

Iterative Boosting Deep Neural Networks for Predicting Click-Through Rate

AMIT LIVNE, ROY DOR, EYAL MAZUZ, TAMAR DIDI, BRACHA SHAPIRA, and LIOR ROKACH, Ben-Gurion University of the Negev

The click-through rate (CTR) reflects the ratio of clicks on a specific item to its total number of views. It has significant impact on websites' advertising revenue. Learning sophisticated models to understand and predict user behavior is essential for maximizing the CTR in recommendation systems. Recent works have suggested new methods that replace the expensive and time-consuming feature engineering process with a variety of deep learning (DL) classifiers capable of capturing complicated patterns from raw data; these methods have shown significant improvement on the CTR prediction task. While DL techniques can learn intricate user behavior patterns, it relies on a vast amount of data and does not perform as well when there is a limited amount of data. We propose XDBoost, a new DL method for capturing complex patterns that requires just a limited amount of raw data. XDBoost is an iterative three-stage neural network model influenced by the traditional machine learning boosting mechanism. XDBoost's components operate sequentially similar to boosting; However, unlike conventional boosting, XDBoost does not sum the predictions generated by its components. Instead, it utilizes these predictions as new artificial features and enhances CTR prediction by retraining the model using these features. Comprehensive experiments conducted to illustrate the effectiveness of XDBoost on two datasets demonstrated its ability to outperform existing state-of-the-art (SOTA) models for CTR prediction.

CCS Concepts: • **Information systems** → *Recommender systems*; • **Theory of computation** → *Boosting*; • **Computing methodologies** → *Neural networks*.

Additional Key Words and Phrases: Click-Through Rate Prediction, Deep Neural Network, Boosting, Recommender Systems

ACM Reference Format:

Amit Livne, Roy Dor, Eyal Mazuz, Tamar Didi, Bracha Shapira, and Lior Rokach. 2020. Iterative Boosting Deep Neural Networks for Predicting Click-Through Rate. 1, 1 (July 2020), 16 pages. <https://doi.org/10.1145/1122445.1122456>

1 INTRODUCTION

Over the past decade, the popularity and use of the Internet have rapidly increased. Each day, the number of users and mobile devices utilizing this technology grows. Surfing the Internet is one of the most commonly performed activities, with users visiting websites via their devices countless times throughout the day. Click-through rate (CTR) is a critical aspect of every web page, and many recommender engines within such sites aim at recommending to visitors the next item in order to maximize CTR. CTR is defined as the ratio of clicks on a specific link to its total number of impressions. Many recommender systems (RSs) suggest a ranked recommendation list to a user. The RS needs to sort the list of recommendations to maximize the chance for positive interaction within the user (i.e., click). RS sort the items within the recommendation list by their estimated CTR. Additionally, in other applications such as online advertising,

Authors' address: Amit Livne, livneam@post.bgu.ac.il; Roy Dor, rdo@post.bgu.ac.il; Eyal Mazuz, mazuze@post.bgu.ac.il; Tamar Didi, tamarin@post.bgu.ac.il; Bracha Shapira, bshapira@bgu.ac.il; Lior Rokach, liorrk@bgu.ac.il, Ben-Gurion University of the Negev, Beer-Sheva, Israel.

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.

© 2020 Association for Computing Machinery.

Manuscript submitted to ACM

Manuscript submitted to ACM

1

improving CTR is important for increasing the revenue. Thus, each CTR estimation is adjusted by the benefit the system receives for each candidate within the recommendation list. Many studies in the recommender systems domain propose various methods for improvement the accuracy of predicting the click-through for e-commerce sites or ads [8, 17, 28].

The authors of [26] introduced factorization machines (FMs) in which features are transformed into embedding space, and pairwise feature interactions are modeled as an inner product of those embedded representations. In [5], the authors suggested extending FMs to capture higher-order feature interactions. Guo et al. [13] presented, DeepFM, A model that combines the power of FM and deep learning to learn both low and high order feature interactions for solving the CTR prediction task. Features play a crucial role in the success of many predictive systems. However, using raw features rarely leads to optimal results. Thus, data scientists spend a lot of effort generating new features to improve predictive models [14].

Due to improvements in computational power, we have recently witnessed the emergence of new methodologies that consider deep learning techniques [13, 16, 20] to solve the CTR prediction task compared to traditional approaches [1, 26]. Moreover, these studies suggest skipping the feature generation phase and including it as part of a DL mechanism that captures the most suitable features automatically. However, these DL algorithms rely on training with a large amount of data to artificially generate features and may not perform as well when using only a small amount of data. To address this limitation, we suggest a new iterative boosting deep neural network (DNN) algorithm, XDBoost, that automatically crafting artificial features using a limited amount of data. Boosting refers to ensemble mechanism that combine several weak learners into a strong learner [9]. Its main idea is training predictors sequentially, where each predictor tries to correct its predecessor. In adaptive boosting (AdaBoost), suggested by [10], a base classifier (i.e., decision tree) is trained first; this classifier is used to make predictions on the training set, and this is followed by increasing the relative weights of misclassified training instances. In [6], the authors presented XGBoost, a scalable tree-based machine learning system, in which the method suggested by [11] was modified, improving the regularized objective. CatBoost, presented by [25], is an innovative boosting algorithm for processing categorical features. Our proposed XDBoost method involves the generation of new features that capture the estimated error distribution using a limited amount of data iteratively. The iterative boosting mechanism aids XDBoost to achieve more accurate CTR prediction.

Our main contributions are summarized as follows:

- (1) We propose a new neural network model, XDBoost, that integrates a boosting mechanism within state of the art (SOTA) DNN to address the CTR prediction task using limited amount of data. Incorporating estimated error via boosting mechanism within DNN allowing XDBoost to improve its CTR predictions.
- (2) We evaluate XDBoost on various datasets, including a public CTR prediction dataset and a proprietary real-world dataset, demonstrating consistent improvement on existing SOTA models for CTR prediction.
- (3) We analyze XDBoost sensitivity to the size of training data and compare the performance of XDBoost to other SOTA models when there is a limited amount of training data available. We show that XDBoost is especially beneficial for scenarios with limited amount of data available.
- (4) We explore XDBoost’s performance in addressing the cold start problem for new items. XDBoost outperforms all of the baselines.

The rest of this paper is structured as follows: Section 2 describes related work, and in Section 3 we present our proposed neural network XDBoost. In Section 4 we describe our evaluation and present the results. Finally, in Section 5 we discuss the results and in Section 6 we provide concluding remarks and discuss future work.

2 RELATED WORK

2.1 CTR Prediction

Development of the Internet and mobile devices in recent years has increased the importance of CTR prediction. As a result, many studies have been performed to try to maximize CTR prediction capabilities. CTR prediction is the task of predicting the probability that user u will click on item i in a given context c . CTR prediction is essential for businesses that rely on the pay-per-click (PPC) model. The two SOTA methods aimed at this task described below are based on a combination of traditional methods and neural network-based methods.

