

BACHELOR INFORMATICA



UNIVERSITY OF AMSTERDAM

Creating an experimental environment for participating in TREC Fair Ranking Track

René Emmaneel

December 3, 2020

Supervisor(s): Fatemeh Sarvi

Signed:

Abstract

The TREC Fair Ranking Track is an annual academic search task with the stated goal of developing a benchmark for evaluating retrieval systems in terms of fairness. To participate we created a framework. [1]

Contents

1	Introduction	7
2	Theoretical background	9
2.1	Ethical aspect	9
2.2	Measuring fairness	9
2.2.1	Measuring author exposure	9
2.2.2	Measuring author relevance	10
2.2.3	Measuring group fairness	10
2.2.4	Measuring relevance	10
2.3	Learning to rank	11
2.3.1	Feature extraction	12
2.3.2	Ranking Support Vector Machine	13
2.4	TREC Fair Ranking Track definition	14
2.4.1	Semantic Scholar Open Research Corpus	14
2.5	15
3	Creating the framework	17
4	Experiments	19
5	Conclusions	21
	Bibliography	23

Introduction

Ranking systems allow items to be seen by the users, where the exposure each item gets is significantly determined by the position in the ranking. Traditionally a ranking system would order items based on their relevance so that the system is of maximum utility to the user. However, modern information access systems often influence both the consumer and the producer of the content that they serve. Examples are hiring platforms where it is important that the system has a high utility for both the employer and the job-seeker. But also more traditional environments such as music, book or video recommendations can be considered two-sided because the indirect matching of users to content creators.

The ranking system allows exposure to the producers, which exposes risks to companies such as Google that hold a lot of power, and has already faced legal challenges such as the European Union antitrust violation fine [2]. While it is unlikely that a universal definition of fairness can be created, there are various attempts to measure unfairness and devise various fair ranking algorithms to ensure that the received attention of a given subject is approximately equal to its deserved attention [1].

To evaluate different fair ranking algorithms, the TREC Fair Ranking track was created. The stated goal was to develop a benchmark for evaluating retrieval systems in terms of fairness, as well as releasing a dataset for benchmarking fair ranking algorithms. The TREC 2020 Track was an academic search task, where a set of academic article abstract and queries were submitted to an academic search engine. The central goal is to provide fair exposure to different groups of authors, where the group definitions can be arbitrary. Various sets of data has been provided, as well as a detailed evaluation protocol. The data of the 2019 track includes the Semantic Scholar Open Corpus, which consists of data describing various papers. It also includes query data, and query sequences.

To participate in the TREC fair ranking track, an environment needs to be build for the entire process. This process includes handling of the input data, implementing the evaluation algorithm, making it modular for various uses, testing of the framework and implementing learning to rank machine learning algorithms. To design and implement the framework the following research question will be answered:

How can a modular framework be constructed to create an experimental environment for participating in the TREC fair ranking track.

TODO: Structure van project

Theoretical background

2.1 Ethical aspect

2.2 Measuring fairness

Fairness can be subjective and differ for various use cases, and as such the TREC fair ranking track has provided a fairness definition for the academic search task. The goal of the fair ranking track is to provide fair exposure to different groups of authors, where the group definition may be arbitrary. Before describing the fairness definition we have to describe the way exposure can influence the discoverability for a certain group, by defining formulas for both the exposure and the relevance an author has.

2.2.1 Measuring author exposure

The exposure an author has is related to positions of the papers he contributed to, where the first ranked paper has the highest exposure and lower ranked papers have lower exposure. The formula used to measure exposure is the Expected Reciprocal Rank metric [3]. This metric describes the idea that a user is more likely to stop browsing if they have seen a highly relevant document. The user will look at each document in order, and will have a chance to stop browsing depending on the relevance of the document. It also implements an abandonment probability γ , which represent the chance the user will examine the next document, meaning that there is a $1 - \gamma$ chance of abandoning the search. The exposure of the document on place i given the Expected Reciprocal Rank metric is as follows.

$$e_i = \gamma^{i-1} \prod_{j=1}^{i-1} (1 - f(r_j))$$

where $f(r_j)$ is a function to transform the relevance of document d_j into a probability of stopping the search after seeing the document. For the fair ranking track the relevance is either 0 or 1, and the function is given as $f(r_d) = 0.7 \cdot r_d$. The continuation probability is given as $\gamma = 0.5$.

