

Semantic Interpretation of Top-N Recommendations

Vito Walter Anelli¹, Tommaso Di Noia¹, Eugenio Di Sciascio¹, Azzura Ragone, and Joseph Trotta

Abstract—Over the years, model-based approaches have shown their effectiveness in computing recommendation lists in different domains and settings. By relying on the computation of latent factors, they can recommend items with a very high level of accuracy. Unfortunately, when moving to the latent space, even if the model embeds content-based information, we miss references to the actual semantics of the recommended item. It makes the interpretation of the recommendation process non-trivial. In this paper, we show how to initialize latent factors in Factorization Machines by using semantic features coming from knowledge graphs to train an interpretable model, which is, in turn, able to provide recommendations with a high level of accuracy. In the presented approach, semantic features are injected into the learning process to retain the original informativeness of the items available in the dataset. By relying on the information encoded in the original knowledge graph, we also propose two metrics to evaluate the semantic accuracy and robustness of knowledge-aware interpretability. An extensive experimental evaluation on six different datasets shows the effectiveness of the interpretable model in terms of both accuracy and diversity of recommendation results and interpretability robustness.

Index Terms—Recommender systems, knowledge graph, interpretable models, factorization machines

1 INTRODUCTION

RESEARCH on transparency and interpretability of predictive models is gaining momentum since the research community is recognizing them as decisive elements in the next generation of recommendation algorithms. When equipped with the interpretability of recommendation results, a system ceases to be just a black-box (transparency) [1], [2], [3], and users are more willing to exploit the predictions extensively [4]. Transparency increases users' trust [5] (also exploiting specific semantic structures [6]), and satisfaction in using the system. For a recommender system, the user's trust is a vital aspect since it also leads to better performance [7]. In a nutshell, we may say that interpretations for recommendation results can be item-based, user-based, or feature-based. Item-based interpretations make use of the shared set of items among users [8]; User-based interpretations rely on sets of most similar users, like in [4]; while Feature-based interpretations may exploit features of recommended items as director, genre, and cast. Among interpretable models, we may distinguish between those based on Content-based Filtering (CBF) approaches and those based on Collaborative filtering (CF) ones. CBF algorithms provide recommendations by exploiting the available content and matching it with a user profile [9]. The use of content features makes the model interpretable. It is worth noticing that these features can be dramatically

different depending on the considered scenario: we may recommend a movie recommendation considering the director, actors, the producer, the genre, whereas we may interpret a book recommendation by the author, the book formats, or the saga. Sometimes, this prevents the straight adoption of a model in a specific knowledge domain. When content is missing, the recommender may rely only on the relationships between users (and the rates they provide to items), collaboratively [10]. Consequently, the interpretation of CF results inevitably reflects the approach adopted by the algorithm. For instance, an item-based and a user-based recommendation could be explained, respectively, as "other users who have experienced A have experienced B" or "similar users have experienced B".

Unfortunately, things change when we adopt more accurate Deep Learning-based [11], Recurrent [12], [13], [14], or model-based algorithms and techniques. Such approaches project items and users in a new vector space of latent features [15], making the final result not directly interpretable. Indeed, it is possible to compute items and users' similarities via latent factor exploitation, but we entirely lose any reference to the original user-item interaction and then to an explicit justification for the recommendations. In the last years, the industry and the research community have proposed many approaches that take advantage of side information to enhance the performance of latent factor models. Side information can refer to items as well as users [16] and can be either structured [17] or semi-structured [18], [19]. Interestingly, in [20] the authors argue about a new generation of knowledge-aware recommendation engines able to exploit information encoded in knowledge graph (\mathcal{KG}) to produce meaningful recommendations.

In this work, we propose a knowledge-aware Hybrid Factorization Machine (kaHFM) to train interpretable models in recommendation scenarios. kaHFM relies on Factorization Machines [21], and it extends them in different

-
- Vito Walter Anelli, Tommaso Di Noia, Eugenio Di Sciascio, and Joseph Trotta are with the Department of Electrical & Information Engineering Polytechnic University of Bari, Bari 70126, Italy. E-mail: {vitowalter.anelli, tommaso.dinoia, eugenio.disciascio, joseph.trotta}@poliba.it.
 - Azzura Ragone is with EY. E-mail: azzurra.ragone@gmail.com.

Manuscript received 1 Mar. 2019; revised 27 May 2020; accepted 13 July 2020.

Date of publication 17 July 2020; date of current version 1 Apr. 2022.

(Corresponding author: Tommaso Di Noia.)

Recommended for acceptance by D. Cai.

Digital Object Identifier no. 10.1109/TKDE.2020.3010215

key aspects making use of the semantic information encoded in a knowledge graph. Moreover, since it is unlikely that an inaccurate model will be adopted, it would be highly beneficial for users to develop accurate recommender systems, built upon an interpretable technique, so that the model can be, in case, exploited to generate explanations. In this sense, we show how kaHFM exploits data coming from knowledge graphs as side information to build a recommender system whose final results are accurate and, at the same time, semantically interpretable. With kaHFM, we build a model in which the meaning of each latent factor is bound to a semantic feature extracted from a knowledge graph. After the model training, we still have an explicit reference to the original semantics of the features describing the items, thus making possible the interpretation of the final results. Furthermore, we show that the explicit mapping of latent features to content-based ones makes it possible to exploit the characteristics of these latter to implement a more effective initialization technique.

We have evaluated kaHFM on six different publicly available datasets. Initially, we have extracted content-based features from data encoded in the DBpedia¹ knowledge graph. Then, we have analyzed the performance of the approach in terms of accuracy, diversity, and novelty of results by exploiting categorical, ontological, and factual features (see Section 2.1). For each of them, we have exploited public mappings to DBpedia. Finally, we have tested the robustness of kaHFM, showing that it ranks important features higher and can regenerate them if we remove them from the original dataset. With kaHFM, we address the following research questions:

- RQ1 Can we develop a model-based recommendation engine whose results are very accurate and, at the same time, interpretable regarding an explicitly stated semantics coming from a knowledge graph?
- RQ2 Can we evaluate whether kaHFM preserves the original semantics of items features after training?
- RQ3 How to measure with an offline evaluation that the proposed model can identify essential features by exploiting their explicit semantics?

We can summarize the main contributions of this paper as:

- Presentation of kaHFM: a framework that exploits data available in knowledge graphs to build semantically interpretable models for recommendation tasks;
- An experimental evaluation based on six different datasets to assess the performance of kaHFM in terms of accuracy, diversity, and novelty of computed results;
- Introduction of two metrics, Semantic Accuracy (SA@K) and robustness (n-Rob@K), to measure the interpretability of a knowledge-aware recommendation engine.

The remainder of the paper is structured as follows: in the next section, we introduce the background technologies behind kaHFM, and then we detail the overall approach in Section 3. The following section is devoted to the introduction of the two metrics we propose to assess the quality of

kaHFM results in terms of interpretability. In Section 5, we describe the experimental setting while in Section 6, we report on related work. Conclusion and future work close the paper.

2 BACKGROUND TECHNOLOGIES

In this section, we briefly recap the leading technologies we have adopted to develop kaHFM. First, we introduce Vector Space Models for recommender systems, and then we give a quick overview of knowledge graphs and their Linked (Open) Data implementation.

Content-based recommender systems rely on the assumption that it is possible to predict the future behavior of users based on their personalized profiles. We can build user profiles by exploiting the characteristics of the items they have liked in the past or some other available side information. Researchers have proposed several approaches that take advantage of side information in different ways: some of them consider tags [22], demographic data [23], or they extract information from collective knowledge bases [24]. Many of the most popular and adopted CBF approaches make use of a Vector Space Model (VSM).

In VSM, we represent users and items in the same vector space through Boolean or weighted vectors. Their respective positions and the distance (or better the proximity) between them provides a measure of how related (or similar) these two entities are. The vector space dimensions may refer to a user or item (or user-item) feature (e.g., the propensity to watch movies in the morning).

The choice of item features may substantially differ depending on their availability and application scenario. Crowd-sourced tags, categorical, ontological, or textual knowledge are just some of the most exploited ones. To sum up, in a CBF approach, we need (i) to get reliable items descriptions, (ii) to find a way to measure the strength of each feature for each item, (iii) to represent users, and finally (iv) to measure similarities.

Regarding the first step, nowadays, we can smoothly get item descriptions from the Web. In particular, thanks to the Linked Open Data initiative, much semantically structured knowledge is publicly available in the form of Linked Data datasets.

2.1 Knowledge Graphs and Linked Data

In 2012, Google has announced its Knowledge Graph² (\mathcal{KG}) as a new tool to improve the identification and retrieval of entities in return to a search query. Most of the knowledge encoded in the Google Knowledge Graph comes from Freebase, which has been a crowd-sourced effort to create a base of facts in every possible knowledge domain.

The original idea of Semantic Web [25] has changed over the years, making possible the creation of a full stack of semantic technologies. More remarkably, these technologies have given birth to the Linking Open Data initiative,³ where a community of researchers and practitioners have devoted

2. <https://googleblog.blogspot.it/2012/05/introducing-knowledge-graph-things-not.html>

3. <http://linkeddata.org>

1. <http://dbpedia.org>

an enormous effort to build publicly available knowledge bases of machine-understandable data.

We can represent a knowledge base that exploits Semantic Web technologies through a graph (knowledge graph) in which entities are linked to each other by binary relations. The Semantic Web community has developed several technologies [26] and languages to manage this kind of knowledge model, e.g., the Resource Description Framework (RDF).

