# Passage Ranking: An Information Retrieval Problem

## NLP Project – Vector

FirstName Surname[†]
Department Name
Institution/University Name
City State Country
email@email.com

FirstName Surname
Department Name
Institution/University Name
City State Country
email@email.com

FirstName Surname
Department Name
Institution/University Name
City State Country
email@email.com

## ABSTRACT

In this report, we explore the task of document relevancy ranking using various evaluation metrics, including BM25, NDCG, and Average Precision. We have based this particular work to implement and test the concept of Information Retrieval (IR) to develop models that return relevant passages given a query. The experiment and evaluation focuses on comparing the performance of three different Machine Learning models: Logistic Regression, LambdaMART, and a Feed-forward Neural Network architecture. We used a dataset of documents and corresponding queries to evaluate the models' effectiveness in ranking the documents by their relevance to the query. Our results indicate that the Feed-forward Neural Network architecture outperforms the other models in terms of all three evaluation metrics. These findings suggest that the Feed-forward Neural Network architecture holds great potential for improving search results and enhancing the user experience in document relevancy ranking. It also shows promising growth in the research of accurate Information Retrieval, Information Extraction and Document Ranking.

## KEYWORDS

Siemese Convolutional Neural Network, LambdaMART, Information Retrieval, BERT, GloVE Embeddings, Average Precision, BM25, NDCG, RankNET, LambdaRank

## PREVIOUS WORK

The task of Information retrieval holds document relevancy and page ranking very important. It has been seen that several metrics of evaluation have been brought up and tested to assess the effectiveness of models given a dataset of query-document pair. The studies show experiments starting from the simplest Machine Learning models such as Logistic Regression and moving on to SVMs, Random Forests, Neural Networks and Transformer-based architectures. In all evaluations, given a large dataset, Neural Network architectures have shown superior performance. The results have shown the n-fold growth in performance of ML models as modern architectures are used.

In [1] the authors propose a duet neural network architecture that composes of one that uses local representation and the other that uses learned distributed representations. The two networks are trained jointly and demonstrate significant results in a Web search ranking task based on queries. [2] uses the setting of Learning Under Privileged Information (LUPI) that uses additional data alongwith the training data. They extend the RankNet model to the LUPI paradigm and demonstrate significant performance growth. A shift towards a sequence-to-sequence pretrained model is introduced in [3] which demonstrates the T5 model against BERT. It is seen that in a data-rich encoder-driven networks, also the sequence-driven networks, do perform well but the T5 model greatly outperforms BERT in a data-poor regime. Authors in [4] train a combination of the Siemese Convolutional Neural Network Encoder with the RankNet ranking model to provide an end-to-end training titled as the ConvRankNet which outperformed all existing feature-based models.

Overall, these studies illustrate the effectiveness of machine learning models, particularly deep learning architectures, in document relevancy ranking. The results demonstrate the potential for improving search results and enhancing the user experience in various information retrieval applications.

## CORPUS DESCRIPTION

### 1. Dataset Structure

The dataset is comprehensive with five distinct files with different format but a single unanimous structure. The first file, test-queries.tsv, is a tab-separated file that contains queries within the test set. Each row in this file is comprised of a query ID (qid) and its corresponding query text.

Additionally, the dataset includes a candidate_passages_top1000.tsv file which contains initial rankings that contain 1000 passages for each of the given queries as it was in the first part of the assignment. This file is formatted as <qid pid query passage>, where qid is the query ID, pid is the ID of the passage retrieved, query is the query text, and passage is the passage text. Figure 1 displays a select few rows from this file.

For the purposes of training and validation, the dataset is also composed of train_data.tsv and validation_data.tsv files. The experimenting students will train their models on the training set and evaluate the model's performance on the validation set.

Within these datasets, an additional relevance column is included indicating the relevance of the passage to the query, which is required during the training and validation processes.

## 2. Feature Generation

### 2.1 Overview

To prepare the text data for our analysis, we first load the pre-trained GloVe embeddings into memory. GloVe embeddings are vector representations of words in a high-dimensional space, which are often used as features in NLP tasks. We then define a function to generate word embeddings for the query and passage texts using the pre-trained embeddings dictionary. The function tokenizes the text and maps each word to its corresponding embedding in the dictionary, and then takes the mean of all the embeddings in the text to obtain a single embedding for the text.