Because the task is to provide fair exposure to authors instead of documents, we have to calculate the cumulative exposure of the documents which each author has contributed to. Given a ranking of documents π , we can calculate the exposure of author a as follows.

$$e_a^\pi = \sum_{i=1}^n (e_i) \cdot I(\pi_i \in \mathcal{D}_a)$$

Where π_i is the document on position i , and \mathcal{D}_a is the documents which include a as an author. Because the final result of the track consist of a series of different rankings, we can calculate the amoritized exposure of author a as,

$$e_a = \sum_{\pi \in \Pi} e_a^\pi$$

where Π is the sequence of all rankings.

2.2.2 Measuring author relevance

The author relevance is simply the sum of the stop probability for each document where a is an author

$$r_a^\pi = \sum_{d \in \mathcal{D}_a} f(r_d)$$

And the amortized relevance for an author is defined as

$$r_a = \sum_{\pi \in \Pi} r_a^\pi$$

2.2.3 Measuring group fairness

We can assume that each author is part of exactly one group. The goal of the fair ranking track to provide fair exposure for each group relatively to the relevance for each group. Let \mathcal{G} be the set of all groups and \mathcal{A}_g be the set of all authors in group g . The group exposure and relevance are defined as,

$$\mathcal{E}_g = \frac{\sum_{a \in \mathcal{A}_g} e_a}{\sum_{g' \in \mathcal{G}} \sum_{a \in \mathcal{A}_{g'}} e_a}$$

$$\mathcal{R}_g = \frac{\sum_{a \in \mathcal{A}_g} r_a}{\sum_{g' \in \mathcal{G}} \sum_{a \in \mathcal{A}_{g'}} r_a}$$

Note that $\sum_{g \in \mathcal{G}} \mathcal{E}_g = 1$ and $\sum_{g \in \mathcal{G}} \mathcal{R}_g = 1$. We can use these metrics to calculate the deviation between the exposure and relevance for each group.

$$\Delta_g = \mathcal{E}_g - \mathcal{R}_g$$

Groups should receive exposure proportional to relevance, meaning that a perfectly fair model provides a deviation of zero for each group. We can compute the fair exposure using the square norm.

$$\Delta = \sqrt{\sum_{g \in \mathcal{G}} \Delta_g^2}$$

2.2.4 Measuring relevance

We measure the relevance of a ranking using the Expected Reciprocal Rank metric, which is the same metric as used in the fairness measurement. Recall that the exposure of document i is defined as,

$$e_i = \gamma^{i-1} \prod_{j=1}^{i-1} (1 - f(r_j))$$

We multiply the exposure of each document with the stop probability to calculate the total utility of a ranking π as follows,

$$u^\pi = \sum_{i=1}^n e_i \cdot f(r_i)$$

The average utility of all rankings, $U = \frac{1}{|\Pi|} \sum_{\pi \in \Pi} u^\pi$, will be the final relevance metric.

