

Ensemble Ranking Model with Multiple Pretraining Strategies for Web Search

Xiaojie Sun*

CAS Key Lab of Network Data
Science and Technology, ICT, CAS
University of Chinese Academy of
Sciences
Beijing, China
sunxiaojie21s@ict.ac.cn

Lulu Yu*

CAS Key Lab of Network Data
Science and Technology, ICT, CAS
University of Chinese Academy of
Sciences
Beijing, China
nothing_0_1@163.com

Yiting Wang*

CAS Key Lab of Network Data
Science and Technology, ICT, CAS
University of Chinese Academy of
Sciences
Beijing, China
wangyiting21s@ict.ac.cn

Keping Bi

CAS Key Lab of Network Data
Science and Technology, ICT, CAS
University of Chinese Academy of
Sciences
Beijing, China
bikeping@ict.ac.cn

Jiafeng Guo

CAS Key Lab of Network Data
Science and Technology, ICT, CAS
University of Chinese Academy of
Sciences
Beijing, China
guojiafeng@ict.ac.cn

ABSTRACT

An effective ranking model usually requires a large amount of training data to learn the relevance between documents and queries. User clicks are often used as training data since they can indicate relevance and are cheap to collect, but they contain substantial bias and noise. There has been some work on mitigating various types of bias in simulated user clicks to train effective learning-to-rank models based on multiple features. However, how to effectively use such methods on large-scale pre-trained models with real-world click data is unknown. To alleviate the data bias in the real world, we incorporate heuristic-based features, refine the ranking objective, add random negatives, and calibrate the propensity calculation in the pre-training stage. Then we fine-tune several pre-trained models and train an ensemble model to aggregate all the predictions from various pre-trained models with human-annotation data in the fine-tuning stage. Our approaches won 3rd place in the “Pre-training for Web Search” task in WSDM Cup 2023 and are 22.6% better than the 4th-ranked team.

KEYWORDS

Pre-training, Ensemble Learning, Negative Sampling

ACM Reference Format:

Xiaojie Sun, Lulu Yu, Yiting Wang, Keping Bi, and Jiafeng Guo. 2023. Ensemble Ranking Model with Multiple Pretraining Strategies for Web Search. In *Proceedings of Make sure to enter the correct conference title from your*

*Equal contribution, and the order is determined randomly.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

WSDM Cup’23, February 27–March 3, 2023, Singapore, Singapore

© 2023 Association for Computing Machinery.
ACM ISBN 978-1-4503-XXXX-X/18/06...\$15.00
<https://doi.org/XXXXXXXX.XXXXXXX>

rights confirmation emai (WSDM Cup’23). ACM, New York, NY, USA, 4 pages.
<https://doi.org/XXXXXXXX.XXXXXXX>

1 INTRODUCTION

It is essential to measure relevance effectively in various search-related scenarios, such as web search and e-commerce search. Standard approaches usually use human-annotated relevance labels as guidance to train a ranker. Although the relevance judgments between query-document pairs have high quality, they require extensive manual efforts and huge costs. State-of-the-art rankers are mainly based on pre-trained language models, which are of large scale and data-hungry. Training such models with limited manual relevance judgments may lead to sub-optimal performance. User clicks could be alternative or supplementary to training such large models since they are cheap to collect and can indicate relevance.

Click data is difficult to use effectively since it contains much noise and various of bias, including position bias, selection bias, trust bias, etc. Based on the large-scale search logs of a popular Chinese web search engine, the WSDM Cup 2023 presents a challenge to alleviate the bias in click data and use it for ranker training to improve retrieval performance. One of the two tasks in the competition is the “Pre-training for Web Search” task. It aims to train a BERT-style ranker with click data and the organizer also provide a human-annotated dataset an unbiased data source that can fine-tune the model. A hidden test set is used to evaluate each team’s model performance at last. This task is a valuable foundation for academic studies on learning an unbiased pre-trained model for effective retrieval.

