

A Broad Learning Approach for Context-Aware Mobile Application Recommendation

Tingting Liang[†], Lifang He[¶], Chun-Ta Lu[‡], Liang Chen[§], Philip S. Yu^{‡*}, Jian Wu[†]

[†] College of Computer Science & Technology, Zhejiang University, China

[¶] Department of Healthcare Policy and Research, Cornell University, NY, USA

[‡] Department of Computer Science, University of Illinois at Chicago, IL, USA

[§] School of Data and Computer Science, Sun Yat-Sen University, China

* Institute for Data Science, Tsinghua University, Beijing, China

{liangtt, wujian2000}@zju.edu.cn, {lifanghescut, jasonclx}@gmail.com, clu29@uic.edu, psyu@cs.uic.edu

Abstract—With the rapid development of mobile apps, the availability of a large number of mobile apps in application stores brings challenge to locate appropriate apps for users. Providing accurate mobile app recommendation for users becomes an imperative task. Conventional approaches mainly focus on learning users’ preferences and app features to predict the user-app ratings. However, most of them did not consider the interactions among the context information of apps. To address this issue, we propose a broad learning approach for Context-Aware app recommendation with Tensor Analysis (CATA). Specifically, we utilize a tensor-based framework to effectively integrate user’s preference, app category information and multi-view features to facilitate the performance of app rating prediction. The multidimensional structure is employed to capture the hidden relationships between multiple app categories with multi-view features. We develop an efficient factorization method which applies Tucker decomposition to learn the full-order interactions within multiple categories and features. Furthermore, we employ a group ℓ_1 -norm regularization to learn the group-wise feature importance of each view with respect to each app category. Experiments on two real-world mobile app datasets demonstrate the effectiveness of the proposed method.

I. INTRODUCTION

The rapid adoption of mobile devices accelerates the proliferation of mobile apps. The number of available apps in the Google Play¹ reached 2.8 million in Mar. 2017, and there have been 2.2 million mobile apps available in the Apple App Store² in Jan. 2017. The surge of mobile apps with diverse functions brings not only great convenience to users but also challenges for discovering appropriate apps. As a consequence, it becomes critical to develop effective approaches of rating prediction for recommending apps for users with accuracy.

There are some recent studies about the mobile apps recommendation, most of which leverage features of apps or users [1] [2] [3]. Karatzoglou et al. [1] proposed a collaborative filtering method for app recommendation by incorporating some contextual features like location, time of day, etc. Liu et al. [2] proposed to incorporate both app functionality and user privacy preference as features and capture the trade-off between them for app recommendation. Most of the previous works only tried one kind of feature or a simple

combination of multiple features and did not consider the complex interactions between those features. Currently there exist many works exploiting multiple views of features in the tasks like recommendation, clustering, etc [4], [5], [6], [7]. In the scenario of app recommendation, the interactions between different views of features are quite important as different views can provide complementary information. For example, assume we have obtained the latent representations for each app from three aspects, i.e., categories, permissions and description text, as shown in Fig. 1. BackCountry Navigator is an app categorized as *Maps&Navigation* and it is mainly used for outdoor navigation which can be inferred from the description text. The permission of getting users’ precise location is acceptable (i.e., a positive value), while the permission of reading SMS is abnormal (i.e., a negative value). It can be found that only the third-order interaction provides a negative result reflecting the unreasonable permission for app function. The category of Instagram merely shows its function for social interaction (i.e., a positive value) and neglects the function of sharing photos (i.e., a negative value). Through the interactions between multiple views, complementary information is provided to show a more sufficient understanding about the app. Obviously, the comprehensive consideration of the features from multiple views would be more insightful on understanding app information and user preference.

To figure out the latent correlations among the context information of apps, we conduct an empirical analysis on the dataset collected from Google Play and discover some important characteristics of mobile apps. For apps in different categories, users’ download behaviors are different. Within some categories like *Maps & Navigation* and *Weather*, users might only download one or two app for a long time use. However, for some categories like *Entertainment*, users are more likely to download many apps in the same category. Generally, the different download behaviors happen because users will consider different reasons (e.g. functions, interface, permissions) to decide whether to download apps for different categories. It can be inferred that users would focus on multiple views of features with different importance for apps in different categories. The analysis of the feature-level correlation between categories shows the similarities

¹Google Play: <https://play.google.com/store/apps>

²Apple App Store: <https://itunes.apple.com/us/genre/ios/id36?mt=8>

	Communication	Map & Navigation	Social ...	Read SMS	Precise Location	Take Pictures	... Message Navigation	Photos	Community ...	#1	#2	#3	
Messenger	1 (+)	0	0	0.80 (+)	0.10 (+)	0.10 (+)	0.80 (+)	0	0.10 (+)	0.10 (+)	3.00	3.00	1.00
BackCountry	0	1 (+)	0	-0.80 (-)	0.60 (+)	0	0	0.80 (+)	0	0	1.60	0.44	-0.16
Instagram	0	0	1 (+)	0.10 (+)	0.10 (+)	-0.80 (-)	0.10 (+)	0	-0.60 (-)	0.80 (+)	0.70	-0.48	-0.18
	Category (C)			Permission (P)			Description Text (D)						

Fig. 1. An example of feature interactions with different orders. The values in the column #1, #2, and #3 represent the summation of first-order, the second-order, and the third-order interactions. $\#1 = C + P + D$, $\#2 = C \times P + C \times D + P \times D$, $\#3 = C \times P \times D$.

between different categories are lower, which implies that the significances of features within a specific view are different for different categories. The category diversities of apps rated by different users are distinct. Some users prefer to download apps from various categories even though the amount of the downloaded apps is small. Based on the analysis, we consider to fuse the user preference, the category information, the features of multiple views, and the complex interactions among them to generate a context-aware category specific app rating prediction model.

In this paper, we propose a broad learning approach for Context-Aware app recommendation with Tensor Analysis (CATA). Specifically, we integrate the interactions among the multiple categories and multiple views of features into a tensor structure through the tensor product of the corresponding feature spaces. The interactions with different orders can fully reflect the complementary relationships, and we use them to predict the user ratings on apps. To effectively learn the full-order interactions³ without physically building a tensor, we further develop an efficient factorization method which employs Tucker decomposition. The Tucker decomposition is applied to factorize the interaction parameters for each order, which can make accurate parameter estimation under sparsity and avoid overfitting. Moreover, we introduce the group ℓ_1 -norm regularization for the global-specific weight matrix to further improve the proposed model.

The main contributions of this paper are as follows:

- We propose a context-aware recommendation approach for mobile apps called CATA that models the interactions with different orders among the multiple categories and multiple views of features as a tensor structure.
- To effectively learn the hidden relationships among the different views of the context information of apps, Tucker decomposition is adopted to factorize the interaction parameters such that the principal components of the latent representations can be retained.
- Empirical studies based on two real world datasets demonstrate the effectiveness of the proposed context-aware recommendation approach.

³Full-order interactions range from the first-order interactions (i.e., single-view features in each category) to the highest-order interactions (i.e., all combinations of features from multiple views and from different categories).

TABLE I
STATISTICS OF THE DATASETS

Dataset	#App	#User	#Feature	#Category	#Rating
Google Play	5460	7165	Text (2574) Permissions (84)	45	67504
Apple App Store	2643	4010	Text (1592)	2	74764

II. DATA ANALYSIS

In this section, we first describe the datasets used for the analysis and experiments. We then provide the statistical characteristics of the employed datasets.

A. Data Description

- **Google Play:** We crawled app’s meta data (e.g., name, category, permissions, description) and user review ratings from its description page in Google Play. We filter users and apps with less than 5 ratings. Each rating record in this dataset is represented in three views, i.e., users, permissions and text. The user view consists of binary feature vectors for user ids which means there is only one non-zero feature in the user view for each rating record. The TF-IDF vector representations of the app permissions and description texts are used as the permission and text view, respectively.
- **Apple’s App Store:** The dataset is offered by [8][9] and consists of the apps in the “Top Free 300” and “Top Paid 300” leaderboards from Feb. 2010 to Sep. 2012, and the related user ratings and review information. As the dataset lacks of the classification information, we use *Free* and *Paid* as two categories, and we remove users and apps with less than 10 ratings. Each rating record in this dataset has two views, i.e., users and text. The user view are constructed using the same way as in the Google Play dataset. The TF-IDF vector representations of the review texts of apps are used as the text view.

