# Detecting Arbitrary Order Beneficial Feature Interactions for Recommender Systems

Yixin Su
yixins1@student.unimelb.edu.au
The University of Melbourne
Australia

Yunxiang Zhao
zhaoyx1993@163.com
Bejing Institute of Biotechnology
China

Sarah Erfani*
sarah.erfani@unimelb.edu.au
The University of Melbourne
Australia

Junhao Gan
junhao.gan@unimelb.edu.au
The University of Melbourne
Australia

Rui Zhang
rayteam@yeah.net
www.ruizhang.info, China

## ABSTRACT

Detecting beneficial feature interactions is essential in recommender systems, and existing approaches achieve this by examining all the possible feature interactions. However, the cost of examining all the possible higher-order feature interactions is prohibitive (exponentially growing with the order increasing). Hence existing approaches only detect limited order (e.g., combinations of up to four features) beneficial feature interactions, which may miss beneficial feature interactions with orders higher than the limitation. In this paper, we propose a hypergraph neural network based model named HIRS. HIRS is the first work that directly generates beneficial feature interactions of arbitrary orders and makes recommendation predictions accordingly. The number of generated feature interactions can be specified to be much smaller than the number of all the possible interactions and hence, our model admits a much lower running time. To achieve an effective algorithm, we exploit three properties of beneficial feature interactions, and propose deep-infomax-based methods to guide the interaction generation. Our experimental results show that HIRS outperforms state-of-the-art algorithms by up to 5% in terms of recommendation accuracy.

## CCS CONCEPTS

• **Information systems** → **Recommender systems**; • **Mathematics of computing** → *Hypergraphs*; • **Computing methodologies** → *Feature selection*.

## KEYWORDS

Recommender Systems, Feature Interactions, Graph Neural Networks, Mutual Information Maximization

*Corresponding Author.

## 1 INTRODUCTION

Features interactions (e.g., co-occurrence of user/item attributes) are essential in providing accurate recommendation predictions [3]. Pairwise feature interactions, as the most basic feature interaction form, have been utilized in many studies [14, 21, 34]. Recently, researchers found that more sophisticated forms of feature interactions (i.e., high-order interactions) help achieve more accurate recommendation predictions [23, 25]. For example, if we observe that male teenagers like Nolan's sci-fi movies, the interactions (order-4) between user gender, user age, movie director, and movie genre provide useful information for recommendation that pairwise feature interactions cannot.

Since the number of all the possible feature interactions grows exponentially as the order increases, considering all the possible high-order feature interactions is prohibitive due to the high complexity. Meanwhile, considering the feature interactions irrelevant to the recommendation result may introduce noise and decrease the prediction accuracy [26, 40]. Therefore, a practical solution to consider high-order feature interactions is to detect a small set of most *beneficial* feature interactions and perform predictions based on the beneficial ones [44]. A set of feature interactions is considered most beneficial if learning from them results in a more accurate prediction than learning from other sets [33]. However, existing interaction detection methods have to examine all the possible feature interactions to find the most beneficial feature interactions [25, 33]. This leaves the complexity issue unsolved and they resort to detecting only limited order (e.g., combinations of up to four features) beneficial feature interactions (Figure 1 *left*). Hence, it is still an urgent yet challenging problem to efficiently detect arbitrary order beneficial feature interactions.

In this work, we propose a method that can efficiently detect beneficial feature interactions of arbitrary orders. Our method learns a neural network that directly generates the beneficial feature interactions (Figure 1 *right*). The number of generated feature interactions can be specified to be much smaller than the number of all the
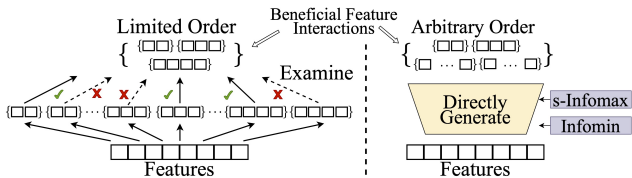
**Figure 1:** *Left*: existing studies detect beneficial feature interactions within limited orders by examining all the possible feature interactions. *Right*: we detect arbitrary order beneficial feature interactions by directly generating them with our deep-infomax-based methods (s-Infomax and Infomin).

possible interactions and hence, our method admits a much lower running time. Then, we achieve recommendation predictions based on the generated feature interactions. Specifically, we propose a hypergraph neural network based model called <u>H</u>ypergraph <u>I</u>nfomax <u>R</u>ecommender <u>S</u>ystem (HIRS). HIRS learns a hypergraph representation for each data sample, with each feature in the data represented as a node and each beneficial feature interaction represented as a hyperedge. The beneficial feature interaction generation is conducted via hyperedge matrix prediction, and the recommendation prediction is conducted via hypergraph classification.

To achieve an effective beneficial feature interaction generation, we exploit three properties of beneficial feature interactions based on information theory:

- *Sufficiency*: all beneficial feature interactions together contain as much relevant information about the true output as possible.
- *Low redundancy*: each beneficial feature interaction contains as less irrelevant information about the true output as possible.
- *Disentanglement*: any two beneficial feature interactions contain as less duplicated relevant information about the true output as possible.

Based on these properties, we leverage Deep Infomax [15] and propose supervised mutual information maximization (s-Infomax) and mutual information minimization (Infomin) methods, together with an $L_0$ activation regularization, to guide the interaction generation.

To effectively learn from the generated beneficial feature interactions for recommendation, we leverage the arbitrary-order correlation reasoning ability of hypergraph neural networks [9]. We propose an interaction-based hypergraph neural network (IHGNN) to perform recommendation prediction based on the generated interactions via hypergraph classification.

We summarize our contributions as follows:

- We propose a hypergraph neural network based model named HIRS. [1]. HIRS is the first work that directly generates a set of beneficial feature interactions of arbitrary orders via hyperedge prediction and then performs recommendation prediction via hypergraph classification.
- We define three properties of beneficial feature interactions to guide the proposed s-Infomax and Infomin methods for beneficial feature interaction generation. We further propose an interaction-based hypergraph neural network (IHGNN), which effectively learns from the generated interactions for recommendation.

---

[1]The implementation of our model can be found at https://github.com/ruizhang-ai/HIRS_Hypergraph_Infomax_Recommender_System.git.

- Experimental results show that HIRS outperforms state-of-the-art feature interaction based models for recommendation. The further evaluations prove our model's ability in beneficial feature interaction generation.

## 2 RELATED WORK

### 2.1 Feature Interaction based Recommendation

Pairwise feature interactions have been widely used for recommendation, such as factorization machine (FM) based models [14, 28, 32] and graph neural network based models [21, 33]. Higher-order feature interactions provide more information about the true output but have not been largely explored due to its high complexity. A multilayer perceptron (MLP) can implicitly model any order feature interactions for recommendation [6, 12, 16], but implicit modeling is ineffective in learning feature interactions [2]. Lian et al. [23] explicitly model all feature interactions under a limited order (e.g., up to four orders). Besides modeling all feature interactions, there are studies that detect a set of beneficial feature interactions for recommendation [25, 33]. However, they have to examine all the possible feature interactions, which leaves the complexity issue unsolved. Different from existing studies, we directly generate beneficial feature interactions without examining all the possible ones.

### 2.2 Graph Neural Networks

Graph Neural Networks (GNNs) are powerful in correlation reasoning between entities [43, 45]. They have been applied in various research domains such as molecular prediction [11], object relation learning [4], and user-item interaction modeling [36, 39]. Recently, GNNs show the power of feature interaction learning for recommendation [21, 22, 33]. However, existing studies only capture pairwise feature interactions in a graph since each edge can only link to two nodes. We are the first to capture arbitrary order feature interactions and propose a Hypergraph Neural Network [9] based model to learn arbitrary order feature interactions uniformly.

### 2.3 Mutual Information Maximization

Mutual information maximization is a critical technique in unsupervised representation learning [24]. Deep Infomax (DIM, Hjelm et al. [15]) leverages a deep neural network-based method [1] to maximize the mutual information between an image's local and global representations. A graph-based DIM is proposed to maximize the mutual information between node representations and graph representations [37]. For the first time, our work applies DIM on hypergraph that maximizes the mutual information between hyperedge representations and graph representations. Contrastive learning is close to DIM but maximizes the mutual information between the two views of the same variable [5]. Tian et al. [35] define properties of good views in contrastive learning. They hypothesize that each view contains all the information of an input variable about the output. We define the properties of beneficial feature interactions hypothesizing that each interaction only contains part of the information about the output. Khosla et al. [17] extend contrastive learning to a supervised setting. Our s-Infomax is similar to [17] but applies the supervised learning between local and global representations.

