

# ILMART: Interpretable Ranking with Constrained LambdaMART

Claudio Lucchese  
claudio.lucchese@unive.it  
Ca' Foscari University of Venice  
Italy

Franco Maria Nardini  
francomaria.nardini@isti.cnr.it  
ISTI-CNR  
Italy

Salvatore Orlando  
orlando@unive.it  
Ca' Foscari University of Venice  
Italy

Raffaele Perego  
raffaele.perego@isti.cnr.it  
ISTI-CNR  
Italy

Alberto Veneri  
alberto.veneri@unive.it  
Ca' Foscari University of Venice  
ISTI-CNR  
Italy

## ABSTRACT

Interpretable Learning to Rank (LtR) is an emerging field within the research area of explainable AI, aiming at developing intelligible and accurate predictive models. While most of the previous research efforts focus on creating post-hoc explanations, in this paper we investigate how to train effective and *intrinsically-interpretable* ranking models. Developing these models is particularly challenging and it also requires finding a trade-off between ranking quality and model complexity. State-of-the-art rankers, made of either large ensembles of trees or several neural layers, exploit in fact an unlimited number of feature interactions making them black boxes. Previous approaches on intrinsically-interpretable ranking models address this issue by avoiding interactions between features thus paying a significant performance drop with respect to full-complexity models. Conversely, ILMART, our novel and interpretable LtR solution based on LambdaMART, is able to train effective and intelligible models by exploiting a limited and controlled number of pairwise feature interactions. Exhaustive and reproducible experiments conducted on three publicly-available LtR datasets show that ILMART outperforms the current state-of-the-art solution for interpretable ranking of a large margin with a gain of nDCG of up to 8%.

## CCS CONCEPTS

• Information systems → Learning to rank; • Computing methodologies → Boosting.

## KEYWORDS

Interpretable Ranking, Interpretable Boosting, LambdaMART

### ACM Reference Format:

Claudio Lucchese, Franco Maria Nardini, Salvatore Orlando, Raffaele Perego, and Alberto Veneri. 2022. ILMART: Interpretable Ranking with Constrained LambdaMART. In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR '22)*, July 11–15, 2022, Madrid, Spain.

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 the author(s) 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](mailto:permissions@acm.org).

SIGIR '22, July 11–15, 2022, Madrid, Spain

© 2022 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 978-1-4503-8732-3/22/07...\$15.00

<https://doi.org/10.1145/3477495.3531840>

2022, Madrid, Spain. ACM, New York, NY, USA, 5 pages. <https://doi.org/10.1145/3477495.3531840>

## 1 INTRODUCTION

Learning to Rank (LtR) is a vast research area focusing on learning effective models for solving ranking tasks. Although being effective, state-of-the-art LtR techniques do not deal with the *interpretability* of the learned model, i.e., the possibility given to humans to easily understand the reasoning behind predictions made by the model. The need for learning trustworthy models in LtR is becoming a topic of increasing interest, which tries to address important concerns about privacy, transparency, and fairness of ranking processes behind many everyday human activities.

For what regards the fairness of ranking algorithms, several solutions have been proposed to create fairer models, where the term *fair* takes different meanings according to the metric employed [2]. However, focusing only on devising a fairer model does not imply to be compliant with new regulations that are emerging. An example of that is *The Right of Explanation*, which allows each European citizen to claim a justification for every algorithm decision made with their data that affects their life [13]. Therefore, for highly-sensitive scenarios new LtR approaches have to be investigated to train ranking models that do not jeopardize interpretability in the name of effectiveness.

In literature, two main lines of work dealing with the interpretation of ranking models are investigated. The first concerns the creation of post-hoc explanations for black-box models, while the second focuses on the design of intrinsically-interpretable models. In the first line, the goal is the definition of techniques that allow, given a ranking model, to generate post-hoc explanations of the decisions made. Here, some of the explanation techniques used in regression or classification tasks have been adapted to the ranking scenario. For example, two recent contributions [17, 19] present two different approaches to adapt LIME [15] to the ranking task. The work by Singhs and Anand focuses only on pointwise LtR techniques and treats the ranking problem as a classification problem, where the probability of a document to be relevant is computed with respect to its score, and a simple linear SVM is used as a local explanation model [17]. The work by Verma and Ganguly evaluates the consistency and correctness of the explanations generated by LIME by varying the sampling strategies used to create the local neighborhood of a document [19]. In the same line, Fernando *et al.* propose an approach to reuse SHAP [12] for the ranking task [5].