Table I shows the basic statistics of the two employed datasets.

B. Characteristics of Google Play dataset

As the Google Play dataset has richer category and feature information, we focus on the analysis for it.

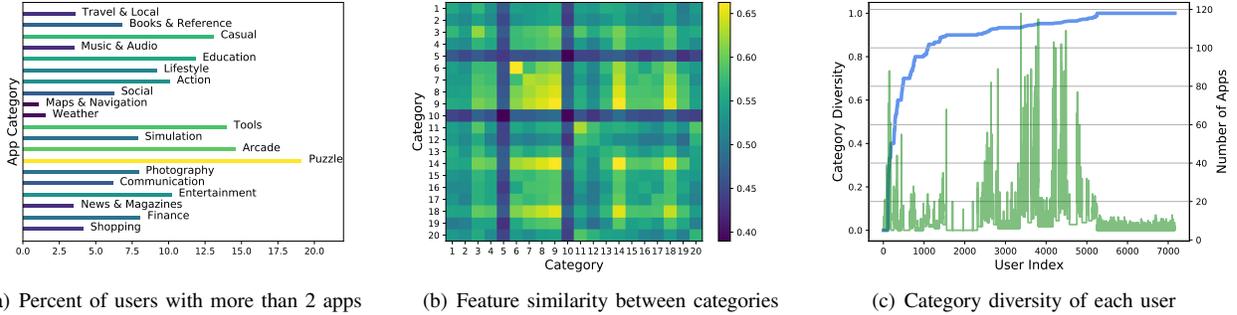


Fig. 2. Characteristics of Google Play dataset. The blue curve and green curve in (c) represent category diversity and app number, respectively.

We calculate the proportions of the users who downloaded more than 2 apps for each category. Due to the space limitation, Fig. 2(a) reports the results of 20 randomly selected app categories. The lower proportion for a category is, the more users only download one or two apps in that category. Generally, apps in the categories with a low proportion can be used for a long time. Taking category *Maps & Navigation* in which only about 1.07% users downloaded more than 2 apps as an example, users usually only download one or two apps (e.g., *Google Map*, *Baidu Map*) in this category as these apps are sufficient to use. The other categories having the low proportions in Fig 2(a) are *Weather*, *Travel& Local*, etc. The categories with high proportions are more general and the apps in them are with more varieties, like *Tools*, *Puzzle*, *Arcade*, etc. For apps in these categories, users might consider different reasons to decide whether to download the apps. For example, users mainly consider if the functions of the apps in category *Tools* will meet their demands. But for apps in category *Maps & Navigation*, compared to the functions, users pay more attention to the features like interface or permissions as they have clearly known the functions. We can learn that users would focus on multiple views of features with different significances for apps in different categories.

After the investigation of relationship between multiple views and categories, we explore the features within a specific view. To investigate the feature-level correlation between categories with respect to a certain view, we calculate the feature-based similarities between each pair of apps from any two categories. Figure 2(b) shows the similarities generated based on app permission feature between any two categories. The category indexes are sorted by the order in Fig. 2(a) (i.e., #1 is *Shopping*). It can be observed that the similarities between two different categories are generally lower than those between the same categories. That is, for apps of different categories, the significance of features within a specific view are distinguishing.

To investigate the relationship between users and categories, we apply the diversity metric widely used for the evaluation of recommender systems [10] to evaluate the category diversity. The category diversity is calculated by $Div(u) = 1 - \frac{\sum_{i,j \in A(u), i \neq j} s(i,j)}{\frac{1}{2}|A(u)|(|A(u)|-1)}$, where $A(u)$ is the set of apps rated by

user u . $s(i, j) = 1$ if app i and j belong to the same category, otherwise, $s(i, j) = 0$. Figure 2(c) shows the category diversity and the number of apps for each user. The green curve presents the number of apps rated by users, and the blue curve is the category diversity of the apps rated by users. The user indexes on the x -axis are sorted by the values of category diversity in an ascending order. The left y -axis shows the value of category diversity and the right y -axis represents the number of apps. It can be found that some users have interactions with many types of apps even though the numbers of apps rated by them are very small while some users rate many apps with few categories. Different users have interactions with categories with different diversities. As discussed above, the importance of features from multiple views and features within a specific view is distinct for different categories. Therefore, for each user, it is critical to model his preference on an app considering the category information and the corresponding relationships with features of multiple views.

Based on the analysis of the relationships among app category, app feature, and user, it requires a recommendation model which can integrate the interactions among the multiple categories, multiple views of features, and users.

III. PRELIMINARIES

In this work, we intend to predict ratings for mobile applications by a tensor-based approach. Before that, we introduce some related concepts and notation in tensor algebra that will be used throughout the paper, and then provide the problem formulation of app rating prediction.

A. Tensor Concepts and Notation

A tensor is a multi-dimensional array which generalizes matrix representation. Each dimension in tensor is called *mode* or *way*. Following prevailing convention, tensors are represented by calligraphic letters, matrices by boldface uppercase letters, vectors by boldfaced lowercase letters, and scalars by lowercase letters. An element of a vector \mathbf{x} , a matrix \mathbf{X} , or a tensor \mathcal{X} is represented by x_i , $x_{i,j}$, $x_{i,j,k}$, etc., depending on the number of modes. All vectors are column vectors unless otherwise specified. For an arbitrary matrix $\mathbf{X} \in \mathbb{R}^{I \times J}$, its i -th row and j -th column vector are represented by \mathbf{x}^i and \mathbf{x}_j ,

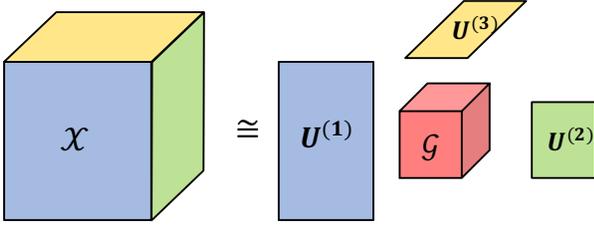


Fig. 3. Tucker decomposition of a third-order weight tensor.

respectively. The outer product of N vectors $\mathbf{x}^{(n)} \in \mathbb{R}^{I_n}$ for all $n \in [1 : N]$ is an N -th-order tensor and defined elementwise as $(\mathbf{x}^{(1)} \circ \dots \circ \mathbf{x}^{(N)})_{i_1, \dots, i_N} = x_{i_1}^{(1)} \dots x_{i_N}^{(N)}$ for $i_n \in [1 : I_n]$. The inner product of two tensors $\mathcal{X}, \mathcal{Y} \in \mathbb{R}^{I_1 \times \dots \times I_N}$ is defined as $\langle \mathcal{X}, \mathcal{Y} \rangle = \sum_{i_1=1}^{I_1} \dots \sum_{i_N=1}^{I_N} x_{i_1, \dots, i_N} y_{i_1, \dots, i_N}$. In particular, for $\mathcal{X} = \mathbf{x}^{(1)} \circ \dots \circ \mathbf{x}^{(N)}$ and $\mathcal{Y} = \mathbf{y}^{(1)} \circ \dots \circ \mathbf{y}^{(N)}$, it holds that

$$\langle \mathcal{X}, \mathcal{Y} \rangle = \prod_{n=1}^N \langle \mathbf{x}^{(n)}, \mathbf{y}^{(n)} \rangle = \prod_{n=1}^N \mathbf{x}^{(n)T} \mathbf{y}^{(n)}. \quad (1)$$

Definitions of Kronecker product, Khatri–Rao product, mode- n product, and Tucker decomposition are given below, which will be applied to build the proposed model.