# 3 PRELIMINARIES

## 3.1 Hypergraph Neural Networks

Different from the edges in a typical graph, each edge in a hypergraph (i.e., hyperedge) can link to an arbitrary number of nodes. A hypergraph $G_H = \langle \mathcal{V}, \mathcal{E} \rangle$ contains a node set $\mathcal{V} = \{i\}_{i=1}^m$, where $m$ is the number of nodes. Each node $i$ is represented as a vector $v_i \in \mathbb{R}^d$ of $d$ dimensions (e.g., node embedding). The hyperedge set $\mathcal{E} = \{e_j\}_{j=1}^k$ contains $k$ edges. Each edge $e_j$ is an $m$-dimensional binary vector ($e_{ji} \in \{0, 1\}$), where $e_{ji} = 1$ if the $j^{th}$ edge links to node $i$, and $e_{ji} = 0$ otherwise. We represent all the node vectors in a graph as a node representation matrix $V \in \mathbb{R}^{m \times d}$ and the edge set as an incidence matrix $E \in \{0, 1\}^{m \times k}$.

In general, Hypergraph Neural Network (HGNN, Feng et al. [9]) first generates a representation $h_j$ for each hyperedge $j$ by aggregating the linked node vectors. Then, each node vector is updated by aggregating the edge representations linked to the node (called node patch representation). The procedure is:

$$h_j = \phi(\{V_i | E_{ij} = 1\}_{i \in \mathcal{V}}), \quad v_i' = \psi(\{h_j | E_{ij} = 1\}_{e_j \in \mathcal{E}}), \quad (1)$$

where $\phi(\cdot)$ and $\psi(\cdot)$ are aggregation functions (e.g., element-wise mean), and $v_i'$ is the patch representation of node $i$. The graph representation $c \in \mathbb{R}^d$ can be generated by further aggregating the node patch representations:

$$c = \eta(\{v_i'\}_{i \in \mathcal{V}}), \quad (2)$$

where $\eta(\cdot)$ is an aggregation function similar to $\phi(\cdot)$ and $\psi(\cdot)$.

## 3.2 Deep Infomax

Given two random variables $A$ and $B$, Mutual Information Neural Estimation (MINE, Belghazi et al. [1]) estimates the mutual information $I(A; B)$ by training a discriminator that distinguishes the samples from their joint distribution $\mathbb{J}$ and from their marginal distribution $\mathbb{M}$. Specifically, MINE uses Donsker-Varadhan representation (DV, Donsker and Varadhan [8]) of KL-divergence as the lower bound to the mutual information:

$$\begin{aligned} I(A; B) = D_{KL}(\mathbb{J} || \mathbb{M}) &\geq \hat{I}_\omega^{DV}(A; B) \\ &= \mathbb{E}_{\mathbb{J}}[T_\omega(a, b)] - \log \mathbb{E}_{\mathbb{M}}[e^{T_\omega(a, b)}], \end{aligned} \quad (3)$$

where $a$ and $b$ are the samples of $A$ and $B$, respectively, and $T_\omega$ is a classifier with parameters $\omega$.

In Deep Infomax [15] that aims to maximize a mutual information, the exact mutual information value is not important. Therefore, $\hat{I}_\omega^{DV}(A; B)$ can be conveniently maximized by optimizing a GAN-like objective function [20, 37]:

$$\mathcal{L} = \mathbb{E}_{\mathbb{J}}[\log T_\omega(a, b)] + \mathbb{E}_{\mathbb{M}}[\log(1 - T_\omega(a, b))]. \quad (4)$$

This objective function can be optimized by training the discriminator $T_\omega$ through a Binary Cross Entropy (BCE) loss [37].

# 4 PROBLEM STATEMENT

Denote by $\mathcal{J}$ the *universe* of features. A *feature-value pair* is a name-value pair, denoted by $(o, w)$, where $o$ is the name of the feature and $w$ is the value of this feature. For example, (*Male*, 1) and (*Female*, 1) mean the categorical features of Male and Female, respectively, where $Male \in \mathcal{J}$ and $Female \in \mathcal{J}$ are considered as two different features, and value 1 means the feature is in the data sample. Let

$\mathcal{D} : \mathcal{X} \times \mathcal{Y}$ be an input-output domain and $D = \{(x_n, y_n)\}_{i=1}^N$ is a set of $N$ input-output training pairs sampled from $\mathcal{D}$.

- $x_n = \{(o, w)\}_{o \in J_n}$ is called a *data sample* consisting of a set of features, where $J_n \subseteq \mathcal{J}$.
- $y_n \in \{0, 1\}$ is the *implicit feedback* (e.g., purchased, clicked) of the user on the item.

**Feature Interaction based Recommendation.** Given $(x_n, y_n)$, a recommendation model that *explicitly* considers feature interactions first selects (either predefined or detected) $k_n$ feature interactions $I_n = \{q_i\}_{i=1}^{k_n}$, where $q_i \subseteq \mathcal{J}_n$ indicates an interaction between the features in $q_i$. Then, a predictive model $F(x_n, I_n) = y_n'$ is designed with a middle state $H = \{h_i\}_{i=1}^{k_n}$, where $h_i \in \mathbb{R}^d$ is the high-level representation of $q_i$, and $y_n'$ is the prediction of $y_n$.

# 5 OUR PROPOSED METHOD

In this section, we first give an overview of our proposed method. Then, we demonstrate HIRS in detail and the empirical risk minimization function of HIRS. Finally, we provide analyses of our proposed model, including i) our proposed hypergraph neural network IHGNN can be easily generated to existing methods for feature interaction modeling; and ii) the interaction generation cost of HIRS.

## 5.1 Model Overview

We start by introducing data representation with hypergraphs. Then we describe the overall structure of HIRS.

*5.1.1 Data Representation with Hypergraphs.* We focus on one input-output pair considering $k$ arbitrary order feature interactions. Existing GNN-based models can only represent pairwise feature interactions [33]. We propose a hypergraph form data representation that is more general in representing arbitrary order feature interactions. A hypergraph is an extension of a standard graph that each edge, i.e., hyperedge, can link to an arbitrary number of nodes. We treat each data sample as a hypergraph, with each node being a feature that occurred in the data sample. Each hyperedge represents an interaction between all the linked features. Formally, for the feature set of each data sample input $x = \{(o_i, w_i)\}_{i=1}^m$, it can be represented as a hypergraph $G_H \langle \mathcal{V}, \mathcal{E} \rangle$. The node set $\mathcal{V} = x$. The hyperedge set $\mathcal{E}$ consists of $k$ feature interactions (to be detected) that are most beneficial to the prediction accuracy (i.e., beneficial feature interactions). $\mathcal{E}$ can be represented as a hyperedge incidence matrix $E \in \{0, 1\}^{m \times k}$.

*5.1.2 The Overall Structure of HIRS.* HIRS contains two components: the first component performs recommendations predictions, and the second component performs s-Infomax and Infomin methods based on the three properties of beneficial feature interactions. The recommendation prediction component of HIRS contains two modules: one is a hyperedge prediction module that performs interaction generation; and the other is the IHGNN module that performs interaction learning. The recommendation prediction component takes a hypergraph without hyperedge as input, i.e., $G_H \langle \mathcal{V}, \emptyset \rangle$. A hyperedge prediction function $f_{gen}(\mathcal{V})$ leverages the node information to generate $k$ beneficial feature interactions as the edge set $\mathcal{E}'$. Then, IHGNN models the hypergraph with

the generated edge set $G_H\langle\mathcal{V}, \mathcal{E}'\rangle$ and outputs the recommendation prediction $y'$. The function combining the two modules is $f_r(\mathcal{V}; \boldsymbol{\rho}) = f_{IHGNN}(G_H\langle\mathcal{V}, f_{gen}(\mathcal{V})\rangle)$, where $\boldsymbol{\rho}$ are all the parameters in $f_{gen}$ and $f_{IHGNN}$. While training, the other component containing s-Infomax and Infomin is conducted on the hyperedge representations ($\boldsymbol{h}$) and the graph representation ($\boldsymbol{c}$) in $f_r$, together with an $L_0$ activation regularization on the predicted hyperedges, to ensure the effectiveness of the interaction generation. Figure 2 demonstrates an overview of HIRS. In the following subsections, we describe these two components in detail.

## 5.2 The Recommendation Prediction Component

The recommendation prediction component (or RPC in short) of HIRS generates beneficial feature interactions via hyperedge prediction and then performs recommendation via IHGNN.

*5.2.1 Hyperedge Prediction.* The hyperedge prediction module $f_{gen}(\mathcal{V})$ generates hyperedges $\mathcal{E}'$ as beneficial feature interactions in an incidence matrix format $E' \in \{0,1\}^{m \times k}$. $E'_{ij}$ equals 1 if the $i^{th}$ feature is in the $j^{th}$ beneficial feature interaction, and 0 otherwise.

Specifically, each node is first mapped into a node vector of $d$ dimensions for hyperedge prediction: $V_i^e = emb^e(o_i) \cdot w_i$, where $emb^e(\cdot)$ maps a feature into a $d$-dimensional embedding[2]. Then, for each node, we concatenate the vector containing the information of the node itself, and the vector summarizing the information of all other nodes in the hypergraph. The concatenated vector will be sent to an MLP to obtain a $k$-dimensional vector $\boldsymbol{u}_i$:

$$\boldsymbol{u}_i = MLP([V_i^e, \sum_{j \neq i} V_j^e]), \tag{5}$$

where $[\cdot]$ is the concatenation operator.