Even though post-hoc explanations can be seen as a viable way to inspect and audit accurate ranking models, they can produce explanations that are not faithful to internal computations of the model, and thus they can be misleading [16]. For this reason, Zhuang *et al.* pioneer the problem of developing intrinsically-interpretable ranking models [20]. In detail, they propose *Neural RankGAM*, an adaptation of GAM [6] to the ranking task. Although Neural RankGAM nicely exploits the interpretability of GAMs, the authors show that the ranking performance achieved is significantly worse than those provided by state-of-the-art black-box models. However, it is worth noticing that in other contexts, such as in regression and classification tasks, other intrinsically interpretable models appear to be as accurate as other black-box models. A notable example of an interpretable and effective solution is Explainable Boosting Machine (EBM), a C++ implementation of the  $\ell_2$  algorithm [9], where each component of the Generalized Additive Models (GAM) are learned through a mixed bagging-boosting procedure [8].

In this paper, we contribute to the development of intrinsically-interpretable ranking models with ILMART, a novel LtR solution based on LambdaMART. ILMART creates effective and intelligible models by encompassing a limited and controlled number of pairwise feature interactions. We evaluate ILMART on three publicly-available LtR datasets. The results of our reproducible experiments show that ILMART outperforms the current state-of-the-art solution for interpretable ranking of a large margin, with a gain in terms of nDCG of up to 8% with respect to Neural RankGAM.

## 2 INTERPRETABLE LAMBDA MART

ILMART produces ensemble-based ranking models having the structure of a GAM, which is considered an interpretable model [8]. Before introducing ILMART, we review the main aspects of GAMs to highlight similarities and differences with respect to the newly proposed technique.

**Generalized Additive Models.** Let  $\mathbf{x} \in \mathbb{R}^3$  be an instance in a 3-dimensional vector space, with  $G_j$  being the value of its  $j$ -th feature. An instance  $\mathbf{x}$  is associated with a target label  $\tilde{y}$  generated by some unknown target function  $\delta(\cdot)$ , i.e.,  $\delta(\mathbf{x}) = \tilde{y}$ . In its basic form, a GAM models the target variable  $\tilde{y}$  as a combination of  $?$  main effects and non-zero interaction effects, as defined in [6]:

$$\delta(\mathbf{x}) = U + \sum_{j=1}^? B_j(G_j) + \sum_{(\beta, \gamma) \in \mathcal{K}} B_{\beta, \gamma}(G_{\beta} \cdot G_{\gamma})$$

Where  $\delta(\mathbf{x}) = [\tilde{y}|\mathbf{x}]$ , with  $\tilde{y}$  following a distribution of the exponential family, and  $\delta(\cdot)$  being the so-called *link function* that is able to describe the relationship between the  $?$  main effects and the interactions effects with the expected value of  $\tilde{y}$ . In addition, the number of non-zero interaction effects is equal to  $?$ , i.e.  $|\{\beta, \gamma: (\beta, \gamma) \in \mathcal{K}\}| = ?$ . Usually, the  $?$  main effects and the interactions are modeled through simple third order splines. From the explainability point of view, the analyst is promptly provided with a plot of each basis function  $B_j$  that clearly describes the relation between each feature and the target variable  $\tilde{y}$ , while  $B_{\beta, \gamma}$  values represent the contribution of the interactions between feature pairs. For classification and regression tasks, (shallow) decision trees instead of splines were investigated in [9].

