

SMR: Medical Knowledge Graph Embedding for Safe Medicine Recommendation

Fan Gong^a, Meng Wang(✉)^{b,c}, Haofen Wang^d, Sen Wang^e, Mengyue Liu^f

^aShanghai Shuguang Hospital Affiliated to Shanghai University of Traditional Chinese Medicine, Pu'an Road, Shanghai, China

^bSchool of Computer Science and Engineering, Southeast University, Nanjing, China

^cKey Laboratory of Computer Network and Information Integration (Southeast University), Ministry of Education, Nanjing, China

^dCollege of Design and Innovation, Tongji University, Shanghai, China

^eThe University of Queensland, Brisbane, Australia

^fSchool of Electronic and Information Engineering, Xi'an Jiaotong University, Xi'an, China

Abstract

Most of the existing medicine recommendation systems that are mainly based on electronic medical records (EMRs) are significantly assisting doctors to make better clinical decisions benefiting both patients and caregivers. Even though the growth of EMRs is at a lighting fast speed in the era of big data, content limitations in EMRs restrain the existed recommendation systems to reflect relevant medical facts, such as drug-drug interactions. Many medical knowledge graphs that contain drug-related information, such as DrugBank, may give hope for the recommendation systems. However, the direct use of these knowledge graphs in systems suffers from robustness caused by the incompleteness of the graphs. To address these challenges, we stand on recent advances in graph embedding learning techniques and propose a novel framework, called Safe Medicine Recommendation (SMR), in this paper. Specifically, SMR first constructs a high-quality heterogeneous graph by bridging EMRs (MIMIC-III) and medical knowledge graphs (ICD-9 ontology and DrugBank). Then, SMR jointly embeds diseases, medicines, patients, and their corresponding relations into a shared lower dimensional space. Finally, SMR uses the embeddings to decompose the medicine recommendation into a link prediction process while considering the patient's diagnoses and adverse drug reactions. Extensive experiments on real datasets are conducted to evaluate the effectiveness of proposed framework.

Keywords: Knowledge Graph, Embeddings, Recommendation System, Drug Safety

1. Introduction

Over the last few years, medicine recommendation systems have been developed to assist doctors in making accurate medicine prescriptions. On the one hand, many researchers [1, 2] adopt rule-based protocols that are defined by the clinical guidelines and the experienced doctors. Constructing, curating, and maintaining these protocols are time-consuming and labor-intensive. Rule-based protocols might be effective

for a general medicine recommendation for a specific diagnosis, but give little help to tailored recommendations for complicated patients. On the other hand, supervised learning algorithms and variations, such as Multi-Instance Multi-label (MIML) learning [3], have been proposed to recommend medicines for patients. Both input features and ground-truth information that are extracted from massive EMRs are trained to obtain a predictive model that outputs multiple labels of the new testing data as medicine recommendations. It is a fact that therapies and treatments in clinical practices are rapidly updated. Unfortunately, supervised learning methods cannot deal with those medicines that are not included in the training phase. Incomplete training data set will be a detriment to the recommendation system performance.

It is reported in [4, 5] that patients with two or more diseases, acute or chronic, often take five or more different medicines simultaneously and have immense health risks. Studies [6, 7] have shown that 3-5% of all in-hospital misused prescriptions blame to ignorances of adverse drug reactions, which is difficult to prohibit even for the highly trained and experienced clinicians. With the assistance from the conventional medicine recommendation systems, clinicians still need to cautiously rule out those recommendations that have potential adverse effects caused by drug-drug interactions. Most of the existing works have largely ignored the exploit of medical facts in medicines, such as drug-drug interactions, which is crucial in medicine recommendation system. One possible reason might be because there is little medical expert knowledge in EMRs. Content limitations in EMRs constrain the systems to barely associate accurate medical facts with the recommended prescriptions, which makes the final recommendation less trustworthy for the complicated patients.