Definition 3.1 (Kronecker Product): The Kronecker product of matrices $\mathbf{X} \in \mathbb{R}^{I \times J}$ and $\mathbf{Y} \in \mathbb{R}^{K \times L}$ is denoted by $\mathbf{X} \otimes \mathbf{Y}$. The result is a matrix of size $(IK) \times (JL)$ and defined by

$$\begin{aligned} \mathbf{X} \otimes \mathbf{Y} &= \begin{bmatrix} x_{1,1}\mathbf{Y} & x_{1,2}\mathbf{Y} & \dots & x_{1,J}\mathbf{Y} \\ \vdots & \vdots & \ddots & \vdots \\ x_{I,1}\mathbf{Y} & x_{I,2}\mathbf{Y} & \dots & x_{I,J}\mathbf{Y} \end{bmatrix} \\ &= [\mathbf{x}_1 \otimes \mathbf{y}_1 \quad \mathbf{x}_1 \otimes \mathbf{y}_2 \quad \dots \quad \mathbf{x}_J \otimes \mathbf{y}_{L-1} \quad \mathbf{x}_J \otimes \mathbf{y}_L]. \end{aligned} \quad (2)$$

Definition 3.2 (Khatri–Rao Product): The Khatri–Rao product of matrices $\mathbf{X} \in \mathbb{R}^{I \times K}$ and $\mathbf{Y} \in \mathbb{R}^{J \times K}$ is denoted by $\mathbf{X} \odot \mathbf{Y}$. The result is a matrix of size $(IJ) \times K$ and defined by

$$\mathbf{X} \odot \mathbf{Y} = [\mathbf{x}_1 \otimes \mathbf{y}_1 \quad \mathbf{x}_2 \otimes \mathbf{y}_2 \quad \dots \quad \mathbf{x}_K \otimes \mathbf{y}_K]. \quad (3)$$

The Khatri–Rao product is the “matching columnwise” Kronecker product.

Definition 3.3 (n -mode Product): The n -mode product of a tensor $\mathcal{X} \in \mathbb{R}^{I_1 \times \dots \times I_N}$ with a matrix $\mathbf{U} \in \mathbb{R}^{J \times I_n}$ denoted by $\mathcal{X} \times_n \mathbf{U}$ is defined as

$$(\mathcal{X} \times_n \mathbf{U})_{i_1, \dots, i_{n-1}, j, i_{n+1}, \dots, i_N} = \sum_{i_n=1}^{I_n} x_{i_1, \dots, i_N} u_{j, i_n}. \quad (4)$$

Figure 3 visualizes the Tucker decomposition of a third-order tensor and table II summarizes the main notations for easy referencing.

TABLE II
LIST OF BASIC SYMBOLS.

Symbol	Definition and description
x	each lowercase letter represents a scale
\mathbf{x}	each boldface lowercase letter represents a vector
\mathbf{X}	each boldface capital letter represents a matrix
\mathcal{X}	each calligraphic letter represents a tensor, set or space
$[1 : N]$	a set of integers in the range of 1 to N inclusively.
$\langle \cdot, \cdot \rangle$	denotes inner product
\circ	denotes outer product
\otimes	denotes Kronecker product
\odot	denotes Khatri–Rao product
\times_n	denotes n -mode product
$\ \cdot\ _F$	denotes Frobenius norm of vector, matrix or tensor

Definition 3.4 (Tucker Decomposition): For a general tensor $\mathcal{X} \in \mathbb{R}^{I_1 \times \dots \times I_N}$, its Tucker decomposition is defined as

$$\begin{aligned} \mathcal{X} &\approx \mathcal{G} \times_1 \mathbf{U}^{(1)} \times_2 \dots \times_N \mathbf{U}^{(N)} \\ &= \sum_{r_1=1}^{R_1} \dots \sum_{r_N=1}^{R_N} g_{r_1, \dots, r_N} \mathbf{u}_{r_1}^{(1)} \circ \dots \circ \mathbf{u}_{r_N}^{(N)} \\ &= \llbracket \mathcal{G}; \mathbf{U}^{(1)}, \dots, \mathbf{U}^{(N)} \rrbracket, \end{aligned} \quad (5)$$

where $\mathbf{U}^{(n)} \in \mathbb{R}^{I_n \times R_n}$ are the factor matrices and can be thought of as the principal components in each mode. $\mathcal{G} \in \mathbb{R}^{R_1 \times \dots \times R_N}$ is called the *core tensor*. $\llbracket \cdot \rrbracket$ is used for short-hand notation.

B. Problem Formulation

Suppose that the scenario of app rating prediction includes a user set \mathcal{U} and mobile app set \mathcal{A} . The numbers of app categories and feature views are C and V . Let N_c be the number of the rating records in the category $c \in [1 : C]$, then the total number of rating records is $N = \sum_c N_c$. Let I_v be the dimensionality of the feature view $v \in [1 : V]$ and $I = \sum_v I_v$.

In this paper, we construct a multi-dimensional tensor to discover the latent interactions among the category information and multi-view features. Each rating record in category c can be represented in V different views, i.e., $\mathbf{x}_c^T = (\mathbf{x}_c^{(1)T}, \dots, \mathbf{x}_c^{(V)T})$, where $\mathbf{x}_c^{(v)} \in \mathbb{R}^{I_v}$ and $\mathbf{x}_c \in \mathbb{R}^I$. Generally, a rating record involves a user, an app, and different types of characteristics of the app. Given a training set of rating records $\mathcal{D} = \{(\mathbf{X}_c^{(1)}, \dots, \mathbf{X}_c^{(V)}, \mathbf{y}_c) | c \in [1 : C]\}$, where $\mathbf{X}_c^{(v)} \in \mathbb{R}^{I_v \times N_c}$ is the feature matrix in the c -th category for v -th view and \mathbf{y}_c is the vector of the rating values of those apps in the c -th category. Our goal is to find a predictive function $f_c : \mathcal{X}_c \rightarrow \mathcal{Y}_c$ for each category that can minimize the expected loss and provide accurate predicted ratings. The regularized objective function to be minimized can be formulated as:

$$\mathcal{H}(\{f_c\}_{c=1}^C) = \sum_{c=1}^C (\mathcal{L}_c(f_c(\{\mathbf{X}_c^{(v)}\}), \mathbf{y}_c)) + \lambda \Omega, \quad (6)$$

where \mathcal{L}_c is the empirical loss in the c -category. Ω is the regularization term and $\lambda > 0$ is the regularization parameter.

\mathcal{L}_c can be rewritten as the average square error of each instance.

$$\begin{aligned}\mathcal{L}_c(f_c(\{\mathbf{X}_c^{(v)}\}), \mathbf{y}_c) &= \frac{1}{N_c} \sum_{n=1}^{N_c} \ell(f_c(\{\mathbf{x}_{c,n}^{(v)}\}), \mathbf{y}_{c,n}) \\ &= \frac{1}{N_c} \sum_{n=1}^{N_c} (f_c(\{\mathbf{x}_{c,n}^{(v)}\}) - \mathbf{y}_{c,n})^2.\end{aligned}\quad (7)$$

IV. PROPOSED METHOD

In this section, we first introduce the context-aware recommendation approach based on tensor analysis (CATA). Then we discuss how to employ Tucker decomposition to learn the proposed model without physically building the tensor.

A. Model for App Rating Prediction

We derive the proposed model from the basic framework of linear analysis. Given a vector of an app rating record $\mathbf{x} \in \mathbb{R}^I$, the basic linear model for the c -th category is written as

$$f_c(\mathbf{x}) = \sum_{i=1}^I w_{c,i} x_i + w_{c,0} = \mathbf{x}^T \mathbf{w}_c + w_{c,0}, \quad (8)$$

where $\mathbf{w}_c \in \mathbb{R}^I$ is the weight vector for the c -th category, and $w_{c,0}$ is the bias factor for adjusting the threshold of the c -th category label assignment.

Let $\mathbf{z} = [1; \mathbf{x}] \in \mathbb{R}^{1+I}$ and $\mathbf{w}_c = [w_{c,0}; \mathbf{w}_c] \in \mathbb{R}^{(1+I)}$, then the bias factor $w_{c,0}$ can be absorbed to \mathbf{w}_c (see [11]). Eq. (8) can thus be rewritten as follows:

$$f_c(\mathbf{x}) = \mathbf{z}^T \mathbf{w}_c. \quad (9)$$

Let $\mathbf{W} \in \mathbb{R}^{(1+I) \times C}$ denote the weight matrix to be learned, whose columns are the vector \mathbf{w}_c . In order to jointly learn multiple linear models for C categories, we introduce a category indicator vector denoted by $\mathbf{e}_c \in \mathbb{R}^C$ to model the second-order interactions between input features and categories. The indicator vector \mathbf{e}_c is defined as

$$\mathbf{e}_c = \underbrace{[0, \dots, 0]_{c-1}}_{c-1}, 1, 0, \dots, 0]^T.$$

Then Eq. (9) can be rewritten as

$$f_c(\mathbf{x}) = \mathbf{z}^T \mathbf{w}_c = \mathbf{z}^T \mathbf{W} \mathbf{e}_c = \langle \mathbf{W}, \mathbf{z} \circ \mathbf{e}_c \rangle. \quad (10)$$

Note that the outer product is used to compute intersections between input features and categories, which consists in the product of all combinations of the variables that define each domain. This data fusion technique provides a good framework to introduce multiple features. When each object is associated with multi-view features, by means of the outer product we can easily extend the above Eq. (10) to the multi-view case and provide a consensus formulation.