The first is factorization machine (FM) [26], a traditional method used to capture interactions between features that became very popular after the Netflix Prize competition [2, 3, 19, 30]. FM has been proven successful in many domains including: computer vision [12] and recommendation systems [26]. However, they often struggle to capture complex patterns and nonlinear interactions, as neural networks often do [4]. The second is a novel artificial neural network (ANN) architecture called Wide & Deep Learning [8] that was created by Google. Wide & Deep is utilizing wide linear models and a DNN, combining the benefits of memorization and generalization to capture more complex patterns and interactions. The Wide & Deep architecture has had significant impact on recent studies and contributed to the creation of new SOTA models. A variant of FM that extended [26], field-aware factorization machines (FFM), proposed by Juan et al. [17], addressed the task of CTR prediction. Feature engineering, a difficult manual task requiring time and domain experts, can significantly improve model performance [8]. Recent studies in the CTR prediction domain suggest using feature extraction methods to avoid manual feature extraction and implicitly generate new features within their models. Guo et al. [13] presented DeepFM to capture both low and high order features from raw features. Inspired by DeepFM and Wide & Deep, Lian et al. [20] replaced the FM layer of DeepFM with a novel cross-network they call the compressed interaction network (CIN), to capture feature interactions. Their method, referred to as xDeepFM, aims to learn certain bounded-degree feature interactions explicitly while learning arbitrary low and high order feature interactions implicitly. Huang et al. [16] used SENET [15] to dynamically evaluate feature importance, combined with bilinear feature interactions (FiBiNET), feeding them to a DNN for prediction.

In this research, we suggest a new method that considers implicit raw feature interactions to capture complex patterns. Additionally, we build on another traditional machine learning concept, boosting.

2.2 Boosting

Boosting is an abstract mechanism for improving the classification by reducing bias and variance usually applied in ensemble methods. A highly accurate prediction rule is found by combining rough and moderately inaccurate rules of thumb, which are called weak learners. A weak learner is defined as a classifier that is only slightly correlated with the true classification. This theory is based on Valiant's probably approximately correct (PAC) learning model [31]. The main variation between many boosting algorithms is their method of weighting training data points and hypotheses, used to create a set of simple rules that have high variance between them.

The most basic boosting algorithm is AdaBoost [10], which has undergone intense theoretical study and empirical testing. AdaBoost uses a set of decision trees called a forest, to perform prediction. It generates one tree at a time, calculates its prediction error, and gives each training instance a different weight. Then, the next tree is generated differently, due to the weight changes of the instances. Newer boosting methods suggested using gradient boosting factorization machines to incorporate a feature selection algorithm with FM [7]. For the CTR prediction task, Ling et al.

[21] suggested an ensemble method using DNN and the gradient boosting decision tree (GBDT) algorithm, and several other studies have suggested methods that use neural network-based FM combined with GBDT [32, 34].

XGBoost [6] is currently considered the SOTA boosting algorithm. It has gained popularity in the machine learning community due to its fast performance which stems from its use of parallelization and hardware optimization. XGBoost is based on GBDT, and it uses regularization by penalizing more complex models to prevent overfitting. It handles different types of sparsity patterns by learning the best missing values depending on training loss, and uses the weighted quantile sketch technique to find the optimal split points effectively. CatBoost [25] is a boosting method that is designed to handle categorical features. CatBoost generates new features from all available categorical feature combinations and utilizes them in the GBDT method for classification. Although most boosting algorithms were designed for ensemble methods, recent work has tried using the boosting concept on neural networks. Roy et al. [29] used boosting techniques to fine-tune a CNN image segmentation class weights in cases where labeled data is limited for solving the multi-class classification task. Notably, they change and update the loss function. Mosca and Magoulas [23] embraced the use of boosting techniques in DNNs by using transfer learning [24] to quickly generate DNN weak learners. However, they build on multiple classifiers and not a **single** classifier as we do.

In our model, we suggest an iterative mechanism similar to gradient boosting, however instead of summing the boosted prediction, we reuse the estimated error prediction as a feature for retraining our model. By using the estimated error prediction as a feature, we are able to fuse DNN with boosting. A detailed description of our model is provided in Section 3.

3 METHOD

We aim to design an iterative boosting DNN architecture for predicting the CTR. Specifically, we present a network that learns the estimated errors during the training phase and incorporates it iteratively as inputs to the network. Knowing the estimated errors and incorporate them within the model allows us to learn the relation between these error values and true labels and therefore improve the CTR predictions.

3.1 Problem Formulation

The input for the CTR prediction task is a set D composed of N quartettes. Each quartette $(u, i, c, Y) \in D$ denotes an interaction event where user u was exposed to an item i while considering contextual side information c regarding this interaction. c includes additional information about u and i in two different representations: categorical fields (i.e., content category or campaign language) and continuous fields (i.e., quality level). $Y \in \{0, 1\}$ is the associated label indicating user click behaviors ($Y=1$ indicates that user u clicked on item i under contextual side information c , $Y=0$ otherwise). The CTR prediction task is building a prediction model to estimate the probability for user u clicking a specific item i in a given context c .

3.2 XDBoost

The proposed method, XDBoost, is based on an iterative three-stage neural network model build on SOTA deep learning classifier (DLC) for CTR Prediction task. XDBoost constructed of a **single** DLC and **several** deep learning regressors (DLRs) that operate sequentially. DLC predictions are determined by the sigmoid function. DLR, is identical SOTA component as DLC. However, DLR aims to *estimate the error produced by the DLC's predictions*. Since the CTR is between the range of zero to one, its error distribution can vary from -1 to +1. Thus, instead of the sigmoid activation used in

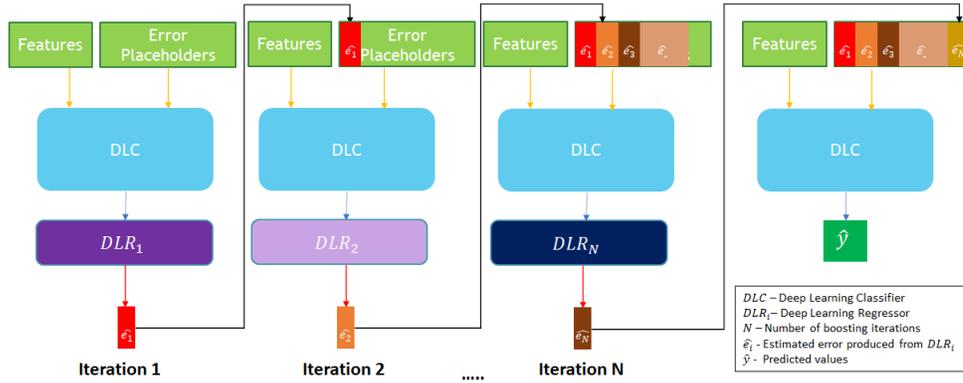


Fig. 1. The architecture of our proposed XDBoost model

DLC, DLR uses hyperbolic tangent activation. XDBoost' DLRs share the same structure. Each DLR is dedicated to a specific iteration of XDBoost.