Once the embeddings are generated, we proceed to normalize the features using a Standard Scaler function. This normalization step is important because it ensures that each feature has zero mean and unit variance, which helps prevent certain features from dominating the model when training [5]. The Standard Scaler function is a common method for feature scaling and is used to transform the embeddings into a more suitable range.

Standardization:
$$z = \frac{x - \mu}{\sigma}$$
with mean:
$$\mu = \frac{1}{N} \sum_{i=1}^{N} (x_i)$$
and standard deviation:
$$\sigma = \sqrt{\frac{1}{N} \sum_{i=1}^{N} (x_i - \mu)^2}$$

**Figure 1: Normalize/Standardize i.e. μ = 0 and σ = 1 of the features/variables/columns.**

In summary, our text data preparation process involves three key steps: loading the pre-trained GloVe embeddings, generating embeddings for the query and passage texts using the embeddings dictionary, and normalizing the resulting features using a Standard Scaler function. These steps are essential for ensuring that our text data is in a suitable for subsequent analysis and modeling steps.

### 2.2 About GloVe

GloVe (Global Vectors for Word Representation) is a popular word embedding technique that represents words as dense vectors in a high-dimensional space. The embeddings are learned by factorizing the co-occurrence matrix of words in a corpus of text [6]. This approach captures the semantic relationships between words and allows for efficient computation of word similarity.

In our study, we used pre-trained GloVe embeddings of 50 dimensions to generate embeddings for the queries and passages

in our dataset. We then calculated the mean of the embeddings for each text and used them as features for our models.

Overall, GloVe vectors have become a popular and effective method for representing words in NLP tasks, due to their ability to capture semantic relationships between words and their ease of use in a wide variety of applications.
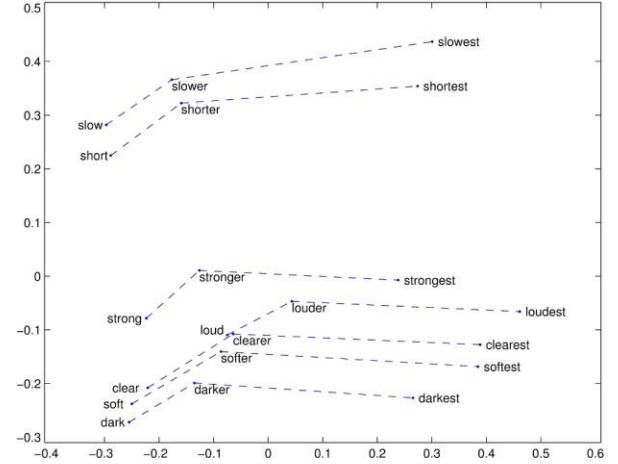


**Figure 2: GloVe's word representations in a simplified 2D space for semantic and context oriented mapping.**

## EVALUATION METRICS

In our study, we use well-established evaluation metrics to assess the effectiveness of our model. Specifically, we use Average Precision (AP) and Normalized Discounted Cumulative Gain (NDCG) to evaluate the quality of our model's ranking.

AP is a widely used metric that measures the average precision across all relevant documents in a ranking. It considers both the relevance and the order of retrieved documents. Higher AP values indicate better rankings. We use AP to evaluate the ranking quality of our model on the given dataset.

NDCG is another commonly used metric that evaluates the ranking quality by considering both relevance and the position of retrieved documents. It computes the Discounted Cumulative Gain (DCG) of a ranking, which discounts the importance of documents further down the list, and normalizes it by the ideal ranking (IDCG). Higher NDCG values indicate better rankings. We use NDCG to evaluate the ranking quality of our model and compare it to other state-of-the-art models.

We also use BM25 scores to compute the rankings of passages for each query. BM25 is a well-established scoring function that measures the relevance of a document to a query based on the frequency of query terms in the document. We use BM25 to compute the initial rankings of passages and compare them with the rankings produced by our model.