It provides a simple graph-based data model to encode knowledge in a structured way through triples in the form $\langle \sigma, \rho, \omega \rangle$. In a knowledge graph, each triple represents the connection $\sigma \xrightarrow{\rho} \omega$ between two nodes, named *subject* (σ) and *object* (ω), through the *relation* (*predicate*) ρ . In an RDF knowledge graph, we usually find different types of encoded information [27]:

- *Factual*. It refers to statements such as *The Matrix was directed by the Wachowskis* or *Melbourne is located in Australia* that describe attributes of an entity;
- *Categorical*. It states something about the subject of an entity. In this direction, the categories of Wikipedia pages are an excellent example. Categories can be used to cluster entities and are often hierarchically organized thus making it possible to define them in a more generic or specific way;
- *Ontological*. It is a more restrictive and formal way to classify entities via a hierarchical structure of classes. Differently from categories, between sub-classes and super-classes, an IS-A (transitive) relation subsists.

If we take a look at the following RDF triples⁴ from the Dbpedia knowledge graph

```
dbr:The_Matrix dbo:director
dbr:The_Wachowski_Brothers.
dbr:The_Matrix dct:subject dbc:Dystopian_
```

films.

```
dbr:The_Matrix rdf:type dbo:Film.
```

we may see that each of them represents one of the types of data mentioned earlier. In the first triple, we state a fact about the movie *The Matrix* (represented by the corresponding URI `dbr:The_Matrix`) saying that it has been directed (`dbo:director`) by the *Wachowski Brothers* (`dbr:The_Wachowski_Brothers`). The second triple encodes categorical information through the predicate `dct:subject` about the same movie. In particular, here we say that it belongs to the category of dystopian films (`dbc:Dystopian_films`). Finally, with the last triple, we classify *The Matrix* as a Film (`dbo:Film`) thanks to the predicate `rdf:type`.

3 KNOWLEDGE-AWARE HYBRID FACTORIZATION MACHINES FOR TOP-N RECOMMENDATION

3.1 Formal Model

Factorization models have proven to be among the best performing approaches as collaborative filtering methods to build a recommender system [21], [28]. Among all the different factorization models, factorization machines propose

4. For the sake of conciseness, we use the CURIE syntax that abbreviates URIs using their namespaces. In this paper, we refer to namespaces as available at <http://prefix.cc>.

x^1	1	0	0	0	...	1	0	0	0	0	...
x^2	1	0	0	0	...	0	1	0	0	0	...
x^3	1	0	0	0	...	0	0	1	0	0	...
x^4	0	1	0	0	...	0	0	1	0	0	...
x^5	0	1	0	0	...	0	0	0	1	0	...
x^6	0	0	1	0	...	1	0	0	0	0	...
x^7	0	0	1	0	...	0	0	1	0	0	...
	U_1	U_2	U_3	U_4	...	I_1	I_2	I_3	I_4	I_5	...
	User					Item					

Fig. 1. A visual representation of X for sparse real valued vectors x^{ui} .

a unified general model to represent most of them. Here we report the definition and the formalization of a factorization model of order 2 for a recommendation problem involving only implicit ratings. Although, we could extend the model to a more expressive representation by taking into account, e.g., demographic and social information [29], multi-criteria [30], and even relations between contexts [31]. For each user $u \in U$ and each item $i \in I$ we build a binary vector $x^{ui} \in \mathbb{R}^{1 \times n}$, with $n = |U| + |I|$, representing the interaction between u and i in the original user-item rating matrix. In this modeling, x^{ui} contains only two 1 values corresponding to u and i while all the other values are set to 0 (see Fig. 1). Additionally, we indicate as m the overall number of interactions between users and items recorded on the platform (more precisely the fraction retained in the training set). Then, we denote with $X \in \mathbb{R}^{m \times n}$ the matrix containing as rows all possible x^{ui} we can generate starting from the original user-item rating matrix as shown in Fig. 1.

We can define the factorization machine (FM) for each vector x as:

$$\hat{y}(x^{ui}) = w_0 + \sum_{j=1}^n w_j \cdot x_j + \sum_{j=1}^n \sum_{j'=j+1}^n x_j \cdot x_{j'} \cdot \sum_{f=1}^k v_{(j,f)} \cdot v_{(j',f)}, \quad (1)$$

where the parameters are: w_0 representing the global bias; w_j giving the importance to every single x_j ; the pair $v_{(j,f)}$ and $v_{(j',f)}$ in $\sum_{f=1}^k v_{(j,f)} \cdot v_{(j',f)}$ measuring the strength of the interaction between each pair of variables x_j and $x_{j'}$. We represent the number of latent factors by k . Researchers usually select the k value at design time.

To make the recommendation results computed by kaHFM semantically interpretable, we want to inject the knowledge encoded within a knowledge-graph in a Factorization Machine. Given a set of features retrieved from a \mathcal{KG} [32], we first bind them to the latent factors. Then, since we address a *Top-N* recommendation problem, we train the model by using the Bayesian Personalized Ranking (BPR) criterion.

In [33], the authors have proposed to encode a Linked Data knowledge graph in a vector space model to develop a CBF recommender system. Given a set of items $I = \{i_1, i_2, \dots, i_N\}$ in a catalog, and their associated triples $\langle i, \rho, \omega \rangle$ in a knowledge graph \mathcal{KG} we may build the set of all possible features as $F = \{ \langle \rho, \omega \rangle \mid \langle i, \rho, \omega \rangle \in \mathcal{KG} \text{ with } i \in I \}$. We can represent each item as a vector of weights $\mathbf{i} = [v_{(i,1)}, \dots, v_{(i,(\rho,\omega))}, \dots, v_{(i,|F|)}]$ where $v_{(i,(\rho,\omega))}$ is the normalized

		- dbc:Space_adventure_films	- dbc:Films_set_in_the_future	- dbc:American_science_fiction_action_films	- dbc:1980s_science_fiction_films	- dbc:Paramount_Pictures_films	- dbc:Midlife_crisis_films	- dbc:American_sequel_films
v_1	0	0.88	0.81	0.7	0	0.60	0.53	...
v_2	1.3	1.12	0.91	0.84	0.65	0.59	0.58	...
v_3	0.5	0	0.71	0	0.28	0.35	0	...
v_4	0	0	0.31	0	0	0	0.6	...
v_5	0	0	0	0	0.18	0	0	...
v_6	0	0.12	0.22	0	0	0	0	...
v_7	1.23	1.03	0.89	0.85	0.56	0.3	0.61	...

Fig. 2. Example of real valued feature vectors for different items v_j .

TF-IDF value for $\langle \rho, \omega \rangle$:

$$v_{(i, \langle \rho, \omega \rangle)} = \frac{|\{\langle \rho, \omega \rangle \mid \langle i, \rho, \omega \rangle \in \mathcal{KG}\}|}{\sqrt{\sum_{\langle \rho, \omega \rangle \in F} |\{\langle \rho, \omega \rangle \mid \langle i, \rho, \omega \rangle \in \mathcal{KG}\}|^2}} \cdot \log \frac{|I|}{|\{j \mid \langle j, \rho, \omega \rangle \in \mathcal{KG} \text{ and } j \in I\}|} \quad (2)$$

Analogously, when we have a set U of users, we may represent users using the features describing the items they have enjoyed in the past. In the following, when no confusion arises, we use f to denote a feature $\langle \rho, \omega \rangle \in F$. Given a user u , we denote with I^u the set of the items enjoyed by u and we have $\mathbf{u} = [v_{(u,1)}, \dots, v_{(u,f)}, \dots, v_{(u,|F|)}]$ with

$$v_{(u,f)} = \frac{\sum_{i \in I^u} v_{(i,f)}}{|\{i \mid i \in I^u \text{ and } v_{(i,f)} \neq 0\}|}$$

Given the vectors \mathbf{u}_j , with $j \in [1 \dots |U|]$, and $\mathbf{i}_{j'}$, with $j' \in [1 \dots |I|]$, we build the matrix $\mathbf{V} \in \mathbb{R}^{n \times |F|}$ (see Fig. 2) where the first $|U|$ rows have a one to one mapping with \mathbf{u}_j while the last ones correspond to $\mathbf{i}_{j'}$. If we go back to Equation (1), we may see that, for each \mathbf{x} , the term $\sum_{j=1}^n \sum_{j'=j+1}^n x_j \cdot x_{j'} \cdot \sum_{f=1}^k v_{(j,f)} \cdot v_{(j',f)}$ is not zero only once, i.e., when both x_j and $x_{j'}$ are equal to 1. In the matrix depicted in Fig. 1, this happens when there is an interaction between a user and an item.

Moreover, the summation $\sum_{f=1}^k v_{(j,f)} \cdot v_{(j',f)}$ represents the dot product between two vectors \mathbf{v}_j and $\mathbf{v}_{j'}$ with a size equal to k . Hence, \mathbf{v}_j represents a latent representation of a user, $\mathbf{v}_{j'}$ that of an item in the same latent space, and we may evaluate the interaction through their dot product.