XDBoost's goal is to derive a learning model that can learn feature interactions in an end-to-end manner without any feature engineering besides raw features. To use XDBoost, we use raw features that includes the user ID and item ID. We create N additional empty features that we call error placeholders. The estimated error distribution of the DLC in iteration i is derived by DLR_i , and it populates the empty placeholder $error_i$. XDBoost includes a single DLC classifier instance and N DLR instances, and can be seen in Figure 1.

Moreover, XDBoost has two properties: the number of boosting iterations defined by N and the error learning rate multiplier factor defined by E_{LR} . A brief description of iteration i consist of three stages is provided below:

- (1) The DLC, an existing SOTA classifier, aims to predict the CTR.
- (2) A new DLR instance DLR_i , learns the classifier's estimated error distribution $error_i$.
- (3) The DLC uses the error estimation provided by $error_i$ as an input, instead of the corresponding empty placeholder feature.

We describe XDBoost's creation, training, and prediction steps in the algorithms that follow. In order to create an instance of XDBoost, we apply algorithm 1.

Algorithm 1: Creating XDBoost

Input: N : number of boosting iterations (scalar), an integer scalar greater than 0.
 E_{LR} : error learning rate multiplier factor, a float scalar bounded by a predefined range $[0, 1]$
 DLC : a SOTA DLC for CTR prediction, a tensorflow model.
Output: $XDBoost$: a new instance of XDBoost

```

1  $XDBoost_N \leftarrow N$ 
2  $XDBoost_{E_{LR}} \leftarrow E_{LR}$ 
3  $XDBoost_{DLC} \leftarrow DLC()$ 
4 for  $i$  in  $XDBoost_N$  do
5    $XDBoost_{DLR_i} \leftarrow DLR_i()$ 
6 return  $XDBoost$ 

```

We initialize the properties of XDBoost in lines 1-3. In line 4, we create a new instance of DLC. Lines 4-5 describe the creation of XDBoost’s DLRs. Last, we return a new untrained instance of *XDBoost* in line 6. XDBoost’s training consists of two steps and performed sequentially as described in Algorithm 2. First, in lines 1-3, we create artificial zero-based features; the number of features is based on the number of boosting iterations declared while initializing XDBoost, defined by $XDBoost_N$. Second, in lines 4-11 we train XDBoost by applying the iterative process. Each iteration i is based on the following three stages. First stage described in lines 5-6: train $XDBoost_{DLC}$ and generate its predictions. Second stage described in lines 7-9: estimating the $XDBoost_{DLC}$ error distribution utilizing $XDBoost_{DLR_i}$. Third stage described in line 10-11: populate artificial feature i and retrain $XDBoost_{DLC}$. Last, in line 12 we return the trained instance of *XDBoost*.

Algorithm 2: Training XDBoost

Input: X_{train} : train set features; $\overrightarrow{y_{train}}$: train set binary target label data
XDBoost: an instance of XDBoost
Output: *XDBoost*: trained boosted instance of XDBoost

- 1 **for** i in $XDBoost_N$ **do**
- 2 $\overrightarrow{error}_i \leftarrow 0$
- 3 $X_{train}[-XDBoost_N + i] \leftarrow \overrightarrow{error}_i$
- 4 **for** i in $XDBoost_N$ **do**
- 5 $XDBoost_{DLC}.fit(X_{train}, \overrightarrow{y_{train}})$
- 6 $\overrightarrow{y}_{DLC} \leftarrow XDBoost_{DLC}.predict(X_{train})$
- 7 $\overrightarrow{error}_{DLC} \leftarrow (\overrightarrow{y_{train}} - \overrightarrow{y}_{DLC})$
- 8 $XDBoost_{DLR_i}.fit(X_{train}, \overrightarrow{error}_{DLC})$
- 9 $\overrightarrow{error}_{train} \leftarrow XDBoost_{DLR_i}.predict(X_{train})$
- 10 $X_{train}[-XDBoost_N + i] \leftarrow (XDBoost_{ELR} * \overrightarrow{error}_{train})$
- 11 $XDBoost_{DLC}.fit(X_{train}, \overrightarrow{y_{train}})$
- 12 **return** *XDBoost*

In contrast to traditional classifiers, generating predictions via XDBoost includes two steps that operate iteratively. In each iteration, the estimated error is predicted in the first step, and the second step populates the artificial features for the test set. Algorithm 3 provides a full description of how XDBoost generates predictions. First step described in line 2: estimating error utilizing $XDBoost_{DLR_i}$. Second step described in line 3: populate artificial feature i in the test set.

Algorithm 3: Predicting via XDBoost

Input: *XDBoost*: XDBoost model; X_{test} : test set features;
Output: $\overrightarrow{y_{test}}$: XDBoost CTR prediction

- 1 **for** i in $XDBoost_N$ **do**
- 2 $\overrightarrow{error}_{test} \leftarrow XDBoost_{DLR_i}.predict(X_{test})$
- 3 $X_{test}[-XDBoost_N + i] \leftarrow (XDBoost_{ELR} * \overrightarrow{error}_{test})$
- 4 $\overrightarrow{y}_{test} = XDBoost_{DLC}.predict(X_{test})$
- 5 **return** $\overrightarrow{y}_{test}$

Last, in line 4-5, we generate *XDBoost* predictions utilizing $XDBoost_{DLC}$ and return its predictions.

4 EXPERIMENTS

In this Section, we describe the extensive experiments conducted to answer the following research questions (RQs):

- (RQ1) Does XDBoost’s boosting mechanism improve the performance of SOTA DLC on raw features? Notably, does XDBoost outperform its base classifier? Specifically, we would like to investigate this question for various portions of the data to examine the effect of small training sets on the performance of XDBoost vs. SOTA DLC.
- (RQ2) Does XDBoost’s boosting mechanism improve the performance of SOTA boosting on raw features? Notably, we would like to investigate the effect of XDBoost performance when examining different portions of the available training set.
- (RQ3) One of the major challenges in recommendation systems is the cold start problem [27]. In classic recommendation systems, the cold start problem occurs when new users or items which may not have any ratings at all are added to the system. How does XDBoost’s ability to generalize and address the cold start problem compare to the abilities of SOTA deep learning algorithms for CTR prediction? How do different portions of the training set influence those results?

The experiments described later in the paper will be conducted in order to address these questions.

4.1 Experimental Settings

4.1.1 *Datasets.* We evaluate the effectiveness of our proposed method on the following datasets:

- **Taboola:** Taboola is an advertising company that provides advertisements, such as the "Around the Web" and "Recommended for You" boxes at the bottom of many online news articles. Taboola provides approximately 450 billion article recommendations each month for more than a billion unique users. The Taboola dataset consists of a sample of 15 days of ad click-through data which is ordered chronologically. It contains click logs with 34 million data instances. Each instance consists of 26 fields which reflect the elements of a single ad impression.
- **Avazu:** The Avazu dataset is widely used in many CTR model evaluation. It consists of several days of ad click-through data which is ordered chronologically. It contains click logs with 40 million data instances. For each instance, there are 24 fields which reflect the elements of a single ad impression. The Avazu dataset is publicly accessible.¹

Table 1 provides a summary of the datasets used in this study.

