# XCrossNet: Feature Structure-Oriented Learning for Click-Through Rate Prediction

Runlong Yu[1], Yuyang Ye[2], Qi Liu[1(✉)], Zihan Wang[3], Chunfeng Yang[4], Yucheng Hu[4], and Enhong Chen[1]

[1] Anhui Province Key Laboratory of Big Data Analysis and Application,
School of Computer Science and Technology,
University of Science and Technology of China, Hefei, China
`yrunl@mail.ustc.edu.cn`, {`qiliuql,cheneh`}`@ustc.edu.cn`
[2] Management Science and Information Systems, Rutgers Business School,
Rutgers University, Newark, USA `yuyang.ye@rutgers.edu`
[3] MOE Key Laboratory of Computational Linguistics,
School of Electronics Engineering and Computer Science,
Peking University, Beijing, China `wzh@stu.pku.edu.cn`
[4] Tencent Inc, Shenzhen, China {`yannisyang,nikohu`}`@tencent.com`

**Abstract.** Click-Through Rate (CTR) prediction is a core task in nowadays commercial recommender systems. Feature crossing, as the mainline of research on CTR prediction, has shown a promising way to enhance predictive performance. Even though various models are able to learn feature interactions without manual feature engineering, they rarely attempt to individually learn representations for different feature structures. In particular, they mainly focus on the modeling of cross sparse features but neglect to specifically represent cross dense features. Motivated by this, we propose a novel Extreme Cross Network, abbreviated XCrossNet, which aims at learning dense and sparse feature interactions in an explicit manner. XCrossNet as a feature structure-oriented model leads to a more expressive representation and a more precise CTR prediction, which is not only explicit and interpretable, but also time-efficient and easy to implement. Experimental studies on Criteo Kaggle dataset show significant improvement of XCrossNet over state-of-the-art models on both effectiveness and efficiency.

## 1   Introduction

Accurate targeting of commercial recommender systems is of great importance, in which Click-Through Rate (CTR) prediction plays a key role. CTR prediction aims to estimate the ratio of clicks to the impression of a recommended item for a user. Therefore, we consider users have negative preferences instead of implicit feedbacks on those un-clicked items [25,26]. In common display advertising systems, advertisers expect lower costs to achieve a higher return on investment. The ad exchange platforms usually trade with advertisers and publishers according to the generalized second price of the maximum effective Cost Per Mille (eCPM). If CTR is overestimated, advertisers could waste campaign

budgets on the useless impression; On the other hand, if CTR is underestimated, advertisers would lose some valuable impressions and the campaigns may under deliver. With multi-billion dollar business on commercial recommendation today, CTR prediction has received growing interest from communities of both academia and industry [3,13,21].

In web-scale commercial recommender systems, the inputs of users' characteristics are in two kinds of structures. The first kind of structure is described by numerical or dense parameters, e.g., "`Age_years=22, Height_cm=165`". Each of such characteristics is formalized as a value associated with a numerical field, while the values are named as dense features. The second kind of structure is described by categorical or sparse parameters, e.g.,"`Gender=Female, Relationship=In love`". Each of such characteristics is formalized as a vector of one-hot encoding associated with a categorical field, while the vectors are named as sparse features. Research shows an important property of recommendation datasets for industrial use cases is the availability of both dense features and sparse features [22]. Thus, Criteo Kaggle dataset[1] is usually regarded as representative of real production use cases. Moreover, the number of dense and sparse features for industrial use cases are often 100s to 1000 with a 50:50 split[2].

Data scientists usually spend much time on interactions of raw features to generate better predictive models. Among these feature interactions, cross features, previously focused more on the product of sparse features, show a promising way to enhance the performance of prediction [15]. Owing to the fact that correct cross features are mostly task-specific and difficult to identify a priori, the crucial challenge is in automatically extracting sophisticated cross features hidden in high-dimensional data. Research on feature crossing as the mainline of CTR prediction has attracted widespread attention in recent years. Shallow models are simple, interpretable, and easy to scale, but limited in expressive ability. Alternatively, deep learning has shown powerful expressive capabilities, nevertheless, deep neural networks (DNNs) require many more parameters than tensor factorization to approximate high-order cross features. Besides, almost all deep models leverage multilayer perceptron (MLP) to learn high-order feature interactions, however, whether plain DNNs indeed effectively represent right functions of cross features remains an open question [10,21].

In addition, most methods neglect to represent cross dense features. There are three major patterns for handling dense features. First, dense features are discarded when crossing features, that is, dense features only participate in the linear part of the model [20]. Second, dense features are directly concatenated with the embeddings of sparse features, which could cause an important feature dimensionality imbalance problem [21]. Third, dense features are converted into sparse features through bucketing, which could introduce hyper-parameters and loss information of dense features [12].