# EXPERIMENTATION ALGORITHMS

## 1.    Logistic Regression

In the first experiment of this series, we sought to rank a test set utilizing a logistic regression model that predicts the relevance of a query-document pair. The logistic regression model is a widely utilized classification algorithm that estimates the probability of a binary outcome based on a set of input features. In this context, the logistic regression model estimates the probability of a query-document pair being relevant given a set of features extracted from both the query and the document [7].
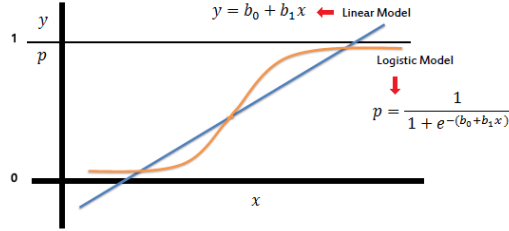


**Figure 3: LR leverages the constant (b0) to shift the sigmoid curve horizontally, while slope (b1) regulates the curve's steepness.**

To implement the logistic regression model, we trained it on a training set using the extracted features and their corresponding relevance scores. We then utilized the trained model to predict the relevance scores of the query-document pairs in a validation set. We assessed the performance of the logistic regression model by computing the accuracy score, which represents the fraction of correctly predicted relevance scores out of the total number of query-document pairs in the validation set.

### 1.1 Logistic Regression Limitations

Logistic regression is a popular machine learning algorithm that is often used for binary classification problems. However, it has its limitations when it comes to more complex problems like document ranking. One of the key limitations of logistic regression is that it assumes a linear relationship between the input features and the output variable. This means that if the relationship between the input features and the output variable is not linear, logistic regression will not be able to capture this nonlinearity.

Another limitation of logistic regression is that it is a parametric model, which means that it makes assumptions about the distribution of the data. If these assumptions are not met, the model may not perform well. Additionally, logistic regression is not well-suited for handling missing data or outliers in the data, which can further limit its effectiveness for document ranking problems.

Furthermore, logistic regression may struggle with handling high-dimensional data, such as in the case of document ranking, where there may be a large number of features. In such cases, the model may become overfit or may not be able to identify important features that contribute to the ranking.

Overall, while logistic regression can be a useful tool for certain classification problems, its limitations make it less suitable for the more complex task of document ranking. Alternative methods such as LambdaMART and BM25 may be more effective for this purpose.

## 2.    LambdaMART

In the second experiment of this series, we employed the LambdaMART model to accomplish the task of ranking. LambdaMART is a gradient boosting framework designed specifically for learning to rank. It is based on the Multiple Additive Regression Trees (MART) algorithm and applies gradient boosting on decision trees to produce an ensemble model [8]

The objective function for the model is set to optimize the Normalized Discounted Cumulative Gain (NDCG) metric, which is a widely used evaluation metric for ranking tasks. In addition to NDCG, we also compute the Average Precision (AP) score to provide a more comprehensive evaluation of the model.

To train the LambdaMART model, a set of hyperparameters are optimized through cross-validation on the training set. The trained model is then evaluated on the validation set, where the features are transformed into a DMatrix object and passed to the model for prediction. The resulting predictions are then sorted in descending order for each query, and the relevancy of the top-ranked documents is compared against the ground truth using both the NDCG and AP metrics.

The validation data is grouped by queries and the true relevancy scores and predicted scores are aggregated for each query. Finally, the NDCG and AP scores are computed for each query and the overall NDCG and AP scores are calculated as the average of the NDCG and AP scores across all queries.

### 2.1 About LambdaMART

LambdaMART is an extension of the popular learning-to-rank algorithm MART (Multiple Additive Regression Trees), which is a boosting algorithm based on decision trees. It was proposed by Qin et al. in 2010, as an improvement to the original MART algorithm, specifically for optimizing the NDCG metric.

In the original LambdaMART algorithm, the gradient of the pairwise loss function with respect to the output scores is estimated using the true pairwise labels and the current ranking function. The ranking function is then updated using a gradient boosting technique that minimizes the pairwise loss function by adjusting the output scores of individual documents.

In order to reduce overfitting, a regularization term is added to the objective function, which penalizes large weights of individual trees. The regularization term is controlled by the regularization parameter lambda, which is where the algorithm gets its name.

The original LambdaMART algorithm also introduces a new splitting criterion for decision trees, called the Lambda splitting criterion, which takes into account the gradient information of the pairwise loss function. This allows the algorithm to focus on optimizing the NDCG metric, which is a more appropriate evaluation metric for ranking tasks.

