

# Implementing of RankNet

Ryan Greenup

February 20, 2021

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Motivation</b>	<b>2</b>
<b>3</b>	<b>Implementation</b>	<b>2</b>
<b>4</b>	<b>Implementation</b>	<b>3</b>
4.1	Neural Network . . . . .	3
4.1.1	The Ranknet Method . . . . .	3
4.1.2	Implementation . . . . .	3
4.2	How to clone . . . . .	4
4.3	Blobs . . . . .	4
4.4	Moons . . . . .	4
4.5	Optimisers . . . . .	4
4.6	Batches . . . . .	4
4.7	Wine . . . . .	4
4.8	Rank Wiki Articles . . . . .	4
<b>5</b>	<b>Difficulties</b>	<b>4</b>
<b>6</b>	<b>Further Research</b>	<b>4</b>
<b>7</b>	<b>Conclusion</b>	<b>5</b>
<b>8</b>	<b>Further Research</b>	<b>5</b>
<b>9</b>	<b>Text and References</b>	<b>5</b>
<b>10</b>	<b>Fractals</b>	<b>5</b>
<b>11</b>	<b>Appendix</b>	<b>5</b>
11.1	Search Engines . . . . .	5
11.1.1	Fuzzy String Match . . . . .	6
# #+TODO: TODO IN-PROGRESS WAITING DONE		

## 1 Introduction

Ranknet is an approach to *Machine-Learned Ranking* (often referred to as “Learning to Rank” [21]) that began development at Microsoft from 2004 onwards [5], although previous work in this area had already been undertaken as early as the 90s (see generally [9, 8, 9, 10, 37]) these earlier models didn’t perform well compared to more modern machine learning techniques [23, §15.5].

Information retrieval is an area that demands effective ranking of queries, although straight-forward tools such as `grep`, relational databases (e.g. `sqlite`, `MariaDB`, `PostgreSQL`) or NoSQL (e.g. `CouchDB`, `MongoDB`) can be used to retrieve documents with matching characters and words, these methods do not perform well in real word tasks across large collections of documents because they do not provide any logic to rank results (see generally [35]).

## 2 Motivation

Search Engines implement more sophisticated techniques to rank results, one such example being TF-IDF weighting [24], well established search engines such as *Apache Lucene* [1] and *Xapian* [14], however, are implementing Machine-Learned Ranking in order to improve results.

This paper hopes to serve as a general introduction to the implementation of the Ranknet technique to facilitate developers of search engines in more modern languages (i.e. Go and Rust) in implementing it. This is important because these more modern languages are more accessible [13] and memory safe [30] than C/C++ respectfully, without significantly impeding performance; this will encourage contributors from more diverse backgrounds and hence improve the quality of profession-specific tooling.

For a non-comprehensive list of actively maintained search engines, see §11.1 of the appendix.

## 3 Implementation

The Ranknet approach is typically implemented using Neural Networks, an early goal of this research was to evaluate the performance of different machine learning algorithms to implement the Ranknet method, this research however is still ongoing.

Further Research to look at the implementation of Ranknet for documents and comparing different approaches to apply the method to on-demand queries is required, although this seems to have been implemented by open source Apache Solr project [25], which may provide guidance for further study. `teamDocFetcherFastDocument`.

Ranking/ is the process of applying machine learning algorithms to ranking problems, it .

A lot of data cannot be clearly categorised or quantified even if there is a capacity to compare different samples, the motivating example is a collection of documents, it might be immediately clear to the reader which documents are more relevant than others, even if the reader would not be able to quantify a “relevance score” for each document.

---

By training a model to identify a more relevant document, a ranking can be applied to the data.

An example of this might be identifying documents in a companies interwiki that are relevant for new employees, by training the model to rank whether one document is more relevant than an other, ultimately an ordered list of documents most relevant for new employees could be created.

## 4 Implementation

This implementation will first apply the approach to a simple data set so as to clearly demonstrate that the approach works, following that the model will be extended to support wider and more complex data types before finally being implemented on a corpus of documents.