Table 1. Dataset Summary

	Taboola	Avazu
Number of records	34M	40M
Number of unique ads	94K	40M
Number of unique users	16M	Unknown
CTR (%)	49.9	15.2
Period of sample (days)	15	11
Number of fields	26	24

¹Avazu dataset: <http://www.kaggle.com/c/avazu-ctr-prediction>

4.1.2 *Evaluation Metrics.* We use two evaluation metrics in our experiments: the area under the ROC curve (**AUC**) and the **Log Loss**. Both metrics are widely used to evaluate CTR prediction performance [13, 16, 22].

- **AUC:** Area under ROC curve is a widely used metric in evaluating classification problems. The upper bound of the AUC is one, and the larger the AUC, the better.
- **Log Loss:** Takes into account the uncertainty of your prediction based on how much it varies from the actual label. The lower bound of the log loss is zero, indicating that the two distributions match perfectly, and a smaller value indicates better performance.

4.1.3 *Baselines.* In order to test the proposed method, and to answer the research questions we conducted a series of offline simulations and compared them to SOTA deep learning algorithms that are designed specifically to address the CTR prediction task. Additionally, because we integrate boosting mechanism in our method, we compare it to other SOTA algorithms that incorporate boosting mechanism. Specifically, we use the following algorithms as baselines:

(1) *SOTA deep learning algorithms:*

- **DeepFM** [13], a DNN model that integrates the architectures of FM and DNNs. It models low order feature interactions like FM and models high order feature interactions like DNNs.
- **xDeepFM** [20], a DNN model suggesting to learn certain bounded-degree feature interactions explicitly combined with low and high order feature interactions implicitly.
- **FiBiNET** [16], a DNN model suggesting a new way to calculate the feature interactions using bilinear function.

(2) *SOTA boosting algorithms:*

- **XGBoost** [6], an open source software library² which provides a gradient boosting model.
- **CatBoost** [25], an open source software library³ which provides a high performance gradient boosting model.

For the implementation of DeepFM, xDeepFM, and FiBiNET we use *deepctr* open source package.⁴

4.1.4 *Train-Validation-Test Split.* In [16, 22], the authors demonstrated their proposed methods on several datasets, using random splits. However, since both of the datasets used in our study 4.1.1 (Taboola and Avazu) are sorted chronologically. We split the data using the timestamp. This split simulates the real-world scenario, as real systems train on available data up to a certain timestamp and predict for the following days. Such split therefore prevents data leakage. We split each dataset as follows: the most recent 20% of the interactions are considered the test set. Of the remaining 80% of the data, the latest 8% is used as a validation set and 72% is used for training.

4.1.5 *Sub-Training Sets.* To address all RQs that relate to the influence of the size of the dataset, i.e. training the model using only a portion of the available training set, we created "sub" training sets from all of the available data while preserving chronologically constant validation and test sets, 8% and 20% of the data, respectively (as described in 4.1.4). We separate all available training data into two components: the last X% is considered a sub-training set, and the remaining data is not in use. An illustration of splitting the data into sub-training sets is provided in Figure 2. In the figure it can be observed that the range of all available training data (72%) consists of two parts: the sub-training set (X%) and the portion that is not in use. X can vary from 1% (indicating that a very small number of instances is used in the training process) to 72% (indicating that all of the available training data is used).

²XGBoost: <https://xgboost.readthedocs.io/en/latest/>

³CatBoost: <https://catboost.ai/>

⁴deepctr: <https://github.com/shenweichen/DeepCTR>

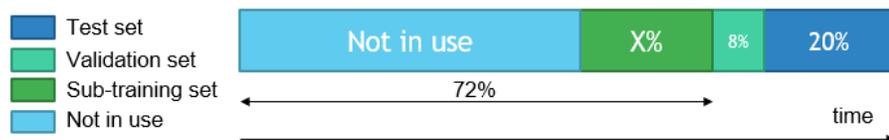


Fig. 2. Splitting data into subsets

We used several different sub-training sets. Specifically, we examine the performance of the following sub-training sets sizes: 1%, 5%, 10%, 20%, 40%, 60% and 72% while preserving the validation and test set constants.

4.1.6 Class Distributions. We conduct our experiments on two datasets, Taboola and Avazu (as described in Section 4.1.1). While the Taboola dataset label data is balanced, the Avazu dataset label data is unbalanced; specifically, its CTR is 15.2%, which creates a situation in which non-click instances are more dominant than click instances. To deal with this imbalance ratio and treat both classes equally, we set a weight for each class; in this way, we are able to mimic a balanced distribution. Thus, we increase the weights of the click class. The weights are determined by the distribution of the training set as follows:

$$ClassWeights = \begin{cases} \text{non-clicks} & 1.0 \\ \text{clicks} & \frac{Count(training\text{-}set_{non\text{-}clicks})}{Count(training\text{-}set_{clicks})} > 1.0 \end{cases} \quad (1)$$

4.1.7 Hyper Parameter tuning. In our experiments, we implement XDBoost with TensorFlow.⁵ The dimension of the embedding layer is set to 64. For the optimization we use the Adam optimizer [18] with a mini-batch size of 1,024, and the learning rate is set to 0.0001. For all deep models, the layer depth is set to three, and all activation functions are ReLU, except for the last activation which is sigmoid. The last activation in each DLR component is tanh. The number of neurons per layer is 128 for the Avazu dataset and 256 for the Taboola dataset. For classification components (i.e., DLC and all baselines), we use binary cross-entropy as the loss. For DLR components, we use the mean absolute error (MAE) as the loss, since MAE is not sensitive to outliers. We conduct our experiments using several RTX 2080 TI GPUs.

4.2 Results

In order to address our research questions, we generate several variants of XDBoost. Each variant is based on a different SOTA deep learning classifier for CTR prediction. Specifically, we suggest the following XDBoost variants: **XDBoost-DeepFM**, **XDBoost-xDeepFM**, and **XDBoost-FiBiNET**.

4.2.1 XDBoost Performance Compared to SOTA DLC Using Different Sub-Training Sets (RQ1). We compare the performance of our boosted algorithms to their non-boosted variants. In addition, the effectiveness of SOTA DLCs is reduced when using a small amount of data, and SOTA boosting algorithms do not need a lot of data to achieve high scores. Therefore, we explore the effectiveness of the suggested variants of XDBoost also compared to SOTA DLC baselines using different sub-training sets of various dataset sizes. Figure 3 and 4 present the AUC and log loss performance on Taboola test set while training on sub-training sets; 72% indicating that all of the available training data is used. As observed, all DLC algorithms using our XDBoost method outperform their base non-boosted classifiers.