With the increasing emergence of knowledge graphs, many world-leading researchers have successfully extracted information from huge volumes of medical databases to build up giant heterogeneous graphs that reflect medical facts of medicines and diseases. For instance, DrugBank [8] is a rich source of medicine information. It contains extensive entities (drugs, drug targets, chemistry, etc.) and relationships (enzymatic pathways, drug-drug interactions, etc.). ICD-9 ontology [9] represents a knowledge base of human diseases and can be used to classify diagnoses of patients. Harnessing well-built medical knowledge graphs in EMRs-based medicine recommendation system might enable the reinvented system to provide appropriate prescriptions for special patients, as well as alerts of possible side effects and serious drug-drug interactions (DDIs).

As shown in Figure1, linking EMRs and medical knowledge graphs to generate a large and high-quality heterogeneous graph is a promising pathway for medicine recommendations in a wider scope, but never easy. Specifically, the newly designed system confronts with the following challenges: **1. Computational Efficiency.** Querying specialized medical entities and relationships based on conventional graph-based algorithms have limitations in portability and scalability. The computational complexity becomes unfeasible when the heterogeneous graph reaches a very large scale. **2. Data Incompleteness.** The medical knowledge graphs also follow the long-tail distribution as same as other types of large-scale knowledge bases. Data incompleteness is another serious problem existing among entities and relationships in such a distribution. For example, since the DDIs is not usually identified in the clinical trial phase, there is a lack of significant DDIs in DrugBank which cannot support a comprehensive pre-

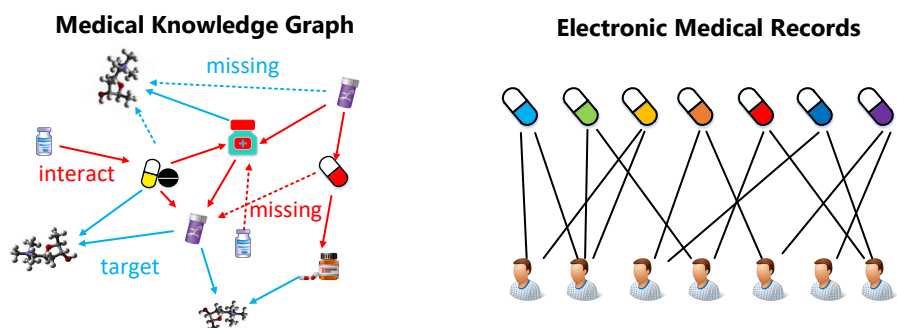


Figure 1: Left part is a medical knowledge graph with missing relationships. Right part is prescription records in EMRs, and each edge indicates that a patient takes a medicine. There is no relationship between medicines in EMRs.

caution to the medication. Last but not least, medicine recommendation suffers from **3. Cold Start**. As conventional systems normally recommend medicines based on the historical records, the pace of recommendation changes cannot keep up with the frequent updates of new therapies and treatments in medical practices. Little information on adverse reactions to the newly updated medicines in historical EMRs or even in well-built knowledge graphs makes the evidence-based recommendation model hardly support the new medicines as updated recommendations.

Taking all the challenges above into account, we propose a novel medicine recommendation framework based on graph embedding techniques, inspired by the idea of link prediction. We name our framework as Safe Medicine Recommendation (SMR) throughout this paper. The recommendation process mainly includes:

1. A large heterogeneous graph is constructed from EMRs and medical knowledge graphs, where the nodes are entities (medicines, diseases, patients), and the edges (or links) represent various relations between entities, such as drug-drug interactions.
2. The different parts of the generated heterogeneous graph (patient-medicine bipartite graph, patient-disease bipartite graph, medicine knowledge graph, disease knowledge graph) are embedded into a shared low-dimension space based on graph-based embedding models. Afterward, a joint learning algorithm is proposed to optimize the integrated graph simultaneously.
3. Based on the learned embeddings, a new patient, represented by the vectors of his/her diagnoses, is modeled as an entity in the disease-patient graph. Recommending medicines for the patient is translated to predict links from the patient to medicines.