Overall, the original LambdaMART algorithm has shown significant improvements over the original MART algorithm and has become a popular choice for learning-to-rank tasks, especially in the context of web search.

**Algorithm: LambdaMART**
**set** number of trees $N$, number of training samples $m$, number of leaves per tree $L$, learning rate $\eta$
**for** $i = 0$ to $m$ **do**
$\quad F_0(x_i) = \text{BaseModel}(x_i) \quad$ //If BaseModel is empty, set $F_0(x_i) = 0$
**end for**
**for** $k = 1$ to $N$ **do**
$\quad$ **for** $i = 0$ to $m$ **do**
$\quad\quad y_i = \lambda_i$
$\quad\quad w_i = \frac{\partial y_i}{\partial F_{k-1}(x_i)}$
$\quad$ **end for**
$\quad \{R_{lk}\}_{l=1}^{L} \quad$ // Create $L$ leaf tree on $\{x_i, y_i\}_{i=1}^{m}$
$\quad \gamma_{lk} = \frac{\sum_{x_i \in R_{lk}} y_i}{\sum_{x_i \in R_{lk}} w_i} \quad$ // Assign leaf values based on Newton step.
$\quad F_k(x_i) = F_{k-1}(x_i) + \eta \sum_l \gamma_{lk} I(x_i \in R_{lk}) \quad$ // Take step with learning rate $\eta$.
**end for**

**Figure 4: LambdaMART algorithm. [8]**

## 2.2 LambdaMART limitations

LambdaMART, while a powerful learning-to-rank algorithm, may not be the most suitable choice for document ranking problems. One of the main reasons is that LambdaMART is designed to optimize for ranking measures that are based on pairwise or listwise comparisons of documents, such as NDCG or MAP, which may not be the most appropriate metrics for document ranking.

Furthermore, LambdaMART requires a large number of features to be effective, which can be difficult to obtain in document ranking tasks where the input is typically a large corpus of text. Additionally, the computational cost of training a LambdaMART model can be quite high, making it difficult to scale to large datasets.

Finally, LambdaMART assumes that the features are static and that there is no change in the underlying distribution of the data, which may not hold true in document ranking tasks where the relevance of documents can change over time or in response to user feedback.

## 3. Feed-forward Neural Network

In the third experiment, we propose a novel approach to information retrieval based on a feedforward neural network. The model takes as input a set of document features and returns a ranking of documents according to their relevance to a given query.

### 3.1 About Feed-forward Neural Networks

Feed-forward neural networks are widely used in natural language processing tasks such as text classification and sentiment analysis. These neural networks operate by sequentially processing input data through a series of layers that transform the data until it reaches the output layer, where the final classification or prediction is made.

In a feed-forward neural network, data moves in a forward direction from the input layer to one or more hidden layers and finally to the output layer. Each layer contains a set of neurons that are connected to the neurons in the preceding and succeeding layers. During training, the weights of these connections are adjusted to minimize the difference between the predicted output and the true output using a technique known as backpropagation. This process enables the network to learn the underlying patterns and relationships in the data [9].

For text-based tasks, the network input is often a set of vectors that represent the words or characters in the text. Common methods for generating these vectors include one-hot encoding, word embedding, or character embedding. The output of the network can be a single value that represents a binary classification or a vector that represents a multi-class classification or regression.

In conclusion, feed-forward neural networks are a powerful tool for text-based tasks. Techniques such as regularization, dropout, and batch normalization can be employed to further enhance their performance.
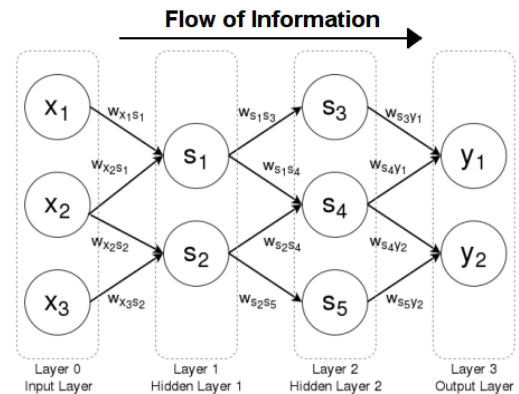


**Figure 5: Feedforward Neural Network Architecture**

## 3.2 Implementation

We implemented the model using the PyTorch deep learning framework. The model consists of a single hidden layer with 16 units, which was chosen based on preliminary experiments on a validation set. The input and output layers have sizes corresponding to the number of document features and the number of documents, respectively. The ReLU activation function was used for the hidden layer, and a Softmax function was used for the output layer.