⁵Tensorflow: <https://www.tensorflow.org>

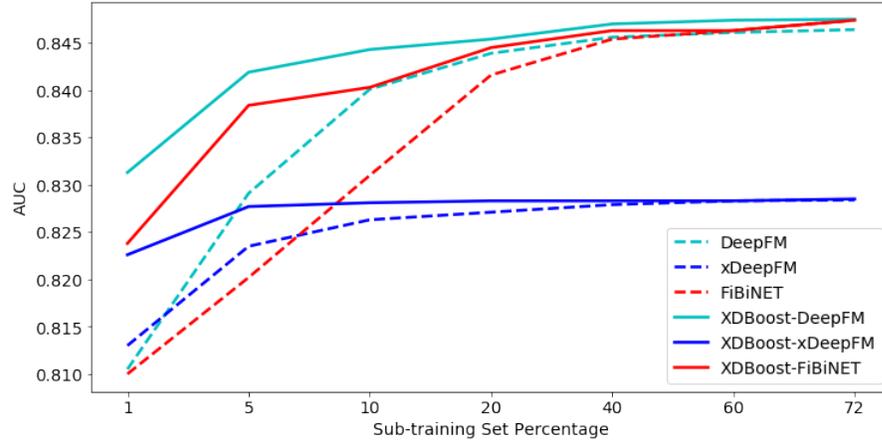


Fig. 3. XDBoost’ variants AUC performance compared to its base classifier – Taboola test set.

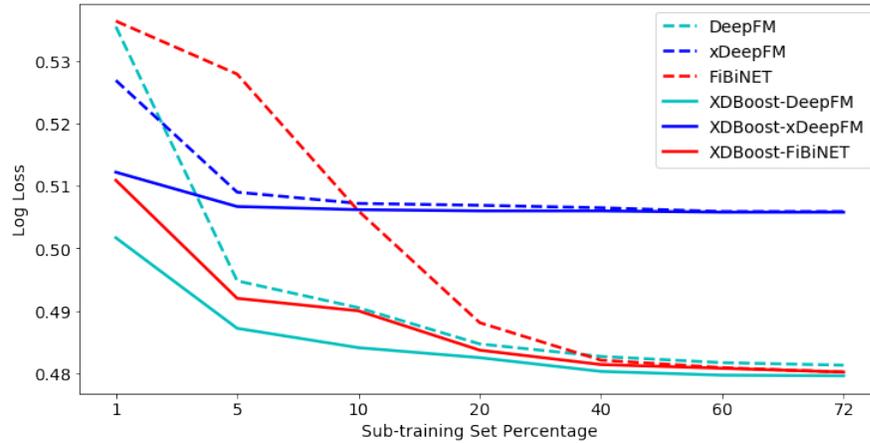


Fig. 4. XDBoost’ Variants log loss performance compared to its base classifier – Taboola test set.

For the Taboola dataset, XDBoost-DeepFM outperforms all baselines in terms of AUC. The improvement ranges between 2.25% to 2.62% for 1% sub-training set compared to other SOTA DLCs. Among the variants of XDBoost, XDBoost-DeepFM results consistently with the best results. When comparing the log loss for different sub-training sets, we can observe trends similar to those of the AUC. For both metrics, XDBoost variants outperform their base classifiers, when using small sub-training sets. However, its impact decreases when using sub-training sets greater than 40% (i.e., improvements vary between 0.15 and 0.13%). Notably, when using 60% and 72% sub-training sets, more than 20M records are used for training. Thus, allowing DLCs to fulfill their potential. Figures 5 and 6 present the AUC and log loss performance on Avazu test set while training on sub-training sets.

When using the Avazu test set, XDBoost-DeepFM outperforms all baselines, in terms of both the AUC and log loss. XDBoost-DeepFM achieves much better results than its base classifier, DeepFM. For instance, obtained an AUC

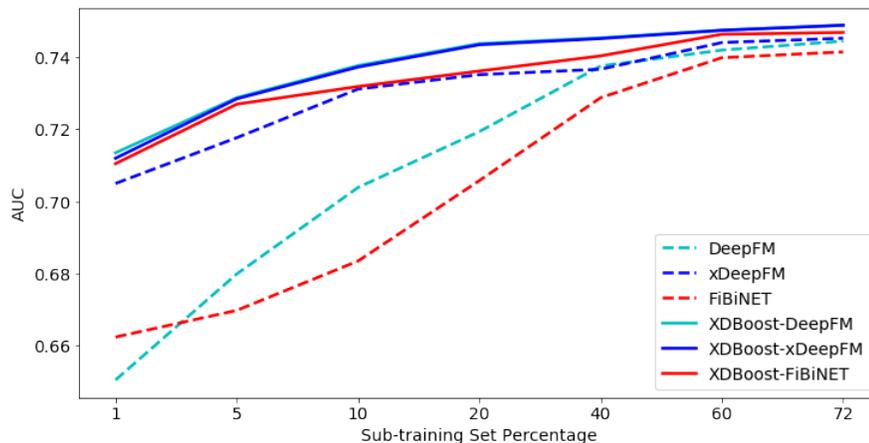


Fig. 5. XDBoost’ Variants AUC performance compared to its base classifier – Avazu test set.

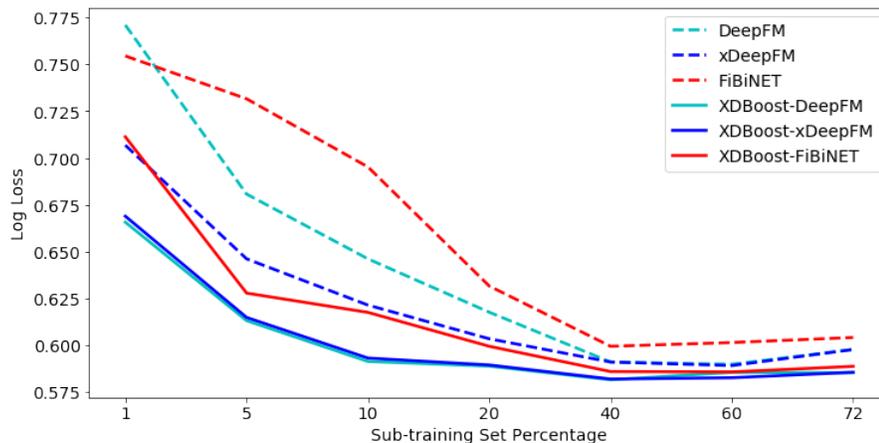


Fig. 6. XDBoost’ Variants log loss AUC performance compared to its base classifier – Taboola test set.

rate of 0.7378, XDBoost-DeepFM needed only 10% of the available data. However, to achieve similar results (0.7375), DeepFM needed to use 40% of the available data. For both metrics, XDBoost variants outperform their base classifiers, when using small sub-training sets. However, similar to the results with the Taboola dataset, their impact decreases when using sub-training sets greater than 40%. In order to determine if our proposed boosting achieves significant improvement over the baselines, we conducted the following statistical tests. We first used the adjusted Friedman test in order to reject the null hypothesis that all classifiers perform the same and then the apply the Bonferroni–Dunn test to examine whether our solution algorithm performs significantly better than existing baselines. The null-hypothesis with a confidence level of 99% for both datasets for all sub-training sets. We proceeded with the Bonferroni-Dunn test and found that *all* the variations of our suggested XDBoost significantly outperform all baselines with 99% confidence level for both datasets when considering sub-training sets smaller than 60%. For larger sub-training sets (i.e., 60% and

72%), XDBOost-DeepFM and XDBOost-FiBiNET statistically outperform all baselines with a 99% confidence level for both dataset as well.