Finally, we map $\boldsymbol{u}_i$ into binary values indicating the predicted hyperedge values $E'_i = \mathcal{B}(\boldsymbol{u}_i)$, where $\mathcal{B}$ is a Bernoulli distribution. The predicted incidence matrix $E' \in \{0,1\}^{m \times k}$ is the concatenation of $E'_i$ of all the $m$ nodes.

However, directly representing $E'$ by binary values results in a non-differentiable problem while training with gradient descent methods [29]. Therefore, we replace the Bernoulli distribution with the hard concrete distribution [27], which is differentiable and can approximate binary values. Specifically, $\boldsymbol{u}_i$ is transferred into hyperedge values of hard concrete distribution:

$$\begin{aligned} z \sim \mathcal{U}(0,1), &\qquad s_{ij} = Sig((\log z - \log(1-z) + \log(u_{ij}))/\tau), \\ \bar{s}_{ij} = s_{ij}(\delta - \gamma) + \gamma, &\qquad E'_{ij} = min(1, max(0, \bar{s}_{ij})), \end{aligned} \tag{6}$$

where $z$ is a uniform distribution, $Sig$ is the Sigmoid function, $\tau$ is a temperature value and $(\gamma, \delta)$ is an interval with $\gamma < 0, \delta > 0$. We set $\gamma = -0.1, \delta = 1.1, \tau = 0.66$ following Shi et al. [29].

*5.2.2 IHGNN.* With the generated edge set $\mathcal{E}'$ ($E'$ in incidence matrix form), we propose an interaction-based hypergraph neural network (IHGNN) to learn $G_H\langle\mathcal{V}, \mathcal{E}'\rangle$.

First, we map each node in $\mathcal{V}$ into a node vector for IHGNN: $V_i^g = emb^g(o_i) \cdot w_i$. Then, different from HGNN [9] that models

---

[2] $emb^e(\cdot)$ ensures a feature has the same embedding in different hypergraphs, and the embedding will be updated while training.

each hyperedge by linearly aggregating the node vectors, we use an MLP (a non-linear method) to model each hyperedge:

$$\boldsymbol{h}_j = f_E(sum(\{V_i^g | E'_{ij} = 1\}_{i \in \mathcal{V}})) \tag{7}$$

where $f_E$ is an MLP, $sum(\cdot)$ is an element-wise sum and $\boldsymbol{h}_j \in \mathbb{R}^d$ is a high-level representation of edge $j$.

The non-linearly modeling ensures IHGNN correctly and effectively models the generated interactions founded by statistical interactions [33], which is theoretically justified in Appendix B. Then, each node representation is updated via aggregating the representations of the corresponding (linked) hyperedges by a function $\psi(\cdot)$ (e.g., element-wise mean). The hypergraph representation $\boldsymbol{c}$ is obtained via further aggregating all the updated node representations by a function $\phi(\cdot)$ that is similar to $\psi(\cdot)$. Finally, we use a readout function $g(\cdot)$ to map $\boldsymbol{c}$ to a scalar value as the hypergraph classification result. The IHGNN function $f_{IHGNN}$ is:

$$f_{IHGNN}(G_H\langle\mathcal{V}, \mathcal{E}'\rangle) = g(\phi(\{\psi(\{\boldsymbol{h}_j | E'_{ij} = 1\}_{j \in \mathcal{E}})\}_{i \in \mathcal{V}})). \tag{8}$$

Inspired by Su et al. [33], we leverage an $L_0$ activation regularization to help detect beneficial feature interactions via sparsifying the generated hyperedge matrix $E'$. However, the $L_0$ activation regularization solely may not lead to a successful interaction detection of arbitrary orders. For example, it may lead to the generated incidence matrix containing only one edge that links to all nodes and the other $k - 1$ edges being empty. As a result, all interaction information assembles in one hyperedge, which is not a beneficial feature interaction. Therefore, we exploit three properties of beneficial feature interactions in the next section to guide the beneficial feature interaction generation.

## 5.3 The s-Infomax and Infomin Component

We first define three properties of beneficial feature interactions (i.e., beneficial properties). Then we describe how we achieve the beneficial properties by s-Infomax and Infomin methods.

*5.3.1 Beneficial Properties.* We formally define the three beneficial properties that beneficial feature interactions in a data sample should try to achieve:

*Sufficiency*: The representations of all beneficial feature interactions together contain as much information of the input about the true output as possible. In information theory, ideally, we aim to reach:

$$I((\boldsymbol{h}_1, \boldsymbol{h}_2, ..., \boldsymbol{h}_k); y) = I(\boldsymbol{x}; y), \tag{9}$$

where $I(\cdot)$ is the mutual information.

*Low Redundancy*: Each beneficial feature interaction representation should contain as less irrelevant information about the true output as possible. Ideally, we aim to reach:

$$\sum_{i=1}^{k} H(\boldsymbol{h}_i | y) = 0, \tag{10}$$

where $H(\boldsymbol{h}_i | y)$ is the entropy of $\boldsymbol{h}_i$ conditioned on $y$.

*Disentanglement*: The representations of any two beneficial feature interactions contain as less duplicated information about the true output as possible. Ideally, we aim to reach:

$$\sum_{i \neq j} I(\boldsymbol{h}_i; \boldsymbol{h}_j; y) = 0. \tag{11}$$

Figure 3 shows the relations between feature interactions and the true output according to the three properties.
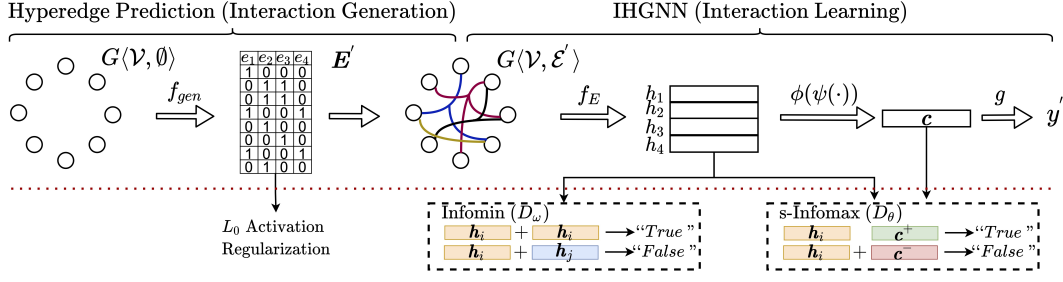
Figure 2: An Overview of HIRS. The part *above* the red dotted line is the recommendation prediction component (RPC). The part *below* the red dotted line is the s-Infomax and Infomin component and the $L_0$ activation regularization, which is only conducted during training.
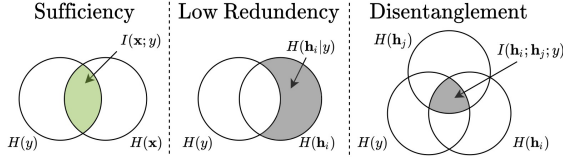


Figure 3: The Venn diagram of the three properties. The *green* part is the information that our generated interactions to retain, and the *gray* part is to discard.

We aim to generate hyperedges that can simultaneously achieve the three properties. In HIRS, all the hyperedge representations $\{h_1, h_2, ..., h_k\}$ are generated by a function with features $x$ as the input, and the output $y$ is predicted through the hyperedge representations, which forms a Markov chain $x \rightarrow \{h_1, h_2, ..., h_k\} \rightarrow y$. Therefore, $I(x; y) \geqslant I((h_1, h_2, ..., h_k); y)$. To achieve the sufficiency property, we just need to maximize $I((h_1, h_2, ..., h_k); y)$ to approach the optimal situation.

Due to the non-negativity of entropy and mutual information, $\sum_{i=1}^{k} H(h_i|y) \geqslant 0$ and $\sum_{i \neq j} I(h_i; h_j; y) \geqslant 0$. To achieve the low redundancy and the disentanglement properties, therefore, we minimize $\sum_{i=1}^{k} H(h_i|y)$ and $\sum_{i \neq j} I(h_i; h_j; y)$.

Combining all the goals together, the objective function is:

$$\max I((h_1, h_2, ..., h_k); y) - \alpha \sum_{i=1}^{k} H(h_i|y) - \beta \sum_{i \neq j} I(h_i; h_j; y), \quad (12)$$

where $\alpha$ and $\beta$ are scalar weights balancing the three properties.

*5.3.2 s-Infomax and Infomin.* Due to the difficulty of directly optimizing Equation (12) [15], we propose supervised mutual information maximization (s-Infomax) and mutual information minimization (Infomin) methods based on Deep Infomax (DIM, Hjelm et al. [15]) as the approximation of Equation (12). The s-Infomax method maximizes the mutual information between each edge representation and the label representation, i.e., $\max I(h_i; y)$. The Infomin method minimizes the mutual information between every pair of edge representations in a hypergraph, i.e., $\min I(h_i; h_j)$. The approximation of Equation (12) is:

$$\max \sum_{i=1}^{k} I(h_i; y) - \sigma \sum_{i \neq j} I(h_i; h_j), \quad (13)$$

where $\sigma$ is a scalar weight.

