



Leibniz
Universität
Hannover



Deconstructing Ranking Abilities of Language Models

Gottfried Wilhelm Leibniz Universität Hannover
Fakultät für Elektrotechnik und Informatik
Institut für verteilte Systeme
Fachgebiet Wissensbasierte Systeme
Forschungszentrum L3S

Thesis by
Fabian Beringer

First examiner: xxxx
Second examiner: xxxx
Advisors: Abhijit Aanand
Jonas Wallat

Hannover, xx.xx.2022

Abstract

...

Contents

| | | |
|----------|--|-----------|
| 1 | Introduction | 5 |
| 1.1 | Motivation | 5 |
| 1.2 | Problem Statement | 6 |
| 1.3 | Contribution | 6 |
| 1.4 | Thesis Outline | 6 |
| 2 | Foundations | 7 |
| 2.1 | Information Retrieval - Ranking Text | 7 |
| 2.1.1 | TF-IDF | 8 |
| 2.1.2 | BM25 | 9 |
| 2.2 | Machine Learning | 9 |
| 2.3 | Deep Learning | 9 |
| 2.3.1 | Deep Neural Networks | 10 |
| 2.3.2 | Optimization | 10 |
| 2.3.3 | Regularization | 11 |
| 2.3.4 | Learning to Rank | 11 |
| 2.4 | Transformer Models | 11 |
| 2.4.1 | Architecture | 11 |
| 2.4.2 | Pre-Training - BERT | 11 |
| 2.5 | Probing | 11 |
| 3 | Previous Work | 12 |
| 4 | Datasets | 13 |
| 4.1 | Probing | 13 |
| 4.1.1 | BM25 Prediction | 13 |
| 4.1.2 | Term Frequency Prediction | 13 |
| 4.1.3 | Named Entity Recognition | 13 |
| 4.1.4 | Semantic Similarity | 13 |
| 4.1.5 | Coreference Resolution | 13 |
| 4.1.6 | Fact Checking | 13 |
| 4.1.7 | Relevance Estimation | 13 |
| 4.2 | Multitask Learning | 13 |

| | | |
|----------|-------------------------------|-----------|
| 5 | Approach | 14 |
| 5.1 | Methodology | 14 |
| 5.2 | Experimental Setup | 14 |
| 5.3 | Evaluation Measures | 14 |
| 5.3.1 | MDL | 14 |
| 5.3.2 | Compression | 14 |
| 5.3.3 | F1 | 14 |
| 5.3.4 | Accuracy | 14 |
| 5.3.5 | Ranking | 14 |
| 6 | Results | 15 |
| 7 | Conclusion | 16 |
| 7.1 | Future Work | 16 |
| | Plagiarism Statement | 17 |
| | List of Figures | 18 |
| | List of Tables | 19 |
| | Bibliography | 20 |

1 Introduction

Due to large successes in computer vision, natural language processing (NLP) and many other fields, machine learning (ML) and in particular deep learning (DL) have gained popularity in the research community and among practitioners. Both as object of research and in real world applications these methods have been widely adopted and shown to be effective for several use cases.

One major field, in which ML and DL has been playing an increasingly important role is information retrieval (IR). Not only has "Neural Information Retrieval" [MC18] been acknowledged as promising field of research within the IR research community, but it is already present in modern web IR systems used by google, amazon, netflix and other major companies, which do heavily rely on ML for delivering relevant content to their users.

While these systems most certainly do not solely operate on raw text, when it comes to finding relevant documents for a given natural language query, the task of text ranking still lies at the heart of the problem.

In 2018 NLP experienced a major breakthrough when [Dev+19]. Core idea pretrain, finetune... While progress in pretraining lms and using for downstream task have been made before this that, this was first time the model achieved sota by large margin on variety of nlp benchmarks. and is often referred to as "the imagenet moment of NLP".

subfields including text ranking of nlp became aware and started to use.

text ranking bert considered breakthrough

- modern ml system often black boxes

- need to understand: mitigate biases, make clearer to use, make better ... we deal with text ranking - advances in nlp, large models, hard to understand

1.1 Motivation

Nowadays, modern approaches for NLP often rely on large language models. These models are first pre-trained on huge amounts of text data and then fine-tuned on a smaller, task specific dataset. Although this approach has shown great effectiveness compared to traditional approaches, due to their size and complexity, these language models are mostly being treated as black boxes.