Suppose that the given rating records are composed by features of V views (denoted as $\{\mathbf{x}^{(v)}\}, v \in [1 : V]$), we can extend Eq. (10) to model the full-order interactions between multi-view features and categories as:

$$f_c(\{\mathbf{x}^{(v)}\}) = \langle \mathcal{W}, \mathbf{z}^{(1)} \circ \dots \circ \mathbf{z}^{(V)} \circ \mathbf{e}_c \rangle, \quad (11)$$

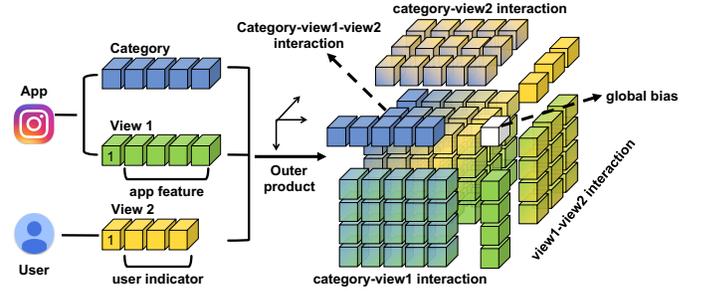


Fig. 4. The context information of each rating record are represented by multiple vectors and combined by the outer product to generate the full-order interactions. The full-order interactions between multiple views and categories are modeled in a tensor structure.

or element-wise as

$$f_c(\{\mathbf{x}^{(v)}\}) = \sum_{s=1}^C \sum_{i_1=0}^{I_1} \dots \sum_{i_V=0}^{I_V} w_{i_1, \dots, i_V, s} (e_{c,s} \prod_{v=1}^V z_{i_v}^{(v)}). \quad (12)$$

Where $\mathbf{z}^{(v)} = [1; \mathbf{x}^{(v)}] \in \mathbb{R}^{(1+I_v)}$ is the input data vector, and $\mathcal{W} = \{w_{i_1, \dots, i_V, s}\} \in \mathbb{R}^{(1+I_1) \times \dots \times (1+I_V) \times C}$ is the weight tensor to be learned, wherein $w_{0, \dots, 0}$ is the global bias, and $w_{i_1, \dots, i_V, s}$ with some indexes satisfying $i_v = 0$ encodes lower-order interactions between views whose $i_{v'} > 0$.

In such a manner, the full-order interactions between multiple views and categories are embedded within the tensor structure, as shown in Fig. 4. However, one drawback might be generated from the model is that not all the categories are fit to the constructed feature tensor and those interactions will be redundant information. Thus, we consider to build a rating predictive function based on both the full-order feature interaction space and the original feature spaces.

Let $\mathcal{Z}_c = \mathbf{z}^{(1)} \circ \dots \circ \mathbf{z}^{(V)} \circ \mathbf{e}_c \in \mathbb{R}^{(1+I_1) \times \dots \times (1+I_V) \times C}$ be the full-order tensor, and $\mathbf{x} = [\mathbf{x}^{(1)}; \dots; \mathbf{x}^{(V)}] \in \mathbb{R}^I$ be the feature vector concatenated by multiple views. We formulate our CATA model as follows:

$$f_c(\{\mathbf{x}^{(v)}\}) = \langle \mathcal{W}, \mathcal{Z}_c \rangle + \mathbf{x}^T \mathbf{d}_c, \quad (13)$$

where $\mathbf{d}_c \in \mathbb{R}^I$ is the category-specific weight vector. For convenience in the following discussion, we denote $\mathbf{D} = [\mathbf{d}_1, \dots, \mathbf{d}_C] \in \mathbb{R}^{I \times C}$.

B. Model Inference

The number of parameters to be estimated in Eq. (13) is $C(\prod_{v=1}^V (1 + I_v) + I)$, which makes it infeasible to directly learning the model. Therefore, we assume that the weight tensor \mathcal{W} can be factorized by Tucker decomposition as

$$\mathcal{W} = \llbracket \mathcal{G}; \Phi, \Theta^{(1)}, \dots, \Theta^{(V)} \rrbracket, \quad (14)$$

where $\mathcal{G} \in \mathbb{R}^{R_0 \times R_1 \times \dots \times R_V}$ is called the core tensor and its entries show the level of interaction between the different components, $\Theta^{(v)} \in \mathbb{R}^{(1+I_v) \times R_v}$ is the shared structure matrix for the v -th view, and $\Phi \in \mathbb{R}^{C \times R_0}$ is the category specific weight matrix.

Then we can transform Eq. (12) into

$$\begin{aligned}
\langle \mathcal{W}, \mathcal{Z}_c \rangle &= \sum_{s=1}^C \sum_{i_1=0}^{I_1} \cdots \sum_{i_V=0}^{I_V} w_{i_1, \dots, i_V, s} (e_{c,s} \prod_{v=1}^V z_{i_v}^{(v)}) \\
&= \sum_{s=1}^C \sum_{i_1=0}^{I_1} \cdots \sum_{i_V=0}^{I_V} \left(\sum_{r_0=1}^{R_0} \cdots \sum_{r_V=1}^{R_V} g_{r_0, \dots, r_V} \phi_{s, r_0} \prod_{v=1}^V \theta_{i_v, r_v}^{(v)} \right) (e_{c,s} \prod_{v=1}^V z_{i_v}^{(v)}) \\
&= \sum_{r_0=1}^{R_0} \cdots \sum_{r_V=1}^{R_V} g_{r_0, \dots, r_V} \left(\sum_{s=1}^C \phi_{s, r_0} e_{c,s} \right) \sum_{i_1=0}^{I_1} \cdots \sum_{i_V=0}^{I_V} \left(\prod_{v=1}^V \theta_{i_v, r_v}^{(v)} z_{i_v}^{(v)} \right) \\
&= \sum_{r_0=1}^{R_0} \cdots \sum_{r_V=1}^{R_V} g_{r_0, \dots, r_V} \langle \boldsymbol{\theta}_{r_1}^{(1)} \circ \cdots \circ \boldsymbol{\theta}_{r_V}^{(V)} \circ \boldsymbol{\phi}_{r_0}, \mathbf{z}^{(1)} \circ \cdots \circ \mathbf{z}^{(V)} \circ \mathbf{e}_c \rangle
\end{aligned} \tag{15}$$

Because $e_{c,s} = 1$ only when $c = s$ and according to Eq. (1), we can further rewrite the equation above into

$$\begin{aligned}
\langle \mathcal{W}, \mathcal{Z}_c \rangle &= \sum_{r_0=1}^{R_0} \cdots \sum_{r_V=1}^{R_V} g_{r_0, \dots, r_V} \phi_{c, r_0} (\mathbf{z}^{(1)T} \boldsymbol{\theta}_{r_1}^{(1)}) \cdots (\mathbf{z}^{(V)T} \boldsymbol{\theta}_{r_V}^{(V)}) \\
&= \mathcal{G} \times_0 \phi^c \times_1 (\mathbf{z}^{(1)T} \boldsymbol{\Theta}^{(1)}) \times_2 \cdots \times_V (\mathbf{z}^{(V)T} \boldsymbol{\Theta}^{(V)}),
\end{aligned} \tag{16}$$

where \times_v is the v -mode product and \times_0 means multiplying the core tensor by the category specific vector ϕ^c . It is worth noting that the first row $\boldsymbol{\theta}^{(v),0}$ within $\boldsymbol{\Theta}^{(v)}$ is associated with the constant value $z_0^{(v)} = 1$ and represents the bias factors of the v -th view. The bias factors make the lower-order interactions active in the rating predictive function.