We prove that Equation (12) reaches the optimal value when Equation (13) reaches the optimal value in Proposition 1. We detail the proof in Appendix A.

**Proposition 1.** *Consider an input-output pair $(x, y)$, where $x$ is a set of features, a function maps $x$ to $y$ in a Markov chain $x \rightarrow \{h_1, h_2, ..., h_k\} \rightarrow y$, where $\{h_1, h_2, ..., h_k\}$ is a set of middle states mapped from $x$, if the $\{h_1, h_2, ..., h_k\}$ makes Equation (13) reach the optimal value, it also makes Equation (12) reach the optimal value.*

Proposition 1 ensures that by optimizing Equation (13), we approximately optimize Equation (12) so that generate feature interactions that fit the three properties.

*5.3.3 The Implementation.* One difficulty of applying s-Infomax is that the label space only has two states, i.e., 0 and 1, which are discrete and are too small to provide enough output information [17]. Hence, we use the hypergraph representations to represent the label space, extending the two states into a continuous high-dimensional space. Following DIM, we sample the *joint distribution* of $I(h_i; y)$ as the pair of $h_i$ and a random hypergraph representation $c^+$ from training samples having the same label $y$. Then, we sample the *marginal distribution* as $h_i$ and a random hypergraph representation $c^-$ from training samples having the opposite label $y$. We use a GAN-based objective function for s-Infomax:

$$\max \sum_i I(h_i; y) = \max \frac{1}{k} \sum_i (\log D_\theta(h_i, c^+) + (1 - \log D_\theta(h_i, c^-))), \quad (14)$$

where $D_\theta \in \mathbb{R}^{2 \times d} \rightarrow \mathbb{R}$ is a discriminator that distinguishes the joint distribution and the marginal distribution.

In Infomin, we minimize the mutual information between different hyperedge representations in a hypergraph. We set the sampled joint distribution as $(h_i, h_j)$ that $i \neq j$, and the marginal distribution as $(h_i, h_i)$. To avoid the overfitting issue that the discriminator in Infomin may trivially compare whether two input vectors are the same, inspired by Gao et al. [10], we use a dropout function $f_a$ to prevent each pair of input vectors from being the same.

To ensure efficiency, we get one sample from the joint distribution and one sample from the marginal distribution for each node. The objective function for Infomin is also a GAN-based objective function but swaps the joint and marginal distributions:

$$\min \sum_{i \neq j} I(h_i; h_j) = \max \frac{1}{k} \sum_i (\log D_\omega(f_a(h_i) f_a(h_i)) + (1 - \log D_\omega(f_a(h_i), f_a(h_j)))). \quad (15)$$

Equation (14) and Equation (15) can be efficiently optimized by training the discriminators $D_\theta$ and $D_\omega$ through Binary Cross-entropy (BCE) loss [37].

Note that the main difference between s-Infomax and Infomin in implementation is the choice of the joint distribution and the marginal distribution that a discriminator to distinguish. The distributions chosen in Infomin are intuitive: we want any two hyperedge representations in a hypergraph to be as different as possible (i.e., have minimum mutual information).

## 5.4 The Empirical Risk Minimization Function

The empirical risk minimization function of HIRS contains the loss function of recommendation predictions, an $L_0$ activation regularization, the loss functions of the s-Infomax discriminator and the Infomin discriminator:

$$\mathcal{R}(\theta, \omega, \rho) = \frac{1}{N} \sum_{n=1}^{N} (\mathcal{L}(f_r((G_H)_n; \rho), y_n) + \lambda_1 \left\| E'_n \right\|_0 \\ + \lambda_2 \mathcal{L}_D(D_\theta) + \lambda_3 \mathcal{L}_D(D_\omega)), \quad (16)$$
$$\rho^*, \theta^*, \omega^* = arg\,min_{\theta, \omega, \rho} \mathcal{R}(\theta, \omega, \rho),$$

where $(G_H)_n = (G_H \langle \mathcal{V}, \emptyset \rangle)_n$, $\lambda_1$, $\lambda_2$ and $\lambda_3$ are weight factors, $\mathcal{L}(\cdot)$ and $\mathcal{L}_D(\cdot)$ are both BCE loss, and $\rho^*$, $\theta^*$, $\omega^*$ are final parameters. Following Shi et al. [29], the $L_0$ activation regularization can be calculated as: $\left\| E'_n \right\|_0 = \sum_{i,j} Sig(\log u_{ij} - \tau \log \frac{-\gamma}{\delta})$.

## 5.5 Generalizing IHGNN to Existing Methods

We show that some existing methods can be represented as an IHGNN in specific settings.

*5.5.1 $L_0$-SIGN as IHGNN.* A hypergraph is an extension of a standard graph. Therefore, GNN-based recommender systems such as $L_0$-SIGN [33] and Fi-GNN [21] can be easily applied by IHGNN by changing the edge matrix into the incidence matrix form and replacing the corresponding feature interaction modeling functions, e.g., dot product by an MLP for $L_0$-SIGN.

*5.5.2 FM as IHGNN.* FM [28] models every pairwise feature interaction by dot product and sums up all the modeling results. We can regard each pairwise feature interaction in FM as a hyperedge linking to two nodes, which forms the $\mathcal{E}^{FM}$ that contains all pairwise node combinations, including all self-loops as the point-wise terms in FM. In addition, FM-based extensions can also be achieved based on the above derivation. For example, NFM [14] can be achieved by replacing the dot product and the $\phi$ by an element-wise production and an MLP, respectively.

*5.5.3 DeepFM as IHGNN.* DeepFM [12] merges the prediction of FM and the prediction of an MLP modeling all the features. We can regard DeepFM as IHGNN by adding a hyperedge linking to all nodes in addition to $\mathcal{E}^{FM}$. Since DeepFM uses different interaction modeling functions in FM and the MLP, the hyperedges are modeled differently according to the edge degree.

## 5.6 Interaction Generation Cost Analysis

HIRS generates $k$ beneficial feature interactions through the hyperedge prediction module $f_{gen}$, which has the time complexity of $O(mkd)$ for a data sample of $m$ features, where $d$ is the embedding dimension. The s-Infomax and the Infomin methods help

| Dataset | #Data | #Users | #Items | #Features |
|---|---|---|---|---|
| MovieLens 1M | 1,149,238 | 5,950 | 3,514 | 6,974 |
| Book-Crossing | 1,050,834 | 4,873 | 53,168 | 43,244 |
| MovieLens 25M | 31,147,118 | 162,541 | 62,423 | 226,122 |

**Table 1: Dataset Statistics.**

the interaction generation and both have the complexity of $O(kd)$. In addition, the IHGNN module models the generated feature interactions for recommendation with a $O(mkd + md)$ complexity. Therefore, the time complexity of HIRS is $O(mkd)$. Existing studies model or examine all the possible interactions by enumerating them $O(2^m d)$ [25] or stacking pairwise interaction modeling layers $O(m^3 d)$ [30, 38] when considering arbitrary orders. Since $k$ can be set to be a much smaller number (e.g., 40) than $m^2$ (e.g., 196 when $m = 14$), HIRS achieves a more efficient interaction modeling. We do not include the complexity of the neural network modules in HIRS and the baselines as they are all linear to $d$. The running time of HIRS and the baselines are detailed in Appendix D.

## 6 EXPERIMENTAL ANALYSIS

In this section, we evaluate the following aspects of our proposed model: (i) the recommendation prediction accuracy of HIRS; (ii) the beneficial feature interactions generation ability of HIRS; (iii) the impact of different components and hyperparameters of HIRS.

## 6.1 Experimental Setting

*6.1.1 Datasets and Baselines.* We evaluate HIRS on three real-world datasets in recommender systems: MovieLens 1M [13], Book-crossing [46], and MovieLens 25M [13]. We summarize the statistics of the datasets in Table 1. Note that although MovieLens 1M and MovieLens 25M are both movie recommendation datasets, they are actually deployed with different features, e.g., MovieLens 25M contains the tag features, while MovieLens 1M does not. In this way, we can evaluate the effectiveness and the robustness of our model when different kinds of features are available in similar scenarios. We compare HIRS with state-of-the-art models considering feature interactions for recommendation. They are FM [19], AFM [42], NFM [14], Fi-GNN [21], $L_0$-SIGN [33], AutoInt [30], DeepFM [12], xDeepFM [23], DCNv2 [38], AutoFIS [25], AFN [7]. More details about the datasets and the baselines are in Appendix C.