4.2.2 XDBOost Performance Compared to SOTA Boosting Algorithms Using Different Sub-Training Sets (RQ2). SOTA boosting algorithms do not need a lot of data to achieve high scores. We explore the effectiveness of the suggested variants of XDBOost compared to SOTA boosting algorithms using different sub-training sets. The results of Taboola and Avazu datasets are presented in Tables 2 and 3 respectively. The best results in each column are denoted in bold. We use the Friedman and Bonferroni-Dunn tests as described in RQ1. Results that are statistically significant ($p < 0.01$) are denoted by an asterisk (*). As observed, all DLC-boosted algorithms outperform significantly the none- DLC boosting algorithms, and the XDBOost-DeepFM performed best.

Table 2. Comparison of SOTA Boosting Algorithms to XDBOost Variants - Taboola Dataset

Models\ Sub-Training	AUC							Log Loss						
	1%	5%	10%	20%	40%	60%	72%	1%	5%	10%	20%	40%	60%	72%
XGBoost	.815	.815	.816	.813	.814	.814	.815	.524	.524	.523	.526	.525	.525	.523
CatBoost	.830	.831	.832	.832	.832	.831	.832	.502	.501	.500	.500	.500	.501	.499
XDBOost-DeepFM	.831*	.842*	.844*	.845*	.847*	.847*	.848*	.500*	.487*	.484*	.482*	.480*	.479*	.479*
XDBOost-xDeepFM	.823*	.828*	.828*	.828*	.828*	.828*	.829*	.512*	.507*	.506*	.506*	.506*	.506*	.506*
XDBOost-FiBiNET	.823*	.838*	.840*	.844*	.846*	.846*	.847*	.511*	.492*	.490*	.484*	.481*	.481*	.480*

Table 3. Comparison of SOTA Boosting Algorithms to XDBOost Variants - Avazu Dataset

Models\ Sub-Training	AUC							Log Loss						
	1%	5%	10%	20%	40%	60%	72%	1%	5%	10%	20%	40%	60%	72%
XGBoost	.676	.698	.696	.696	.696	.699	.701	.727	.684	.681	.679	.647	.634	.629
CatBoost	.678	.706	.712	.727	.728	.728	.726	.725	.671	.693	.653	.622	.627	.620
XDBOost-DeepFM	.714*	.729*	.738*	.744*	.746*	.748*	.749*	.666*	.613*	.591*	.588*	.582*	.583*	.585*
XDBOost-xDeepFM	.712*	.728*	.737*	.743*	.745*	.747*	.749	.669*	.615*	.593*	.589*	.582*	.584*	.586
XDBOost-FiBiNET	.711*	.727*	.732*	.736*	.741*	.747*	.747*	.711*	.628*	.618*	.599*	.586*	.586*	.589*

When using the Taboola dataset, XDBOost-DeepFM outperforms all baselines for all sub-training sets. As we assumed, traditional tree-based boosting algorithms (i.e., XGBoost and CatBoost) perform well when using a small amount of data. However, they cannot model complex nonlinear interactions like DNNs.

When using the Avazu dataset, XDBOost-DeepFM is much more effective compared to XGBoost and CatBoost. XDBOost-DeepFM outperforms all baselines by at least 5.08% in terms of the AUC when training on 1% sub-training set.

4.2.3 Cold Start (RQ3). To examine the cold start handling of our method we used only the Taboola Dataset. In the Avazu dataset, each instance represents a unique ad ID (item) without any information regarding the user ID. Thus, we can conclude that this dataset is, by definition, a user- cold start dataset, thus, we could not test the effect of the portion of the cold start users on results. However, this is not the case for the Taboola dataset 1. In this RQ we focus on cold start problem for new ads. Notably, we investigate the effect of XDBOost performance when examining different portions of sub-training sets. A cold start ad is defined as an ad that was not appearing the the training set. Thus, we filter from the test set items which appear in the training set, resulting in a smaller test set with fewer records within the originally test set. We use several different sub-training sets, and therefore each sub-training require a corresponding smaller test set.

For example, when training our model on the 1% sub-training set, we filter the ads that appear in the 1% sub-training set from the original test set (i.e., 6M records) resulting in a smaller test set (i.e., 186K). Then we evaluated every model on the smaller test set. When using greater sub-training sets, the amount of records within the corresponding filtered test set decrease dramatically. We conducted several experiments on the Taboola dataset to explore the cold start problem with different sub-training sets. We compared all baselines to our suggested XDBoost-DeepFM, which yields the best results in RQ1 and RQ2. Table 4 present the results regarding new ads. The best results in each column are denoted in bold. Results that are statistically significant compared to *all* baseline with significance level of 95% and 99% are denoted by an asterisk (*) and two asterisks (**) respectively. We use the Friedman and Bonferroni-Dunn tests as described previously. Moreover, we mark the significant level of XDBoost-DeepFM compared to all baselines with color in Figure 7.

Table 4. Comparison of all baselines algorithms to XDBoost-DeepFM addressing the cold start when Examining New Ads.

Models\ Sub-Training	AUC							Log Loss						
	1%	5%	10%	20%	40%	60%	72%	1%	5%	10%	20%	40%	60%	72%
XGBoost	.780	.764	.755	.752	.740	.747	.745	.561	.573	.579	.581	.590	.585	.588
CatBoost	.793	.786	.781	.776	.769	.774	.774	.544	.549	.550	.557	.563	.558	.561
DeepFM	.771	.763	.776	.785	.793	.799	.797	.577	.576	.561	.549	.538	.534	.534
xDeepFM	.758	.765	.763	.762	.754	.759	.760	.582	.574	.575	.576	.580	.577	.578
FiBiNET	.782	.791	.791	.797	.795	.799	.799	.560	.537	.541	.539	.537	.531	.531
XDBoost-DeepFM	.798**	.804**	.799*	.800*	.797*	.800	.801	.541**	.532**	.536*	.534*	.535	.530	.530

We can observe from Table 4 that XDBoost-DeepFM outperforms all other SOTA models in all cases, in terms of both the AUC and log loss. While XDBoost excels for all sizes of sub-training sets, it has less impact sub-training set of size 40% and greater. As we can observe from Figure 7 comparing to SOTA boosting algorithms, XDBoost-DeepFM significant level remains 99% for all sub-training sets. In contrast, the significant confidence level of XDBoost-DeepFM decrease as the size of sub-training sets increase comparing to SOTA DLCs. Notably, XDBoost-DeepFM is not significant only when using 60% & 72% sub-training sets (i.e., more than 20M records for training) and only by comparing it to DeepFM and FiBiNET.