Existing studies on unbiased learning mainly focus on learning-to-rank models, that typically learn an unbiased ranking function of various types of features such as relevance, recency, quality, and popularity using clicks as learning objectives. Semantic matching signals are often included as features in such models so the underlying model cannot be trained in an end-to-end manner. It is not clear whether existing unbiased learning methods can alleviate bias

from click data and achieve good performance when training an end-to-end ranker based on a much larger pre-trained model.

We investigate several representative unbiased learning methods in the task of pre-training for web search and find that their benefit is not as large as in the learning-to-rank setting. We also observe that fine-tuning the BERT-style ranker with unbiased human-annotated data will improve the performance a lot. To mitigate multiple biases in the click data that harm model training and boost retrieval quality, we adopt multiple strategies in our runs, which are shown to be effective. Our solution has two stages: 1) pre-training the model with the click data, and 2) fine-tuning and ensemble based on the human-annotated data. In the first stage, we incorporate heuristic-based features, refine the ranking objective, add random negatives, and calibrate the propensity calculation. In the second stage, we augment training samples by duplicating high-frequency samples and learn an ensemble model using the predictions of our multiple runs and heuristics as features. We won 3rd place in the competition and our performance is 22.6% higher than the team ranked 4th place.

2 DEBIASED PRE-TRAINING

In this section, we first describe the model architecture we use, and then we introduce our learning objective and inverse propensity weighting schemes.

2.1 Model Architecture

To leverage more matching signals, we include heuristic-based features such as BM25 [3] and query likelihood with multiple smoothing methods [6] in our model, which is as same as the procedure describe in [5]. As shown in Figure 1, our model has a wide and deep architecture that learns semantic matching using deep neural networks and incorporates much more exact matching features (e.g., TF, IDF, TF-IDF, BM25, DIR [6], etc.) with a shallow network. We adopt the well-known cross encoder, multi-layer bidirectional Transformer encoder blocks [4] (Hidden State=768, Attention Head=12, Layer=3 or 12), to encode the concatenated query and document text. Then the vector of the [CLS] token in the output of the encoder layers captures the interactions between the query and document terms. We project the heuristic-based dense features to hidden vectors and concatenate them with the embedding of [CLS]. Then MLP layers are used to map the concatenated vectors to a final score.

2.2 Model Training

We introduce the strategies we use in model training, including label refinement and negative sample selection, and then we introduce the two loss functions we use.

Label Refinement. We find that training an end-to-end ranker based on the pre-trained model on the click data, whether with or without unbiased learning, cannot outperform heuristic-based features such as BM25. So, instead of using clicks as labels, we obtain the learning targets by incorporating a well-performing heuristic-based feature with the clicks according to Figure 2. In Figure 2, c_{ij} represents whether the j th document in the document list given query q_i is clicked, f_{ij} is the corresponding feature value, and finally, we perform a Softmax operation on all y_{ij} to get \tilde{y}_{ij} so that all labels sum up to one.

Negative Sample Selection. Only the top 10 documents are recorded in the click logs while the evaluation data includes the top 30 and other documents ranked in the top 900 with an interval of 30 (e.g., 30,60,90,...,990). Since top results are usually relevant and user clicks cannot cover all the relevant documents. Training a ranker with the click data will inevitably lead to sub-optimal performance on the test set. To let the model obtain the ability to discriminate relevant documents from irrelevant ones, besides the non-clicked documents in the top 10, we include random samples as negatives during training as well.

Loss Function. As previous studies show that listwise ranking loss has superior performance compared to pointwise and pairwise training, we train our model with a listwise loss based on attention allocation[1].

$$L_{listwise} = - \sum_{i=1}^N \sum_{j=1}^K [y_{ij} \frac{\exp(x_{ij})}{\sum_{k=1}^K \exp(x_{ik})}] \quad (1)$$