---

[1] https://labs.criteo.com/2013/12/download-terabyte-click-logs/
[2] The statistics for the dense and sparse features and proportion are based on the survey outcome conducted in December 2019 with the MLPerf Advisory Board.

Based on all these observations, we propose a novel **Extreme Cross Network (XCrossNet)**, to represent feature structure-oriented interactions. Modeling with XCrossNet consists of three stages. In the Feature Crossing stage, we separately propose a cross layer for crossing dense features and a product layer for crossing sparse features. In the Feature Concatenation stage, cross dense features and cross sparse features interact through a concatenate layer and a cross layer. Lastly, in the Feature Selection stage, we employ an MLP for capturing non-linear interactions and their relative importance. Experimental results on Criteo Kaggle dataset demonstrate the superior performance of XCrossNet over the state-of-the-art baselines.

## 2    Related Work

Studies on CTR prediction can be categorized into five classes which will be respectively introduced below.

**Generalized linear models.** Logistic Regression (LR) models such as FTRL are widely used in CTR prediction for their simplicity and efficiency [9,16]. Ling Yan et al. argue that LR cannot capture nonlinear feature interactions and propose Coupled Group Lasso (CGL) to solve it [24]. Human efforts are usually needed for LR models. Gradient Boosting Decision Tree (GBDT) is a method to automatically do feature engineering and search interactions [4], then the transformed feature interactions can be fed into LR. In practice, tree-based models are more suitable for dense features but not for sparse features.

**Quadratic polynomial mappings and Factorization Machines.** Poly2 enumerates all pairwise feature interactions to avoid feature engineering which works well on dense features [2]. For sparse features, Factorization Machine (FM) and its variants project each feature into a low-dimensional vector and model cross features by inner product [20]. SFM introduces Laplace distribution to model the parameters and better fit the sparse data with a higher ratio of zero elements [17]. FFM enables each feature to have multiple latent vectors to interact with features from different fields [8]. As FM and its variants can only model order-2nd cross features. An efficient algorithm Higher-Order FM (HOFM) for training arbitrary-order cross features was proposed by introducing the ANOVA kernel [1]. As reported in [23], HOFM achieves marginal improvement over FM whereas using many more parameters and only its low-order (usually less than 5) form can be practically used.

**Implicit deep learning models.** As deep learning has shown promising representation capabilities, several models use deep learning to improve FM. Attention FM (AFM) enhances the importance of different order-2nd cross features via attention networks [23]. Neural FM (NFM) stacks deep neural networks on top of the output of the order-2nd cross features to model higher-order cross features [6]. FNN uses FM to pre-train low-order features and then feeds feature embeddings into an MLP [27]. In contrast, DSL uses MLP to pre-train high-order non-linear features and then feeds them with basis features into an FM layer [7]. Moreover, CCPM uses convolutional layers to explore local-global

dependencies of cross features [14]. IPNN (also known as PNN) feeds the interaction results of the FM layer and feature embeddings into an MLP [18]. PIN introduces a micro-network for each pair of fields to model pairwise cross features [19]. FGCNN combines a CNN and MLP to generate new features for feature augmentation [11]. However, all these approaches learn the high-order cross features in an implicit manner, therefore lack good model explainability.

**Wide&Deep based models.** Jianxun Lian et al. argue that implicit deep learning models focus more on high-order cross features but capture little low-order cross features [10]. To overcome this problem, there has been proposing a hybrid network structure, namely Wide&Deep, which combines a shallow component and a deep component with the purpose of learning both memorization and generalization [3]. Wide&Deep framework revolutionizes the development of CTR prediction, and attracts industry partners a lot from the beginning. As for the first Wide&Deep model proposed by Google [3], it combines a linear model (wide part) and DNN, while the input of the wide part still relies on feature engineering. Later on, DeepFM uses an FM layer to replace the wide component. Deep&Cross [21] and xDeepFM [10] take outer product of features at the bit- and vector-wise level respectively. However, xDeepFM uses so many parameters that great challenges are posed to identify important cross features in the huge combination space.

**AutoML based models.** There exist some pre-trained approaches using AutoML techniques to deal with cross features. AutoCross is proposed to search over subsets of candidate features to identify effective interactions [15]. This requires training the whole model to evaluate the selected feature interactions, but the candidate sets are incredibly many. AutoGroup treats the selection process of high-order feature interactions as a structural optimization problem, and solves it with Neural Architecture Search [12]. It achieves state-of-the-art performance on various datasets, but is too complex to be applied in industrial applications.