*6.1.2 Experimental Setup.* In the experiments, we use element-wise mean as the linear aggregation function for both $\psi(\cdot)$ and $\phi(\cdot)$. $g(\cdot)$ is a linear regression function. The MLPs for $f_{gen}(\cdot)$ and $f_E$ consist of a fully connected layer of size 64 with a ReLU activation. Both $D_\theta$ and $D_\omega$ are Bilinear models. The dropout probability of $f_a$ is 0.1. The node embedding size for interaction modeling and edge prediction is 64. The hyperedge number $k$ is 40. We set $\lambda_1 = 0.02$, $\lambda_2 = 1$, $\lambda_3 = 0.1$, and Adam [18] as the optimization algorithm. We randomly split each dataset into training, validation, and test sets with a proportion of 70%, 15%, and 15%, respectively. We evaluate the performance of HIRS and the baseline models using the popular ranking metrics Recall@10, Recall@20, NDCG@10, and NDCG@20.

## 6.2 Model Performance

Table 2 shows the performance of our model and the baseline models. The results are the average of 10 runs. The best results are in

| | HO | BD | MovieLens 1M | | | | Book-Crossing | | | | MovieLens 25M | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | R@10 | R@20 | N@10 | N@20 | R@10 | R@20 | N@10 | N@20 | R@10 | R@20 | N@10 | N@20 |
| FM | | | 0.8527 | 0.8996 | 0.8521 | 0.8845 | 0.7802 | 0.8834 | 0.8123 | 0.8539 | 0.8331 | 0.8972 | 0.9535 | 0.8791 |
| AFM | | | 0.8738 | 0.9176 | 0.8734 | 0.8994 | 0.8084 | 0.8892 | 0.8473 | 0.8779 | 0.8484 | 0.9082 | 0.8618 | 0.8845 |
| NFM | | | 0.8979 | 0.9318 | 0.8925 | 0.9084 | 0.8625 | 0.9174 | 0.8492 | 0.8843 | 0.8604 | 0.9170 | 0.8620 | 0.8903 |
| Fi-GNN | | | 0.9060 | 0.9341 | 0.9091 | 0.9226 | 0.8818 | 0.9262 | 0.8689 | 0.8861 | 0.8692 | 0.9124 | 0.8667 | 0.8905 |
| $L_0$-SIGN | | √ | <u>0.9092</u> | <u>0.9386</u> | <u>0.9176</u> | <u>0.9268</u> | 0.8817 | 0.9266 | <u>0.8723</u> | 0.8914 | 0.8729 | 0.9261 | 0.8704 | 0.8915 |
| AutoInt | √ | | 0.8991 | 0.9318 | 0.9077 | 0.9118 | 0.8739 | 0.9263 | 0.8636 | 0.8894 | 0.8707 | 0.9233 | 0.8692 | 0.8899 |
| DeepFM | √ | | 0.8975 | 0.9308 | 0.9074 | 0.9128 | 0.8765 | 0.9245 | 0.8648 | 0.8800 | 0.8667 | 0.9213 | 0.8650 | 0.8883 |
| xDeepFM | √ | | 0.9049 | 0.9342 | 0.9134 | 0.9114 | 0.8791 | 0.9249 | 0.8692 | 0.8826 | 0.8673 | 0.9226 | 0.8696 | 0.8902 |
| DCNv2 | √ | | 0.9014 | 0.9311 | 0.9079 | 0.9104 | 0.8773 | 0.9234 | 0.8671 | 0.8858 | 0.8678 | 0.9214 | 0.8667 | 0.8890 |
| AutoFIS | √ | √ | 0.9028 | 0.9317 | 0.9096 | 0.9208 | 0.8804 | 0.9227 | 0.8706 | 0.8887 | 0.8702 | 0.9228 | 0.8712 | 0.8922 |
| AFN | √ | √ | 0.9032 | 0.9333 | 0.9112 | 0.9232 | <u>0.8821</u> | <u>0.9285</u> | 0.8710 | <u>0.8917</u> | <u>0.8743</u> | <u>0.9271</u> | <u>0.8737</u> | <u>0.8938</u> |
| HIRS | √ | √ | **0.9545** | **0.9663** | **0.9464** | **0.9500** | **0.9279** | **0.9584** | **0.9081** | **0.9211** | **0.9223** | **0.9545** | **0.8846** | **0.9011** |
| *Improv.* | | | 4.98% | 2.95% | 3.13% | 2.50% | 5.19% | 3.22% | 4.10% | 3.30% | 5.49% | 2.96% | 1.26% | 0.81% |
| *p-value* | | | 0.5% | 0.5% | 0.5% | 0.5% | 0.5% | 0.5% | 0.5% | 0.5% | 0.5% | 0.5% | 1.4% | 3.7% |

**Table 2: Comparing the prediction performance of HIRS with the baselines. HO and BD indicate whether the model can consider high-order feature interactions and perform beneficial feature interaction detection, respectively. R@j refers to Recall, and N@j refers to NDCG.**
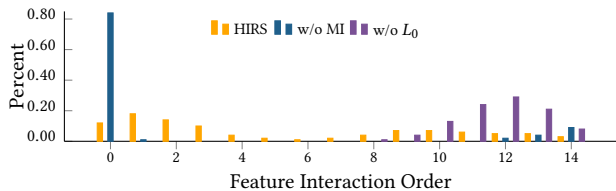


**Figure 4: Comparing the distribution of the orders of the generated feature interactions from HIRS, HIRS without s-Infomax and Infomin, and HIRS without $L_0$ regularization for MovieLens 1M.**
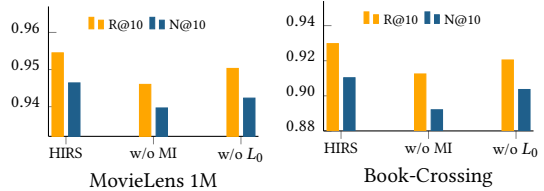


**Figure 5: Comparing the performance of HIRS, HIRS without s-Infomax and Infomin, and HIRS without $L_0$ regularization.**

bold, and the best baseline results are underlined. The rows *Improv* (standing for Improvements) and *p-value* show the improvement and statistical significance test results (through Wilcoxon signed-rank test [41]) of HIRS to the best baseline results, respectively.

We observe that: (i) Our model outperforms all the baselines on the three datasets (up to 5.49%), and the significance test results are all lower than the threshold of 5%, which show that the improvements are significant. (ii) For the models that only consider pairwise feature interactions, GNN-based models (e.g., Fi-GNN and $L_0$-SIGN) perform better than other baselines (e.g., NFM). This shows the ability of GNNs in interaction modeling. (iii) Models considering higher-order feature interactions (e.g., xDeepFM and AFN) perform better than those only consider pairwise interactions (e.g., NFM and Fi-GNN). Thus, considering high-order feature interactions is beneficial for providing better prediction accuracy. (iv) Detecting beneficial feature interactions is vital in removing noise interactions and improving prediction accuracy (e.g., $L_0$-SIGN v.s. Fi-GNN). (v) Our model combines GNN-based interaction modeling, higher-order interaction modeling, and beneficial feature interaction detection, and achieves the best performance.

## 6.3 Effectiveness of the Beneficial Properties

To ensure our model can detect beneficial feature interactions that fit the beneficial properties (sufficiency, low-redundancy, and disentanglement), we propose the s-Infomax and Infomin module and

the $L_0$ activation regularization. This section, we evaluate the influence of them on the beneficial feature interaction detection and the prediction accuracy.

We run two variants of HIRS: (i) without s-Infomax and Infomin (*w/o MI*), and (ii) without the $L_0$ activation regularization (*w/o $L_0$*). Figure 4 shows the distribution of the generated feature interactions on three settings (the full model and the two variants). We observe that when using the variant *w/o MI*, most of the generated hyperedges are empty (order 0), and a few hyperedges connect to most features (order 12-14). It proves the importance of the three properties on arbitrary order beneficial interaction generation. When using the variant *w/o $L_0$*, most of the generated interactions have relative higher orders than using our full model. This proves that $L_0$ regularization helps ease the issue of mixed interaction in each hyperedge. Next, we show the prediction accuracies of the two variants and the full model in Figure 5. We show the results of Recall@10 and NDCG@10. Similar conclusions can be drawn from Recall@20 and NDCG@20 (the same as the remaining experiments). We can see that both *w/o MI* and *w/o $L_0$* reduce the prediction accuracy. This proves that our full model successfully generates more beneficial feature interactions than the two variants and hence gains the best performance.