**Interpretable LambdaMART.** We now introduce ILMART, our novel method based on LambdaMART to learn ranking models that are both accurate and interpretable. ILMART presents similarities with GAM, and builds the prediction  $\hat{y}$  of an unknown target variable  $\tilde{y}$  as follows:

$$\hat{y} = \sum_{j \in \mathcal{J}} g_j(G_j) + \sum_{(\beta, \gamma) \in \mathcal{K}} g_{\beta, \gamma}(G_{\beta} \cdot G_{\gamma}) \quad (1)$$

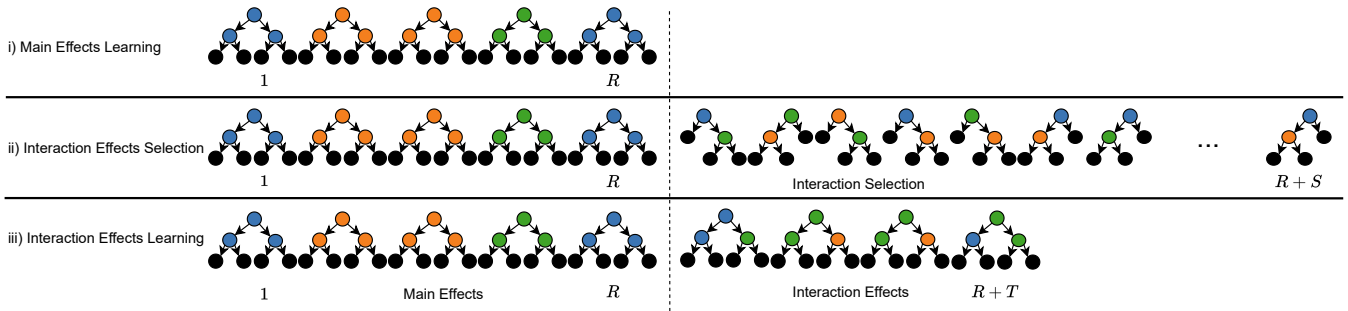
where  $g_j(G_j)$  and  $g_{\beta, \gamma}(G_{\beta} \cdot G_{\gamma})$  are ensembles of trees. Each tree  $\ell \in g_j(G_j)$  models a *main effect* in which the only feature allowed for each split is  $G_j$ .

On the other hand, each tree  $\ell \in g_{\beta, \gamma}(G_{\beta} \cdot G_{\gamma})$  models the *interaction effects*, where the only features allowed in a split of tree  $\ell$  are either  $G_{\beta}$  or  $G_{\gamma}$ . Similarly to the notation used to describe a GAM,  $?$  denotes the number of distinct main effects  $G_j$ , whereas  $?$  denotes the number of distinct interaction effects  $(G_{\beta} \cdot G_{\gamma})$ .

To learn the ensemble of trees modeled by Eq. (1), we adopt a modified version of the LambdaMART boosting algorithm. Unlike GAM, we do not make any assumption on the distribution of  $\tilde{y}$  to guide the learning process. The learning strategy to generate the main effects  $\ell \in g_j(G_j)$  and the interaction effects  $\ell \in g_{\beta, \gamma}(G_{\beta} \cdot G_{\gamma})$  in Equation 1 is made of three steps depicted in Fig. 1: *i)* Main Effects Learning, which learns a set of trees, each one working on a main effect (single feature); *ii)* Interaction Effects Selection, which selects the top- $?$  most important interactions effects (feature pairs); *iii)* Interaction Effects Learning, which learns a set of trees, each one exploiting one of the interaction effects selected by the previous step.

*Main Effects Learning.* In the initial step, illustrated on top of Fig. 1, we learn an ensemble of  $?$  trees modeling the main effects. To this end, we constrain the LambdaMART boosting procedure to use a single feature per tree in the first  $?$  boosting rounds. In other words, at a given boosting round of LambdaMART, when feature  $G_j$  is chosen for the root's split of a tree  $\ell$ , we force the algorithm to use  $G_j$  for all the further splits in  $\ell$ , until a stopping criterion for tree growth is reached. LambdaMART stops its boosting rounds when no improvements are observed on the validation set. Eventually, we obtain an ensemble of  $?$  trees only using for their splits a set  $\mathcal{J}$  of  $?$  features,  $? \leq 3$  and  $? \leq ?$ , modeling the main effects of Equation 1. Thus,  $?$  is not an input parameter of ILMART, but it is the number of distinct features used by the ensemble of  $?$  trees modeling the main effects only. We name this initial model ILMART. It is worth noting that, as a side effect of this first step, we achieve the additional effect of potentially reducing the cardinality of the features, from 3 to  $?$ .