Using Eq. (16) to replace the first term in Eq. (13), the rating predictive function can be represented by

$$f_c(\mathbf{x}^{(v)}) = \mathcal{G} \times_0 \phi^c \times_1 (\mathbf{z}^{(1)T} \boldsymbol{\Theta}^{(1)}) \times_2 \cdots \times_V (\mathbf{z}^{(V)T} \boldsymbol{\Theta}^{(V)}) + \mathbf{x}^T \mathbf{d}_c. \tag{17}$$

The whole framework of the proposed CATA method is illustrated in Fig. 5.

C. Model Estimation

We propose to learn the app rating prediction model CATA by minimizing the following regularized empirical risk:

$$\begin{aligned}
\min \mathcal{H}(\boldsymbol{\Phi}, \{\boldsymbol{\Theta}^{(v)}\}, \mathcal{G}, \mathbf{D}) &= \sum_{c=1}^C \mathcal{L}_c(f_c(\{\mathbf{X}_c^{(v)}\}), \mathbf{y}_c) + \\
&\alpha \Omega_\alpha(\boldsymbol{\Phi}, \{\boldsymbol{\Theta}^{(v)}\}, \mathcal{G}) + \beta \Omega_\beta(\mathbf{D}).
\end{aligned} \tag{18}$$

The regularization term Ω_α and Ω_β can be set as Frobenius norm, $\ell_{2,1}$ norm, or other structural regularization. In this paper, we adopt the alternating block coordinate descent approach for the optimization of the given objective function. The whole learning procedure is summarized in Algorithm 1.

With all other parameters fixed, the minimization over $\boldsymbol{\Theta}^{(v)}$ consists of learning the parameters $\boldsymbol{\Theta}^{(v)}$ by a regularization method, and the partial derivative of \mathcal{H} w.r.t. $\boldsymbol{\Theta}^{(v)}$ is given by

$$\frac{\partial \mathcal{H}}{\partial \boldsymbol{\Theta}^{(v)}} = \sum_{c=1}^C \frac{\partial \mathcal{L}_c}{\partial f_c} \frac{\partial f_c}{\partial \boldsymbol{\Theta}^{(v)}} + \alpha \frac{\partial \Omega_\alpha(\boldsymbol{\Theta}^{(v)})}{\partial \boldsymbol{\Theta}^{(v)}} \tag{19}$$

where $\frac{\partial \mathcal{L}_c}{\partial f_c} = \frac{1}{N_c} \left[\frac{\partial \ell_{c,1}}{\partial f_c}, \dots, \frac{\partial \ell_{c,N_c}}{\partial f_c} \right]^T \in \mathbb{R}^{N_c}$ and $\frac{\partial \ell_{c,n}}{\partial f_c} = 2(f_c - \mathbf{y}_{c,n})$ for $n \in [1 : N_c]$.

For convenience, we let $\boldsymbol{\pi} \in \mathbb{R}^{1 \times (R_1 \cdots R_V)}$ denote the Kronecker product in a reverse order from V to 1 $\prod_{v=V}^1 \otimes (\mathbf{z}^{(v)T} \boldsymbol{\Theta}^{(v)})$ and $\boldsymbol{\pi}^{(-v)} \in \mathbb{R}^{1 \times (R_1 \cdots R_{v-1} R_{v+1} \cdots R_V)}$ denote $\prod_{v'=V, v' \neq v}^1 \otimes (\mathbf{z}^{(v')T} \boldsymbol{\Theta}^{(v')})$. Let $\boldsymbol{\Pi} = [\boldsymbol{\pi}_1, \dots, \boldsymbol{\pi}_N]^T$ and $\boldsymbol{\Pi}^{(-v)} = [\boldsymbol{\pi}_1^{(-v)}, \dots, \boldsymbol{\pi}_N^{(-v)}]^T$. Then we have that

$$\begin{aligned}
\frac{\partial \mathcal{L}_c}{\partial f_c} \frac{\partial f_c}{\partial \boldsymbol{\Theta}^{(v)}} &= \frac{1}{N_c} \sum_{n=1}^{N_c} \frac{\partial \ell_{c,n}}{\partial f_{c,n}} \frac{\partial f_{c,n}}{\partial \boldsymbol{\Theta}^{(v)}} \\
&= \sum_{n=1}^{N_c} \mathbf{z}_{c,n}^{(v)} \left(\frac{1}{N_c} \frac{\partial \ell_{c,n}}{\partial f_{c,n}} \right) (\boldsymbol{\pi}^{(-v)} \otimes \phi^c) \mathbf{G}_{(v)}^T \\
&= \mathbf{z}_c^{(v)} \left((\boldsymbol{\Pi}^{(-v)})^T \odot \left(\frac{\partial \mathcal{L}_c}{\partial f_c} \phi^c \right)^T \right)^T \mathbf{G}_{(v)}^T,
\end{aligned} \tag{20}$$

where $\mathbf{G}_{(n)}$ the n -mode matricization of tensor \mathcal{G} .

With all other parameters fixed, the minimization over $\boldsymbol{\Phi}$ consists of learning each parameter component ϕ^c independently. The partial derivative of \mathcal{H} w.r.t. $\boldsymbol{\Phi}$ is given by

$$\frac{\partial \mathcal{H}}{\partial \boldsymbol{\Phi}} = \left[\frac{\partial \mathcal{L}_1}{\partial f_1} \frac{\partial f_1}{\partial \phi^1}; \dots; \frac{\partial \mathcal{L}_C}{\partial f_C} \frac{\partial f_C}{\partial \phi^C} \right] + \alpha \frac{\partial \Omega_\alpha(\boldsymbol{\Phi})}{\partial \boldsymbol{\Phi}}. \tag{21}$$

Following the derivation in Eq.(20), we have that

$$\frac{\partial \mathcal{L}_c}{\partial f_c} \frac{\partial f_c}{\partial \phi^c} = \left(\frac{\partial \mathcal{L}_c}{\partial f_c} \right)^T \boldsymbol{\Pi}_c \mathbf{G}_{(0)}^T \tag{22}$$

By keeping all other parameters fixed, we can get the partial derivative of \mathcal{H} w.r.t. the core tensor \mathcal{G} as follows,

$$\frac{\partial \mathcal{H}}{\partial \mathcal{G}} = \sum_{c=1}^C \frac{\partial \mathcal{L}_c}{\partial f_c} \frac{\partial f_c}{\partial \mathcal{G}} + \alpha \frac{\partial \Omega_\alpha(\mathcal{G})}{\partial \mathcal{G}} \tag{23}$$

Following the derivation in Eq.(20), we have that

$$\begin{aligned}
\left(\frac{\partial \mathcal{L}_c}{\partial f_c} \frac{\partial f_c}{\partial \mathcal{G}} \right)_{(1 \times (R_0 \cdots R_V))} &= \sum_{n=1}^{N_c} \left(\frac{1}{N_c} \frac{\partial \ell_{c,n}}{\partial f_{c,n}} \right) (\boldsymbol{\pi} \otimes \phi^c) \\
&= \left(\frac{\partial \mathcal{L}_c}{\partial f_c} \right)^T (\boldsymbol{\Pi}_c \otimes \phi^c)
\end{aligned} \tag{24}$$

By keeping all other parameters fixed, we can get the partial derivative of \mathcal{H} w.r.t. the core tensor \mathbf{D} as follows,

$$\frac{\partial \mathcal{H}}{\partial \mathbf{D}} = \left[\mathbf{X}_1 \frac{\partial \mathcal{L}_1}{\partial f_1}; \dots; \mathbf{X}_C \frac{\partial \mathcal{L}_C}{\partial f_C} \right] + \beta \frac{\partial \Omega_\beta(\mathbf{D})}{\partial \mathbf{D}}, \tag{25}$$

where $\mathbf{X}_c = [\mathbf{X}_c^{(1)}; \dots; \mathbf{X}_c^{(V)}] \in \mathbb{R}^{I \times N_c}$ is the concatenated feature matrix for the c -th category.

D. Group ℓ_1 -Norm

As mentioned in section IV-C, the regularization terms can be Frobenius norm, $\ell_{2,1}$ norm, or other structural regularization, here we present a proper regularization term for parameter \mathbf{D} to further improve the performance of rating prediction.