## 6.4 Ablation Study on RPC Modules

We evaluate the impact of two key functions, the hyperedge detection function $f_{gen}$ and the non-linear interaction modeling function
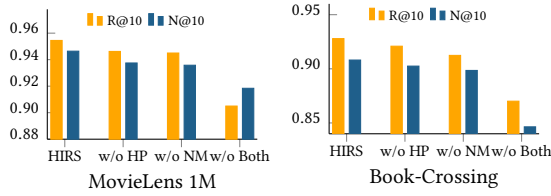
Figure 6: Comparing the recommendation performance of HIRS, HIRS without hyperedge prediction, HIRS without non-linear interaction modeling, and HIRS without both.
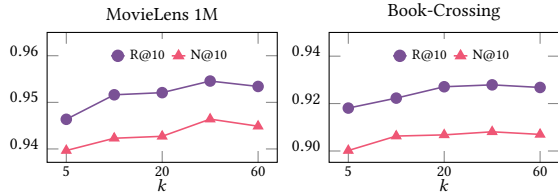


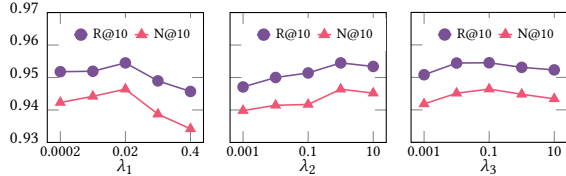Figure 7: Comparing the recommendation performance of HIRS with different hyperedge numbers.



Figure 8: Comparing the performance of different weight values of the $L_0$ activation regularization, s-Infomax, and Infomin for Movie-Lens 1M.

$f_E$, in the two modules (the hyperedge prediction module and the IHGNN module) of the recommendation prediction component (RPC), respectively. To evaluate the impact of $f_{gen}$, we remove $f_{gen}$ and use a fully connected matrix as the hyperedge matrix, i.e., $E' = 1$. We name this variant as *w/o HP*. To evaluate $f_E$ in HIRS, we replace $f_E$ by a linear function, which simply remove the ReLU activation in $f_E$. We name this variant as *w/o NM*.

We summarize the results in Figure 6. The variant *w/o HP* achieves an inferior performance. In this situation, the hyperedge does not indicate a feature interaction, but all the features. Therefore, the interaction modeling function does not model beneficial feature interactions but the mix of all the feature information, which inevitably contains the noise information that decreases the prediction accuracy. The variant *w/o NM* also achieves an inferior performance. This is because the non-linear interaction modeling helps correctly model the generated feature interactions and hence improves the accuracy. This is consistent with the claim in Su et al. [33]. We also show the results of the variant with both modifications (*w/o Both*). The worst performance it achieves is not surprising since the whole model under this setting can be regarded as a linear mapping from feature embeddings to the final prediction.

## 6.5 Parameter Study

We evaluate the impact of the number of feature interactions ($k$), the weight values for $L_0$ activation regularization ($\lambda_1$), s-Infomax ($\lambda_2$), and Infomin ($\lambda_3$) in Equation (16) on the prediction accuracy.

Figure 7 shows the performance of HIRS using different numbers of edges ($k$). The prediction accuracy increases when the $k$ increases from 5 to 40, but flattens over 40. This is because that when the edge number is small, some of the beneficial feature interactions are discarded or mixed in one hyperedge. However, the number of edge slots is enough to contain beneficial feature interactions when $k > 40$. Next, we evaluate different $\lambda_1, \lambda_2, \lambda_3$ values in Equation (16) on the MovieLens 1m dataset. As shown in Figure 8, every weight has a sweat point that leads to the best performance. For the $L_0$ regularization weight ($\lambda_1$), it balances the sparsity and the useful information retraining of the beneficial feature interactions. s-Infomax and Infomin balance the three properties. The sufficiency and the low redundancy are achieved by s-Infomax, and the disentanglement is achieved by Infomin. For example, a high s-Infomax weight ensures sufficiency and low redundancy, but may harm the disentanglement. Only when the three properties achieve simultaneously we can successfully generate beneficial feature interactions, and leverage them for accurate recommendations.

## 6.6 Case Study on the Detected Interactions

We do a case study to visualize the detected feature interactions. Figure 9 shows the detected feature interactions (in an incidence matrix) of a data sample in the MovieLens 1M dataset, where empty ones were removed. Each column is a feature interaction with corresponding features colored. We sort the feature interactions by order and set them with different colors based on their order, e.g., dark blue for second order feature interactions and purple for 13 order feature interactions (the number of colored cells in each column indicates the order of the corresponding interaction). From the figure, we observe that: (i) Feature interactions provide potential explanations for the recommendation prediction results. For example, interaction 13 shows the interaction between age and comedy genre, which indicates that users aged 25-35 may prefer comedy movies. (ii) Features such as gender and age participate in many feature interactions, which shows the importance of these features in the prediction. (iii) User id (*e.g., User_212*) and movie id (e.g., *Naked Gun*) usually appears in high order feature interactions (e.g., interaction 30-34). This indicates that the id features can be used to infer complex preferences or properties. (iv) There are a few duplicated feature interactions, e.g., interactions 13 and 17. It means that these interactions are important. Having them duplicated adds the weights of these interactions in recommendation predictions.

Due to the space limitation, we present some experimental results for the Book-Crossing dataset and the MovieLens 25M dataset in Appendix E, which are consistent with demonstrated results in the above sections.

## 7 CONCLUSION

We proposed a method for detecting arbitrary order beneficial feature interactions for recommendation. The method represents each data sample as a hypergraph and uses a hypergraph neural network-based model named HIRS for recommendation. HIRS achieves superior results by directly generating beneficial feature interactions via hyperedge prediction and performing recommendation predictions via hypergraph classification. We define three properties of beneficial feature interactions, and propose s-Infomax and Infomin
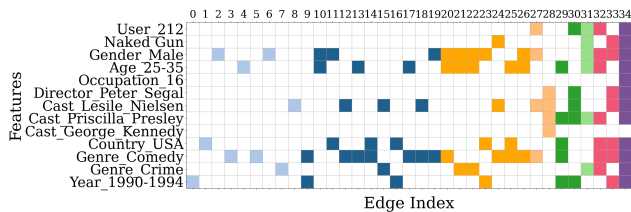
**Figure 9: The detected feature interactions (in an incidence matrix) of a data sample in the MovieLens 1M dataset. Feature interactions are sorted by their orders and columns with the same color represent the interactions of the same order.**

methods for effective interaction generation. In the future work, we will explore the possibility to automatically determine the best edge number, and to extend s-Infomax and Infomin methods to other domains such as anomaly detection.