The simple list-wise loss function is shown in the Equation (1), where N is the total number of queries, K is the length of the list, that is, the number of documents containing user feedback that is the longest saved in the click log. On this basis, since the documents displayed to the user already have different degrees of relevance, in order to make the model see more negative samples, we randomly add a specified number of random negative samples with different content of queries from the same training batch to the list, and set its label to zero. In addition, considering the documents after the last clicked one may not actually be observed by the user, we replace these documents with random negative samples to alleviate the false negative problem. To encourage the model learned with the heuristic-click-based labels to outperform the original heuristic features, we also tried to relax the listwise matching constraint to pairwise so that the model has more freedom to behave differently from the target label. We attempt to build the order and use the pairwise loss function to optimize (see Section 3.1 for details). The priority order between two documents is mainly constructed based on the following relationships:

- Documents that are clicked are more relevant than documents that are not clicked.
- Documents with high feature values are more relevant than documents with low values.
- Documents shown to users outperform random negative samples.

Then the relative order can be obtained and we optimize the pre-training model using this objective.

2.3 Inverse Propensity Weighting

Previous studies have shown the effectiveness of the Dual Learning Algorithm (DLA) in the task of unbiased learning to rank. However, we find that the propensity weights learned by DLA from the click data are not decreasing from the first to the tenth position, which is inconsistent with existing research. Also, since we conduct label refinement and its impact on propensity learning is unknown, we use static propensity calculated from click ratios for model training. The output of model variants trained with DLA and click-ratio-based propensities are used as features in the final ensemble model

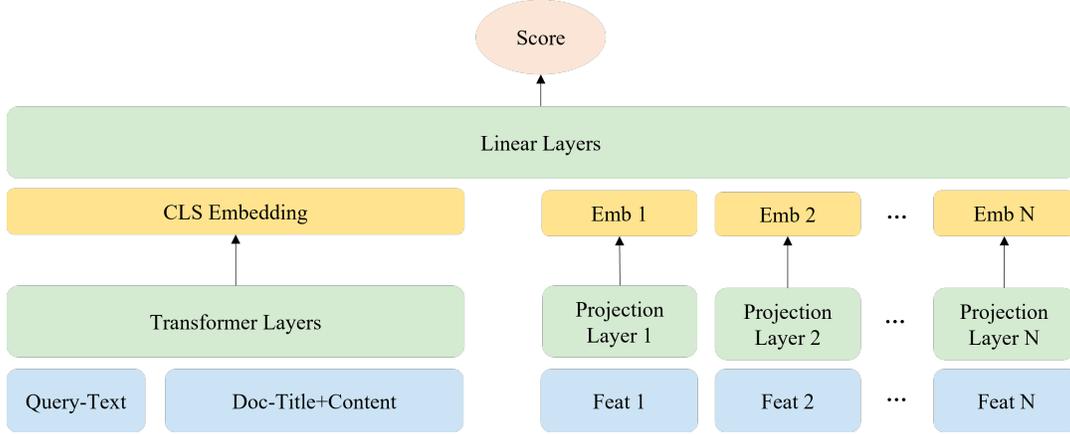
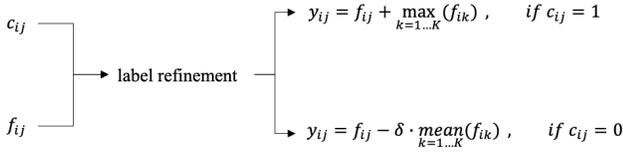


Figure 1: Model architecture.

Figure 2: The generation of the label y_{ij} of each document in the corresponding list for each query q_i .

(in Section 3.4). The two inverse propensity weighting mechanisms are as follows:

- **DLA.** Using the dual learning algorithm, the propensity model and the ranking model are jointly optimized. It is worth mentioning that we fixed the propensity weight value of additive negative samples to 1, which is slightly different from the original DLA.
- **Click-ratio-based propensity.** We use static inverse propensity weights according to the click ratio on different ranking positions of the whole training set in [7] as the following:

$$pw_i = \left(\frac{cr_1}{cr_i}\right)^\alpha, i = 1, \dots, 10. \quad (2)$$

Here, α is used to control the relative size of propensity weights and we fix it as 0.25. The fixed inverse propensity weights are set to 1, 1.19, 1.44, 1.58, 1.89, 1.95, 2.12, 2.26, and 2.51 for the top-10 ranked documents.