First, to inject the knowledge coming from \mathcal{KG} into kaHFM, we keep Equation (1) and we set $k = |F|$. In other words, we impose the number of latent factors equal to the number of features describing all the items in our catalog. We want to stress that we aim not to represent each feature through a latent vector, but to associate each factor to a semantic feature. Finally, the latent vectors will be composed only of semantic features. Hence, we initialize the parameters \mathbf{v}_j and $\mathbf{v}_{j'}$ with their corresponding rows from \mathbf{V} which in turn represent respectively \mathbf{u}_j and $\mathbf{i}_{j'}$. In this way, we try to identify each latent factor with a

TABLE 1
Top-10 Features Computed by kaHFM for the Movie "Star Trek II - The Wrath of Khan"

kaHFM	TF-IDF	Predicate	Object
1.3669	0.2584	dct:subject	dbc:Space_adventure_films
1.1252	0.2730	dct:subject	dbc:Films_set_in_the_future
0.9133	0.2355	dct:subject	dbc:American_science_fiction_action_films
0.8485	0.3190	dct:subject	dbc:1980s_science_fiction_films
0.6529	0.1549	dct:subject	dbc:Paramount_Pictures_films
0.5989	0.3468	dct:subject	dbc:Midlife_crisis_films
0.5940	0.1797	dct:subject	dbc:American_sequel_films
0.5862	0.2661	dct:subject	dbc:Film_scores_by_James_Horner
0.5634	0.2502	dct:subject	dbc:Films_shot_in_San_Francisco
0.5583	0.1999	dct:subject	dbc:1980s_action_thriller_films

TABLE 2
Top-10 Features Computed by kaHFM for the Movie "Star Trek - First Contact"

kaHFM	TF-IDF	Predicate	Object
1.2434	0.2858	dct:subject	dbc:Space_adventure_films
1.0355	0.3020	dct:subject	dbc:Films_set_in_the_future
0.8956	0.2605	dct:subject	dbc:American_science_fiction_action_films
0.8951	0.3451	dct:subject	dbc:Android_robot_films
0.7338	0.3105	dct:subject	dbc:Time_travel_films
0.6665	0.2701	dct:subject	dbc:Film_scores_by_Jerry_Goldsmith
0.6581	0.2205	dct:subject	dbc:1990s_action_films
0.6561	0.2279	dct:subject	dbc:1990s_science_fiction_films
0.6118	0.1988	dct:subject	dbc:American_sequel_films
0.5649	0.1713	dct:subject	dbc:Paramount_Pictures_films

corresponding explicit feature. The intuition is that after the training phase, the resulting matrix $\hat{\mathbf{V}}$ still refers to the original features. However, it contains better values for $v_{(j,f)}$ and $v_{(j',f)}$ that take into account also the latent interactions between users, items and features. It is noteworthy that after the training phase \mathbf{u}_j and $\mathbf{i}_{j'}$ (corresponding to $v_{(j,f)}$ and $v_{(j',f)}$ in \mathbf{V}) contain non-zero values also for features that are not originally in the description of the user u or of the item i .

In Tables 1 and 2 we show examples for values after the training (in the column kaHFM) together with the original TF-IDF ones computed for two movies from the Yahoo! Movies⁵ dataset.

3.2 Optimization

We can readily train Factorization Machines to reduce the prediction error via gradient descent methods, alternating least-squares (ALS), and MCMC. Since we have formulated our problem as a *Top-N* recommendation task, it is desirable to train kaHFM using a learning to rank approach like Bayesian Personalized Ranking Criterion (BPR) [34]. The BPR criterion optimizes using a stochastic gradient descent algorithm on a set D_S of triples (u, i, j) , with $i \in I^u$ and $j \notin I^u$, selected through random sampling from a uniform distribution. We may formulate BPR optimization criterion as:

$$\begin{aligned} \text{BPR} - \text{OPT} &= \sum_{(u,i,j) \in D_S} \ln \sigma(\hat{x}_{uij}) - \lambda_{\Theta} \|\Theta\|^2 \\ &= \sum_{(u,i,j) \in D_S} \ln \sigma(\hat{y}(\mathbf{x}^{ui}) - \hat{y}(\mathbf{x}^{uj})) - \lambda_{\Theta} \|\Theta\|^2. \end{aligned} \quad (3)$$

5. http://research.yahoo.com/Academic_Relations

In this formulation, $\sigma(\cdot)$ is a sigmoid function, and we may define the update step as:

$$\Theta \leftarrow \Theta + \alpha \left(\frac{e^{-\hat{x}_{uij}}}{1 + e^{-\hat{x}_{uij}}} \cdot \frac{\partial}{\partial \Theta} \hat{x}_{uij} + \lambda \Theta \right), \quad (4)$$

where α is the chosen learning rate. Since we are in an implicit feedback setting, we may assume that there is only an instance for the pair user-item. Hence, in our model, we can derive \hat{x}_{uij} as:

$$\hat{x}_{uij} = \hat{y}(\mathbf{x}^{ui}) - \hat{y}(\mathbf{x}^{uj}) = w_i x_i - w_j x_j + \sum_{f \in F} x_u x_i v_{(u,f)} v_{(i,f)} - x_u x_j v_{(u,f)} v_{(j,f)}. \quad (5)$$

For kaHFM, we can efficiently compute it by estimating the partial derivatives (to update the factorized parameters) for the only active entities involved in the transactions, w_i , w_j , v_u , v_i , and v_j :

$$\frac{\partial}{\partial \Theta} \hat{x}_{uij} = \begin{cases} 1, & \text{if } \theta = w_i, \\ -1, & \text{if } \theta = w_j, \\ v_{(u,f)}, & \text{if } \theta = v_{(i,f)}, \\ -v_{(u,f)}, & \text{if } \theta = v_{(j,f)}, \\ (v_{(i,f)} - v_{(j,f)}), & \text{if } \theta = v_{(u,f)}, \\ 0, & \text{otherwise} \end{cases} \quad (6)$$

Applying Equation (6) in Equation (4), we may iteratively update the model parameters to maximize the BPR-OPT criterion. The algorithm updates sequentially each sampled triple and continues the training until it reaches the provided number of iterations. Although we consider a high number of factors (features), the approach preserves the linear computational efficiency of BPR. Indeed, we perform each iteration in $\mathcal{O}(m)$ steps, where m is the overall number of transactions in the training set.

3.3 Personalized Recommendation

Once the training phase returns the optimal model parameters, the item recommendation step can take place. We extract the items vectors \mathbf{v}_j from the matrix \mathbf{V} , with the associated optimal values. Afterward, we use them to implement an Item-kNN recommendation approach.⁶ Hence, we measure similarities between each pair of items i and j by evaluating the cosine similarity of their corresponding vectors in \mathbf{V} :

$$cs(i, j) = \frac{\mathbf{v}_i \cdot \mathbf{v}_j}{\|\mathbf{v}_i\| \cdot \|\mathbf{v}_j\|}.$$

Given a set of neighbors for the item i , denoted as N^i , such that $i \notin I^u$ and a user u we may predict the score assigned by u to i as:

$$score(u, i) = \frac{\sum_{j \in N^i \cap I^u} cs(i, j)}{\sum_{j \in N^i} cs(i, j)}. \quad (7)$$

6. Please note that this choice leaves the item bias out of the final score. To provide a fair comparison with kNN methods, we have chosen to exploit a standard similarity metric, renouncing to it. Going beyond standard cosine similarity, we may compute a custom similarity by considering, e.g., the product between w_i and w_j : $\frac{w_i \cdot w_j}{\|w_i\| \|w_j\|} + \frac{w_i \cdot w_j}{\max(w_i, w_j)^2}$.

4 SEMANTIC ACCURACY AND GENERATIVE ROBUSTNESS OF kaHFM

The proposed approach lets us keep the explicit meaning of the “latent” factors computed via a factorization machine, thus making possible an interpretation of the recommended results. In the following, we propose an automated offline procedure to assess that kaHFM preserves the semantics of the features represented in \mathbf{V} after the training phase. We refer to the values computed by the procedure as *Semantic Accuracy*. A different but related aspect is that of evaluating whether kaHFM assigns a higher value to essential features. We refer to this behavior as *Robustness*. Interestingly, we may evaluate both *Semantic Accuracy* and *Robustness* in an offline setting.

4.1 Semantic Accuracy

The main idea behind Semantic Accuracy is to evaluate, given an item i , how well kaHFM can return its original features in the top-K list \mathbf{v}_i . In other words, given the set of features of i represented by $F^i = \{f_1^i, \dots, f_m^i, \dots, f_M^i\}$, with $F^i \subseteq F$, we check if the values in \mathbf{v}_i corresponding to $f_m^i \in F^i$ are higher than those corresponding to $f \notin F^i$. For the set of M features initially describing i , we see how many of them appear in the set $top(\mathbf{v}_i, M)$ representing the top- M features in \mathbf{v}_i . Then, we normalize this number by the size of F^i and average on all the items within the catalog I .

$$\text{Semantic Accuracy (SA@M)} = \frac{\sum_{i \in I} \frac{|top(\mathbf{v}_i, M) \cap F^i|}{|F^i|}}{|I|}.$$

In many practical scenarios, we may have $|F| \gg M$. Hence, we might also be interested in measuring the accuracy for different sizes of the top list. Since item descriptions may consist of a different number of features, the size of the top list could be a function of the original size of the description. Thus, we have measured SA@ nM with $n \in \{1, 2, 3, 4, 5, \dots\}$ and evaluate the number of features in F^i available in the $Top\text{-}n \cdot M$ elements of \mathbf{v}_i .

$$\text{SA@nM} = \frac{\sum_{i \in I} \frac{|top(\mathbf{v}_i, n \cdot M) \cap F^i|}{|F^i|}}{|I|}$$

4.2 Robustness

SA@ nM may be very useful to understand if kaHFM assigns weights according to the original description of item i . However, we still do not know if a high value in \mathbf{v}_i means that the corresponding feature is vital to define i . In other words, are we sure that kaHFM promotes meaningful features for i ?

To measure the “meaningfulness” for a given feature, we suppose that a particular feature $\langle \rho, \omega \rangle$ may be useful to describe an item i . Nevertheless, the corresponding triple, $\langle i, \rho, \omega \rangle$, is not represented in the knowledge graph. In case kaHFM was robust in generating weights for unknown features, it should discover the importance of that feature and modify its value to make it enter the $Top\text{-}K$ features in \mathbf{v}_i .

Starting from this observation, the idea to measure robustness is then to “forget” a triple involving i and check if kaHFM can generate it.