### 4.1 Neural Network

Neural Networks [31]

The Ranknet method is typically implemented using a Neural Network, although other machine learning techniques can also be used [5, p. 1], Neural Networks are essentially a collection of different regression models that are fed into one another to create a non-linear classifier, a loss function is used to measure the performance of the model with respect to the parameters (e.g. RMSE <sup>1</sup> or BCE <sup>2</sup>) and the parameters are adjusted so as to reduce this error by using the *Gradient Descent Technique* (although there are other optimisation algorithms such as RMSProp and AdaGrad [26] that can be shown to perform better see [2]). The specifics of Neural Networks are beyond the scope of this paper (see [12] or more generally [31]).

#### 4.1.1 The Ranknet Method

The Ranknet method is concerned with a value  $p_{ij}$  that measures the probability that an observation  $i$  is ranked higher than an observation  $j$ .

A Neural Network ( $n$ ) is trained to return a value  $s_k$  from a feature vector  $\mathbf{X}_k$ :

$$n(\mathbf{X}_i) = s_i \quad \exists k$$

So as to minimise the error of:

$$p_{ij} = \frac{1}{1 + e^{\sigma \cdot (s_i - s_j)}} \quad \exists \sigma \in \mathbb{R}$$

#### 4.1.2 Implementation

The first step is to create a simple data set and design a neural network that can classify that data set, this can then be extended.

---

<sup>1</sup>RMSE Root Mean Square Error

<sup>2</sup>BCE Binary Cross Entropy

### 4.2 How to clone

How can the reader clone this onto there machine?

put on the summer repo then provide instructions to clone this working example onto there machine to try it out.

### 4.3 Blobs

### 4.4 Moons

### 4.5 Optimisers

### 4.6 Batches

### 4.7 Wine

### 4.8 Rank Wiki Articles

## 5 Difficulties

- Don't use torch
  - Do it by hand first because it can be hard to see if the correct weights are being updated sensibly, making debugging very difficult.
  - R or Julia would be easier because counting from 0 get's pretty confusing when dealing with  $\{1, 0\}$ ,  $\{-1, 0, 1\}$ .
- Don't use misclassification rate to measure whether the ranking
  - In hindsight this is obvious, but at the time misclassification was a tempting metric because of it's interpretability

was correct

Very difficult to see if the model is working

- A continuous function will still produce an ordered pattern in the ranking of results, even if the model hasn't been trained, so visualising isn't helpful either.
- Implement it on a data set that already has order, obfuscate the order and then contrast the results
  - or use a measurement
- Plot the loss function of the training data live, the model is slow to train and waiting for it to develop was a massive time drain.

## 6 Further Research

It is still not clear how the performance of Ranknet compares to traditional approaches implemented by search engines (see §11.1), further study would ideally:

- 
- Write a program to query a corpus of documents using an existing search engine.
    - Or possibly just implement TF-IDF weighting in order to remove variables.
  - Extend the program to implement machine learned ranking
  - Measure and contrast the performance of the two models to see whether there are any significant improvements.

This could be implemented with TREC datasets [34] using a cumulated-gain cost function [16] as demonstrated in previous work [35].

## 7 Conclusion

## 8 Further Research

- Apply this to documents to get a sorted list.
- The "Quicksort" algorithm likely needs a random pivot to be efficient [32]

## 9 Text and References

Fractals are complex shapes that often occur from natural processes, in this report we hope to investigate the emergence of patterns and complex structures from natural phenomena. We begin with an investigation into fractals and the concept of dimension and then discuss links between fractal patterns and natural processes.

This is a Reference [33] and another [27] and yet another [4].

## 10 Fractals

Images are shown in figure .