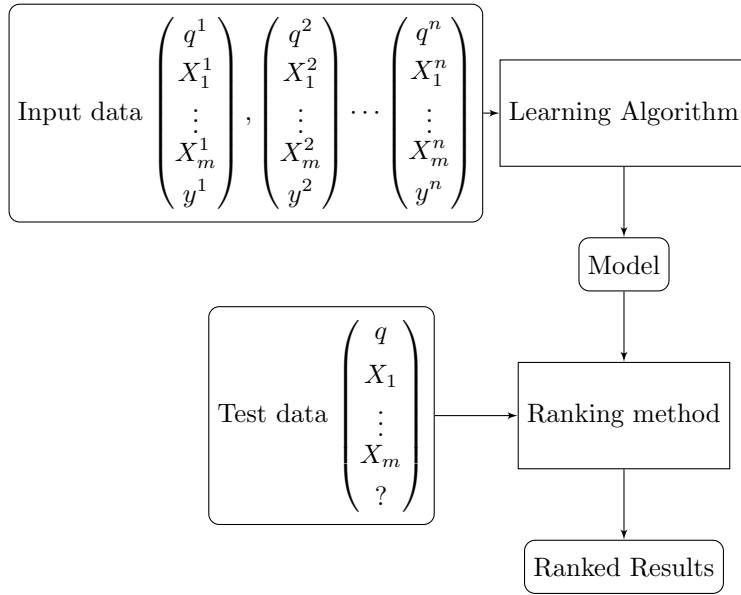


Figure 2.1: Learning to rank framework

2.3 Learning to rank

Before looking at a way to create a model for learning to fairly rank, we will first propose a model for the standard learning to rank model on the TREC fair ranking track data. Learning to rank, when applied to document retrieval, is a task as follows. In training a set of queries and documents for each query is provided, as is the relevance judgement for each document. Using this data, a ranking function is created such that the model can predict a ranked list of documents for a query. The two major approaches for creating a ranking function is either a non-learning approach or a learning to rank approach. Learning approaches for information retrieval (IR) have been widely studied and shown to be useful for ranking articles. [4]

A learning to rank framework is systematically represented in Figure 2.1. As can be seen in the figure, a training set is needed to train the ranking model. The training data consist typically of n training queries, and m documents for each query. Each document is represented by a feature vector, which is a list of values calculated before training based on the document and query. Each query also has a relevance judgement y , which can for example be the order in which the documents are ranked, or can be a vector representing the relevance for each document based on the query. After training a model has been created, which can be used to rank a set of given documents based on the query.

There are three main learning to rank algorithms. The pointwise approach, which learns a function for directly predicting the relevance for each document for a query. The pairwise approach, which learns a function that takes in a pair of documents and a query and tries to predict the document that is most relevant for the given query. Lastly there is the listwise approach, which learns a function to directly rank a set of input documents based on a query.

In this thesis we have decided to create a pairwise ranking algorithm. This decision was made based on the training data available in the TREC fair ranking track. As described in [section], the training data consist of queries and a certain amount of documents per query, with a boolean relevance score for each document. By matching documents that are relevant to a certain query with documents that are irrelevant to the same query, we can easily create document pairs for the pairwise approach. By creating these document pairs, we can reduce the problem to a classification problem to decide which pair of the document is more relevant. We will first discuss

ID	Feature Description
1	$\sum_{q_i \in q \cap d} TF(q_i, d)$ in title
2	$\sum_{q_i \in q} IDF(q_i)$ in title
3	$\sum_{q_i \in q \cap d} TF(q_i, d) \cdot IDF(Q_i)$ in title

Table 2.1: Features used in the rankSVM model

the various features extracted for this model, after which we will discuss the learning algorithm itself.

2.3.1 Feature extraction

In learning to rank, each query-document pair is represented by a feature vector, and each feature gives an indication to the relevance of the document to the query. For example one feature could be the frequency of the query words in the abstract of the document. In this section we will discuss the various types of features, and which features we have used in constructing the feature vectors for the learning to rank framework.

Features can be categorized in 3 distinct categories. The query-document category means that the feature is dependant on both the feature and document, examples are term frequency and BM25. The other two categories are only dependant on the query or the document respectively.

The different features as used in the rankSVM model for the TREC fair ranking track are described in table 2.1. The explanation about these features are listed below.

- Features (A-B) are the term frequency (TF) features. The features will count the amount of occurrences of each term in the query in the respective field on the document.
- Features (C-D) are the inverse document frequency (IDF) features. The inverse document frequency function is a measure about how much information a given word provides. The function is defined as follows.

$$IDF(t) = \log \frac{N}{n(t)}$$

with N being the total number of documents, and n(t) being the amount of documents with the term t

- Features (E-F) are the term frequency - inverse document frequency (TF-IDF) features. For every term of the query, the TF and IDF of the specific term is multiplied. This feature is correlated by a high frequency of a term, which has a low document frequency.
- Features (G-H) are the BM25 features. This is a TF-IDF-like retrieval function, first developed for TREC-3 [5]. The function used is defined as follows.

$$BM25(q, d) = \sum_{q_i \in q \cap d} IDF(q_i) \cdot \frac{TF(q_i, d) \cdot (k + 1)}{TF(q_i, d) + k \cdot (1 - b + b \cdot \frac{length(d)}{avglength()})}$$

where length(d) is the amount of words in document d, and avglength() is the average amount of words in all documents. k and b are free variables, which can be set to various values depending on the dataset, but are set to k=1.2 and b=0.75, which experiments have shown to be reasonable values [6].