## 3   Extreme Cross Network (XCrossNet)

In this section, we will introduce the problem statement and describe the details of Extreme Cross Network (XCrossNet) in the following three steps: Feature Crossing, Feature Concatenation, and Feature Selection. The complete XCrossNet model is depicted in Figure 1.

### 3.1   Problem Statement

In web-scale commercial recommender systems, the inputs of users' characteristics are in two kinds of structures. The first kind of structure is described by numerical or dense parameters, denoted as $D$. The second kind of structure is described by categorical or sparse parameters, denoted as $S$. Suppose that the dataset for training consists of $n$ instances $([D; S], y)$, where $D = [D_1, D_2, \cdots, D_M]$ indicates dense features including $M$ numerical fields, and $S = [S_1, S_2, \cdots, S_N]$ indicates sparse features including $N$ categorical fields, and
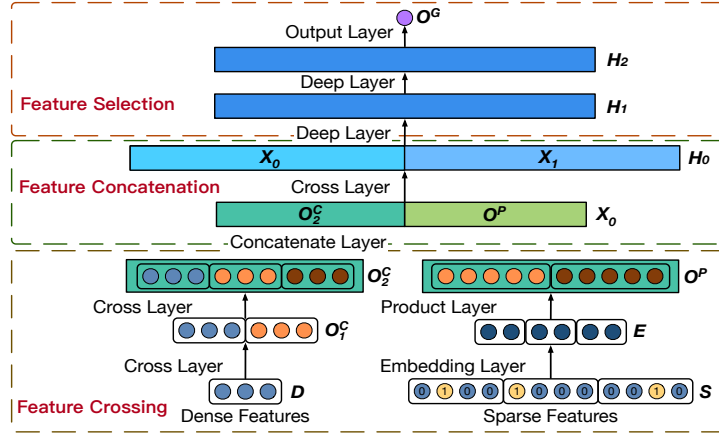
**Fig. 1.** The structure of XCrossNet.

$y \in \{0, 1\}$ indicates the user's click behaviors ($y = 1$ means the user clicked the item, and $y = 0$ otherwise). The task of CTR prediction is to build a prediction model $\hat{y} = pCTR\_Model([\boldsymbol{D}; \boldsymbol{S}])$ to estimate the ratio of clicks to impressions of a given feature context.

### 3.2 Feature Crossing

A cross feature is defined as a synthetic feature formed by multiplying (crossing) two features. Crossing combinations of features can provide predictive abilities beyond what those features can provide individually. Based on the definition, cross features can be generalized to high-order cases. If we consider individual features as order-1st features, an order-$k$th cross feature is formed by multiplying $k$ individual features.

**Cross layers on dense features.** First we introduce a novel cross layer for crossing dense features (see in Fig. 2). Cross layers have the following formula:

$$
\begin{aligned}
\boldsymbol{C_1} &= \boldsymbol{D} \cdot \boldsymbol{D}^{\mathsf{T}} \cdot \boldsymbol{W_{C,0}} + \boldsymbol{b_{C,0}}, \quad \boldsymbol{O_1^C} = [\boldsymbol{D}; \boldsymbol{C_1}], \\
\boldsymbol{C_{l+1}} &= \boldsymbol{D} \cdot \boldsymbol{C_l}^{\mathsf{T}} \cdot \boldsymbol{W_{C,l}} + \boldsymbol{b_{C,l}}, \quad \boldsymbol{O_{l+1}^C} = [\boldsymbol{O_l^C}; \boldsymbol{C_{l+1}}],
\end{aligned}
\tag{1}
$$

where $\boldsymbol{D} \in \mathbb{R}^M$ indicates the input dense features, and $\boldsymbol{C_l} \in \mathbb{R}^M$ is a column vector denoting the order-$(l + 1)$th cross features. Later we prove how $\boldsymbol{C_l}$ expresses multivariate polynomials of degree $(l + 1)$ after weighted mapping. $\boldsymbol{W_{C,l}}, \boldsymbol{b_{C,l}} \in \mathbb{R}^M$ are the weight and bias parameters respectively, and $\boldsymbol{O_l^C}, \boldsymbol{O_{l+1}^C}$ denote the outputs from the $l$-th and the $(l + 1)$-th cross layers.