## 11 Appendix

### 11.1 Search Engines

There are many open source search engines available , a cursory review found the following popular projects:

- [Zettair](#) (C) [15]
- [Apache lucene/Solr](#) (Java) [1]
  - Implemented by [DocFetcher](#) [7]
- [Sphinx](#) (C++) [38]
- [Xapian](#) (C++) [29]

- Implemented by [Recoll](#) [17]

More Modern Search engines include:

- [LunrJS](#) (JS) [28]
- [Bleve Search](#) (Go) [24]
- [Riot](#) (Go) [36]
- [Tantivy](#) (Rust) [6]
- [SimSearch](#) (Rust) [22]

### 11.1.1 Fuzzy String Match

Somewhat related are programs that rank string similarity, such programs don't tend to perform well on documents however (so for example these would be effective to filter document titles but would not be useful for querying documents):

- [fzf](#) [3]
- [fzy](#) [11]
- [peco](#) [20]
- [Skim](#) [39]
- [go-fuzzyfinder](#) [19]
- [Swiper](#) [18]

## References

- [1] Apache Software Foundation. *Learning To Rank | Apache Solr Reference Guide 6.6*. Solr Reference Guide. 2017. URL: [https://lucene.apache.org/solr/guide/6\\_6/learning-to-rank.html](https://lucene.apache.org/solr/guide/6_6/learning-to-rank.html) (visited on 02/20/2021) (cit. on pp. 2, 5).
- [2] Vitaly Bushaev. *Understanding RMSprop — Faster Neural Network Learning*. Medium. URL: <https://towardsdatascience.com/understanding-rmsprop-faster-neural-network-learning-62e116fcf29a> (visited on 02/16/2021) (cit. on p. 3).
- [3] Junegunn Choi. *Junegunn/Fzf*. Feb. 20, 2021. URL: <https://github.com/junegunn/fzf> (visited on 02/20/2021) (cit. on p. 6).
- [4] Christopher Burges. *From RankNet to LambdaRank to LambdaMART: An Overview (MSR-TR-2010-82)*. Jan. 1, 2010. URL: <https://www.microsoft.com/en-us/research/uploads/prod/2016/02/MSR-TR-2010-82.pdf> (cit. on p. 5).
- [5] Christopher Burges. *RankNet: A Ranking Retrospective*. Microsoft Research. July 7, 2015. URL: <https://www.microsoft.com/en-us/research/blog/ranknet-a-ranking-retrospective/> (visited on 02/15/2021) (cit. on pp. 2, 3).
- [6] Clement Renault, Marin, and Quentin de Quelen. *Meilisearch/MeiliSearch*. MeiliSearch, Feb. 20, 2021. URL: <https://github.com/meilisearch/MeiliSearch> (visited on 02/20/2021) (cit. on p. 6).
- [7] Docfetcher Development Team. *DocFetcher - Fast Document Search*. URL: <http://docfetcher.sourceforge.net/en/index.html> (visited on 02/15/2021) (cit. on p. 5).
- [8] Norbert Fuhr. “Optimum Polynomial Retrieval Functions Based on the Probability Ranking Principle”. In: *ACM Transactions on Information Systems* 7.3 (July 1, 1989), pp. 183–204. ISSN: 1046-8188. DOI: [10.1145/65943.65944](https://doi.org/10.1145/65943.65944). URL: <https://doi.org/10.1145/65943.65944> (visited on 02/15/2021) (cit. on p. 2).
- [9] Norbert Fuhr. “Probabilistic Models in Information Retrieval”. In: *The Computer Journal* 35.3 (June 1, 1992), pp. 243–255. ISSN: 0010-4620. DOI: [10.1093/comjnl/35.3.243](https://doi.org/10.1093/comjnl/35.3.243). URL: <https://doi.org/10.1093/comjnl/35.3.243> (visited on 02/15/2021) (cit. on p. 2).
- [10] Fredric C. Gey. “Inferring Probability of Relevance Using the Method of Logistic Regression”. In: *Proceedings of the 17th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*. SIGIR ’94. Berlin, Heidelberg: Springer-Verlag, Aug. 1, 1994, pp. 222–231. ISBN: 978-0-387-19889-7 (cit. on p. 2).
- [11] John Hawthorn. *Jhawthorn/Fzy*. Feb. 20, 2021. URL: <https://github.com/jhawthorn/fzy> (visited on 02/20/2021) (cit. on p. 6).
- [12] HMKCode. *Backpropagation Step by Step*. URL: <https://hmkcode.com/ai/backpropagation-step-by-step/> (visited on 02/16/2021) (cit. on p. 3).
- [13] Daniel Hunt. “A Thesis Submitted in Partial Fulfillment of the Requirements for the Degree of Master of Science in Computer Science”. In: (), p. 30 (cit. on p. 2).