## ACKNOWLEDGMENTS

## REFERENCES

[1] Mohamed Ishmael Belghazi, Aristide Baratin, Sai Rajeshwar, Sherjil Ozair, Yoshua Bengio, Aaron Courville, and Devon Hjelm. 2018. Mutual Information Neural Estimation. In *ICML*. 531–540.
[2] Alex Beutel, Paul Covington, Sagar Jain, Can Xu, Jia Li, Vince Gatto, and Ed H Chi. 2018. Latent Cross: Making Use of Context in Recurrent Recommender Systems. In *WSDM*. 46–54.
[3] Mathieu Blondel, Masakazu Ishihata, Akinori Fujino, and Naonori Ueda. 2016. Polynomial Networks and Factorization Machines: New Insights and Efficient Training Algorithms. In *ICML*. 850–858.
[4] Michael B Chang, Tomer Ullman, Antonio Torralba, and Joshua B Tenenbaum. 2016. A compositional object-based approach to learning physical dynamics. In *ICLR*. 1–14.
[5] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. 2020. A Simple Framework for Contrastive Learning of Visual Representations. In *ICML*. 1597–1607.
[6] 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 *RecSys*. 7–10.
[7] Weiyu Cheng, Yanyan Shen, and Linpeng Huang. 2020. Adaptive Factorization Network: Learning Adaptive-order Feature Interactions. In *AAAI*. 3609–3616.
[8] Monroe D Donsker and SR Srinivasa Varadhan. 1983. Asymptotic Evaluation of Certain Markov Process Expectations for Large Time. IV. *Communications on Pure and Applied Mathematics* 36, 2 (1983), 183–212.
[9] Yifan Feng, Haoxuan You, Zizhao Zhang, Rongrong Ji, and Yue Gao. 2019. Hypergraph Neural Networks. In *AAAI*. 3558–3565.
[10] Tianyu Gao, Xingcheng Yao, and Danqi Chen. 2021. SimCSE: Simple Contrastive Learning of Sentence Embeddings. In *EMNLP*. 6894–6910.
[11] Justin Gilmer, Samuel S Schoenholz, Patrick F Riley, Oriol Vinyals, and George E Dahl. 2017. Neural Message Passing for Quantum Chemistry. In *ICML*. 1263–1272.
[12] Huifeng Guo, Ruiming Tang, Yunming Ye, Zhenguo Li, and Xiuqiang He. 2017. DeepFM: A Factorization-machine based Neural Network for CTR Prediction. In *IJCAI*. 1725–1731.
[13] F Maxwell Harper and Joseph A Konstan. 2015. The Movielens Datasets: History and Context. *TIIS* (2015), 1–19.
[14] Xiangnan He and Tat-Seng Chua. 2017. Neural Factorization Machines for Sparse Predictive Analytics. In *SIGIR*. 355–364.
[15] R Devon Hjelm, Alex Fedorov, Samuel Lavoie-Marchildon, Karan Grewal, Phil Bachman, Adam Trischler, and Yoshua Bengio. 2019. Learning Deep Representations by Mutual Information Estimation and Maximization. In *ICLR*. 1–14.
[16] Kurt Hornik, Maxwell Stinchcombe, and Halbert White. 1989. Multilayer Feedforward Networks are Universal Approximators. *Neural networks* 2, 5 (1989), 359–366.
[17] Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola, Aaron Maschinot, Ce Liu, and Dilip Krishnan. 2020. Supervised Contrastive Learning. In *NeurIPS*. 18661–18673.
[18] Diederik P Kingma and Jimmy Ba. 2015. Adam: A Method for Stochastic Optimization. In *ICLR*. 1–15.
[19] Yehuda Koren. 2008. Factorization Meets the Neighborhood: A Multifaceted Collaborative Filtering Model. In *SIGKDD*. 426–434.
[20] Maosen Li, Siheng Chen, Ya Zhang, and Ivor W Tsang. 2020. Graph Cross Networks with Vertex Infomax Pooling. In *NeurIPS*. 14093–14105.
[21] Zekun Li, Zeyu Cui, Shu Wu, Xiaoyu Zhang, and Liang Wang. 2019. Fi-GNN: Modeling Feature Interactions via Graph Neural Networks for CTR Prediction. In *CIKM*. 539–548.
[22] Zekun Li, Shu Wu, Zeyu Cui, and Xiaoyu Zhang. 2021. GraphFM: Graph Factorization Machines for Feature Interaction Modeling. (2021).
[23] 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 *SIGKDD*. 1754–1763.
[24] Ralph Linsker. 1988. Self-organization in a Perceptual Network. *Computer* 21, 3 (1988), 105–117.
[25] Bin Liu, Chenxu Zhu, Guilin Li, Weinan Zhang, Jincai Lai, Ruiming Tang, Xiuqiang He, Zhenguo Li, and Yong Yu. 2020. AutoFIS: Automatic Feature Interaction Selection in Factorization Models for Click-Through Rate Prediction. In *SIGKDD*. 2636–2645.
[26] Christos Louizos, Max Welling, and Diederik P Kingma. 2018. Learning Sparse Neural Networks through L0 Regularization. In *ICLR*. 1–11.
[27] Chris J Maddison, Andriy Mnih, and Yee Whye Teh. 2017. The Concrete Distribution: A Continuous Relaxation of Discrete Random Variables. In *ICLR*. 1–12.
[28] Steffen Rendle. 2010. Factorization Machines. In *ICDM*. 995–1000.
[29] Chengzhi Shi, Ben Glocker, and Daniel C Castro. 2019. PVAE: Learning Disentangled Representations with Intrinsic Dimension via Approximated L0 Regularization. In *MLR*. 1–6.
[30] Weiping Song, Chence Shi, Zhiping Xiao, Zhijian Duan, Yewen Xu, Ming Zhang, and Jian Tang. 2019. AutoInt: Automatic Feature Interaction Learning via Self-attentive Neural Networks. In *CIKM*. 1161–1170.
[31] Daria Sorokina, Rich Caruana, Mirek Riedewald, and Daniel Fink. 2008. Detecting Statistical Interactions with Additive Groves of Trees. In *ICML*. 1000–1007.
[32] Yixin Su, Sarah Monazam Erfani, and Rui Zhang. 2019. MMF: Attribute Interpretable Collaborative Filtering. In *IJCNN*. 1–8.
[33] Yixin Su, Rui Zhang, Sarah Erfani, and Zhenghua Xu. 2021. Detecting Beneficial Feature Interactions for Recommender Systems. In *AAAI*. 4357–4365.
[34] Yixin Su, Rui Zhang, Sarah M. Erfani, and Junhao Gan. 2021. Neural Graph Matching based Collaborative Filtering. In *SIGIR*. 849–858.
[35] Yonglong Tian, Chen Sun, Ben Poole, Dilip Krishnan, Cordelia Schmid, and Phillip Isola. 2020. What Makes for Good Views for Contrastive Learning?. In *NeurIPS*. 4116–4126.
[36] Rianne van den Berg, Thomas N. Kipf, and Max Welling. 2017. Graph Convolutional Matrix Completion. arXiv:1706.02263 [stat.ML]
[37] Petar Veličković, William Fedus, William L Hamilton, Pietro Liò, Yoshua Bengio, and R Devon Hjelm. 2019. Deep Graph Infomax. In *ICLR*. 1–17.
[38] Ruoxi Wang, Rakesh Shivanna, Derek Cheng, Sagar Jain, Dong Lin, Lichan Hong, and Ed Chi. 2021. DCN V2: Improved Deep & Cross Network and Practical Lessons for Web-scale Learning to Rank Systems. In *WWW*. 1785–1797.
[39] Xiang Wang, Xiangnan He, Meng Wang, Fuli Feng, and Tat-Seng Chua. 2019. Neural Graph Collaborative Filtering. In *SIGIR*. 165–174.
[40] Xiaojie Wang, Rui Zhang, Yu Sun, and Jianzhong Qi. 2019. Doubly Robust Joint Learning for Recommendation on Data Missing Not at Random. In *ICML*. 6638–6647.
[41] Frank Wilcoxon. 1992. Individual Comparisons by Ranking Methods. In *Breakthroughs in Statistics*. Springer, 196–202.
[42] Jun Xiao, Hao Ye, Xiangnan He, Hanwang Zhang, Fei Wu, and Tat-Seng Chua. 2017. Attentional Factorization Machines: Learning the Weight of Feature Interactions via Attention Networks. In *IJCAI*. 3119–3125.
[43] Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. 2019. How Powerful are Graph Neural Networks?. In *ICLR*. 1–13.
[44] Zilin Zeng, Hongjun Zhang, Rui Zhang, and Chengxiang Yin. 2015. A Novel Feature Selection Method Considering Feature Interaction. *Pattern Recognition* 48, 8 (2015), 2656–2666.
[45] Yunxiang Zhao, Jianzhong Qi, Qingwei Liu, and Rui Zhang. 2021. WGCN: Graph Convolutional Networks with Weighted Structural Features. In *SIGIR*. 624–633.
[46] Cai-Nicolas Ziegler, Sean M McNee, Joseph A Konstan, and Georg Lausen. 2005. Improving Recommendation Lists Through Topic Diversification. In *WWW*. 22–32.

# A PROOF OF PROPOSITION 1

**Proposition 1.** *Consider an input-output pair $(x, y)$, where $x$ is a set of features, a function maps $x$ to $y$ in a Markov chain $x \rightarrow \{h_1, h_2, ..., h_k\} \rightarrow y$, where $\{h_1, h_2, ..., h_k\}$ is a set of middle states mapped from $x$, if the $\{h_1, h_2, ..., h_k\}$ makes Equation (13) reach the optimal value, it also makes Equation (12) reach the optimal value.*

PROOF. Due to the non-negativity of mutual information, we have $I(h_i; h_j) \geqslant 0$. $\sum_{i \neq j} I(h_i; h_j)$ reaches the smallest value when $I(h_i; h_j) = 0$ for all $i \neq j$. That is, each pair of hyperedge representations are independent. Therefore, we have

$$I((h_1, h_2, ..., h_k); y) = \sum_{i=1}^{k} I(h_i; y). \qquad (17)$$

It means that when $\sum_{i \neq j} I(h_i; h_j)$ reaches the smallest value, maximizing $\sum_{i=1}^{k} I(h_i; y)$ equals to maximizing $I((h_1, h_2, ..., h_k); y)$. Next, since

$$0 \leqslant I(h_i; h_j; y) \leqslant I(h_i; h_j), \qquad (18)$$

$\sum_{i \neq j} I(h_i; h_j; y) = 0$ when $\sum_{i \neq j} I(h_i; h_j) = 0$.

Finally, since $I(h_i; y) = H(h_i) - H(h_i|y)$, maximizing $I(h_i; y)$ tries to reduce $H(h_i|y)$ and in the optimal situation, $H(h_i|y) = 0$.

In summary, when the formulation in Equation (13) reaches the optimal situation, we have maximized $I((h_1, h_2, ..., h_k); y)$, and have $I(h_i; h_j; y) = 0$ and $H(h_i|y) = 0$, which make Equation (12) reach the optimal situation. □

# B IHGNN SATISFYING THE STATISTICAL INTERACTION PRINCIPLE

We first give the definition of high-order statistical interaction following Sorokina et al. [31].

**Definition 1. (Statistical Interaction)** Consider a function $f(\cdot)$ with input variables $x = \{x_1, x_2, ..., x_m\}$. Then, a set of variables $I \subseteq x$ is a statistical interaction of function $f(\cdot)$ if and only if there does **not** exist a set of functions $f_{\backslash i}(\cdot)$, where $x_i \in I$ and $f_{\backslash i}(\cdot)$ is not a function of $x_i$, such that

$$f(x) = \sum_{x_i \in I} f_{\backslash i}(x_1, x_{i-1}, x_{i+1}, ..., x_m) \qquad (19)$$

Then, we prove that IGHNN models a set of features as a statistical interaction if they form a hyperedge in the input hypergraph in Proposition 2.