We denote $\boldsymbol{\alpha} = [\alpha_1, \cdots, \alpha_M]$. If our proposed cross layer expresses any cross features of order-$(l+1)$th, it could approximate to any multivariate polynomials
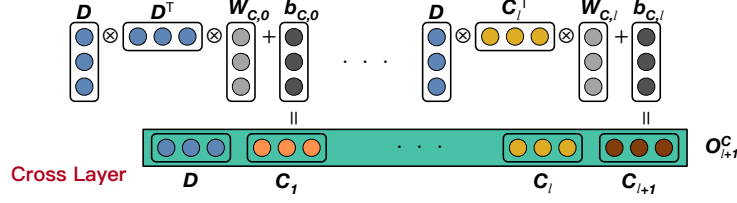
**Fig. 2.** The structure of cross layers.

of degree $(l+1)$, denoted as $P_{l+1}(\boldsymbol{D})$:

$$P_{l+1}(\boldsymbol{D}) = \left\{ \sum_{\boldsymbol{\alpha}} W_{\boldsymbol{\alpha}} D_1^{\alpha_1} D_2^{\alpha_2} \cdots D_M^{\alpha_M} \,\Big|\, |\boldsymbol{\alpha}| = l+1 \right\}, \tag{2}$$

where $|\boldsymbol{\alpha}| = \sum_{i=1}^{M} \alpha_i$. For simplicity, here we use $\boldsymbol{W^i} = [W_1^i, W_2^i, \cdots, W_M^i]$ to denote the original subscript of $\boldsymbol{W_{C,i}}$. We study the coefficient $\hat{W}_{\boldsymbol{\alpha}}$ given by $\boldsymbol{C_l^T} \cdot \boldsymbol{W^l}$ from cross layers, since it constitutes the output $\boldsymbol{O_{l+1}^C}$ from the $(l+1)$-th cross layer. Besides, the following derivations do not include bias terms. Then:

$$\begin{aligned} \boldsymbol{C_l^T} \cdot \boldsymbol{W^l} &= \left(\boldsymbol{C_{l-1}^T} \cdot \boldsymbol{W^{l-1}}\right) \cdot \left(\boldsymbol{D^T} \cdot \boldsymbol{W^l}\right) = \prod_{i=0}^{l} \boldsymbol{D^T} \cdot \boldsymbol{W^i} \\ &= \prod_{i=0}^{l} [D_1, D_2, \cdots, D_M]^T \cdot [W_1^i, W_2^i, \cdots, W_M^i]. \end{aligned} \tag{3}$$

Afterwards, let $\boldsymbol{I}$ denotes the multi-index vectors of orders $[0, 1, \cdots, l]$, and $I_j$ denotes the order of field $j$. Clearly $\boldsymbol{C_l^T} \cdot \boldsymbol{W^l}$ from cross layers approaches the coefficient $\hat{W}_{\boldsymbol{\alpha}}$ as:

$$\hat{W}_{\boldsymbol{\alpha}} = \sum_{k=1}^{M} \sum_{|\boldsymbol{I}|=\alpha_k} \prod_{j=1}^{M} W_j^{I_j}. \tag{4}$$

With $\boldsymbol{C_l^T} \cdot \boldsymbol{W^l}$ approximate to multivariate polynomials of degree $(l+1)$, the output $\boldsymbol{O_{l+1}^C}$ from the $(l+1)$-th cross layer that includes all cross features to order-$(l+1)$th could approximate polynomials in the following class:

$$P_{l+1}(\boldsymbol{D}) = \left\{ \sum_{\boldsymbol{\alpha}} W_{\boldsymbol{\alpha}} D_1^{\alpha_1} D_2^{\alpha_2} \cdots D_M^{\alpha_M} \,\Big|\, 0 \le |\boldsymbol{\alpha}| \le l+1 \right\}. \tag{5}$$

**Embedding and product layers on sparse features.** Here we introduce the embedding layer and product layer for crossing sparse features (see in Fig. 3). As sparse features $\boldsymbol{S}$ are represented as vectors of one-hot encoding of high-dimensional spaces, we employ an embedding layer to transform these one-hot encoding vectors into dense vectors $\boldsymbol{E}$ as:

$$\begin{aligned} \boldsymbol{E} &= [\boldsymbol{E_1}, \cdots, \boldsymbol{E_i}, \cdots, \boldsymbol{E_N}], \\ \boldsymbol{E_i} &= \text{embed}(\boldsymbol{S_i}), \ \left(\boldsymbol{E_i} \in \mathbb{R}^K, i = 1, \cdots, N\right) \end{aligned} \tag{6}$$

where $\boldsymbol{S_i}$ indicates the input sparse feature of field $i$, $K$ denotes the embedding size, and $\boldsymbol{E_i}$ denotes the feature embedding of field $i$.
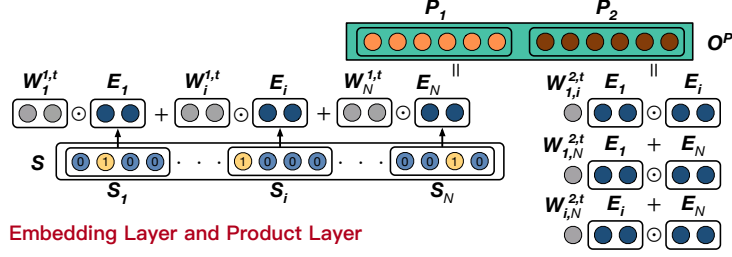
**Fig. 3.** The structure of embedding layer and product layer.