## REFERENCES

---

- [14] James Aylett. *GSoCProjectIdeas/LearningtoRankStabilisation – Xapian*. Oct. 20, 2019. URL: <https://trac.xapian.org/wiki/GSoCProjectIdeas/LearningtoRankStabilisation> (visited on 02/20/2021) (cit. on p. 2).
- [15] Dr Rolf Jansen. *Cyclaero/Zettair*. Dec. 22, 2020. URL: <https://github.com/cyclaero/zettair> (visited on 02/20/2021) (cit. on p. 5).
- [16] Kalervo Järvelin and Jaana Kekäläinen. “Cumulated Gain-Based Evaluation of IR Techniques”. In: *ACM Transactions on Information Systems* 20.4 (Oct. 1, 2002), pp. 422–446. ISSN: 1046-8188. DOI: [10.1145/582415.582418](https://doi.org/10.1145/582415.582418). URL: <http://doi.org/10.1145/582415.582418> (visited on 02/20/2021) (cit. on p. 5).
- [17] Jean-Francois Dockes. *Recoll User Manual*. URL: <https://www.lesbonscomptes.com/recoll/usermanual/usermanual.html> (visited on 02/15/2021) (cit. on p. 6).
- [18] Oleh Krehel. *Abo-Abo/Swiper*. Feb. 20, 2021. URL: <https://github.com/abo-abo/swiper> (visited on 02/20/2021) (cit. on p. 6).
- [19] ktr. *Ktr0731/Go-Fuzzyfinder*. Feb. 16, 2021. URL: <https://github.com/ktr0731/go-fuzzyfinder> (visited on 02/20/2021) (cit. on p. 6).
- [20] lestrrat. *Peco/Peco*. peco, Feb. 18, 2021. URL: <https://github.com/peco/peco> (visited on 02/20/2021) (cit. on p. 6).
- [21] Tie-Yan Liu. “Learning to Rank for Information Retrieval”. In: *Foundations and Trends® in Information Retrieval* 3.3 (June 26, 2009), pp. 225–331. ISSN: 1554-0669, 1554-0677. DOI: [10.1561/15000000016](https://doi.org/10.1561/15000000016). URL: <https://www.nowpublishers.com/article/Details/INR-016> (visited on 02/15/2021) (cit. on p. 2).
- [22] Andy Lok. *Andylokandy/Simsearch-Rs*. Feb. 15, 2021. URL: <https://github.com/andylokandy/simsearch-rs> (visited on 02/20/2021) (cit. on p. 6).
- [23] Christopher D. Manning, Prabhakar Raghavan, and Hinrich Schütze. *Introduction to Information Retrieval*. New York: Cambridge University Press, 2008. 482 pp. ISBN: 978-0-521-86571-5 (cit. on p. 2).
- [24] Marty Schoch and Abhinav Dangeti. *Bleve Search Documentation*. URL: <http://blevesearch.com/> (visited on 02/20/2021) (cit. on pp. 2, 6).
- [25] 21 Nov 2017 Michael A. Alcorn Michael Alcorn. *An Introduction to Machine-Learned Ranking in Apache Solr*. Opensource.com. URL: <https://opensource.com/article/17/11/learning-rank-apache-solr> (visited on 02/15/2021) (cit. on p. 2).
- [26] Mahesh Chandra Mukkamala and Matthias Hein. “Variants of RMSProp and Adagrad with Logarithmic Regret Bounds”. In: *Proceedings of the 34th International Conference on Machine Learning - Volume 70*. ICML’17. Sydney, NSW, Australia: JMLR.org, Aug. 6, 2017, pp. 2545–2553 (cit. on p. 3).
- [27] Olympia Nicodemi, Melissa A. Sutherland, and Gary W. Towsley. *An Introduction to Abstract Algebra with Notes to the Future Teacher*. Upper Saddle River, NJ: Pearson Prentice Hall, 2007. 436 pp. ISBN: 978-0-13-101963-8 (cit. on p. 5).
- [28] Oliver Nightingale. *Olivernn/Lunr.js*. Feb. 20, 2021. URL: <https://github.com/olivernn/lunr.js> (visited on 02/20/2021) (cit. on p. 6).