For the original feature spaces, *i.e.*, the second term in Eq. (13), the feature of a specific view might be more or less discriminative for different app categories. For instance, the description information is more useful for the distinguishing

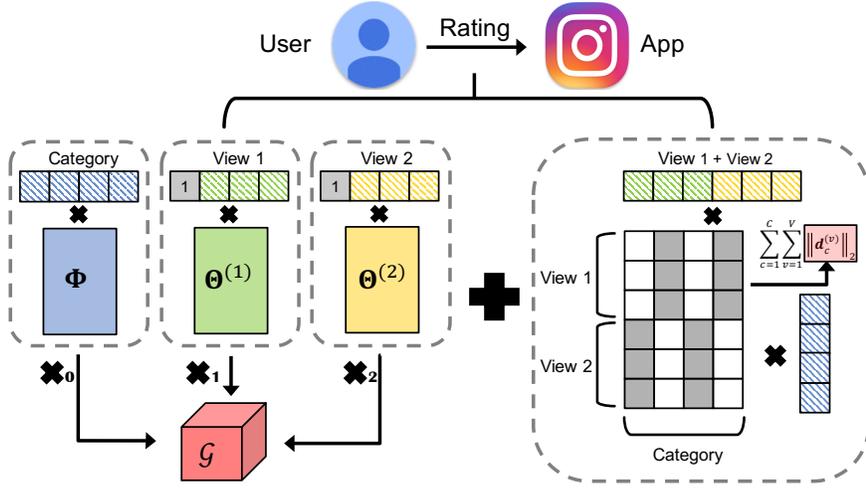


Fig. 5. Framework of the proposed CATA method. The left part is the first term in Eq. (17) which represents the approximation of the full-order feature interactions through Tucker decomposition. The right part is the second term in Eq. (17) which involves the category-specific weight matrix and the group ℓ_1 -norm.

Algorithm 1: Learning CATA Model

Input: Training data \mathcal{D} , number of factors R , regularization parameter α, β , and learning rate η

Output: Model parameters $\{\Theta^{(v)}\}, \Phi, \mathcal{G}, \mathbf{D}$

- 1 Initialize $\{\Theta^{(v)}\}, \Phi, \mathcal{G}, \mathbf{D} \sim \mathcal{N}(0, \sigma)$.
 - 2 **repeat**
 - 3 Fixing $\{\Theta^{(v)}\}, \mathcal{G}$, and \mathbf{D} , update Φ
 - 4 **for** $v = 1 : V$ **do**
 - 5 Fixing $\{\Theta^{(v')}\}_{v' \neq v}, \Phi, \mathcal{G}$, and \mathbf{D} , update $\Theta^{(v)}$
 - 6 Fixing $\{\Theta^{(v)}\}, \Phi$, and \mathbf{D} , update \mathcal{G}
 - 7 Fixing $\{\Theta^{(v)}\}, \Phi$, and \mathcal{G} , update \mathbf{D}
 - 8 **until convergence**;
-

of apps in the *Lifestyle* category than that of apps in *Map & Navigation* category. It is mainly because *Lifestyle* is a broad cluster and the functionality which could be extracted from the description text of each app in it are very different from each other. Consider this, we introduce group ℓ_1 -norm (G_1 -norm, for short) for regularization, which is defined as $\|\mathbf{D}\|_{G_1} = \sum_{c=1}^C \sum_{v=1}^V \|\mathbf{d}_c^{(v)}\|_2$ [12]. The G_1 -norm applies ℓ_2 -norm within each view and ℓ_1 -norm between views, so it can enforce the sparsity between different views. It means that if a specific view of features are not significant for the apps in a certain category, the weights with very small values will be assigned to them for the corresponding category. The G_1 -norm can further improve the performance of app rating prediction as it captures the global relationships between views. The right part of Fig. 5 simply shows the category-specific weight matrix as an illustration. The elements with gray color have large values. It can be found that the G_1 -norm effectively emphasizes the view-wise weight learning corresponding to each category.

V. EXPERIMENTS

In this section, we will verify the effectiveness of the proposed method by conducting a series experiments compared to five well known baselines.

A. Experimental Setup

After the filtering for the Google Play dataset, we first select the top 20 categories with the most apps and then filter users and apps with less than 5 ratings. We obtain 3065 apps and 3895 users with 36791 rating records. The numbers of permissions and text tokens are 83 and 1762, respectively.

We randomly select $K\%$ ($K = 60, 70, 80$), 10%, and 10% of the rating records in each categories as training set, validation set, and testing set. The parameters of all the baselines are set to the optimal values. For the proposed methods, all the dimensions of the core tensor are set as 5, and the learning rate is set $\eta = 0.1$. The maximum numbers of iterations are set as 400. Grid searching is employed to select the optimal regularization parameters for all the comparison methods. Each experiment is repeated for 5 times, and the mean and standard deviation of each metric in both dataset are reported in the next subsection.

We use the Mean Absolute Error (MAE) and Root Mean Square Error (RMSE) [13] to evaluate the performance of the proposed approach and the other compared methods. A smaller MAE or RMSE means the better performance.

B. Compared Methods

In order to demonstrate the effectiveness of the proposed CATA approach, we compare the following methods.

- **PMF.** It is the Probabilistic Matrix Factorization proposed by Salakhutdinov and Minh [14], and the method is widely used for rating prediction tasks.
- **MTFL.** It is the Multi-Task Feature Learning algorithm [12] which is a multivariate regression model with group ℓ_1 -norm.

TABLE III
PERFORMANCE COMPARISON ON THE GOOGLE PLAY DATASET. THE BEST RESULTS ARE LISTED IN BOLD.

Training	Metrics	PMF	MTFL	FM	MVM	MFM	CATA	CATA-G
60%	MAE	0.9597±0.0157	0.9272±0.0247	0.8964±0.0131	0.8735±0.0299	0.8761±0.0177	0.7650±0.0571	0.7679±0.0152
	RMSE	1.3015±0.0234	1.4616±0.0305	1.2133±0.0206	1.2122±0.0353	1.2021±0.0171	1.1720±0.0402	1.1586±0.0209
70%	MAE	0.9463±0.0083	0.8889±0.0053	0.8786±0.0065	0.8496±0.0070	0.8660±0.0350	0.7911±0.0442	0.7864±0.0141
	RMSE	1.2834±0.0090	1.4257±0.0117	1.1981±0.0096	1.1836±0.0089	1.1959±0.0122	1.1656±0.0151	1.1626±0.0115
80%	MAE	0.9369±0.0108	0.8575±0.0201	0.8568±0.0112	0.8384±0.0134	0.8439±0.0266	0.7826±0.0274	0.7765±0.0206
	RMSE	1.2745±0.0157	1.3904±0.0260	1.1785±0.0131	1.1741±0.0113	1.1815±0.0147	1.1502±0.0211	1.1419±0.0154

TABLE IV
PERFORMANCE COMPARISON ON THE APPLE APP STORE DATASET. THE BEST RESULTS ARE LISTED IN BOLD.

Training	Metrics	PMF	MTFL	FM	MVM	MFM	CATA	CATA-G
60%	MAE	1.0609±0.0062	0.9856±0.0221	0.9426±0.0038	0.9463±0.0098	0.9311±0.0113	0.9342±0.0087	0.9271±0.0052
	RMSE	1.3180±0.0064	1.3064±0.0231	1.2890±0.0069	1.2717±0.0191	1.2422±0.0070	1.2396±0.0059	1.2377±0.0060
70%	MAE	1.0551±0.0056	0.9856±0.0188	0.9345±0.0091	0.9409±0.0135	0.9198±0.0164	0.9246±0.0081	0.9247±0.0077
	RMSE	1.3046±0.0091	1.3062±0.0219	1.2745±0.0123	1.2526±0.0157	1.2312±0.0135	1.2257±0.0126	1.2262±0.0112
80%	MAE	1.0541±0.0052	0.9842±0.0130	0.9404±0.0057	0.9427±0.0080	0.9526±0.0206	0.9325±0.0019	0.9288±0.0057
	RMSE	1.3065±0.0033	1.2955±0.0153	1.2800±0.0090	1.2493±0.0141	1.2372±0.0083	1.2317±0.0066	1.2286±0.0079

- **FM.** It is the Factorization Machine [15] that explores pairwise interactions between all features without view segmentation. We implement the FM by concatenating the category indicator vector and all the feature vectors as the input feature vector.
- **MVM.** It is the Multi-view Machine [16] that models the features from multiple views as a tensor structure to explore the full-order interactions between them.
- **MFM.** It is the Multilinear Factorization Machines [17] that learns task-specific feature map and the task-view shared multilinear structures from full-order interactions by applying a joint factorization.
- **CATA.** It is the proposed rating prediction model in this paper that effectively integrates user’s preference, app category and features of multiple views and applies Tucker decomposition to learn the full-order interactions.
- **CATA-G.** It is the variation of the proposed CATA that uses group ℓ_1 -norm for the category-specific weight matrix \mathbf{D} .