3 FINE-TUNING

In this section, we describe the strategies we use in the fine-tuning stage including efficient design of ranking loss, negative sample selection, sample augmentation, and ensemble.

3.1 Ranking Loss

The document list under a query contains 5 classes: {bad, fair, good, excellent, perfect}. The most common way to deal with multi-level relevance labels is to transform them into positive and negative classes, thus using a point-wise loss function. Essentially, the point-wise method is to approximate the ranking problem to a regression

problem, but the ranking task does not pursue accurate scoring, and relative scoring is acceptable. At the same time, the training of the model will be dominated by queries with a large number of labeled documents. Therefore, a pairwise method for modeling the relative relationship between positive and negative samples is needed. The pairwise method organizes a sample as $\langle q, d^+, d^- \rangle$, which means that d^+ is more relevant to q than d^- .

$$L_{pairwise} = - \sum_{i=1}^N \left[\frac{\exp(x_i^+)}{\exp(x_i^+) + \exp(x_i^-)} \right] \quad (3)$$

In order to improve the retrieval performance of the model, as shown in Equation (4), we further introduce a certain number of negative samples into the pair-wise loss function formula, so that the model can better improve the relative score of positive samples. T is the number of negative samples for each query plus 1, so a sample becomes $\langle q, d_1^-, d_2^-, \dots, d_{T-1}^-, d_T^+ \rangle$.

$$L = - \sum_{i=1}^N \left[\frac{\exp(x_{iT}^+)}{\sum_j \exp(x_{ij})} \right] \quad (4)$$

3.2 Training Labels

Since the expert annotation dataset contains multi-level labels, we degenerate the 5 levels into two when using the pointwise and pairwise loss function due to its simplicity. We think that if a document is marked as "perfect" (label=4), "excellent" (label=3), or "good" (label=2), then the document is a positive sample, otherwise, the document is a negative sample. d^+ in Equations (3) and (4) is sampled from all positive samples of a query, and d^- is selected from its negative samples. We also tried to train the model with the listwise loss according to Equation (1) but got no gain. It is worth mentioning that probably we have not tuned the model sufficiently and more fine-grained exploration of the multi-grade label and more effective loss functions can be used to fine-tune in the future.

3.3 Sample Augmentation

Unlike click data in the search log, the cost of manual labeling is high and thus may lead to a limited number of queries in the expert

annotated dataset. Therefore, it is promising to explore practical methods to enhance the effect of core samples in the fine-tuning stage. The queries in the human annotation data are de-duplicated and we find that there exist much more tail queries than head queries in the provided dataset. The numbers of high, mid, and low-frequency queries are 1092, 1820, and 1789 in the annotation dataset, respectively. Table 1 illustrates the distribution of the relevance labels for different frequency queries. In this paper, we consider high-frequency queries as the core samples in this stage due to the following reasons.

It is known that head queries are easier so long tail queries are more important to differentiate the ability of a search engine. Positive documents contribute to model training a lot since most queries contain only a few perfect documents. We find that head queries have more positive documents since there are more relevant resources in the annotation set and previous user clicks can help search engines rank documents for such queries. We train the model more on the head queries to encourage the model to learn common matching patterns sufficiently.

Table 1: Distribution of relevance labels.

Grade	Label	Ratio of label		
		High	Mid	Low
Perfect	4	0.49%	0.15%	0.02%
Excellent	3	12.99%	8.00%	2.71%
Good	2	35.06%	31.32%	21.33%
Fair	1	15.96%	9.40%	5.16%
Bad	0	35.50%	51.13%	70.78%

Considering the above factors, we amplified the effect of the high-frequency samples which is more important and accurately annotated through duplicating these samples and feeding them together with other samples during the fine-tuning stage. This strategy enhances the performance significantly and we leave the exploration of other augmentation strategies to the future.