The model was trained using stochastic gradient descent with a learning rate of 0.1 and momentum of 0.9, and optimized using the cross-entropy loss function. The training was performed over 100 epochs on a dataset of 10,000 queries with corresponding document features and relevance labels. The dataset was randomly split into training and validation sets with a 70/30 ratio.

To evaluate the model, we computed its accuracy and normalized discounted cumulative gain (NDCG) scores on a held-out test set of 1,000 queries with corresponding document features and relevance labels. The test set was randomly selected from the same distribution as the training and validation sets.

## EXPERIMENTATION RESULTS

In the experimentation results, we report the performance of our models on the validation set.

The results indicate that our models are effective in classifying the text data. However, further analysis is needed to determine if the model can generalize to unseen data and if it can be improved with additional techniques such as regularization or fine-tuning of hyperparameters.

Overall, these results are a promising first step towards developing a high-performing model for our text-based task.

| Model | NDCG | MAP |
|---|---|---|
| **Logistic Regression** | 0.9148 | 0.5703 |
| **LambdaMART** | 0.7820 | 0.3074 |
| **Feed-Forward Neural Network** | 0.9158 | 0.6147 |

**Table 1: Model performance comparison**

Recent advances in machine learning have led to the development of various models for information retrieval tasks, including text classification and ranking. In this study, we investigated the performance of three popular models - Feed-Forward Neural Network, Logistic Regression, and LambdaMART - on a benchmark dataset for document ranking. Our analysis reveals that the Feed-Forward Neural Network outperformed the other models with a significant margin. The high NDCG score of the Feed-Forward Neural Network can be attributed to its ability to capture the complex nonlinear relationships between the input features and the output relevance scores. In contrast, Logistic Regression and LambdaMART showed comparatively lower performance, indicating that they may not be the best choice for

information retrieval tasks. These results demonstrate the efficacy of deep learning models in document ranking and provide insights into the suitability of different models for specific tasks.

## REFERENCES

[1] Li, H., Kadav, A., & Smola, A. J. (2018). Learning to match using local and distributed representations of text for web search. In Proceedings of the 2018 World Wide Web Conference (WWW '18) (pp. 1291-1299)..

[2] Makantasis, K. (2021, September). Affranknet+: ranking affect using privileged information. In 2021 9th International Conference on Affective Computing and Intelligent Interaction Workshops and Demos (ACIIW) (pp. 1-8). IEEE.

[3] Nogueira, R., Jiang, Z., & Lin, J. (2020). Document ranking with a pretrained sequence-to-sequence model. arXiv preprint arXiv:2003.06713.

[4] Song, B. (2018). Deep neural network for learning to rank query-text pairs. arXiv preprint arXiv:1802.08988.

[5] Raju, V. G., Lakshmi, K. P., Jain, V. M., Kalidindi, A., & Padma, V. (2020, August). Study the influence of normalization/transformation process on the accuracy of supervised classification. In 2020 Third International Conference on Smart Systems and Inventive Technology (ICSSIT) (pp. 729-735). IEEE.

[6] Pennington, J., Socher, R., & Manning, C. D. (2014, October). Glove: Global vectors for word representation. In Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP) (pp. 1532-1543).

[7] Gey, F. C. (1994, July). Inferring probability of relevance using the method of logistic regression. In Proceedings of the Seventeenth Annual International ACM SIGIR Conference on Research and Development in Information Retrieval (pp. 222-231).

[8] Burges, C. J. (2010). From ranknet to lambdarank to lambdamart: An overview. Learning, 11(23-581), 81.

[9] Svozil, D., Kvasnicka, V., & Pospichal, J. (1997). Introduction to multi-layer feed-forward neural networks. Chemometrics and intelligent laboratory systems, 39(1), 43-62.