## REFERENCES

---

- [29] Olly Betts and Richard Boulton. *Xapian/Xapian*. Xapian, Feb. 19, 2021. URL: <https://github.com/xapian/xapian> (visited on 02/20/2021) (cit. on p. 5).
- [30] Jeffrey M. Perkel. “Why Scientists Are Turning to Rust”. In: *Nature* 588.7836 (7836 Dec. 1, 2020), pp. 185–186. DOI: [10.1038/d41586-020-03382-2](https://doi.org/10.1038/d41586-020-03382-2). URL: <https://www.nature.com/articles/d41586-020-03382-2> (visited on 02/20/2021) (cit. on p. 2).
- [31] Philip Picton. *Neural Networks*. Basingstoke, Hampshire ; New York: Palgrave, 1994. 195 pp. ISBN: 978-0-333-94899-6 (cit. on p. 3).
- [32] Tim Roughgarden, director. *Quicksort Overview*. Jan. 28, 2017. URL: [https://www.youtube.com/watch?v=ETolcpLN7kk&list=PLEAYkSg4uSQ37A6\\_NrUnTHEKp6EkAxTma&index=25](https://www.youtube.com/watch?v=ETolcpLN7kk&list=PLEAYkSg4uSQ37A6_NrUnTHEKp6EkAxTma&index=25) (visited on 02/15/2021) (cit. on p. 5).
- [33] Enmei Tu et al. *A Graph-Based Semi-Supervised k Nearest-Neighbor Method for Nonlinear Manifold Distributed Data Classification*. June 3, 2016. arXiv: [1606.00985](https://arxiv.org/abs/1606.00985) [cs, stat]. URL: <http://arxiv.org/abs/1606.00985> (visited on 12/02/2020) (cit. on p. 5).
- [34] US National Institute of Standards and Technology. *Text REtrieval Conference (TREC) Home Page*. URL: <https://trec.nist.gov/> (visited on 02/20/2021) (cit. on p. 5).
- [35] Vik Singh. *A Comparison of Open Source Search Engines*. Vik’s Blog. July 6, 2009. URL: <https://partyondata.com/2009/07/06/a-comparison-of-open-source-search-engines-and-indexing-twitter/> (visited on 02/20/2021) (cit. on pp. 2, 5).
- [36] vz. *Go-Ego/Riot*. ego, Feb. 20, 2021. URL: <https://github.com/go-ego/riot> (visited on 02/20/2021) (cit. on p. 6).
- [37] S. K.M. Wong and Y. Y. Yao. “Linear Structure in Information Retrieval”. In: *Proceedings of the 11th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*. SIGIR ’88. New York, NY, USA: Association for Computing Machinery, May 1, 1988, pp. 219–232. ISBN: 978-2-7061-0309-4. DOI: [10.1145/62437.62452](https://doi.org/10.1145/62437.62452). URL: <https://doi.org/10.1145/62437.62452> (visited on 02/14/2021) (cit. on p. 2).
- [38] Yuri Schapov, Andrew Aksyonoff, and Ilya Kuznetsov. *Sphinxsearch/Sphinx*. Sphinx Technologies Inc, Feb. 18, 2021. URL: <https://github.com/sphinxsearch/sphinx> (visited on 02/20/2021) (cit. on p. 5).
- [39] Jinzhou Zhang. *Lotabout/Skim*. Feb. 19, 2021. URL: <https://github.com/lotabout/skim> (visited on 02/20/2021) (cit. on p. 6).