To implement such a process, we may proceed by following these steps:

- we train kaHFM thus obtaining optimal values v_i for all the features in F^i ;
- we identify the feature $f_{MAX}^i \in F^i$ with the highest value in v_i ;
- we train the model again initializing $f_{MAX}^i = 0$, and we compute v_i' .

Hence, if $f_{MAX}^i \in \text{top}(v_i', M)$, we may say that kaHFM shows high robustness in identifying important features.

Given a catalog I , we may then define the *Robustness for 1 removed feature @M* (1-Rob@M) as the number of items for which $f_{MAX}^i \in \text{top}(v_i', M)$ divided by the size of I .

$$1 - \text{Rob@M} = \frac{\sum_{i \in I} |\{i \mid f_{MAX}^i \in \text{top}(v_i', M)\}|}{|I|}.$$

Similarly to $SA@nM$, we may define 1-Rob@nM.

5 EXPERIMENTAL EVALUATION

Datasets. To provide an answer to RQ1, we have evaluated the performance of our method on six datasets belonging to three different domains (Music, Books, and Movies). The `Last.fm` dataset corresponds to user-artist plays on Last.fm online music system released during *HETRec 2011*⁷ *Workshop*. It contains social networking, tagging, artists, and music listening information from a set of 2,000 users. `LibraryThing` represents books' ratings collected in the `LibraryThing` website⁸ community. It contains social networking, tagging, and rating information on a [1..10] scale. `Yahoo!Movies` (*Yahoo! Webscope dataset ydata-ymovies-user-movie-ratings-content-v1_0*)⁹ contains movie ratings generated on Yahoo! Movies up to November 2003. It provides content, rating, and demographic information on a [1..5] scale, and mappings to `MovieLens` and `EachMovie` datasets. `Facebook Movies`, `Facebook Music`, and `Facebook Books` datasets have been released for the Linked Open Data challenge co-located with *ESWC 2015*,¹⁰ and they refer to movies, music, and book domains, respectively. Only implicit feedback is available for these datasets, but for each item, they have provided a link to `DBpedia`. To map items in `Last.fm` and `LibraryThing` to `DBpedia` resources, we have exploited a freely available mapping.¹¹ We have extracted all the updated items-features mappings (`Yahoo!Movies`, `LibraryThing`, `Last.fm`, `Facebook Movies`, `Facebook Music` and `Facebook Books`), and we have made them publicly available.¹² We have shown datasets statistics in Table 3.

Evaluation Protocol and Experimental Setting. We have chosen "All Unrated Items"¹³ [35] protocol to compare different algorithms. In All Unrated Items, we consider as candidate items all the items not yet rated by each user. We have split the dataset using Hold-Out 80-20 splitting strategy, retaining for every user the 80 percent of their ratings in the

TABLE 3
Datasets Statistics

Dataset	#Users	#Items	#Transactions	#Features	Sparsity
<code>Yahoo!Movies</code>	4,000	2,626	69,846	988,734	99.34%
<code>LibraryThing</code>	7,223	11,695	410,210	183,182	99.51%
<code>Last.fm</code>	1,375	8,289	60,701	434,817	99.47%
<code>Facebook Music</code>	52,068	5,749	1,374,994	345,452	99.54%
<code>Facebook Movies</code>	32,143	3,901	689,561	180,573	99.45%
<code>Facebook Books</code>	1,398	2,933	18,978	111,401	99.53%

training set and the remaining 20 percent in the test set. Moreover, whenever timestamps associated with every transaction was available, we have performed a temporal split [36], [37]. We have tested our approach against the most related content-based and collaborative algorithms in terms of Accuracy, Diversity, and Novelty [38]. We have compared kaHFM¹³ against a canonical 2-degree Factorization Machine (see Section 3.1) optimized via BPR (BPR-FM).

Moreover, since we exploit items similarity in the last step of our approach (see Equation (7)), we have compared kaHFM against an *Attribute Based Item-kNN* (ABItem-kNN) algorithm, that represents each item as a vector of weights, computed through a TF-IDF model. In ABItem-kNN, we have computed the attributes via Equation (2). For the sake of completeness we have compared kaHFM also against a pure Item-kNN, that is an item-based implementation of the k-nearest neighbors algorithm. Items in the neighborhood are then used to predict a score for each user-item pair. For all the methods that consider neighbors we have considered the neighbors in the range [10, ... 100]. Regarding BPR parameters, we have considered the *learning rate* varying in [0.005, ..., 0.5], the number of iterations in [10, ..., 100], and the number of factors in [20, ..., 100]. The *bias regularization*, *user regularization*, *positive item regularization*, and *negative item regularization* have been set respectively to 0, $\frac{1}{20}$, $\frac{1}{20}$, and $\frac{1}{200}$ of the learning rate, while we have adopted a sampler "without replacement" to sample the triples, as suggested by authors [34]. We have compared kaHFM also against the corresponding User-based nearest neighbor scheme (User-kNN), and `MostPopular`, a simple baseline that shows high performance in specific scenarios [39]. In our context, we have considered mandatory to also compare against a pure knowledge-graph content-based baseline based on the Vector Space Model (VSM) [33]. Finally, to compare our method against state-of-art algorithms, we have chosen VAE-CF and BPR-SSLIM. In VAE-CF [40], the authors train a non-linear probabilistic model taking advantage of Bayesian inference to estimate the parameters. BPR-SSLIM is the BPR optimized variant of *SSLIM* (specifically *cSSLIM* [41]), a well-known hybrid recommender system. It is based on a popular sparse linear model, *SLIM*, that efficiently learn a joint user-feature vector space exploiting users' implicit feedback.

5.1 Features Extraction

The feature extraction is one of the most sensitive steps in our approach. A wrong feature selection may result in noisy data or the lack of some crucial features. We may logically

13. Implementation available at <https://github.com/sisinflab/HybridFactorizationMachines/>

7. <http://ir.ii.uam.es/hetrec2011/>

8. <https://www.librarything.com/>

9. http://research.yahoo.com/Academic_Relations

10. <https://2015.eswc-conferences.org/program/semwebeval.html>

11. <https://github.com/sisinflab/LODrecsys-datasets>

12. <https://github.com/sisinflab/LinkedDatasets/>

TABLE 4
Filtering Thresholds

Dataset	Threshold
Yahoo!Movies	99.62
LibraryThing	99.91
Last.fm	99.88
Facebook Music	99.83
Facebook Movies	99.74
Facebook Books	99.66

split the preprocessing into three steps: (i) “*Extraction*”, in which we retrieve data from the DBpedia knowledge graph, (ii) “*Selection*” where we select the features involved in the specific experiment, and (iii) “*Filtering*” in which we remove uninformative features [32].

Extraction. Thanks to the publicly available mappings, all the items from the datasets come with a DBpedia link. Exploiting this reference, we retrieve all the $\langle \rho, \omega \rangle$ pairs. We have excluded some noisy features (based on the following predicates) already in this early extraction. In particular we have removed: owl:sameAs, dbo:thumbnail, prov:wasDerivedFrom, foaf:depiction, foaf:isPrimaryTopicOf. The rationale is that they give us only a little information about the nature of the entity in the specific knowledge base (e.g., the links between DBpedia and Wikipedia pages) or they are links to multimedia data or even external datasets. After this cleaning step, we have indexed all the features, saved them separately, and associated them with the item *id*.

Selection. We have performed our experiments with three different settings to analyze the impact of the different kinds of features on the recommendation accuracy and diversity. We have chosen the features as they are present in all the different domains and because of their factual, categorical, or ontological meaning:

- *Categorical Setting (CS):* We have selected only the features containing the property `dcterms:subject`.
- *Ontological Setting (OS):* In this case, the only features we have considered is `rdf:type`.
- *Factual Setting (FS):* We have considered all the features but those selected in OS.

Filtering. This last step corresponds to the removal of irrelevant features, that bring little value to the recommendation task, but, at the same time, pose scalability issues. We have performed the preprocessing phase following [32] and [42] with a single threshold. We have reported the corresponding thresholds (*tm* [32], or, equivalently *p* [42] for missing values) for each dataset in Table 4.

We have discarded features with more than *tm* missing values. For a fair comparison, we have used the same features for Attribute-based Item-kNN (ABItem-kNN) and kaHFM. We have depicted the characteristics of each dataset (varying the setting) in Table 5.

5.2 Accuracy, Diversity and Novelty With kaHFM

To evaluate our approach, we have measured accuracy, novelty, and aggregate diversity metrics. The considered accuracy metrics are Precision@N (*Prec@N*) and Normalized

TABLE 5
Considered Features in the Different Settings

Datasets	Categorical Setting		Ontological Setting		Factual Setting	
	Total	Selected	Total	Selected	Total	Selected
Yahoo!Movies	26155	747	38699	1240	950035	3186
LibraryThing	9443	1169	14585	1934	168597	5826
Last.fm	16422	1315	30734	3032	404083	9413
Facebook Music	15016	1057	27988	2531	317464	7881
Facebook Movies	8843	1103	13828	1848	166745	5427
Facebook Books	6231	263	9881	592	101520	1315

Discounted Cumulative Gain (*nDCG@N*). *nDCG@N* measures the usefulness of a suggested item, considering its relevance and its position in the recommendation list. Hence, we have computed it only for datasets that provide explicit ratings (i.e., LibraryThing and Yahoo!Movies). Since the recommended results may vary in length depending on the user, cumulative gain at each position is normalized across users. *EPC@N* (Expected Popularity Complement) is used to measure novelty, or more precisely, the ability of the algorithm to select items that belong to the long tail. Finally, we have measured diversity through *Catalog Coverage* (aggregate diversity in Top-N list, *AD@N*), *Gini Index* (*Gini@N*), and *Shannon entropy* (*SE@N*). In particular, *Catalog Coverage* denotes the ability of a system in selecting as many elements as possible from the whole catalog while *Gini index* and *Shannon entropy* are used to measure the distributional inequality with different approaches. Both accuracy and novelty metrics have been computed by averaging their values per-user. To compute those metrics we have used the implementation provided by the RankSys¹⁴ framework. We have performed the evaluation considering Top-10 ([39], [43], [44]) recommendations for all the datasets. When a rating score was available, we have adopted the *Threshold-based relevant items* condition [45] to take into account only relevant items. A relevance threshold of 4/5 and 8/10 has been set for Yahoo!Movies and LibraryThing, respectively. Finally, we have exploited *Prec@N*, or *nDCG@N* (where available), to select the best models.