Afterwards, we can propose a product layer for cross sparse features. First, we donate order-2nd cross sparse features as $\boldsymbol{P_2}$, and order-1st sparse features as $\boldsymbol{P_1}$, thus the output of product layer is $\boldsymbol{O^P} = [\boldsymbol{P_1}; \boldsymbol{P_2}]$.

The cross feature of two sparse features of field $i$ and field $j$ equals the inner product of two embedding vectors as $\langle \boldsymbol{E_i}, \boldsymbol{E_j} \rangle$. Intuitively, we expect cross features to be vectors, so we concatenate the weighted sums of inner products to formulate order-2nd cross features as:

$$\boldsymbol{P_2} = [P_2^1, \cdots, P_2^t, \cdots, P_2^T], \tag{7}$$

where $T$ is the size of product layer, and $\boldsymbol{P_2}$ is a $T$ dimensional vector, of each dimension $P_2^t$ denotes a weighted sum of inner products of two sparse features. Thus, we have $P_2^t = \sum_{i=1}^{N} \sum_{j=1}^{N} W_{i,j}^{2,t} \langle \boldsymbol{E_i}, \boldsymbol{E_j} \rangle$. We assume that the weighted parameter $W_{i,j}^{2,t} = \Theta_i^t \cdot \Theta_j^t$ for reduction, so $P_2^t$ can be given as:

$$P_2^t = \sum_{i=1}^{N} \sum_{j=1}^{N} \Theta_i^t \cdot \Theta_j^t \langle \boldsymbol{E_i}, \boldsymbol{E_j} \rangle = \left\langle \sum_{i=1}^{N} \Theta_i^t \cdot \boldsymbol{E_i}, \sum_{j=1}^{N} \Theta_j^t \cdot \boldsymbol{E_j} \right\rangle. \tag{8}$$

The feature vector of order-1st features has a similar formula as follows:

$$\boldsymbol{P_1} = [P_1^1, \cdots, P_1^t, \cdots, P_1^T], \tag{9}$$

where $\boldsymbol{P_1}$ is a $T$ dimensional vector, of each dimension $P_1^t$ denotes a weighted sum of sparse features. The weighted feature can be expressed as inner product $\langle \boldsymbol{W_i^{1,t}}, \boldsymbol{E_i} \rangle$. Thus, we have $P_1^t = \sum_{i=1}^{N} \langle \boldsymbol{W_i^{1,t}}, \boldsymbol{E_i} \rangle$.

### 3.3 Feature Concatenation

In the Feature Concatenation stage, in order to learn feature interactions of different structures, cross dense features $\boldsymbol{O^C}$ and cross sparse features $\boldsymbol{O^P}$ are concatenated as a vector through a concatenate layer, then the concatenated feature vector is fed into a cross layer, which can be expressed as:

$$\begin{aligned}
\boldsymbol{X_0} &= [\boldsymbol{O^C}; \boldsymbol{O^P}], \\
\boldsymbol{X_1} &= \boldsymbol{X_0} \cdot \boldsymbol{X_0^\mathsf{T}} \cdot \boldsymbol{W_{X,0}} + \boldsymbol{b_{X,0}}, \quad \boldsymbol{H_0} = [\boldsymbol{X_0}; \boldsymbol{X_1}],
\end{aligned} \tag{10}$$

where $\boldsymbol{X_0}$ denotes the concatenated feature of cross dense features and cross sparse features, $\boldsymbol{X_1}$ denotes the cross features between two kinds of feature structures, $\boldsymbol{H_0}$ denotes the output from this cross layer, and $\boldsymbol{W_{X,0}}, \boldsymbol{b_{X,0}}$ are the weight and bias parameters of this cross layer.