- Features (G-H) are the language model information retrieval (LMIR) features. LMIR is a statistical language model, with the goal of predicting the probability of the document's language model generating the terms of the query [7]. The way LMIR aims to find this

probability is by combining the various probabilities of a certain term in the query being generated by the document's model ($P(t | d)$), which gives the equation as follows.

$$P(q | d) = \prod_{i=1}^M P(t_i | d)$$

where t_1, \dots, t_M are the terms in query q and $P(t | d)$ is the document's language model. There are many variants of LMIR, which mainly differ in the use of the document's language model, of which we use three different methods as described by Zhai and Lafferty [8].

- Features (G-H) are the Jelinek Mercer LMIR features. The language model is based on a combination of the percentage of occurrence of a term in the given document, as well as the percentage of occurrence of a term in all documents combined. The document's language model can then be constructed as follows.

$$P(t_i | d) = (1 - \lambda) \frac{TF(t_i, d)}{LEN(d)} + \lambda \frac{TOT_TF(t_i, C)}{TOT_LEN(C)}$$

where $TOT_TF(t_i, C)$ is the term frequency of term i in the entire corpus, $TOT_LEN(C)$ is the total amount of terms in the corpus and λ is the smoothing factor, which is set to 0.1 in our framework.

- Features (G-H) are the Dirichlet LMIR features.

$$P(t_i | d) = \frac{TF(t_i, d) + \mu TOT_TF(t_i, C)}{TOT_LEN(C) + \mu}$$

with μ is 2000

- Features (G-H) are the Absolute discount LMIR features.

$$P(t_i | d) = \frac{\max(TF(t_i, d) - \delta, 0)}{TOT_TF(t_i, C)} + \sigma \frac{TOT_TF(t_i, C)}{TOT_LEN(C)}, \text{ with } \sigma = \delta \frac{UNIQUE(d)}{LEN(d)}$$

with μ is 0.7

2.3.2 Ranking Support Vector Machine

Support Vector Machine (SVM) algorithm is a supervised learning model with the goal of categorizing data represented as points in space [9]. Given a set of training examples, each being labeled in one of the two possible categories, the SVM algorithm builds a model to categorize new examples in one of the two categories. The training examples are represented as a combination of a label and a feature vector, and the model output is represented as the maximum-margin hyperplane that divides the points of one category and the other, so that the distance from the hyperplane and the closest point from either group is maximized. Utilizing the SVM algorithm to categorize documents as relevant or not has been used to create information retrieval models before [10]. However the main goal of a ranking system is to predict the relative order of documents based on a query, while the above described pointwise approach merely learns the degree to which a document is relevant based on the query. We instead use an algorithm that directly takes the goal of ordering more relevant documents higher into account.

Ranking SVM is a variant of the above described SVM algorithm. The goal of the algorithm is to learn a model which, given two documents and a query, learns which document is more relevant [11]. The main benefit of the pairwise approach is that creating a model directly creates a way to sort documents by relevance by utilizing the model. The model that will be created now is purely classification based, as in the goal of the model is to predict which document is more relevant, and not by how much.

2.4 TREC Fair Ranking Track definition

2.4.1 Semantic Scholar Open Research Corpus

The Semantic Scholar Open Research Corpus is a large archive containing metadata of research papers [12]. It contains 47GB data in JSON format, containing information about papers in many different disciplines. A large dataset is needed for sufficient training sample size. The following data is available for most papers.

- **S2 Paper ID**
- **DOI**
- **Title**
- **Year**
- **Abstract**
- **Venue name**
- **Authors** (resolved to author IDs)
- **Inbound citations** (resolved to S2 paper IDs)
- **outbound citations** (resolved to S2 paper IDs)

For the 2019 TREC Fair ranking Track there was no semantic corpus subset available, however a training set was provided. Using this training set we extract the semantic corpus for the used documents using the available API, which gives the data for a specific document given the S2 Paper ID.