Tables 6, and 7 show the results of our experiments regarding accuracy and diversity. In the tables, we highlight in **bold** the best result while we underline the second one. We have denoted statistically significant results (adopting a Student’s paired t-test with a 0.05 level) with a † mark.

LibraryThing experiments show that our approach outperforms the competing algorithms for all the considered metrics and settings. Instead, MostPopular, and then User-kNN show the worst results. If we observe BPR-FM aggregate diversity, we may notice relatively high values. This may suggest that BPR-FM is struggling to reach convergence.

Yahoo!Movies experiments show that in all three settings, our method is the most accurate, while it is the second one regarding diversity. Moreover, it shows similar values to the best performing one: ABItem-kNN. It is particularly noteworthy the VSM performance regarding accuracy. It behaves in three different ways depending on the considered setting. It is evidently due to the descriptions of items in the dataset, joint with the prediction formula. In this

14. <http://ranksys.org/>

sense, ABItem-kNN performance confirms that applying an item-based scheme to features vectors leads to better results.

In Last.fm, VAE-CF is the best performing method regarding the accuracy, followed by BPR-SSLIM. However, both show low values concerning diversity. On the other hand, kaHFM shows a good trade-off between accuracy and diversity.

Even in Facebook Movies, kaHFM shows to be the most precise algorithm for all the considered settings. Regarding diversity, it is only the second-best after ABItem-kNN. However, kaHFM almost doubles the accuracy of ABItem-kNN, preserving high diversity values. On the other hand, BPR-SSLIM shows good accuracy results, but it provides less diversified and tailored recommendations.

Finally, Facebook Music and Facebook Books show a similar trend. kaHFM shows the best accuracy performance in all the experiments but the Facebook Music Factual Setting, in which it is still very close to VSM. Furthermore, for the two datasets, ABItem-kNN and kaHFM show the best diversity values, providing more personalized recommendation lists. More, we may notice that VAE-CF and BPR-SSLIM show excellent accuracy performance in Facebook Books, but even there, they provide less diversified recommendations.

Let us discuss the baselines more related to our approach. We have compared kaHFM against ABItem-kNN to check whether the collaborative trained features may lead to better similarity values. This hypothesis seems to be confirmed since kaHFM always beats ABItem-kNN regarding accuracy. Moreover, we want to check if a knowledge-graph-based initialization of latent factors may improve the performance of Factorization Machines. kaHFM beats BPR-FM 16 times over 18. This behavior may happen since the random initialization takes a while to drive the Factorization Machine to reach competitive performance. However, we have specifically investigated this aspect in Section 5.3. Lastly, we are interested in checking whether collaborative trained features lead to better accuracy results than a purely informativeness-based Vector Space Model (in detail its knowledge-aware variant). It seems to be confirmed through experiments since kaHFM beats VSM 17 times over 18. Although we have measured very interesting results, it is worth mentioning that accuracy and interpretability come at a cost. Indeed, public DBpedia mappings are not always available. In this sense, the designer should ask either domain experts to describe the items semantically or provide a mapping. Moreover, the method could not recommend not-mapped items, thus forcing to work on dataset linking (toward DBpedia or other public datasets). Finally, the quality itself of the knowledge base may have an impact. Nevertheless, this aspect deserves a specific investigation that is outside the scope of this work.

5.3 Evolution of kaHFM Performance Over Iterations

In this experiment, we try to assess some crucial aspects. First, we want to know whether the knowledge-graph-based initialization is generally better than standard random initialization. Second, we desire to investigate if kaHFM can reach convergence sooner than BPR-FM. Third, we want to assess whether kaHFM training generally improves the original representation of features (for a Top-N recommendation task). For the experimental evaluation, we have analyzed

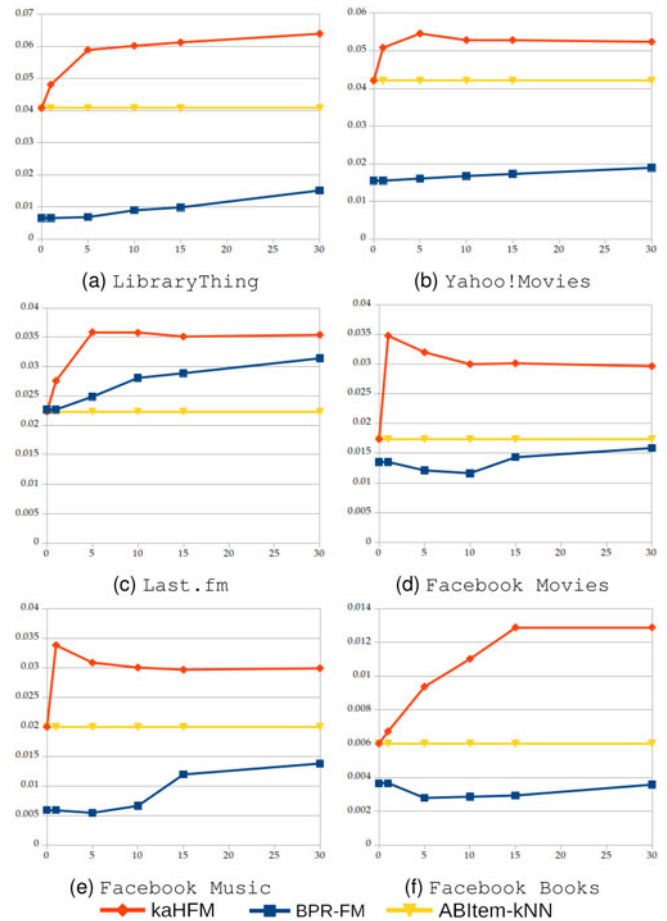


Fig. 3. Precision@10 varying # iterations 0, 1, 5, 10, 15, 30.

kaHFM, BPR-FM, and ABItem-kNN performance. Since we need to compare them fairly, we have considered the same underlying model. In detail, we have used the same parameters adopted for kaHFM and the same number of hidden factors (see the “Selected” column in Table 5). Regarding BPR parameters, we have set *learning rate*, *bias regularization*, *user regularization*, *positive item regularization*, and *negative item regularization* to 0.05, 0, 0.0025, 0.0025, and 0.00025, respectively. For the sake of reproducibility, the BPR parameters correspond to mymedialite¹⁵ implementation.

Then, we have produced recommendations for the Categorical Setting of the six datasets considering 0,1,5,10,15,30 iterations. Results considering the other settings are available online.¹⁶ In Fig. 3, the plots show the evolution of Precision@10 for the considered iterations. First, we may focus on the difference between BPR-FM and ABItem-kNN. In almost all cases, ABItem-kNN Precision is higher than the BPR-FM starting value. It denotes that generally, the semantic features initialized with TF-IDF perform better than random. Of course, this cannot be a conclusive remark, because BPR associates an implicit meaning to the single factors during training. We may only affirm that a TF-IDF initialization puts features weights closer to a local minimum. Then, we may compare the behavior of kaHFM and

15. <http://www.mymedialite.net/>

16. <https://github.com/sisinflab/papers-results/tree/master/kaHfm-results/>

TABLE 8
Semantics Accuracy Results for Different Values of M. F.A.
Denotes the Feature Average Number per Item

Semantics Accuracy	@M	@2M	@3M	@4M	@5M	F.A.
Yahoo!Movies	0.847	0.863	0.865	0.868	0.873	12.143
LibraryThing	0.960	0.996	0.998	0.999	0.999	3.820
Last.fm	0.960	0.987	0.991	0.994	0.995	6.615
Facebook Music	0.892	0.948	0.962	0.970	0.974	7.113
Facebook Movies	0.864	0.883	0.889	0.894	0.899	12.856
Facebook Books	0.995	1	1	1	1	3.133

BPR-FM. The six pairs of curves show that kaHFM moves faster toward convergence. Finally, we may compare kaHFM and ABItem-kNN. All the plots seem to suggest that our approach improves the original item representation. However, it only implies that it leads to a Precision improvement, without any considerations about the preservation of the semantics.

Given the obtained results we may say that the answer to RQ1 is positive when adopting kaHFM.

5.4 Semantic Accuracy

The previous experiments have shown the effectiveness of kaHFM in terms of accuracy, diversity, and novelty. In practical terms, we have shown that: (i) content initialization generally leads to better performance with kaHFM, (ii) the obtained items vectors are better fine-tuned than the original ones for a *Top-N* item recommendation task, (iii) results may depend on the kind of features we extract from the Knowledge Graph. However, we still do not know whether kaHFM preserves the original semantics of the features after the training (as we want to assess by posing RQ2). In Section 4.1, we have introduced Semantics Accuracy ($SA@nM$) as a metric to automatically check if item feature values reflect the actual meaning of that feature. Thus, we have measured $SA@nM$ with $n \in \{1, 2, 3, 4, 5\}$ and $M = 10$, and we have evaluated the number of ground features available in the *Top-nM* elements of \mathbf{v}_i for each of the six datasets.