*Interaction Effects Selection.* In the second step, we still exploit the LambdaMART boosting procedure to *select* the top- $?$  interaction effects, as illustrated in the middle of Fig. 1. Specifically, we continue the learning of the ILMART model obtained so far, by enabling the interactions between all possible pairs of  $?$  features identified by the previous step, according to the interaction heredity principle [4]. This is accomplished by constraining LambdaMART to fit trees with only 3 leaves (and 2 splits), where the 2 splits use a pair of distinct



**Figure 1: ILMART learning process.** Each color corresponds to a different feature used in the splits, and nodes colored in black are leaves. The trees for the main and interaction effects are represented with the same depth for illustration purposes only.

features. The boosting round is stopped when the size of the set  $\mathcal{K}$ , composed of the distinct feature pairs used by the ensemble of  $\ell$  trees, is exactly  $\binom{M}{2}$ , where  $\ell \leq \frac{M^2}{2}$ . The  $\ell$  trees are generated solely for the purpose of identifying the interaction effects and they are discarded afterwards.

Note that for selecting the main  $\ell$  interaction effects, we rely on the capability of LambdaMART in identifying, at each boosting round, the features and the splits with the highest contributions in minimizing the loss function.

Thus, the interaction effects selected by the boosting procedure are likely to be generated in order of importance.

*Interaction effects learning.* Finally, the last step sketched on the bottom of Fig. 1 consists in adding to the initial ILMART model  $\ell$  new trees learned by constraining LambdaMART to use only the top- $\ell$  feature pairs identified in the previous step. The number  $\ell$  of new learned trees is not an input of algorithm, but it is chosen on the basis of the validation set. The final model, obtained by adding to ILMART the  $\ell$  trees working on the  $\ell$  interaction effects is named ILMART $_{\ell}$ .

### 3 EXPERIMENTAL EVALUATION

We experimentally assess the performance of ILMART and ILMART $_{\ell}$  on three public datasets for LTR, namely ISTEELLA-S, WEB30K, and YAHOO. The ISTEELLA-S dataset [11] includes 33,018 queries with an average of 103 documents per query. Each document-query pair is represented by 220 features. The WEB30K dataset [14] is composed of more than 30,000 queries (in 5 folds), with an average of 120 documents per query and 136 features per query-document pair. The YAHOO dataset [3] is composed of two sets, including document-query pairs with 699 features. The feature vectors of the three datasets are labeled with relevance judgments ranging from 0 (irrelevant) to 4 (perfectly relevant). In our experiments, we use train/validation and test splits from fold 1 of the WEB30K dataset and data from “set 1” of the YAHOO dataset.

**Competitors.** We compare the performance of ILMART and ILMART $_{\ell}$  with the following competitors:

- Neural RankGAM (NRGAM), a neural-based approach to learn GAMs for ranking. The approach exploits standalone neural networks to instantiate sub-models for each individual feature. It is the current state-of-the-art technique for learning interpretable models for ranking [20].

- EBM, an efficient implementation of “ $\ell_2$ ” using a mixed bagging-boosting procedure [9]. In this case the learned model is used as a pointwise ranker, since the implementations available are only made to solve regression and classification tasks. We denote with EBM $_{\ell}$  the EBM model using interaction effects.

**Metrics.** We measure the performance in terms of normalized Discounted Cumulative Gain (nDCG) at three different cutoffs, i.e.,  $\{1 \cdot 5 \cdot 10\}$ . We compute the nDCG metric by employing exponential weighing of the relevance [1]. By default, queries with missing relevant documents have been assigned a nDCG = 1.0. The statistical significance is computed with a two-sided Fisher’s randomization test [18] with significance level  $\alpha < 0.05$ . The test is computed using the RankEval library [10].