3.4 Ensemble Strategy

A simple but effective way is adopted in the ensemble stage. Taking the extracted features and the output scores of previous pre-trained and fine-tuned models as input, we train a Gradient Boosting Decision Tree model using Light Gradient Boosting Machine (LightGBM) [2] in the validation set with LambdaRank objective, and the process can be formulated as follows:

$$f^* = \arg \min_f \mathbb{E}_r L(r, f(x)) \quad (5)$$

The core hyper-parameters are num_leaves, max_depth, learning_rate, and num_iteration. We tried different combinations of these parameters and we tuned the hyper-parameters on a small proportion of the validation set, which is described in Section 4. We tried normalization to the input but got no gain.

4 EXPERIMENTS

Considering that queries with fewer than 10 retrieved documents may be too specific and not likely to have truly relevant results, we filter out the query-document pair with no clicks and those

queries with less than 10 candidate documents. We use the AdamW optimizer. We use $\delta = 2$ and $\tau = 0.1$ in the stage of label refinement. Our method is implemented in PaddlePaddle and we use the official model for initialization. In the fine-tuning stage, the provided validation set (i.e. human-annotation set) is randomly divided into two parts using the query-id, and we use 80% for training and 20% for validating. We select the best model based on the performance of the 20% validation set. Table 2 shows the top 5 teams and their final scores in this task. Our solution got 9.83 for the online dcg@10 and won 3rd place.

Table 2: Top 5 scores of the competition. Our team won 3rd in the public score leaderboard.

Rank	Team Name	DCG@10
1	Tencent Search	12.16525
2	THUIR	10.04097
3	Cannot Retrieve	9.83148
4	Accepted	8.02173
5	DisTime	7.30951

5 CONCLUSION

In this paper, we detail our winning solution to the task of pre-training for web search in WSDM Cup 2023. We use the Transformer model as the backbone, with the combination of negative sample augmentation and target modification in the pre-training stage. A pairwise ranking loss, key sample augmentation, and ensemble strategy are used in the fine-tuning stage. Our solution achieved the DCG@10 score of 9.83148 and finally, we ranked 3rd on the leaderboard.

ACKNOWLEDGMENTS

This work was supported by the Lenovo-CAS Joint Lab Youth Scientist Project. Any opinions, findings, conclusions, and recommendations expressed in this material are those of the authors and do not necessarily reflect those of the sponsors.

REFERENCES

- [1] Qingyao Ai, Keping Bi, Jiafeng Guo, and W Bruce Croft. 2018. Learning a deep listwise context model for ranking refinement. In *The 41st international ACM SIGIR conference on research & development in information retrieval*. 135–144.
- [2] Guolin Ke, Qi Meng, Thomas Finley, Taifeng Wang, Wei Chen, Weidong Ma, Qiwei Ye, and Tie-Yan Liu. 2017. Lightgbm: A highly efficient gradient boosting decision tree. *Advances in neural information processing systems* 30 (2017).
- [3] Stephen E Robertson and Steve Walker. 1994. Some simple effective approximations to the 2-poisson model for probabilistic weighted retrieval. In *SIGIR'94: Proceedings of the Seventeenth Annual International ACM-SIGIR Conference on Research and Development in Information Retrieval*. Springer, 232–241.
- [4] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems* 30 (2017).
- [5] Lulu Yu, Yiting Wang, Xiaojie Sun, Keping Bi, and Jiafeng Guo. 2023. Feature-Enhanced Network with Hybrid Debiasing Strategies for Unbiased Learning to Rank. *arXiv preprint arXiv:2302.07530* (2023).
- [6] Chengxiang Zhai and John Lafferty. 2004. A study of smoothing methods for language models applied to information retrieval. *ACM Transactions on Information Systems (TOIS)* 22, 2 (2004), 179–214.
- [7] Lixin Zou, Haitao Mao, Xiaokai Chu, Jiliang Tang, Wenwen Ye, Shuaiqiang Wang, and Dawei Yin. 2022. A large scale search dataset for unbiased learning to rank. *arXiv preprint arXiv:2207.03051* (2022).