This is where several issues arise.

goal shine more light on how probing specifically targeted to IR use knowledge to improve

1.2 Problem Statement

I am still alive. - shine light on knowledge dist - research question

1.3 Contribution

1.4 Thesis Outline

2 Foundations

2.1 Information Retrieval - Ranking Text

Ranking is an integral part of the information retrieval (IR) process. Typically the general IR problem can be formulated as follows: A user with a need for information expresses this information need through formulation of a query. Now given the query and a collection of documents, the IR system's task is to provide a list of documents that satisfy the user's information need. Further, the retrieved documents should be sorted by relevance w.r.t. the user's information need in descending order, i.e. the documents considered most relevant should be at the top of the list.

While from this formulation only, the task might appear simple, there are several caveats to look out for when it comes to ranking. For instance, there is no restriction on the structure of the query. While we might expect a natural language question like "What color are giraffes?" a user might decide to enter a keyword query like "giraffes color". The same applies to documents: Depending on the corpus we are dealing with, the documents might be raw text, structured text like html or even another type of media e.g. image, audio or a combination thereof.

Another possible issue is a mismatch in information need of the user and the corresponding query. Even if we find a perfect ordering of documents with respect to the query, we can not know for certain that the query actually reflects the user's information need. The user might not even know exactly what they're looking for until discovery through an iterative process, i.e. the information need is fuzzy and can not be specified through an exact query from the beginning on.

Further, a query might require additional context information in order for an IR system to find relevant documents. For example, depending on the time at which a query is prompted, the correct answer might change: "Who is president?" should return a different set of documents, as soon as a new president has been elected. Also, since not specified further, it is up to interpretation which country's president the user is interested in and might depend on their location. In addition, even the corpus might not be static either and change or grow over time, e.g. web search has to deal with an ever-growing corpus: the internet.

While this list of issues is not comprehensive, at this point the complexity of the ranking problem should have become apparent.

Because this work focuses on the ranking of text in the context of web search, we will now give a formal definition with that scenario in mind:

Given a set of $|Q|$ natural language queries $Q = \{q_i\}_{i=1}^{|Q|}$ and a collection of $|D|$ documents $D = \{d_i\}_{i=1}^{|D|}$, we want to find a scoring function $s : Q \times D \rightarrow \mathbb{R}$, such that for any query $q \in Q$ and documents $d, d' \in D$, it holds true that $s(q, d) > s(q, d')$ if d is more relevant w.r.t q than d' .

To give the reader a more concrete idea and as we are going to build upon it throughout this work, we will now discuss two traditional approaches to text retrieval which, unlike neural retrieval, are based on exact matching, meaning query and document terms are compared directly.

2.1.1 TF-IDF

Term Frequency - Inverted Document Frequency weighting (TF-IDF), is a classical ranking approach that, given a query, assigns a relevance score to each document based on two assumptions:

1. A document is relevant if terms from the query appear in it often.
2. A document is relevant if the terms shared with the query are also rare in the collection.

From these assumptions, two metrics are derived:

1. Term-Frequency

$$w_{t,d} = \begin{cases} 1 + \log \text{tf}_{t,d} & \text{if } \text{tf}_{t,d} > 0 \\ 0 & \text{otherwise} \end{cases} \quad (2.1)$$

where $\text{tf}_{t,d}$ is the count of term t in document d . The logarithmic scaling is motivated by the idea that a document does not linearly become more relevant by the number of terms in it: A document containing a term 10 times more often doesn't necessarily mean it is 10 times more important, e.g. the document might just be very long document and contain more words in general. Note that this is just one possible normalization scheme out of many.

2. Inverted Document Frequency

$$\text{idf}(t, d) = \log \frac{|D|}{\text{df}_t} \quad (2.2)$$

where df_t counts the number of documents that a term occurs in over the full corpus. This way, terms that occur less frequent relative to the corpus size will receive a high IDF score and those that are more frequent a lower score.

To compute TF-IDF we can simply sum over the product of TF and IDF for each term in the query to produce a relevance score:

$$\text{score}(q, d) = \sum_{t \in q} w_{t,d} \times \text{idf}_t \quad (2.3)$$

Alternatively, vector space idf vector

2.1.2 BM25

2.2 Machine Learning

Machine learning can be described as a set of statistical methods, for automatically recognizing and extracting patterns from data. Typically we can distinguish between two main types of machine learning: Supervised learning and unsupervised learning.