Models\ Sub-Training	New Ads							Legend		
	1%	5%	10%	20%	40%	60%	72%	$P_{value} \leq 0.01$	$P_{value} \leq 0.05$	Not Significant
XGBoost	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green
CatBoost	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green
DeepFM	Green	Green	Green	Green	Green	Red	Red	Green	Green	Green
xDeepFM	Green	Green	Green	Green	Green	Red	Red	Green	Green	Green
FiBiNET	Green	Green	Green	Green	Green	Red	Red	Green	Green	Green

Fig. 7. Significant level conducted by Bonferroni-Dunn test of XDBoost-DeepFM compared to each baseline for cold-start items.

5 DISCUSSION

Our experiments show that the XDBoost method has improved the performance of SOTA DLCs and SOTA boosting models on two datasets for both AUC and log loss (RQ1 and RQ2) for the CTR prediction task. On both datasets, we can

see the impact of XDBoost is increasing when using smaller sub-training sets. Thus, it can be useful, especially when the available training data is limited. The effectiveness of SOTA DLCs is reduced when using a small amount of data. In contrast, SOTA boosting algorithms do not need a lot of data to achieve high scores. XDBoost exploits the traits of both. It does not require a lot of data to achieve high scores; it is improving when using a limited amount of data. Thus, its superiority over DLCs is mostly significant in small training data scenarios which is a very common challenge since in many cases training data is expensive or difficult to collect [33]. When compared to SOTA tree-based boosting algorithms (RQ2) on balance dataset (i.e., Taboola), we can observe that these algorithms' performance barely improve even while considering small sub-training sets. In contrast, XDBoost keeps improving with greater sub-training sets. On the other hand, while observing results on unbalanced data (i.e., Avazu dataset) improvement can be seen for both SOTA tree-based boosting algorithms for bigger training sets. We assume the improvement of CatBoost is significant (7% improvement between 1% to 72% sub-training sets) on the Avazu dataset compared to XGBoost (3.3% improvement) and XDBoosted models (5% improvement) because the Avazu dataset consists of binary and categorical features only, hence it is more suitable for CatBoost that is superior in handling such data. XDBoost can significantly increase DNN models performance, especially when training data is limited. Thus, we recommend adopting the XDBoost architecture when using neural networks and the available data for training is limited in order to maximize the model's performance. We have demonstrated XDBoost on three examples of DNN algorithms. However, it can be applied easily to other DNNs as well. When examining the results it is important to note that a small improvement in the offline AUC is likely to lead to a significant increase in online CTR. As reported in [8], compared with logistic regression, Wide & Deep improves the offline AUC by 0.275%, and the improvement of the online CTR is 3.9%. Thus, even a minor improvements can be greatly beneficial.

6 CONCLUSIONS

In this paper, we propose XDBoost, an iterative boosted DNN architecture for predicting the CTR relying on raw features only. XDBoost learns the estimated error rate during the training phase and incorporate it iteratively as an input to the network. This allows learning of the relation between these error values and the true label. Extensive experiments conducted on two large-scale datasets show consistent improvement over existing SOTA models for CTR prediction. XDBoost's effectiveness compared to other SOTA methods increases when there is a limited amount of data. While exploring the cold start problem for new items, XDBoost outperformed all other baselines in terms of both the AUC and log loss. In the future, we plan to mix different SOTA regressors in the XDBoost architecture. Additionally, we aim to generalize XDBoost in order to solve other recommendation tasks, such as ranking and rating prediction.

REFERENCES

- [1] K Bauman, A Kornetova, V Topinskiy, and D Leshiner. 2010. CTR prediction based on click statistic. *Machine Learning in Online Advertising*, 8–13.
- [2] Robert M Bell and Yehuda Koren. 2007. Lessons from the Netflix prize challenge. *Acm Sigkdd Explorations Newsletter* 9, 2 (2007), 75–79.
- [3] James Bennett, Charles Elkan, Bing Liu, Padhraic Smyth, and Domonkos Tikk. 2007. KDD Cup and Workshop 2007. *SIGKDD Explor. Newsl.* 9, 2, 51–52. <https://doi.org/10.1145/1345448.1345459>
- [4] Christopher Bishop. 2006. *Pattern Recognition and Machine Learning*. Springer. <https://www.microsoft.com/en-us/research/publication/pattern-recognition-machine-learning/>
- [5] Mathieu Blondel, Akinori Fujino, Naonori Ueda, and Masakazu Ishihata. 2016. Higher-Order Factorization Machines. In *Proceedings of the 30th International Conference on Neural Information Processing Systems (Barcelona, Spain) (NIPS'16)*. Curran Associates Inc., Red Hook, NY, USA, 3359–3367.
- [6] Tianqi Chen and Carlos Guestrin. 2016. XGBoost: A Scalable Tree Boosting System. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (San Francisco, California, USA) (KDD '16)*. Association for Computing Machinery, New York, NY, USA, 785–794. <https://doi.org/10.1145/2939672.2939785>