**Proposition 2.** *Consider an input-output pair $(x, y)$. Let $G\langle V, \emptyset \rangle$ be the input of HIRS, where $V = x$, the IHGNN function (Equation (8)) flags a statistical interaction among all the features in $q \subseteq V$ if $q = \{a_i | E'_{ik} = 1\}_{i \in V}$, where $E'_k$ is a hyperedge predicted by $f_{gen}$ (Equation (5)).*

PROOF. We prove by contradiction: assume there is a $q = \{c_i | E'_{ik} = 1\}_{i \in V}$ but the IHGNN function does not flag a statistical interaction on all the features in $q$.

In HIRS, the edge modeling function $f_E$ is a neural network, which is a non-linear function. Without losing generality, we set both the aggregation function $\phi$ and $\psi$ as the element-wise mean function and $g$ is a linear regression function. Therefore, we can regard $f_{IGHNN}$ as:

$$f_{IGHNN}(G_H \langle V, \mathcal{E}' \rangle) = \sum_{j \in \mathcal{E}'} \alpha_j f_E(sum(\{V_i^g | E'_{ij} = 1\}_{i \in V})), \qquad (20)$$

where $\alpha_j$ are some scalar values.

According to Definition 1, if $f_{IHGNN}$ does not flag a statistical interaction on all the features in $q$, we can rewrite Equation (20) as a linear combination of a set of functions that each function does not rely on at least one of the features in $q$. That is:

$$f_{IHGNN}(G_H \langle V, \mathcal{E}' \rangle) = \sum_{i \in q} f_{\backslash i}(\{V_j^g\}_{j \in V \backslash \{i\}}), \qquad (21)$$

However, since all features in $q$ form a hyperedge in $\mathcal{E}'$, Equation (20) contains a component $f_E(sum(\{V_i^g | i \in q\}_{i \in V}))$, which non-linearly models all features in $q$. This component does not belong to any part of the right hand side of Equation (21). Therefore, Equation (21) does not hold, which is a contradiction of the assumption. □

Note that according to the definition of statistical interaction, if $f_{IHGNN}$ flags a statistical interaction on all features in $q$, $f_{IHGNN}$ flags a statistical interaction on any $q_s \subseteq q$. This is the reason that we use Infomin to force each hyperedge to contain different feature interactions and prevent multiple feature interactions from merging into one hyperedge.

# C DATASETS AND BASELINES

*Datasets.* We study three real-world datasets to evaluate our model: **MovieLens 1M** [13] contains users' ratings on movies. Each data sample contains a user and a movie with their corresponding attributes as features. We further collect movies' other features, such as directors and casts from IMDB to enrich the datasets. **Book-crossing** [46] contains users' implicit and explicit ratings of books. Each data sample contains a user and a book with their corresponding features. The reprocessed words in the book titles are also regarded as features of the book. **MovieLens 25M** [13] contains movie attribute information and the tags (e.g., gunfight, dragon) that users gave to the movies. We regard the attributes genres, years, and the 10 most relevant tags for movie as features.

We transfer the explicit ratings to implicit feedback. We regard the ratings greater than 3 as positive ratings for MovieLens 1M and MovieLens 25M, and regard all rated explicit ratings as positive ratings for Book-crossing due to its sparsity. Then, we randomly select the same number of negative samples equal to the number of positive samples for each user.

*Baselines.* We compare our model with eleven state-of-the-art models that consider feature interactions, including:

**FM** [19] models every pairwise feature interactions by dot product and sum up all the modeling results as the prediction. **AFM** [42] calculates an attention value for each feature interaction on top of FM. **NFM** [14] models each pairwise feature interaction and aggregates all the modeling results by an MLP. **Fi-GNN** [21] represents each data sample as a feature graph and models all pairwise interactions in a self-attention based GNN. $L_0$-**SIGN** [33] represents each data sample as a feature graph. It detects beneficial pairwise feature interactions and uses a GNN to model the detected ones for recommendation. **AutoInt** [30] explicitly models all pairwise feature interactions using a multi-head self-attentive neural network and aggregate all the modeling results as the prediction. **DeepFM** [12] combines the results of an MLP that implicitly models high-order

| Method | Order | | | | | | |
|--------|:-----:|:-----:|:-----:|:-----:|:---:|:------:|:------:|
| | 2 | 3 | 4 | 5 | ... | 13 | 14 |
| AutoInt | 5.48 | 6.26 | 10.96 | 12.59 | ... | 23.34 | 24.27 |
| DCNv2 | 9.40 | 14.09 | 18.79 | 22.71 | ... | 46.19 | 50.11 |
| xDeepFM | 18.79 | 39.93 | 59.51 | 78.30 | ... | 229.42 | 248.21 |
| HIRS ($k$=40) | | | | 15.16 | | | |
| HIRS ($k$=20) | | | | 11.38 | | | |

**Table 3: Comparing the time (in seconds) of different models for one epoch on the MovieLens 1M dataset. All the models are run on the same machine equipped with a CPU: Intel(R) Core(TM) i7-9700K @ 3.60GHz, and a GPU: GeForce RTX 2080Ti.**

feature interactions and a FM that models pairwise feature interactions. **xDeepFM** [23] is an extension of DeepFM that models high-order feature interactions in explicit way. **DCNv2** [38] leverages deep and cross network to implicitly and explicitly model feature interactions. **AutoFIS** [25] detects limited high-order beneficial feature interactions by setting a gate for each feature interaction. **AFN** [7] leverages logarithmic neural network to learn high-order feature interactions adaptively from data. For all baselines, we use the same MLP settings in all baselines (if use) as our interaction modeling function in HIRS for a fair comparison. In terms of he highest order for baselines that explicit model high-order feature interactions, we follow the order settings in their original papers.

## D RUNNING TIME ANALYSIS

We have theoretically analyzed the efficiency advantage of HIRS over other baselines that explicitly consider high-order feature interactions. In this section, we experimentally evaluate the running time of HIRS and these baselines. Specifically, we run HIRS and baselines on the MovieLens 1M dataset and record the time consuming they used in each epoch while training. For the baselines, we record the running time of setting different order limitations. We then compare with HIRS that consider arbitrary order feature interactions. Table 3 demonstrates the time consuming results. Note that since the running environment of our model and baseline models are not exact the same, the time comparison results reported are not completely accurate. For example, some baseline models directly call packages for compiling, which is more time efficient. However, we can still gain useful information from the results.

From the table, we can see that the running time of our model approximately equals to baseline models modeling low order feature interactions, which is consistent with our theoretical analysis in Section 5.6. Meanwhile, the running time for baseline models increase with the highest order increasing, while HIRS has constant running time. The results proves the efficiency of our model in modeling arbitrary order feature interactions, while baseline models have much higher running time when considering up to arbitrary orders (order 14). In addition, when we detect fewer feature interactions for each data sample, e.g., $k = 20$, the running time is further decreased comparing to that of $k = 40$.

## E ADDITIONAL EXPERIMENTAL RESULTS

In this section, we list additional results (in Figures 11-12) that are not shown in our paper due to the space limit. These results have consistent trends to the results reported in the paper, which provide further evidence on the robustness of our model.
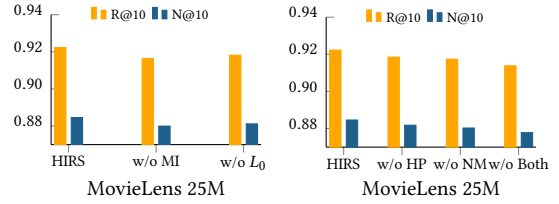


**Figure 10:** *Left*: Additional results for Figure 5 on the MovieLens 25M dataset. *Right*: Additional results for Figure 6 on the MovieLens 25M dataset.



**Figure 11: Additional results for Figure 4. Comparing the distribution of the generated feature interaction orders from HIRS, HIRS without s-Infomax and Infomin, and HIRS without $L_0$ regularization on the Book-Crossing and the MovieLens 25M datasets.**
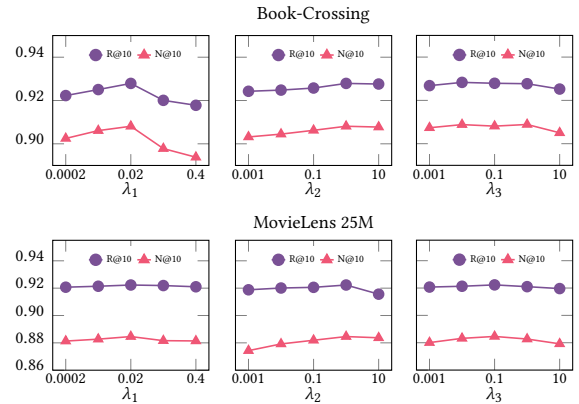


**Figure 12: Additional results for Figure 8. Comparing the performance of different weight values of $L_0$ activation regularization, s-Infomax, and Infomin on the Book-Crossing and the MovieLens 25M dataset.**