**Implementation and training settings.** We implement ILMART as an extension of the LightGBM library<sup>1</sup> [7]. The extension works by adding the possibility to constraint the boosting procedure to use only a limited number of features per tree. The source code used in our experiments along with the models trained on the three public datasets is publicly available online<sup>2</sup>.

The training of ILMART and ILMART $_{\ell}$  optimizes nDCG@10. We perform hyper-parameter tuning by varying the number of leaves in the trees in  $\{32 \cdot 64 \cdot 128\}$  and the learning rate in  $\{0.001 \cdot 0.01 \cdot 0.1\}$ . Early stopping, i.e., stopping boosting if no improvement on the validation set is observed for 100 consecutive rounds, is used for training in step one (main effects learning) and three (interaction effects learning). The optimal ILMART models learned with the settings above are made of 914, 895 and 1,969 trees for the WEB30K, YAHOO, and the ISTEELLA-S datasets, respectively. On the same datasets, the three optimal ILMART $_{\ell}$  models are made of 1,365, 1,215 and 3,024 trees, respectively. Neural RankGAM is also trained by optimizing nDCG@10. We use the public implementation of NRGAM available in Keras<sup>3</sup> and made available by the authors of the original paper [20]. We train NRGAM by using the same settings reported by the authors. In addition, we employ a batch size of 128, and we train the networks with 3,000 epochs for WEB30K, 700 for YAHOO and 1,000 for ISTEELLA-S, with early stopping set to 100 epochs. We also learn EBM by optimizing nDCG@10. The performance of EBM is fine-tuned by varying the number of outer-bags in  $\{5, 10, 15\}$  for

<sup>1</sup><https://github.com/microsoft/LightGBM>

<sup>2</sup><https://github.com/veneres/ilmart>

<sup>3</sup><https://github.com/tensorflow/ranking/issues/202>

**Table 1: Performance comparison in terms of nDCG. Statistically significant improvements w.r.t. NRGAM are marked with an asterisk (\*). Best results are reported in bold.**

Dataset	Method	nDCG				?
		@1	@5	@10		
WEB30K	NRGAM	44.80	43.72	45.73	136	0
	EBM	41.70	43.14	45.66	136	0
	EBM <sub>g</sub>	46.51*	45.96*	48.01*	136	50
	ILMART	45.17	45.00*	47.05*	79	0
	ILMART <sub>g</sub>	<b>48.94*</b>	<b>47.76*</b>	<b>49.55*</b>	79	46
YAHOO	NRGAM	67.92	70.13	75.07	699	0
	EBM	70.33*	72.77*	77.29*	699	0
	EBM <sub>g</sub>	70.22*	72.85*	77.36*	699	50
	ILMART	70.57*	72.64*	77.21*	153	0
	ILMART <sub>g</sub>	<b>70.80*</b>	<b>73.20*</b>	<b>77.57*</b>	153	48
ISTELLA-S	NRGAM	63.37	63.04	69.18	220	0
	EBM	57.75	59.67	66.90	220	0
	EBM <sub>g</sub>	62.29	63.05	69.71	220	50
	ILMART	66.99*	66.29*	72.36*	106	0
	ILMART <sub>g</sub>	<b>69.06*</b>	<b>68.22*</b>	<b>74.04*</b>	106	50

the two versions with and without interactions. In our experiments, we employ a publicly-available implementation of EBM<sup>4</sup>. For what regards the two strategies encompassing interaction effects, i.e., ILMART<sub>g</sub> and EBM<sub>g</sub>, we limit the number of maximum interactions allowed to 50.

**Experimental results.** Table 1 reports the results of the evaluation of ILMART, ILMART<sub>g</sub> and their state-of-the-art competitors. For all the datasets considered, results show that ILMART achieves a statistically significant improvement in terms of nDCG for almost all the cutoffs except the case of nDCG@1 on WEB30K. More importantly, ILMART<sub>g</sub> always achieves the best performance in terms of nDCG that always result in a statistically significant improvement over NRGAM. The best improvements of ILMART<sub>g</sub> over *all* the competitors is obtained on ISTELLA-S, where nDCG@10 improves of up to 6% w.r.t. EBM<sub>g</sub>. On WEB30K, we observe that ILMART<sub>g</sub> improves nDCG@10 of up to 8% w.r.t NRGAM, although the best competitor is, in this case, EBM<sub>g</sub>, whose performance still lags behind ILMART<sub>g</sub>. On WEB30K, EBM<sub>g</sub> also achieves statistically significant improvements in terms of nDCG with respect to NRGAM. Even though EBM<sub>g</sub> is used as a pointwise ranker, these differences confirm the efficacy of tree-based approaches in Ltr and the need to include pairwise interactions to learn more effective ranking functions.