- [7] Chen Cheng, Fen Xia, Tong Zhang, Irwin King, and Michael R. Lyu. 2014. Gradient Boosting Factorization Machines. In *Proceedings of the 8th ACM Conference on Recommender Systems* (Foster City, Silicon Valley, California, USA) (*RecSys '14*). Association for Computing Machinery, New York, NY, USA, 265–272. <https://doi.org/10.1145/2645710.2645730>
- [8] Heng-Tze Cheng, Levent Koc, Jeremiah Harmsen, Tal Shaked, Tushar Chandra, Hrishi Aradhye, Glen Anderson, Greg Corrado, Wei Chai, Mustafa Ispir, Rohan Anil, Zakaria Haque, Lichan Hong, Vihan Jain, Xiaobing Liu, and Hemal Shah. 2016. Wide & Deep Learning for Recommender Systems. In *Proceedings of the 1st Workshop on Deep Learning for Recommender Systems* (Boston, MA, USA) (*DLRS 2016*). Association for Computing Machinery, New York, NY, USA, 7–10. <https://doi.org/10.1145/2988450.2988454>
- [9] Harris Drucker, Corinna Cortes, Lawrence D Jackel, Yann LeCun, and Vladimir Vapnik. 1994. Boosting and other ensemble methods. *Neural Computation* 6, 6 (1994), 1289–1301.
- [10] Yoav Freund and Robert E. Schapire. 1995. A decision-theoretic generalization of on-line learning and an application to boosting. In *Computational Learning Theory*, Paul Vitányi (Ed.). Springer Berlin Heidelberg, Berlin, Heidelberg, 23–37.
- [11] Jerome Friedman, Trevor Hastie, Robert Tibshirani, et al. 2000. Additive logistic regression: a statistical view of boosting (with discussion and a rejoinder by the authors). *The annals of statistics* 28, 2 (2000), 337–407.
- [12] David Guillamet and Jordi Vitria. 2002. Non-negative Matrix Factorization for Face Recognition. In *Topics in Artificial Intelligence*, M. Teresa Escrig, Francisco Toledo, and Elisabet Golobardes (Eds.). Springer Berlin Heidelberg, Berlin, Heidelberg, 336–344.
- [13] Huifeng Guo, Ruiming Tang, Yunming Ye, Zhenguo Li, and Xiuqiang He. 2017. deepFM: a factorization-machine based neural network for CTR prediction. *arXiv preprint arXiv:1703.04247* 24, 3 (2017), 262–290.
- [14] Xinran He, Junfeng Pan, Ou Jin, Tianbing Xu, Bo Liu, Tao Xu, Yanxin Shi, Antoine Atallah, Ralf Herbrich, Stuart Bowers, and Joaquin Quiñero Candela. 2014. Practical Lessons from Predicting Clicks on Ads at Facebook. In *Proceedings of the Eighth International Workshop on Data Mining for Online Advertising* (New York, NY, USA) (*ADKDD'14*). Association for Computing Machinery, New York, NY, USA, 1–9. <https://doi.org/10.1145/2648584.2648589>
- [15] Jie Hu, Li Shen, and Gang Sun. 2018. Squeeze-and-Excitation Networks. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- [16] Tongwen Huang, Zhiqi Zhang, and Junlin Zhang. 2019. FiBiNET: Combining Feature Importance and Bilinear Feature Interaction for Click-through Rate Prediction. In *Proceedings of the 13th ACM Conference on Recommender Systems* (Copenhagen, Denmark) (*RecSys '19*). Association for Computing Machinery, New York, NY, USA, 169–177. <https://doi.org/10.1145/3298689.3347043>
- [17] Yuchin Juan, Yong Zhuang, Wei-Sheng Chin, and Chih-Jen Lin. 2016. Field-Aware Factorization Machines for CTR Prediction. In *Proceedings of the 10th ACM Conference on Recommender Systems* (Boston, Massachusetts, USA) (*RecSys '16*). Association for Computing Machinery, New York, NY, USA, 43–50. <https://doi.org/10.1145/2959100.2959134>
- [18] Diederik Kingma and Jimmy Ba. 2014. Adam: A Method for Stochastic Optimization. *International Conference on Learning Representations* (2014).
- [19] Yehuda Koren. 2009. The bellkor solution to the netflix grand prize. *Netflix prize documentation* 81, 2009 (2009), 1–10.
- [20] Jianxun Lian, Xiaohuan Zhou, Fuzheng Zhang, Zhongxia Chen, Xing Xie, and Guangzhong Sun. 2018. XDeepFM: Combining Explicit and Implicit Feature Interactions for Recommender Systems. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining* (London, United Kingdom) (*KDD '18*). Association for Computing Machinery, New York, NY, USA, 1754–1763. <https://doi.org/10.1145/3219819.3220023>
- [21] Xiaoliang Ling, Weiwei Deng, Chen Gu, Hucheng Zhou, Cui Li, and Feng Sun. 2017. Model Ensemble for Click Prediction in Bing Search Ads. In *Proceedings of the 26th International Conference on World Wide Web Companion* (Perth, Australia) (*WWW '17 Companion*). International World Wide Web Conferences Steering Committee, Republic and Canton of Geneva, CHE, 689–698. <https://doi.org/10.1145/3041021.3054192>
- [22] Bin Liu, Ruiming Tang, Yingzhi Chen, Jinkai Yu, Huifeng Guo, and Yuzhou Zhang. 2019. Feature Generation by Convolutional Neural Network for Click-Through Rate Prediction. In *The World Wide Web Conference* (San Francisco, CA, USA) (*WWW '19*). Association for Computing Machinery, New York, NY, USA, 1119–1129. <https://doi.org/10.1145/3308558.3313497>
- [23] Alan Mosca and George D Magoulas. 2017. Deep Incremental Boosting.
- [24] Sinno Jialin Pan and Qiang Yang. 2009. A survey on transfer learning. *IEEE Transactions on knowledge and data engineering* 22, 10 (2009), 1345–1359.
- [25] Liudmila Prokhorenkova, Gleb Gusev, Aleksandr Vorobev, Anna Veronika Dorogush, and Andrey Gulin. 2018. CatBoost: Unbiased Boosting with Categorical Features. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems* (Montréal, Canada) (*NIPS'18*). Curran Associates Inc., Red Hook, NY, USA, 6639–6649.
- [26] Steffen Rendle. 2010. Factorization Machines. In *Proceedings of the 2010 IEEE International Conference on Data Mining (ICDM '10)*. IEEE Computer Society, USA, 995–1000. <https://doi.org/10.1109/ICDM.2010.127>
- [27] Francesco Ricci, Lior Rokach, and Bracha Shapira. 2015. Recommender Systems: Introduction and Challenges. In *Recommender Systems Handbook*, Francesco Ricci, Lior Rokach, and Bracha Shapira (Eds.). Springer US, Boston, MA, 1–34. https://doi.org/10.1007/978-1-4899-7637-6_1
- [28] Matthew Richardson, Ewa Dominowska, and Robert Ragno. 2007. Predicting clicks: estimating the click-through rate for new ads. In *Proceedings of the 16th international conference on World Wide Web*. 521–530.
- [29] Abhijit Guha Roy, Sailesh Conjeti, Debodoot Sheet, Amin Katouzian, Nassir Navab, and Christian Wachinger. 2017. Error Corrective Boosting for Learning Fully Convolutional Networks with Limited Data. In *Medical Image Computing and Computer Assisted Intervention - MICCAI 2017*, Maxime Descoteaux, Lena Maier-Hein, Alfred Franz, Pierre Jannin, D. Louis Collins, and Simon Duchesne (Eds.). Springer International Publishing, Cham, 231–239.

- [30] Gábor Takács, István Pilászy, Bottyán Németh, and Domonkos Tikk. 2008. Matrix Factorization and Neighbor Based Algorithms for the Netflix Prize Problem. In *Proceedings of the 2008 ACM Conference on Recommender Systems (Lausanne, Switzerland) (RecSys '08)*. Association for Computing Machinery, New York, NY, USA, 267–274. <https://doi.org/10.1145/1454008.1454049>
- [31] Leslie G Valiant. 1984. A theory of the learnable. *Commun. ACM* 27, 11 (1984), 1134–1142.
- [32] Xiaochen Wang, Gang Hu, Haoyang Lin, and Jiayu Sun. 2019. A Novel Ensemble Approach for Click-Through Rate Prediction Based on Factorization Machines and Gradient Boosting Decision Trees. In *Web and Big Data*, Jie Shao, Man Lung Yiu, Masashi Toyoda, Dongxiang Zhang, Wei Wang, and Bin Cui (Eds.). Springer International Publishing, Cham, 152–162.
- [33] Karl Weiss, Taghi M Khoshgoftaar, and DingDing Wang. 2016. A survey of transfer learning. *Journal of Big data* 3, 1 (2016), 9.
- [34] Feng Zhou, Hua Yin, Lizhang Zhan, Huafei Li, Yeliang Fan, and Liu Jiang. 2018. A Novel Ensemble Strategy Combining Gradient Boosted Decision Trees and Factorization Machine Based Neural Network for Clicks Prediction. In *2018 International Conference on Big Data and Artificial Intelligence (BDAI)*. IEEE, 29–33. <https://doi.org/10.1109/BDAL.2018.8546685>