Table 8 shows the results for the Categorical setting of all the different datasets. In general, the results we obtain are noteworthy. As an example, we may examine the worst one to describe better the meaning of the values we obtain. In Yahoo!Movies Categorical setting, 747 different features compose each item vector (see Table 5). After the training phase, on average, more than ten (equal to $0.847 \cdot 12.143$) over twelve features (last column in Table 8) correspond to the original features. Consequently, kaHFM has computed almost the same best features starting from hundreds of them. Even then, the obtained results may provide an affirmative answer to RQ2.

5.5 Generative Robustness

In Section 4.2, we have introduced a procedure to measure the capability of kaHFM to compute meaningful features. Here, we have computed $1\text{-Rob}@nM$ for the six adopted datasets, and we have depicted the results in Table 9.

Even here, we focus on the CS setting. For a better understanding, we may start by focusing on Yahoo!Movies,

TABLE 9
1-Robustness for Different Values of M. Column F.A. Denotes
the Feature Average Number per Item

1-Robustness	@M	@2M	@3M	@4M	@5M	F.A.
Yahoo!Movies	0.487	0.645	0.713	0.756	0.793	12.143
LibraryThing	0.275	0.481	0.554	0.597	0.632	3.820
Last.fm	0.125	0.281	0.346	0.394	0.430	6.615
Facebook Music	0.714	0.893	0.935	0.955	0.966	7.113
Facebook Movies	0.821	0.945	0.970	0.980	0.984	12.856
Facebook Books	0.315	0.516	0.605	0.682	0.745	3.133

which seemingly shows a bad behavior. As mentioned earlier, Table 8 shows that kaHFM has guessed ten on twelve different features for Yahoo!Movies. In this experiment, we eliminate a feature, thus making kaHFM to guess an average of nine over twelve features. What we assess is whether kaHFM can guess the removed feature in the remaining three slots. Results in Table 9 show that kaHFM puts that single removed feature in one of the three slots, the 48.7 percent of the times choosing among 747 overall features. We believe the example may help to appreciate even more the results on Facebook Music and Facebook Movies. For the remaining datasets, Table 8 shows that there are no free slots. Thus, after removing a feature, kaHFM has only one missing slot to fill with the right feature. Let us focus on Facebook Books. It contains 263 different features in the item vector (see Table 5) and a low average number of features per item in the \mathcal{KG} (3.133). Results show that kaHFM fills correctly the missing slot 31 percent of the cases. Overall, results show that kaHFM promotes essential features, as asked by RQ3.

6 RELATED WORK

The core of our model is a Factorization Machines (FM) model [21]. Nowadays FMs are the most widely used factorization models because they offer a number of advantages w.r.t. other latent factors models such as SVD++ [46], PITF [47], FPMC [48]. First, FMs are designed for a generic prediction task while the others tackle specific tasks. Moreover, a FM with n features is equivalent to a linear model with $n + n \cdot (n - 1) / 2$ features and parameters can be estimated accurately even in high data sparsity scenarios. Nevertheless, several improvements have been proposed for FMs. For instance Neural Factorization Machines [49] have been developed to fix the inability of classical FMs to capture non linear structure of real-world data. Furthermore, Attentional Factorization Machines [50] have been proposed that use an attention network to learn the importance of feature interactions. The factorization models have also been adopted to feed Active Learning-based Recommender Systems [51]. In detail, the authors propose the original idea of Dynamic Active Learning Budget to distribute the limited active learning resources. Usually only top recommended items are provided to the user. For this reason, ranking has become a much more important task than rating prediction [52]. This has led to *Learning to Rank* algorithms that can be further categorized in Point-wise [53], Pair-wise [34] and List-wise [44]. In particular, Pair-wise approaches are usually considered as a good trade-off between ordering

performances and computational complexity. Among this class of algorithms, Bayesian Personalized Ranking (BPR) [34] is one of the most widely adopted. It exploits stochastic gradient descent algorithm to learn the relative order between positive and negative items (see Section 3.2). BPR can be applied to Matrix Factorization and Factorization Machines (as in our work and in [54]). One of the first attempts to overcome the interpretability problem for matrix factorization is the Explicit Factor Model (EFM) [18], [55]. Products' features and users' opinions are extracted with phrase-level sentiment analysis from users' reviews to feed a matrix factorization framework.

After that, we have observed a few improvements to EFM to deal with temporal dynamics [56] and to use tensor factorization [19]. They also adopt pair-wise learning to rank approach (in their case BPR on a three-fold user-item-feature space). A further advance in MF-based explainable recommendation models is Explainable Matrix Factorization (EMF) [57] in which they exploit a neighborhood model to generate explanations. Similarly, in [8], the authors propose an explainable Restricted Boltzmann Machine model. It learns a network model (with an additional visible layer) that takes into account a degree of explainability. In [58], the authors compute recommendations by generating and ranking personalized explanations in the form of explanation chains. OCuLaR [59] provides interpretable recommendations from positive examples based on the detection of co-clusters between users (clients) and items (products). In [60], the authors propose a Multi-Level Attraction Model (MLAM) in which they build two attraction models for cast and story. It provides the interpretability of the model in terms of attractiveness of the Sentence level, Word level, and Cast member. In [61], the authors train a matrix factorization model to complete the $U \times I$ matrix. They then use the complete (approximated) rating matrix to compute a set of association rules that explain the obtained recommendations. Among the works that exploit content information to produce explainable recommendations, Tagsplanations [62] is worth to mention. Community tags feed it, and it exploits a relevance measure to weight tags considering items, and user preferences. Furthermore, also demographic-based recommendation explanations have been inspected [23], to recommend items for specific types (age, location, gender) of users.

Among model-based recommender systems, [63] proposes to reconstruct the user-item matrix using a non-convex function to approximate the rank in the *Top-N* recommendation setting. A particularly interesting work is [64], where authors have proposed to overpass the classic limitations of similarity-based methods, taking advantage of kernels. It is worth noticing that the techniques in the latter may drive an extension of kaHFM.

7 CONCLUSION AND FUTURE WORK

In this work, we have proposed an interpretable method, kaHFM, in which we bind the meaning of latent factors for a Factorization machine to data coming from a knowledge graph. We have evaluated kaHFM on six different publicly available datasets, and we have compared it against state-of-the-art algorithms.

Results show that our approach generally outperforms the other approaches concerning accuracy, diversity, and novelty, considering different sets of semantics-aware features. We have shown that the generated recommendation lists are more precise and personalized, and they select more items from the long tail. In particular, we have considered Ontological, Categorical, and Factual information extracted from a freely available knowledge graph. Summing up, the experimental evaluation shows that: (RQ1) the learned model shows competitive performance regarding the accuracy, novelty, and diversity and, at the same time, is effectively interpretable; (RQ2) the estimated features show to be semantically meaningful; (RQ3) the model is robust to re-generate essential features after removing them.

In the future, we are interested in testing kaHFM in different scenarios, other than recommender systems. Moreover, we are already working on some model improvements. In detail, we are interested in considering other metrics for relevance that may precisely fit particular scenarios. In fact, this is possible since the method itself is agnostic to the specific adopted measure. Furthermore, it would be useful to exploit kaHFM to provide suggestions to knowledge graphs maintainers for adding relevant missing features to the knowledge base. In this sense, we would like to evaluate our approach in knowledge graph completion tasks.

ACKNOWLEDGMENTS

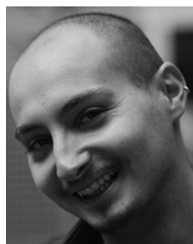
The authors with Polytechnic University of Bari acknowledge partial support of the following projects: Innonetwork CONTACT, PON ARS01_00821 FLET4.0, FINCONS Smart Digital Solutions for the Creative Industry, PON OK-INSAD.

REFERENCES

- [1] R. R. Sinha and K. Swearingen, "The role of transparency in recommender systems," in *Proc. Extended Abstr. Conf. Hum. Factors Comput. Syst.*, 2002, pp. 830–831. [Online]. Available: <http://doi.acm.org/10.1145/506443.506619>
- [2] N. Tintarev and J. Masthoff, "Designing and evaluating explanations for recommender systems," in *Recommender Systems Handbook*, Berlin, Germany: Springer, 2011, pp. 479–510. [Online]. Available: https://doi.org/10.1007/978-0-387-85820-3_15
- [3] M. Zanker, "The influence of knowledgeable explanations on users' perception of a recommender system," in *Proc. 6th ACM Conf. Recommender Syst.*, 2012, pp. 269–272. [Online]. Available: <http://doi.acm.org/10.1145/2365952.2366011>
- [4] J. L. Herlocker, J. A. Konstan, and J. Riedl, "Explaining collaborative filtering recommendations," in *Proc. ACM Conf. Comput. Supported Cooperative Work*, 2000, pp. 241–250. [Online]. Available: <http://doi.acm.org/10.1145/358916.358995>
- [5] R. Falcone, A. Sapienza, and C. Castelfranchi, "The relevance of categories for trusting information sources," *ACM Trans. Internet Techn.*, vol. 15, no. 4, pp. 13:1–13:21, 2015. [Online]. Available: <http://doi.acm.org/10.1145/2803175>
- [6] N. Drawel, H. Qu, J. Bentahar, and E. Shakhshuki, "Specification and automatic verification of trust-based multi-agent systems," *Future Gener. Comput. Syst.*, vol. 107, pp. 1047–1060, 2020. [Online]. Available: <https://doi.org/10.1016/j.future.2018.01.040>
- [7] X. Li, H. Fang, Q. Yang, and J. Zhang, "Who is your best friend?: Ranking social network friends according to trust relationship," in *Proc. 26th Conf. User Model. Adaptation Personalization*, 2018, pp. 301–309. [Online]. Available: <http://doi.acm.org/10.1145/3209219.3209243>
- [8] B. Abdollahi and O. Nasraoui, "Explainable restricted boltzmann machines for collaborative filtering," *CoRR*, vol. abs/1606.07129, 2016. [Online]. Available: <http://arxiv.org/abs/1606.07129>

