

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

Leif Azzopardi¹, Yashar Mosfeghi², Martin Halvey¹,
Krisztian Balog³, Juan Manuel Fernandez Luna⁴, Charlie Hull⁵

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

² University of Glasgow *Yashar.Mosfeghi@glasgow.ac.uk*

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

⁴ University of Granada *jmfluna@decsai.ugr.es*

⁵ Flax *charlie@flax.co.uk*

September 21, 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 Lucene in action in industry, in teaching and learning environments, and evaluation forums. (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 and the resources created in teaching and learning IR evaluation. 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 and raising numerous issues. However, by adopting and contributing to most widely used and supported Open Source IR toolkit, it was clear that 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. 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 Solr¹.

Introduction Talk: Why are we here?

Leif Azzopardi, University of Strathclyde: Leif explained how after attending the lively Reproducibility workshop 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 inclusive and welcoming to such a large community of search practitioners as we could, and have perhaps failed to transfer our knowledge 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 Lucene based tools? 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 like Lucene. 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.

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 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 at reluctant to adopt open source, but nowadays it has been 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

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

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, 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 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 Manuel Fernandez Luna (University of Granada) : In this talk, we shall describe how the University of Granada is supporting teaching and learning Information Retrieval (TLIR) discipline across different courses in the Computer Science studies (degree and masters), and how this is done by means of the Lucene API. Later we shall present our thoughts about how Lucene could be used in TLIR context and present some proposals for improving the Lucene experience, both for students and lecturers.

Black Boxes are Harmful

Sauparna “Rup” Palchowdhury (NIST) : Having seen students and practitioners in the IR community grapple with abstruse documentation that accompany search systems, Rup wanted to direct some attention to Lucene’s internals and demonstrate how to do IR experiments using it. Rup argued that it is important to ensure that a search system that you have built works “correctly” - and one way to do this is by evaluating its output on test-collections like those from TREC. He further argues that using such systems as black boxes can mislead students because they learn little from their results. Rup points out that this is because there are many possible points of failure when implementing and configuring a search system (document corpus problems, checksum mismatches, wrong document-query-rel triplets, one parameter with many meanings, switches not mutually exclusive, bugs in algorithms, different naming conventions). This means it is very difficult to reproduce experiments and results reported in research papers.

As a more concrete example he described how Lucene conceptualizes the scoring of a document:

$$score(q, d) = coord-factor(q, d) \times query-boost(q) \times \frac{V(q) \times V(d)}{|V(q)|} \times doc-len-norm(d) \times doc-boost(d) \quad (1)$$

Rup focused in on implementations of some of the core retrieval algorithms, i.e. TFxIDF and BM25 - and described some of the numerous variations that exist, and how these can lead to quite different levels of performance.

His summarised a number of key issues that have to be considered when using search systems in the context of evaluation and reproducibility. These included have a focus on standardizing how we describe and naming variables, methods and functions, and using a well defined notation, creating sharable experimental artifacts, and providing implementations that are tracable to a source.

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 Lucene/Solr et al 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, rather than data. His working goal behind this project was to automate the process. He described how this motivated the use of Learning To Rank, a technique that enables the automatic tuning of a information retrieval

Table 1: Attempt at encoding the picture

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
trec_eval	Measure the performance	
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	Customize to implement retrieval algorithms without the Lucene scorer i.e. a simple model assuming term independence

system by applying machine learning when estimating parameters. He points out that sophisticated models can make more nuanced ranking decisions than a traditional ranking function when tuned in such a manner. 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 then how the Solr plugin works. At Bloomberg, they have integrated a learning to rank component directly into Solr (and released the code as Open Source), enabling others to easily build their own Learning To Rank systems and access the rich matching features readily available in Solr.

3 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.

3.1 Teaching and Learning

Juanma

Krisztian

Martin

How do you go about teaching IR? What level?

What kinds of things do you need/want from such resources?

How do we see Lucene fitting in? Benefits to students?

3.2 Challenges and Opportunities

Some challenges from charlie - lack of good real-world test data for researchers. Companies need to provide this - many open source search companies are small and therefore find it hard to support apprenticeships, internships etc. or access funding such as KTPs - Lucene community can be hard to enter - learning curve can be steep, this is a very large and complex project

3.3 Feedback

4 Resources

As part of the workshop 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.
- 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.

In the repository, 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 - the idea being that students would be able to vary the indexing and then see the effect on performance.
 - Implementation of other retrieval models: inheriting from Lucene's BM25Similarity Class, BM25 for Long documents was implemented BM25L[], 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 [].
 - Rather than scoring through Lucene's mechanics, others attempted to implement BM25 by directly accessing the inverted index - again to provide a comparison in terms of both efficiency and effectiveness for scoring queries.
 - A QueryExpansionRetrievalApp
 - Hacking the innerloop
 - Additional Examples on how to access and work with Lucene's index and
-

5 Summary

6 Acknowledgments

We thank the European Science Foundation / ELIAS Network for funding the workshop **add number**. We would also like to thank our speakers as well as Bloomberg, FlaxSearch, Lucid-Works and the University of Strathclyde. Finally, we would like to thank all the participants for their contributions to the workshops and hackathon. Also, thanks to Manisha, Guido and Casper for their offline contributions.