In the case of supervised learning, we have a set of training instances $X = \{x_i\}_{i=0}^N$ and corresponding labels $Y = \{y_i\}_{i=0}^N$, assigning a certain characteristic to each data point. For example, this characteristic might be a probability distribution over a set of classes or a regression score. Now given the training data and labels, the goal is to find a hypothesis that explains the data, such that for unseen data points $x' \notin X$, the labels $y' \notin Y$ can be inferred automatically. If each y_i represents one or more categories from a fixed set of classes $C = \{c_i\}_{i=1}^{|C|}$, this is called a classification problem.

In contrast, in unsupervised learning there is no access to any labels whatsoever. Characteristics of the data need to be learned solely from the data X itself. Examples for this include clustering where X is clustered into groups, representation learning which usually tries to find vector representations for X , as well as dimensionality reduction that, if each X is already a vector, tries to compress them into more compact but still informative representations.

That being said, the separating lines between supervised and unsupervised learning are blurry. Especially with the emergence of semi-supervised approaches and "end2end" representation learning, modern ML methods often integrate parts of both.

2.3 Deep Learning

Deep learning is a subfield of ML that makes use of a class of models called Deep Neural Networks (DNN). Typically DNNs find application in the supervised learning scenario and are often used for classification tasks. In the following we explain the basic mechanisms of DNNs and common approaches to train them.

2.3.1 Deep Neural Networks

In essence, a Deep Neural Network (DNN) is a function approximator that applies a series of non-linear transformations to its inputs, in order to produce an output. In the simplest form, an input vector x is multiplied by a single weight matrix, a bias vector is added and the resulting vector is passed through a non-linear activation function σ .

$$\text{FFN}(x) = \sigma(Wx + b) \quad (2.4)$$

where W and b are both learnable parameters. This model is commonly referred to as single layer feed-forward neural network (FFN) or single layer perceptron.

Because a single layer FFN is limited to problems that require linear separation when it comes to classification, in order to learn more complex decision boundaries, multiple layers can be applied in sequence.

$$\begin{aligned} h^{(0)} &= \sigma^{(0)}(W^{(0)}x + b^{(0)}) \\ h^{(1)} &= \sigma^{(1)}(W^{(1)}h^{(0)} + b^{(1)}) \\ &\vdots \\ y^{\text{out}} &= \sigma^{(N)}(W^{(N)}h^{(N-1)} + b^{(N)}) \end{aligned} \quad (2.5)$$

2.3.2 Optimization

The arguably most common way for optimizing a neural network is the gradient descent algorithm and its variants.

Gradient Descent

batch and stochastic

Mini-Batch Gradient Descent

Backpropagation

Momentum

Adam

2.3.3 Regularization

Weight Decay

Dropout

2.3.4 Learning to Rank

2.4 Transformer Models

2.4.1 Architecture

2.4.2 Pre-Training - BERT

2.5 Probing

3 Previous Work

4 Datasets

4.1 Probing

To assess the distribution of knowledge across layers of a BERT model, we design a set of tasks that require information on different abstraction levels in order to be solved.

For each task, we generate a dataset by sampling instances from the MSMARCO validation and test set and automatically infer targets. In the following we provide a list of all probing tasks and details on the corresponding dataset generation process.

4.1.1 BM25 Prediction

4.1.2 Term Frequency Prediction

4.1.3 Named Entity Recognition

4.1.4 Semantic Similarity

4.1.5 Coreference Resolution

4.1.6 Fact Checking

4.1.7 Relevance Estimation

4.2 Multitask Learning

5 Approach

5.1 Methodology

5.2 Experimental Setup

5.3 Evaluation Measures

5.3.1 MDL

5.3.2 Compression

5.3.3 F1

5.3.4 Accuracy

5.3.5 Ranking

MAP

MRR

NDCG

Precision

6 Results

7 Conclusion

7.1 Future Work

Plagiarism Statement

I hereby confirm that this thesis is my own work and that I have documented all sources used.

Hannover, xx.xx.2022

(Fabian Beringer)

List of Figures

List of Tables

Bibliography

- [Dev+19] Jacob Devlin et al. “BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding”. In: *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*. Minneapolis, Minnesota: Association for Computational Linguistics, June 2019, pp. 4171–4186. URL: <https://www.aclweb.org/anthology/N19-1423>.
- [MC18] Bhaskar Mitra and Nick Craswell. “An Introduction to Neural Information Retrieval”. In: *Foundations and Trends® in Information Retrieval* 13.1 (2018), pp. 1–126. URL: <https://www.microsoft.com/en-us/research/publication/introduction-neural-information-retrieval/>.