### 3.4   Feature Selection

In the Feature Selection stage, we employ an MLP to capture non-linear interactions and the relative importance of cross features. The deep layers and the output layer respectively have the following formula:

$$\begin{aligned}
\boldsymbol{H_i} &= \mathrm{ReLU}(\boldsymbol{W_{H,i-1}} \cdot \boldsymbol{H_{i-1}} + \boldsymbol{b_{H,i-1}}), \\
O^G &= \mathrm{Sigmoid}(\boldsymbol{W_{H,i}} \cdot \boldsymbol{H_i} + \boldsymbol{b_{H,i}}),
\end{aligned} \tag{11}$$

where $\boldsymbol{H_i}, \boldsymbol{H_{i-1}}$ are hidden layers, $\mathrm{ReLU}(\cdot)$ and $\mathrm{Sigmoid}(\cdot)$ are activation functions, $\boldsymbol{W_{H,i}}, \boldsymbol{W_{H,i-1}}$ are weights, and $\boldsymbol{b_{H,i}}, \boldsymbol{b_{H,i-1}}$ are biases, and $O^G$ is the output result.

For CTR prediction, the loss function is the Logloss as follows:

$$\boldsymbol{\mathcal{L}} = -\frac{1}{n} \sum_{i=1}^{n} y_i \log(O^G) + (1 - y_i) \log(1 - O^G), \tag{12}$$

where $n$ is the total number of training instances. The optimization process is to minimize the following objective function:

$$\boldsymbol{\mathcal{J}} = \boldsymbol{\mathcal{L}} + \lambda ||\boldsymbol{\Theta}||, \tag{13}$$

where $\lambda$ denotes the regularization term, and $\boldsymbol{\Theta}$ denotes the set of learning parameters, including cross layers, embedding layer, product layer, deep layers and output layer.

## 4   Experiments

In this section, extensive experiments are conducted to answer the following research questions[3]:

**RQ1:** How does XCrossNet perform compared with the state-of-the-art CTR prediction models?
**RQ2:** How does the feature dimensionality imbalance impact CTR prediction?
**RQ3:** How do hyper-parameter settings impact the performance of XCrossNet?

### 4.1   Experimental Setup

**Dataset.** Experiments are conducted on Criteo Kaggle dataset, which is from a world-wide famous Demand-Side Platforms. Criteo Kaggle dataset contains one month of $45,840,617$ ad click instances. It has 13 integer feature fields and 26 categorical feature fields. We select 7 consecutive days of samples as the training set while the next one day for evaluation.

---

[3] We release the source code at https://github.com/bigdata-ustc/XCrossNet/.

| Model | AUC(%) | Logloss |
|---|---|---|
| LR | 78.00 | 0.5631 |
| GBDT | 78.62 | 0.5560 |
| FM | 79.09 | 0.5500 |
| AFM | 79.13 | 0.5517 |
| FFM | 79.80 | 0.5438 |
| CCPM | 79.55 | 0.5469 |
| Wide&Deep | 79.77 | 0.5446 |
| Cross | 78.70 | 0.5550 |
| Deep&Cross | 79.76 | 0.5445 |
| FNN | 79.87 | 0.5428 |
| DeepFM | 79.91 | 0.5423 |
| IPNN | 80.13 | 0.5399 |
| PIN | 80.18 | 0.5394 |
| CIN | 78.81 | 0.5538 |
| xDeepFM | 80.06 | 0.5408 |
| FGCNN | 80.22 | 0.5389 |
| AutoGroup | 80.28 | 0.5384 |
| **XCrossNet** | **80.68** | **0.5339** |

**Table 1.** Performance comparison of different CTR prediction models.



**Fig. 4.** Training time comparison of different CTR prediction models.



**Fig. 5.** Impact of feature dimensionality imbalance.

**Baselines.** As aforementioned, we use following highly related state-of-the-art models as baselines: **LR** [9], **GBDT** [4], **FM** [20], **AFM** [23], **FFM** [8], **CCPM** [14], **Wide&Deep** [3], **Deep&Cross** [21] and its shallow part **Cross** network, **FNN** [27], **DeepFM** [5], **IPNN** [18], **PIN** [19], **xDeepFM** [10] and its shallow part **CIN**, **FGCNN** [11], and **AutoGroup** [12].

**Hyper-parameter settings.** For model optimization, we use Adam with a mini-batch size of 4096, and the learning rate is set as 0.001. We use the L2 regularization with $\lambda = 0.0001$ for all neural network models. For Wide&Deep, Deep&Cross, FNN, DeepFM, IPNN, PIN, xDeepFM, and XCrossNet, the numbers of neurons per deep layer are 400, and the depths of deep layers are set as 2. For our XCrossNet, the number of cross layers on dense features is set as $l=4$. In the main experiments, we set the embedding size for all models be a fixed value of 20.