- [9] H. S. M. Cramer *et al.*, "The effects of transparency on trust in and acceptance of a content-based art recommender," *User Model. User-Adapted Interact.*, vol. 18, no. 5, pp. 455–496, 2008. [Online]. Available: <https://doi.org/10.1007/s11257-008-9051-3>
- [10] V. W. Anelli, T. D. Noia, E. D. Sciascio, A. Ragone, and J. Trotta, "The importance of being dissimilar in recommendation," in *Proc. 34th ACM/SIGAPP Symp. Appl. Comput.*, 2019, pp. 816–821. [Online]. Available: <https://doi.org/10.1145/3297280.3297360>
- [11] S. Chakraborty *et al.*, "Interpretability of deep learning models: A survey of results," in *Proc. IEEE SmartWorld Ubiquitous Intell. Comput. Adv. Trusted Computed Scalable Comput. Commun. Cloud Big Data Comput. Internet People Smart City Innov.*, 2017, pp. 1–6. [Online]. Available: <https://doi.org/10.1109/UIC-ATC.2017.8397411>
- [12] Y. Zhu *et al.*, "Query-based interactive recommendation by meta-path and adapted attention-gru," in *Proc. 28th ACM Int. Conf. Inf. Knowl. Manage.*, 2019, pp. 2585–2593. [Online]. Available: <https://doi.org/10.1145/3357384.3357805>
- [13] Y. Zhu *et al.*, "A brand-level ranking system with the customized attention-gru model," in *Proc. 27th Int. Joint Conf. Artif. Intell.*, 2018, pp. 3947–3953. [Online]. Available: <https://doi.org/10.24963/ijcai.2018/549>
- [14] Y. Zhu *et al.*, "What to do next: Modeling user behaviors by time-LSTM," in *Proc. 26th Int. Joint Conf. Artif. Intell.*, 2017, pp. 3602–3608. [Online]. Available: <https://doi.org/10.24963/ijcai.2017/504>
- [15] Y. Koren, R. M. Bell, and C. Volinsky, "Matrix factorization techniques for recommender systems," *IEEE Comput.*, vol. 42, no. 8, pp. 30–37, 2009. [Online]. Available: <https://doi.org/10.1109/MC.2009.263>
- [16] X. Wang, X. He, F. Feng, L. Nie, and T. Chua, "TEM: tree-enhanced embedding model for explainable recommendation," in *Proc. World Wide Web Conf.*, 2018, pp. 1543–1552. [Online]. Available: <http://doi.acm.org/10.1145/3178876.3186066>
- [17] Z. Sun, J. Yang, J. Zhang, A. Bozzon, L. Huang, and C. Xu, "Recurrent knowledge graph embedding for effective recommendation," in *Proc. 12th ACM Conf. Recommender Syst.*, 2018, pp. 297–305. [Online]. Available: <http://doi.acm.org/10.1145/3240323.3240361>
- [18] Y. Zhang, G. Lai, M. Zhang, Y. Zhang, Y. Liu, and S. Ma, "Explicit factor models for explainable recommendation based on phrase-level sentiment analysis," in *Proc. 37th Int. ACM SIGIR Conf. Res. Develop. Inf. Retrieval*, 2014, pp. 83–92. [Online]. Available: <http://doi.acm.org/10.1145/2600428.2609579>
- [19] X. Chen, Z. Qin, Y. Zhang, and T. Xu, "Learning to rank features for recommendation over multiple categories," in *Proc. 39th Int. ACM SIGIR Conf. Res. Develop. Inf. Retrieval*, 2016, pp. 305–314. [Online]. Available: <http://doi.acm.org/10.1145/2911451.2911549>
- [20] Y. Zhang and X. Chen, "Explainable recommendation: A survey and new perspectives," *CoRR*, vol. abs/1804.11192, 2018. [Online]. Available: <http://arxiv.org/abs/1804.11192>
- [21] S. Rendle, "Factorization machines," in *Proc. 10th IEEE Int. Conf. Data Mining*, 2010, pp. 995–1000. [Online]. Available: <https://doi.org/10.1109/ICDM.2010.127>
- [22] Y. Zhu, Z. Guan, S. Tan, H. Liu, D. Cai, and X. He, "Heterogeneous hypergraph embedding for document recommendation," *Neurocomputing*, vol. 216, pp. 150–162, 2016. [Online]. Available: <https://doi.org/10.1016/j.neucom.2016.07.030>
- [23] W. X. Zhao, S. Li, Y. He, L. Wang, J. Wen, and X. Li, "Exploring demographic information in social media for product recommendation," *Knowl. Inf. Syst.*, vol. 49, no. 1, pp. 61–89, 2016. [Online]. Available: <https://doi.org/10.1007/s10115-015-0897-5>
- [24] T. di Noia, R. Mirizzi, V. C. Ostuni, and D. Romito, "Exploiting the web of data in model-based recommender systems," in *Proc. 6th ACM Conf. Recommender Syst.*, 2012, pp. 253–256. [Online]. Available: <http://doi.acm.org/10.1145/2365952.2366007>
- [25] T. Berners-Lee, J. Hendler, and O. Lassila, "The Semantic Web," *Sci. Amer.*, vol. 284, no. 5, pp. 34–43, May 2001. [Online]. Available: <http://www.sciam.com/article.cfm?articleID=00048144-10D2-1C70-84A9809EC%588EF21>
- [26] V. W. Anelli, A. Cali, T. D. Noia, M. Palmonari, and A. Ragone, "Exposing open street map in the linked data cloud," in *Proc. 29th Int. Conf. Ind. Eng. Appl. Appl. Intell. Syst.*, 2016, pp. 344–355. [Online]. Available: https://doi.org/10.1007/978-3-319-42007-3_29
- [27] V. W. Anelli, T. D. Noia, E. D. Sciascio, A. Ragone, and J. Trotta, "How to make latent factors interpretable by feeding factorization machines with knowledge graphs," in *Proc. 18th Int. Semantic Web Conf.*, 2019, pp. 38–56. [Online]. Available: https://doi.org/10.1007/978-3-030-30793-6_3
- [28] V. W. Anelli, T. D. Noia, E. D. Sciascio, C. Pomo, and A. Ragone, "On the discriminative power of hyper-parameters in cross-validation and how to choose them," in *Proc. 13th ACM Conf. Recommender Syst.*, 2019, pp. 447–451. [Online]. Available: <https://doi.org/10.1145/3298689.3347010>
- [29] G. Adomavicius and Y. Kwon, "Context-aware recommender systems," in *Recommender Systems Handbook*. Berlin, Germany: Springer, 2015, pp. 191–226. [Online]. Available: https://doi.org/10.1007/978-1-4899-7637-6_6
- [30] G. Adomavicius and Y. Kwon, "Multi-criteria recommender systems," in *Recommender Systems Handbook*. Berlin, Germany: Springer, 2015, pp. 847–880. [Online]. Available: https://doi.org/10.1007/978-1-4899-7637-6_25
- [31] Y. Zheng, B. Mobasher, and R. D. Burke, "Incorporating context correlation into context-aware matrix factorization," in *Proc. IJCAI Joint Workshop CPR+ITWP, 24th Int. Joint Conf. Artif. Intell.*, 2015, pp. 21–27. [Online]. Available: <http://ceur-ws.org/Vol-1440/Paper5.pdf>
- [32] T. di Noia, C. Magarelli, A. Maurino, M. Palmonari, and A. Rula, "Using ontology-based data summarization to develop semantics-aware recommender systems," in *Proc. 15th Int. Conf. ESWC*, 2018, pp. 128–144. [Online]. Available: https://doi.org/10.1007/978-3-319-93417-4_9
- [33] T. di Noia, R. Mirizzi, V. C. Ostuni, D. Romito, and M. Zanker, "Linked open data to support content-based recommender systems," in *Proc. 8th Int. Conf. Semantic Syst.*, 2012, pp. 1–8. [Online]. Available: <http://doi.acm.org/10.1145/2362499.2362501>
- [34] S. Rendle, C. Freudenthaler, Z. Gantner, and L. Schmidt-Thieme, "BPR: Bayesian personalized ranking from implicit feedback," in *Proc. 25th Conf. Uncertainty Artif. Intell.*, 2009, pp. 452–461. [Online]. Available: https://dlsppit.org/uai/displayArticleDetails.jsp?mmnu=1&smnu=2&article%_id=1630&proceeding_id=25
- [35] H. Steck, "Evaluation of recommendations: Rating-prediction and ranking," in *Proc. 7th ACM Conf. Recommender Syst.*, 2013, pp. 213–220.
- [36] A. Gunawardana and G. Shani, "Evaluating recommender systems," in Berlin, Germany: Springer, *Recommender Systems Handbook*, 2015, pp. 265–308. [Online]. Available: https://doi.org/10.1007/978-1-4899-7637-6_8
- [37] V. W. Anelli, T. D. Noia, E. D. Sciascio, A. Ragone, and J. Trotta, "Local popularity and time in top-n recommendation," in *Proc. 41st Eur. Conf. IR Res.*, 2019, pp. 861–868. [Online]. Available: https://doi.org/10.1007/978-3-030-15712-8_63
- [38] V. W. Anelli, V. Bellini, T. D. Noia, W. L. Bruna, P. Tomeo, and E. D. Sciascio, "An analysis on time- and session-aware diversification in recommender systems," in *Proc. 25th Conf. User Model. Adaptation Personalization*, 2017, pp. 270–274. [Online]. Available: <https://doi.org/10.1145/3079628.3079703>
- [39] P. Cremonesi, Y. Koren, and R. Turrin, "Performance of recommender algorithms on top-n recommendation tasks," in *Proc. 4th ACM Conf. Recommender Syst.*, 2010, pp. 39–46. [Online]. Available: <http://doi.acm.org/10.1145/1864708.1864721>
- [40] D. Liang, R. G. Krishnan, M. D. Hoffman, and T. Jebara, "Variational autoencoders for collaborative filtering," in *Proc. World Wide Web Conf.*, 2018, pp. 689–698. [Online]. Available: <https://doi.org/10.1145/3178876.3186150>
- [41] X. Ning and G. Karypis, "Sparse linear methods with side information for top-n recommendations," in *Proc. 6th ACM Conf. Recommender Syst.*, 2012, pp. 155–162. [Online]. Available: <https://doi.org/10.1145/2365952.2365983>
- [42] H. Paulheim and J. Fürnkranz, "Unsupervised generation of data mining features from linked open data," in *Proc. 2nd Int. Conf. Web Intell. Mining Semantics*, 2012, pp. 31:1–31:12. [Online]. Available: <http://doi.acm.org/10.1145/2254129.2254168>
- [43] Z. Cao, T. Qin, T. Liu, M. Tsai, and H. Li, "Learning to rank: From pairwise approach to listwise approach," in *Proc. 24th Int. Conf. Mach. Learn.*, 2007, pp. 129–136. [Online]. Available: <http://doi.acm.org/10.1145/1273496.1273513>
- [44] Y. Shi, M. Larson, and A. Hanjalic, "List-wise learning to rank with matrix factorization for collaborative filtering," in *Proc. 4th ACM Conf. Recommender Syst.*, 2010, pp. 269–272. [Online]. Available: <http://doi.acm.org/10.1145/1864708.1864764>
- [45] P. G. Campos, F. Diez, and I. Cantador, "Time-aware recommender systems: A comprehensive survey and analysis of existing evaluation protocols," *User Model. User-Adapted Interact.*, vol. 24, no. 1-2, pp. 67–119, 2014. [Online]. Available: <https://doi.org/10.1007/s11257-012-9136-x>