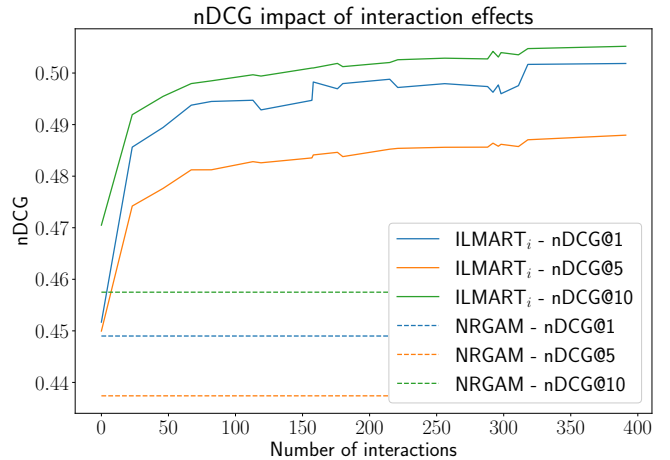
Furthermore, it is worth highlighting that our proposed techniques perform a feature selection during the *main effects learning* phase, while the competitors leave the feature selection to the final user. This is an important difference between ILMART and its competitors as a high number of features used do not allow to easily understand the global behavior of the model. For example, the YAHOO dataset comes with a large number of features, i.e., 699. ILMART and ILMART<sub>g</sub> effectively reduce the initial set of 699 features to 153 (column ? in Table 1) with, at the same time, a significant

<sup>4</sup><https://interpret.ml/docs/ebm.html>

increase of the effectiveness of the learned model. In addition, we point out that eventually ILMART<sub>g</sub> may have a number of interaction effects smaller than 50 due to the training of the new trees which may not include all the pairs selected.

We now investigate how the number of interactions added to ILMART<sub>g</sub> can influence its prediction accuracy. Fig. 2 compares the results of ILMART<sub>g</sub> with the ones of NRGAM by reporting their results in terms of nDCG. We observe that the ranking accuracy of ILMART<sub>g</sub> increases steeply when the first interaction effects are added; then, the gain obtained by adding new interaction trees saturates. This means that ILMART<sub>g</sub> is able to significantly improve the accuracy of an interpretable Ltr model by adding a small number of interaction effects, thus without hindering the overall interpretability of the model. Even though the addition of interactions effects improves a lot the accuracy of the model, we acknowledge that the difference from the full-complexity model LambdaMART is still marked. For example, on WEB30K an optimized black-box LambdaMART model achieves a nDCG@10 close to 52%.

Finally, to show how the main and interaction effects can be explored, Fig. 3 presents three main effects and one interaction effects learned with ILMART<sub>g</sub> trained over WEB30K. The three main effects selected are the ones associated with the features with highest feature importance, i.e., PageRank (PR), Query-url Click Count (QUCC), and Language Model for Information Retrieval (LMIR), while the interaction effect between URL Click Count (UCC) and QUCC is the one with the highest average contribution for each value of ( $G_g \cdot G_g$ ). The 2D plots and the heatmap are built by aggregating all the trees using the same features, i.e., by representing the contribution of each  $g_g(G_g)$  and each  $g_{gg}(G_g \cdot G_g)$ . From the plots, the analyst can easily understand the exact behavior of the model at varying values of each feature. For example, taking into account the interaction effect, the higher the value of UCC and QUCC the larger is the positive contribution of  $g_{gg}(G_g \cdot G_g)$ .



**Figure 2: nDCG by varying the number of interaction effects in ILMART<sub>g</sub> on the WEB30K dataset.**