C. Performance Comparison

In this subsection, we present the performance comparisons between the proposed CATA methods and the baselines with respect to two metrics, *i.e.*, MAE and RMSE.

Table III and Table IV show the performance of all the prediction methods on the Google Play and Apple App Store datasets. We can find that the proposed approach consistently outperforms the other baselines on both datasets in almost all cases. It demonstrates the superiority of the context-aware prediction approach which utilizes higher-order decomposition to learn the full-order interactions. It can be observed that CATA-G method performs better than CATA overall, which indicates that the employed group ℓ_1 -norm can effectively improve the rating prediction accuracy by enforcing the sparsity between different views of features.

It is not surprising that PMF has poor performance in both datasets since it doesn’t employ any other features of

apps. MTFL also performs badly mainly because it ignores the segmentation of feature views and the interactions between the multiple views of features. Compared to MTFL, the improvement achieved by FM illustrates the necessity of feature interactions. Both MVM and MFM outperform FM, especially for the Google Play dataset, and the results generated by them are competitive with each other. It is mainly because that MVM and MFM consider full-order interactions including the higher-order feature interactions and global bias. However, it is crucial for predicting app ratings to distinguish different categories. The propose CATA methods achieve the best performance because of the consideration of category-specific multi-view feature interactions. The application of Tucker decomposition effectively facilitates the performance as it permits the interactions within each modality [18] while the CP decomposition used in MFM does not.

Comparing the two datasets, it can be found that the superiority of the proposed approach is more significant for Google Play dataset. It is mainly caused by the fewer categories and feature views in the dataset of Apple’s App Store. The fewer categories might not sufficiently discriminate the important features in the specific category, while the fewer views of features would lead to the lack of some interaction information between features from different views. Moreover, the two categories in Apple App Store, *i.e.*, “Free” and “Paid”, do not have their own unique characteristics and are not easy to differentiate from each other. Nonetheless, even with very limited context information, the CATA methods still outperform the baselines in almost all cases.

D. Impact of Feature Views

In order to explore the impact of the feature views for the proposed CATA-G method, we conduct experiments based on different numbers of views. As each rating record in Apple App Store dataset only has two views, *i.e.*, user and review text, we use Google Play dataset for the experiments. Figure 6 shows the prediction performance of the CATA-G

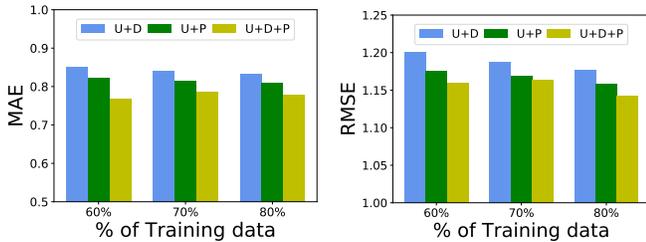


Fig. 6. Performance of the CATA-G method with different feature views on Google Play dataset.

method with two and three feature views. Note that U, D, and P respectively denote user, description text, and permission. It can be observed that the CATA-G method consistently performs best when incorporating three views of features, which benefits from the complementary information generated by the interactions among the various views of features. It indicates that the incorporation of multiple views of features can effectively improve the accuracy of rating prediction for apps. Consider the results produced by two views, we can find that the adoption of permission brings better results than that of description text. This is probably because the features extracted from description text are more sophisticated and higher dimensional, which may provide redundant information.

E. Category-Specific Performance

In this section, we further analyze the performance of the proposed method for each category based on Google Play dataset. Figure 7 shows the MAE and RMSE values of the proposed method and the top 2 baseline methods in each category. The category indexes on the x -axis is sorted by the numbers of rating records within the categories in an ascending order. We can find that MFM performs better than MVM in the categories with few training instances. It indicates that when few instances are available, the method incorporating the category information can improve the performance as it explores the information from other complementary information. The performance of the proposed CATA-G method is the worst with few instances, as CATA-G has more model parameters to learn and requires more instances. For the categories with more instances, CATA-G makes significant improvements and outperforms the other two methods. Compared with MFM, the superiority of CATA-G is the application of Tucker decomposition which can effectively retain the principal components of the weight tensor. Another interesting observation in Fig. 7 is that CATA-G makes the top 5 improvements for category *Simulation*, *Action*, *Casual*, *Arcade*, and *Puzzle* (*i.e.*, #11, #12, #15, #16, and #19). The apps in the five categories are game apps, and it means the features of them are more complicated as each game app has its specific theme setting. Therefore, the proposed method has a greater ability to discriminate the importance of each feature in a complicated feature sets.

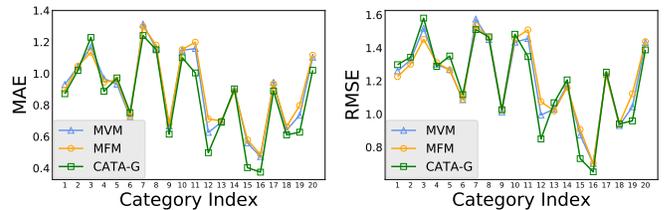


Fig. 7. Category-specific performance of MVM, MFM, and CATA-G on Google Play dataset. The category indexes on the x -axis is sorted by the numbers of rating records within the categories in an ascending order.

F. Sensitivity Analysis

There are two hyper-parameters (*i.e.*, α and β) in the proposed CATA approach. They are used to control the trade-off between the empirical loss and the prior knowledge encoded by the regularizations. To learn the impacts of the two hyper-parameters on the performance of app rating prediction, we run the proposed approach with different values for α and β on the two datasets. From Fig. 8, we can observe that the performance is stable for most pairs of the two hyper-parameters. For each dataset, the effects of the two hyper-parameters on MAE and RMSE are different. For the Google Play dataset, Figs. 8(a) and (b) show that the unstable and worse MAE and RMSE are produced when given a larger α (*i.e.*, $\alpha = 1$) or a smaller β (*i.e.*, in the range from 10^{-2} to 10). And the best performance is achieved by the relatively large value of β (*i.e.*, in the range from 10^2 to 10^5) with $\alpha = 0.1$. Figure 8 (c) and (d) report the results of the Apple App Store dataset, from which we can find that the performance is more stable than that of Google Play. Similarly, when the value of α is larger or the value of β is smaller, the MAE and RMSE are relatively higher. The best performance of the MAE is achieved by $\alpha = 10^{-3}$ while the RMSE is much lower when the value of α is set 10^{-2} . The value of β is in the range from 10^0 to 10^5 . For both datasets, the best performance is generated by a larger β and the larger β means the model hyper-parameters for the category-specific weight matrix \mathbf{D} can be small. It indicates that the part of full-order interactions among multiple categories and multiple views of features is much more important.

VI. RELATED WORK

To the best of our knowledge, this is the first work to consider mining the full-order interactions among app context information with tensor analysis to facilitate mobile app recommendation. From the conceptual perspective, two topics can be seen as closely related to this work: mobile app recommendation, tensor factorization and its applications. We give a short overview of these areas and distinguish from other existing methods.