### 4.2   Overall Performance (RQ1)

Table 4 summarizes the performance of all compared methods on Criteo Kaggle datasets, while the training time on Tesla K80 GPUs is shown in Fig. 4 for comparison of efficiency. From the experimental results, we have the following key observations: Firstly, most neural network models outperform linear models (i.e.,
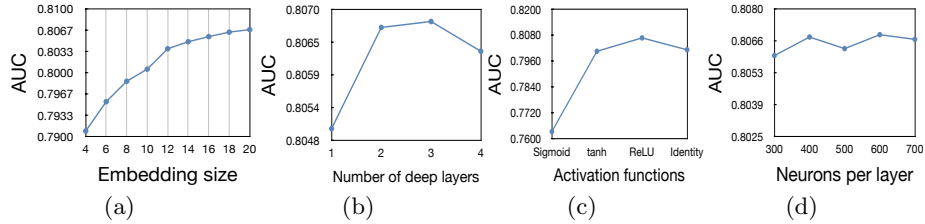
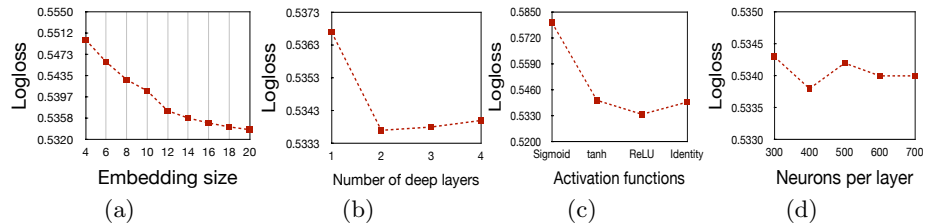**Fig. 6.** Impact of network hyper-parameters on AUC performance.



**Fig. 7.** Impact of network hyper-parameters on Logloss performance.

LR), tree-based models (i.e., GBDT), and FM variants (i.e., FM, FFM, AFM), which indicates MLP can learn non-linear feature interactions and endow better expressive ability. Meanwhile, comparing IPNN, PIN with FNN, Wide&Deep based models, we find that explicitly modeling low-order feature interactions can simplify the training of MLP and boost the performance. Secondly, XCrossNet achieves the best performance. Statistically, XCrossNet significantly outperforms the best baseline in terms of AUC and Logloss on $p$-value $< 0.05$ level, which indicates feature structure-oriented learning can provide better predictive abilities. Thirdly, from the training time comparison, we can observe XCrossNet is very efficient, especially compared to field-aware models, mainly because these models further allow each feature to learn several vectors where each vector is associated with a field, which leads to huge parameter consumption and time consumption.

### 4.3   Feature Dimensionality Imbalance Study (RQ2)

In XCrossNet, we denote $\frac{\dim(O^C)}{\dim(O^P)} \Big/ \frac{M}{N}$ as the balance index of dimensions of dense and sparse features. Noted that, the dimension of cross dense features $O^C$ equals $M \cdot l$, increasing with the depth of cross layers. As for Criteo Kaggle dataset, $M = 13$ and $N = 26$, we set the depths of cross layers from 1 to 8, while the corresponding dimensions of cross dense features are from 13 to 104. Experimental results are shown in Fig. 5 in terms of AUC. We can observe that increasing the depth of cross layers benefits XCrossNet to achieve stable improvements on AUC performance, mainly because the higher dimensions of cross

dense features are able to boost the balance index, which results in relatively balanced impacts of dense and sparse features on prediction.

### 4.4   Hyper-parameter Study (RQ3)

We study the impact of hyper-parameters of XCrossNet, including (1) embedding size; (2) number of deep layers; (3) activation function; (4) neurons per layer. Figures 6a and 7a demonstrate the impact of embedding size. We can observe that model performance boosts steadily when the embedding size increase from 4 to 20. Even with very low embedding sizes, XCrossNet still has comparable performance to some popular Wide&Deep based models with high embedding size. Specifically, XCrossNet achieves AUC$>$ 0.800 and Logloss$<$ 0.541 with embedding size set as 10, which is even better than DeepFM with embedding size set as 20. Figures 6b and 7b demonstrate the impact of the number of deep layers. The model performance boosts with the depth of MLP at the beginning. However, it starts to degrade when the depth of MLP is set to greater than 3. As shown in Figures 6c and 7c, ReLU is indeed more appropriate for hidden neurons of deep layers compared with different activation functions. As shown in Figures 6d and 7d, model performance barely boosts as the number of neurons per layer increasing from 300 to 700. We consider 400 is a more suitable setting to avoid the model being overfitting.