Using this data, we can extract the corpus into smaller csv files containing the metadata of each paper and author. The TREC fair ranking track provided the following 3 csv files, however we also created functionality to extract the following information ourself given the Semantic Scholar Open Research Corpus. The paper metadata is as follows.

paper_metadata.csv:

- **S2 Paper ID**
- **Title**
- **Year**
- **Venue**
- **Number of inbound citations**

Using the same data, we can also extract metadata information for each author of the papers. During the extraction, we count for each author the amount of papers they worked on and the amount of citations each papers has. We use this information for calculating the i10 and H_index scores, as explained below.

author_metadata.csv:

- **Author ID**
- **Name**
- **Citation count:** sum of all citations of each published paper
- **Paper count**

- **i10:** amount of publications with atleast 10 citations
- **H_index:** The maximum value of h such that the given author has published h papers that have each been cited at least h times.
- **H_class:** The class is 'H' if the author has a H_index of atleast 10, and the class is 'L' if the H_index is less than 10

To combine the papers and authors, we need a third linker file, which combined the paper ID with each corresponding author ID, and the position in which they were attributed.

- **S2 Paper ID**
- **Author ID**
- **Position**

2.5

Creating the framework

Suspendisse vitae elit. Aliquam arcu neque, ornare in, ullamcorper quis, commodo eu, libero. Fusce sagittis erat at erat tristique mollis. Maecenas sapien libero, molestie et, lobortis in, sodales eget, dui. Morbi ultrices rutrum lorem. Nam elementum ullamcorper leo. Morbi dui. Aliquam sagittis. Nunc placerat. Pellentesque tristique sodales est. Maecenas imperdiet lacinia velit. Cras non urna. Morbi eros pede, suscipit ac, varius vel, egestas non, eros. Praesent malesuada, diam id pretium elementum, eros sem dictum tortor, vel consectetur odio sem sed wisi.

CHAPTER 4

Experiments

CHAPTER 5

Conclusions

Bibliography

- [1] Asia J. Biega, Krishna P. Gummadi, and Gerhard Weikum. Equity of attention: Amortizing individual fairness in rankings. *CoRR*, abs/1805.01788, 2018.
- [2] Mark Scott. Google fined record \$2.7 billion in e.u. antitrust ruling. *New York Times*, 2017.
- [3] Olivier Chapelle, Donald Metzler, Ya Zhang, and Pierre Grinspan. Expected reciprocal rank for graded relevance. In *Proceedings of the 18th ACM Conference on Information and Knowledge Management*, CIKM '09, page 621–630, New York, NY, USA, 2009. Association for Computing Machinery.
- [4] Tao Qin, Tie-Yan Liu, Jun Xu, and Hang Li. Letor: A benchmark collection for research on learning to rank for information retrieval. 2016.
- [5] S.E. Robertson, S. Walker, S. Jones, M.M. Hancock-Beaulieu, and M. Gatford. Okapi at trec-3. pages 109–126, 1996.
- [6] Hinrich Schütze Christopher D. Manning, Prabhakar Raghavan. An introduction to information retrieval. *Cambridge University Press*, page 270, 2009.
- [7] Tao Tao, Xuanhui Wang, Qiaozhu Mei, and ChengXiang Zhai. Language model information retrieval with document expansion. In *Proceedings of the Human Language Technology Conference of the NAACL, Main Conference*, pages 407–414, New York City, USA, June 2006. Association for Computational Linguistics.
- [8] Chengxiang Zhai and John Lafferty. A study of smoothing methods for language models applied to ad hoc information retrieval. In *Proceedings of the 24th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, SIGIR '01, page 334–342, New York, NY, USA, 2001. Association for Computing Machinery.
- [9] Corinna Cortes and Vladimir Vapnik. Support-vector networks. In *Machine Learning*, pages 273–297, 1995.
- [10] Ramesh Nallapati. Discriminative models for information retrieval. In *Proceedings of the 27th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, SIGIR '04, page 64–71, New York, NY, USA, 2004. Association for Computing Machinery.
- [11] Thorsten Joachims. Optimizing search engines using clickthrough data. *KDD '02*, page 133–142, New York, NY, USA, 2002. Association for Computing Machinery.
- [12] Waleed Ammar, Dirk Groeneveld, Chandra Bhagavatula, Iz Beltagy, Miles Crawford, Doug Downey, Jason Dunkelberger, Ahmed Elgohary, Sergey Feldman, Vu Ha, Rodney Kinney, Sebastian Kohlmeier, Kyle Lo, Tyler Murray, Hsu-Han Ooi, Matthew Peters, Joanna Power, Sam Skjonsberg, Lucy Lu Wang, Chris Wilhelm, Zheng Yuan, Madeleine van Zuylen, and Oren Etzioni. Construction of the literature graph in semantic scholar. In *NAACL*, 2018.