍGUEZ-HERNÁNDEZ, M. Teaching and learning information retrieval based on a visual and interactive tool: sulair. In *4Th International Conference on Education and New Learning Technologies (EDULEARN)* (2012), pp. 6634–6642.
 - [5] LUCENE. <https://lucene.apache.org>.
 - [6] LV, Y., AND ZHAI, C. When documents are very long, bm25 fails! In *Proceedings of the 34th International ACM SIGIR Conference on Research and Development in Information Retrieval* (2011), SIGIR '11, pp. 1103–1104.
 - [7] MACDONALD, C., MCCREADIE, R., SANTOS, R. L., AND OUNIS, I. From puppy to maturity: Experiences in developing terrier. *Proc. of OSIR at SIGIR* (2012), 60–63.
 - [8] MELUCCI, M. *Information retrieval. Metodi e modelli per i motori di ricerca*. Informatica: Nuova serie. Franco Angeli, 2013.
 - [9] METZLER, D., AND KURLAND, O. Experimental methods for information retrieval. In *Proceedings of the 35th International ACM SIGIR Conference on Research and Development in Information Retrieval* (2012), SIGIR '12, pp. 1185–1186.
 - [10] OGILVIE, P., AND CALLAN, J. P. Experiments using the lemur toolkit. In *TREC* (2001), vol. 10, pp. 103–108.
 - [11] OUNIS, I., AMATI, G., PLACHOURAS, V., HE, B., MACDONALD, C., AND JOHNSON, D. Terrier information retrieval platform. In *Proceedings of the 27th European Conference on Advances in Information Retrieval Research* (Berlin, Heidelberg, 2005), ECIR'05, Springer-Verlag, pp. 517–519.
 - [12] PYLUCENE. <http://lucene.apache.org/pylucene/>.
 - [13] STROHMAN, T., METZLER, D., TURTLE, H., AND CROFT, W. B. Indri: a language-model based search engine for complex queries. Tech. rep., in *Proceedings of the International Conference on Intelligent Analysis*, 2005.
 - [14] ZOBEL, J., WILLIAMS, H., SCHOLER, F., YIANNIS, J., AND HEIN, S. The zettair search engine. *Search Engine Group, RMIT University, Melbourne, Australia* (2004).
-