## 5   Conclusion

Due to the fact that previous work rarely attempts to individually learn representations for different feature structures, this paper presented a novel feature structure-oriented learning model, namely Extreme Cross Network (XCrossNet), for improving CTR prediction in recommender systems. A XCrossNet model starts with a Feature Crossing stage, followed by a Feature Concatenation stage and a Feature Selection stage. The main contribution of our approach is to represent dense and sparse feature interactions in an explicit and efficient way. Empirical studies verified the effectiveness of our model on Criteo Kaggle dataset.

## References

1. Blondel, M., et al.: Higher-order factorization machines. In: NeurIPS. pp. 3351–3359 (2016)
2. Chang, Y., et al.: Training and testing low-degree polynomial data mappings via linear SVM. Journal of Machine Learning Research (JMLR) **11**, 1471–1490 (2010)

3. Cheng, H.T., et al.: Wide & deep learning for recommender systems. In: 1st Workshop on Deep Learning for Recommender Systems. pp. 7–10 (2016)
4. Friedman, J.H.: Greedy function approximation: a gradient boosting machine. Annals of Statistics pp. 1189–1232 (2001)
5. Guo, H., et al.: Deepfm: A factorization-machine based neural network for CTR prediction. In: IJCAI. pp. 1725–1731 (2017)
6. He, X., Chua, T.S.: Neural factorization machines for sparse predictive analytics. In: SIGIR. pp. 355–364 (2017)
7. Huang, Z., et al.: An ad ctr prediction method based on feature learning of deep and shallow layers. In: CIKM. pp. 2119–2122 (2017)
8. Juan, Y., et al.: Field-aware factorization machines for ctr prediction. In: RecSys. pp. 43–50 (2016)
9. Lee, K., et al.: Estimating conversion rate in display advertising from past performance data. In: SIGKDD. pp. 768–776. ACM (2012)
10. Lian, J., et al.: xdeepfm: Combining explicit and implicit feature interactions for recommender systems. In: SIGKDD. pp. 1754–1763 (2018)
11. Liu, B., et al.: Feature generation by convolutional neural network for click-through rate prediction. In: WWW. pp. 1119–1129 (2019)
12. Liu, B., et al.: Autogroup: Automatic feature grouping for modelling explicit high-order feature interactions in CTR prediction. In: SIGIR. pp. 199–208. ACM (2020)
13. Liu, Q., et al.: Personalized travel package recommendation. In: ICDM. pp. 407–416. IEEE (2011)
14. Liu, Q., et al.: A convolutional click prediction model. In: CIKM. pp. 1743–1746 (2015)
15. Luo, Y., et al.: Autocross: Automatic feature crossing for tabular data in real-world applications. In: SIGKDD. pp. 1936–1945 (2019)
16. McMahan, H.B., et al.: Ad click prediction: a view from the trenches. In: SIGKDD. pp. 1222–1230. ACM (2013)
17. Pan, Z., et al.: Sparse factorization machines for click-through rate prediction. In: ICDM. pp. 400–409. IEEE (2016)
18. Qu, Y., et al.: Product-based neural networks for user response prediction. In: ICDM. pp. 1149–1154. IEEE (2016)
19. Qu, Y., et al.: Product-based neural networks for user response prediction over multi-field categorical data. ACM Transactions on Information Systems (ACM TOIS) **37**(1), 1–35 (2018)
20. Rendle, S.: Factorization machines. In: ICDM. pp. 995–1000. IEEE (2010)
21. Wang, R., et al.: Deep & cross network for ad click predictions. In: ADKDD, pp. 1–7 (2017)
22. Wu, C.J., et al.: Developing a recommendation benchmark for mlperf training and inference. arXiv preprint arXiv:2003.07336 (2020)
23. Xiao, J., et al.: Attentional factorization machines: learning the weight of feature interactions via attention networks. In: IJCAI. pp. 3119–3125 (2017)
24. Yan, L., et al.: Coupled group lasso for web-scale CTR prediction in display advertising. In: ICML. vol. 32, pp. 802–810 (2014)
25. Yu, R., et al.: Collaborative list-and-pairwise filtering from implicit feedback. IEEE Transactions on Knowledge and Data Engineering (IEEE TKDE) . https://doi.org/10.1109/TKDE.2020.3016732
26. Yu, R., et al.: Multiple pairwise ranking with implicit feedback. In: CIKM. pp. 1727–1730. ACM (2018)
27. Zhang, W., et al.: Deep learning over multi-field categorical data. In: European Conference on Information Retrieval. pp. 45–57. Springer (2016)