- [46] Y. Koren, "Factorization meets the neighborhood: A multifaceted collaborative filtering model," in *Proc. 14th ACM SIGKDD Int. Conf. Knowl. Discov. Data Mining*, 2008, pp. 426–434. [Online]. Available: <http://doi.acm.org/10.1145/1401890.1401944>
- [47] S. Rendle and L. Schmidt-Thieme, "Pairwise interaction tensor factorization for personalized tag recommendation," in *Proc. 3rd Int. Conf. Web Search Web Data Mining*, 2010, pp. 81–90. [Online]. Available: <http://doi.acm.org/10.1145/1718487.1718498>
- [48] S. Rendle, C. Freudenthaler, and L. Schmidt-Thieme, "Factorizing personalized markov chains for next-basket recommendation," in *Proc. 19th Int. Conf. World Wide Web*, 2010, pp. 811–820. [Online]. Available: <http://doi.acm.org/10.1145/1772690.1772773>
- [49] X. He and T. Chua, "Neural factorization machines for sparse predictive analytics," in *Proc. 40th Int. ACM SIGIR Conf. Res. Develop. IR*, 2017, pp. 355–364.
- [50] J. Xiao, H. Ye, X. He, H. Zhang, F. Wu, and T. Chua, "Attentional factorization machines: Learning the weight of feature interactions via attention networks," in *Proc. 26th Int. Joint Conf. Artif. Intell.*, 2017, pp. 3119–3125.
- [51] Y. Zhu *et al.*, "Addressing the item cold-start problem by attribute-driven active learning," *IEEE Trans. Knowl. Data Eng.*, vol. 32, no. 4, pp. 631–644, Apr. 2020. [Online]. Available: <https://doi.org/10.1109/TKDE.2019.2891530>
- [52] S. M. McNee, J. Riedl, and J. A. Konstan, "Being accurate is not enough: How accuracy metrics have hurt recommender systems," in *Proc. Extended Abstracts Conf. Hum. Factors Comput. Syst.*, 2006, pp. 1097–1101. [Online]. Available: <http://doi.acm.org/10.1145/1125451.1125659>
- [53] Y. Koren and J. Sill, "Ordrec: An ordinal model for predicting personalized item rating distributions," in *Proc. ACM Conf. Recommender Syst.*, 2011, pp. 117–124. [Online]. Available: <http://doi.acm.org/10.1145/2043932.2043956>
- [54] I. Bayer, "fastFM: A library for factorization machines," *J. Mach. Learn. Res.*, vol. 17, pp. 184:1–184:5, 2016. [Online]. Available: <http://jmlr.org/papers/v17/15-355.html>
- [55] Y. Zhang, "Explainable recommendation: Theory and applications," *CoRR*, vol. abs/1708.06409, 2017. [Online]. Available: <http://arxiv.org/abs/1708.06409>
- [56] Y. Zhang *et al.*, "Daily-aware personalized recommendation based on feature-level time series analysis," in *Proc. 24th Int. Conf. World Wide Web*, 2015, pp. 1373–1383.
- [57] B. Abdollahi and O. Nasraoui, "Explainable matrix factorization for collaborative filtering," in *Proc. 25th Int. Conf. World Wide Web WWW Companion Vol.*, 2016, pp. 5–6. [Online]. Available: <http://doi.acm.org/10.1145/2872518.2889405>
- [58] A. Rana and D. Bridge, "Explanation chains: Recommendations by explanation," in *Proc. Poster Track 11th ACM Conf. Recommender Syst.*, 2017, pp. 187–195. [Online]. Available: http://ceur-ws.org/Vol-1905/recsys2017_poster4.pdf
- [59] M. Vlachos, C. Duenner, R. Heckel, V. G. Vassiliadis, T. Parnell, and K. Atasu, "Addressing interpretability and cold-start in matrix factorization for recommender systems," *IEEE Trans. Knowl. Data Eng.*, vol. 31, no. 7, pp. 1253–1266, Jul. 2019.
- [60] L. Hu, S. Jian, L. Cao, and Q. Chen, "Interpretable recommendation via attraction modeling: Learning multilevel attractiveness over multimodal movie contents," in *Proc. 27th Int. Joint Conf. Artif. Intell.*, 2018, pp. 3400–3406. [Online]. Available: <https://doi.org/10.24963/ijcai.2018/472>
- [61] G. Peake and J. Wang, "Explanation mining: Post hoc interpretability of latent factor models for recommendation systems," in *Proc. 24th ACM SIGKDD Int. Conf. Knowl. Discov. Data Mining*, 2018, pp. 2060–2069. [Online]. Available: <http://doi.acm.org/10.1145/3219819.3220072>
- [62] J. Vig, S. Sen, and J. Riedl, "Tagsplanations: Explaining recommendations using tags," in *Proc. 14th Int. Conf. Intell. User Interfaces*, 2009, pp. 47–56. [Online]. Available: <http://doi.acm.org/10.1145/1502650.1502661>
- [63] Z. Kang, C. Peng, and Q. Cheng, "Top-N recommender system via matrix completion," in *Proc. 13th AAAI Conf. Artif. Intell.*, 2016, pp. 179–185. [Online]. Available: <http://www.aaai.org/ocs/index.php/AAAI/AAAI16/paper/view/11824>
- [64] Z. Kang, C. Peng, and Q. Cheng, "Kernel-driven similarity learning," *Neurocomputing*, vol. 267, pp. 210–219, 2017. [Online]. Available: <https://doi.org/10.1016/j.neucom.2017.06.005>



Vito Walter Anelli is currently working toward the PhD degree in electrical and information engineering at the Polytechnic University of Bari, Italy. His current research interests fall in the areas of recommender systems, knowledge representation, and user modeling. He has published his works in national and international workshops and conferences as well as in international journals, and served as reviewer in venues of the aforementioned fields, including prestigious journals such as *Knowledge-Based Systems*, and *IEEE Transactions on Knowledge and Data Engineering*.



Tommaso Di Noia is full professor of information and data management with the Polytechnic University of Bari, Italy. Over the years, he has authored several papers in international journals and prominent conferences in his research areas. His current research interests include recommender systems and machine learning, knowledge graphs and Semantic Web technologies, personalized information access, and preference representation and reasoning.



Azzurra Ragone is currently Innovation and Transformation lead at EY. Her research interests fall in the areas of fairness in machine learning, recommender systems, and service design. She worked as a researcher for several years at the Polytechnic of Bari, the University of Trento, the University of Michigan (USA) and the University of Milan Bicocca. She has published her research works in more than 80 national and international workshops and conferences as well as in international journals.



Eugenio Di Sciascio received the MS degree (Hons.) from the University of Bari, Bari, Italy, and the PhD degree from the Polytechnic University of Bari, Bari. He has been an assistant professor with the University of Lecce and an associate professor with the Polytechnic University of Bari. He is currently a full professor of information technology engineering with the Polytechnic University of Bari, where he leads the research group of the SisInf Lab, the Information Systems Laboratory. His current research interests include multimedia information retrieval, knowledge representation and e-commerce. He is involved in several national and European research projects related to his research interests. He was a recipient of the Best Paper Award of the ICEC-2004, IEEE CEC-EEE-2006, ICEC-2007, SEMAPRO-2010, and ICWE-2010 conferences for his coauthored papers.



Joseph Trotta received the MS degree in computer engineering from the Polytechnic University of Bari. He is currently a research assistant at the Polytechnic University of Bari. His current research interests include recommender systems, and user modeling. He has published his works in national and international workshops and conferences.

▷ For more information on this or any other computing topic, please visit our Digital Library at www.computer.org/csdl.