Mobile App Recommendation has drawn an increasing number of attentions as an effective way to alleviate information overload in app market. Most of the existing works are trying to leverage one or more kinds of features to improve the recommendation performance. Yan *et al.* [19] developed the AppJoy system that recommends mobile apps by based

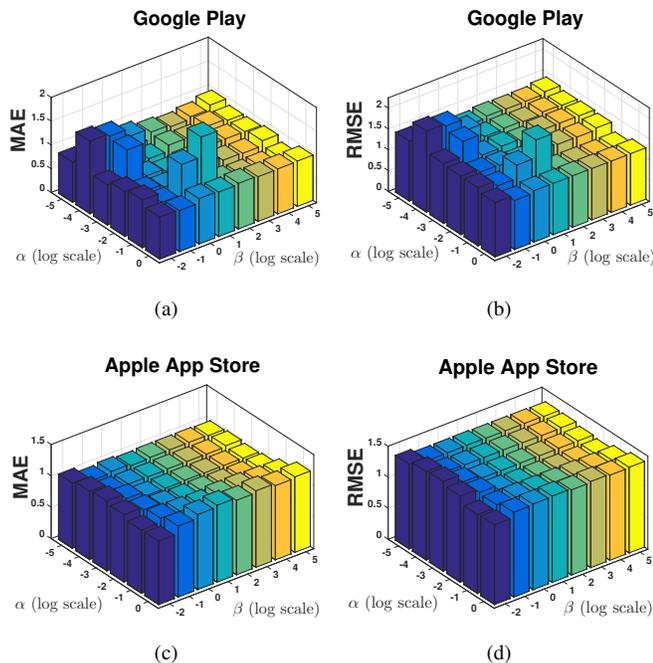


Fig. 8. Sensitivity analysis of hyper-parameters.

on the analysis of users' usage records. In [20], Yin *et al.* applies users' view/download sequences to mine the actual value and tempting value of apps, which are used to build a recommendation model considering the contest between apps. Features from the other sources are incorporated in some works. For instance, to address the cold-start problem, Lin *et al.* [21] proposed to apply app followers' features collected from Twitter to model the app and estimate which users may like the app. Zhu *et al.* [3] presented a method to evaluate the security risks of apps and proposed a flexible app recommendation approach combining both apps' popularity and users' security preferences through the modern portfolio theory. A recommendation model which can capture the trade-off between app functionality and user privacy preference was proposed in [2]. However, most of these works did not consider the complex interactions among the features of different views. In this work, we propose to model the interactions as a tensor structure and leverage tensor factorization to learn the latent relationships.

Tensor Factorization and Applications Tensor factorization is a method to divide a tensor in multidimensionality into many smaller parts. A comprehensive survey on tensor factorization can be found in [22]. Two well-known methods in this area are CANDECOMP/PARAFAC (CP) factorization and Tucker factorization. Both of them can be considered as higher-order generalization of Singular value decomposition (SVD) and Principle Component Analysis (PCA). These methods are used to decompose tensor data into simpler form, containing better features and intrinsic multi-way structures. CP factorization has been frequently investigated in the multi-view learning literature because of its simplicity. Specifically,

[23] first introduced to use the outer product operator to fuse multi-view features in the tensor structure and proposed a CP factorization based multi-view feature selection method. Later, [16] extended this approach to consider the full-order interactions between features, and proposed a CP factorization based multi-view machine (MVM) for multi-view prediction problems. Recently, [17] extended the MVM method to deal with multi-task multi-view prediction problems and proposed a CP factorization based multilinear factorization machine. However, to the best of our knowledge, none of the studies explored Tucker decomposition in the scenario of multi-view learning. Tucker decomposition is more general than the CP decomposition, and permits the interactions within each mode while the CP decomposition does not [18]. This paper gives an application of Tucker decomposition into multi-view learning task.

VII. CONCLUSIONS

In this paper, we propose a context-aware recommendation approach based on tensor analysis (CATA) for mobile apps. The proposed CATA approach models the interactions among the multiple categories and multiple views of features of apps as a tensor structure. CATA applies the Tucker decomposition to collectively learn the category-specific features and the latent relationships integrated in the full-order interactions without physically building the tensor. To further improve the performance of app recommendation, we present a group ℓ_1 norm regularization for the global category-specific weight matrix. Extensive experiments based on two real-world app datasets demonstrate the effectiveness of the proposed CATA approach.

REFERENCES

- [1] A. Karatzoglou, L. Baltrunas, K. Church, and M. Böhmer, "Climbing the app wall: enabling mobile app discovery through context-aware recommendations," in *CIKM*. ACM, 2012, pp. 2527–2530.
- [2] B. Liu, D. Kong, L. Cen, N. Z. Gong, H. Jin, and H. Xiong, "Personalized mobile app recommendation: Reconciling app functionality and user privacy preference," in *WSDM*. ACM, 2015, pp. 315–324.
- [3] H. Zhu, H. Xiong, Y. Ge, and E. Chen, "Mobile app recommendations with security and privacy awareness," in *SIGKDD*. ACM, 2014, pp. 951–960.
- [4] C.-T. Lu, S. Xie, W. Shao, L. He, and S. Y. Philip, "Item recommendation for emerging online businesses," in *IJCAI*, 2016, pp. 3797–3803.
- [5] W. Shao, L. He, and S. Y. Philip, "Clustering on multi-source incomplete data via tensor modeling and factorization," in *PAKDD*. Springer, 2015, pp. 485–497.
- [6] W. Shao, L. He, C.-T. Lu, X. Wei, and S. Y. Philip, "Online unsupervised multi-view feature selection," in *ICDM*. IEEE, 2016, pp. 1203–1208.
- [7] W. Shao, L. He, and S. Y. Philip, "Multiple incomplete views clustering via weighted nonnegative matrix factorization with $\ell_{\{2, 1\}}$ regularization," in *ECML/PKDD*. Springer, 2015, pp. 318–334.
- [8] H. Zhu, H. Xiong, Y. Ge, and E. Chen, "Discovery of ranking fraud for mobile apps," *IEEE Transactions on Knowledge and Data Engineering*, vol. 27, no. 1, pp. 74–87, 2015.
- [9] H. Zhu, C. Liu, Y. Ge, H. Xiong, and E. Chen, "Popularity modeling for mobile apps: A sequential approach," *IEEE Transactions on Cybernetics*, vol. 45, no. 7, pp. 1303–1314, 2015.
- [10] F. Ricci, L. Rokach, and B. Shapira, *Introduction to recommender systems handbook*. Springer, 2011.
- [11] S. Shalev-Shwartz and S. Ben-David, *Understanding machine learning: From theory to algorithms*. Cambridge University Press, 2014.

- [12] H. Wang, F. Nie, H. Huang, J. Yan, S. Kim, S. Risacher, A. Saykin, and L. Shen, "High-order multi-task feature learning to identify longitudinal phenotypic markers for alzheimer's disease progression prediction," in *NIPS*, 2012, pp. 1277–1285.
- [13] A. Gunawardana and G. Shani, "A survey of accuracy evaluation metrics of recommendation tasks," *Journal of Machine Learning Research*, vol. 10, no. Dec, pp. 2935–2962, 2009.
- [14] R. Salakhutdinov and A. Mnih, "Probabilistic matrix factorization." in *NIPS*, vol. 1, no. 1, 2007, pp. 2–1.
- [15] S. Rendle, "Factorization machines," in *ICDM*. IEEE, 2010, pp. 995–1000.
- [16] B. Cao, H. Zhou, G. Li, and P. S. Yu, "Multi-view machines," in *WSDM*. ACM, 2016, pp. 427–436.
- [17] C.-T. Lu, L. He, W. Shao, B. Cao, and P. S. Yu, "Multilinear factorization machines for multi-task multi-view learning," in *WSDM*. ACM, 2017, pp. 701–709.
- [18] A. Cichocki, M. Mørup, P. Smaragdis, W. Wang, and R. Zdunek, "Advances in nonnegative matrix and tensor factorization," *Computational Intelligence and Neuroscience: CIN*, 2008.
- [19] B. Yan and G. Chen, "Appjoy: personalized mobile application discovery," in *MobiSys*. ACM, 2011, pp. 113–126.
- [20] P. Yin, P. Luo, W.-C. Lee, and M. Wang, "App recommendation: a contest between satisfaction and temptation," in *WSDM*. ACM, 2013, pp. 395–404.
- [21] J. Lin, K. Sugiyama, M.-Y. Kan, and T.-S. Chua, "Addressing cold-start in app recommendation: latent user models constructed from twitter followers," in *SIGIR*. ACM, 2013, pp. 283–292.
- [22] T. G. Kolda and B. W. Bader, "Tensor decompositions and applications," *SIAM review*, vol. 51, no. 3, pp. 455–500, 2009.
- [23] B. Cao, L. He, X. Kong, P. S. Yu, Z. Hao, and A. B. Ragin, "Tensor-based multi-view feature selection with applications to brain diseases," in *ICDM*. IEEE, 2014, pp. 40–49.