The primary contributions of this work are summarized as follows:

1. We have developed graph-based embedding models to learn the effective representations of patients, diseases, and medicines in a shared low-dimension space. The representation of medicines enables the proposed framework can even effectively recommend newly emerged medicines for patients, which distinguishes most of the existing works.

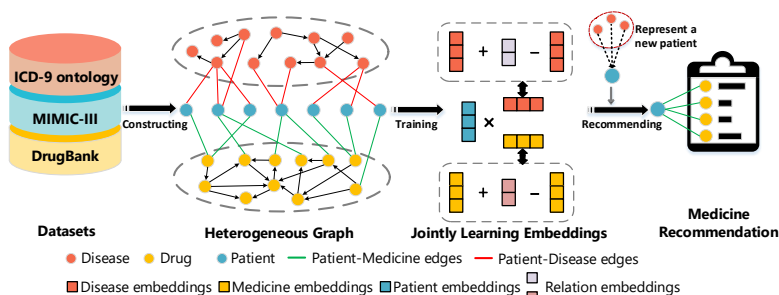


Figure 2: Overview of our framework. Patient-disease and patient-medicine graphs are all bipartite graphs, while disease and medicine graphs are general graphs. Patient, disease, and medicine are encoded into a low dimensional metric by graph-based methods. Diverse factors can be connected through patients.

2. To recommend safe medicines for new patients, we propose a novel method for modeling a patient based on the learned graph embeddings and make a safe recommendation by minimizing the potential adverse drug reactions.
3. We have conducted extensive experiments on large real-world datasets (MIMIC-III, DrugBank, and ICD-9 ontology) to evaluate the effectiveness of our framework. The experimental results have shown that the proposed framework outperforms all the compared methods.
4. To our best knowledge, we firstly propose a framework to conduct Safe Medicine Recommendation (SMR) and formulate it as a link prediction problem. The implementation generates a high-quality heterogeneous graph in which relationships among patients, diseases, and medicines can be unveiled in a wider scope.

The remainder of this paper is organized as follows: Section 2 details our proposed framework SMR. Section 3 reports the experimental results and Section 4 reviews related work. Section 5 presents the conclusions and future work.

2. The Proposed Framework

In this section, we will first describe the notations and formulate medicine recommendation problem, and then present graph embedding models and how to use learned embeddings to recommend safe medicines for patients.

2.1. Problem Formulation

Before we focus on the medicine recommendation problem, we first briefly introduce the important notations employed in the remainder of this paper. Table 1 also summarizes them.

The medical knowledge graph describes the medical entities collected from the integrated sources, as well as relationships among these entities. For instance, a triple (*glucocorticoid*, *adverse interaction*, *aspirin*) indicates that there is a relationship *adverse interaction* from *glucocorticoid* to *aspirin* in DrugBank. We define the medical knowledge graph as follow.

Table 1: Notations.

Variable	Interpretation
\mathcal{N}, \mathcal{R}	the set of entities and relations
(h, r, t)	a triple in knowledge graph
\mathbf{h}, \mathbf{t}	both are k dimensional vector embeddings
\mathbf{r}	a d dimensional relation embedding
\mathbf{H}_r	a $k \times d$ dimensional projection matrix
p, m, d	a patient, medicine, disease
$\mathbf{p}, \mathbf{m}, \mathbf{d}$	a k dimensional vector of a patient, a medicine or a disease
$\mathbb{R}^k, \mathbb{R}^d$	k and d dimensional latent space of entities and relations

Definition 1 (Medical Knowledge Graph). The medical knowledge graph $G = (\mathcal{N}, \mathcal{R})$ is a set of triples in the form (h, r, t) , where \mathcal{N} is a set of entities, \mathcal{R} is a set of relations, $h, t \in \mathcal{N}$ and $r \in \mathcal{R}$.

To capture the co-relationships of patients, diseases, and medicines in EMRs, we define the patient-disease, patient-medicine bipartite graphs as follow.

Definition 2 (Patient-Medicine Bipartite Graph). The patient-medicine bipartite graph is denoted as $G_{pm} = (\mathcal{P} \cup \mathcal{M}, \mathcal{E}_{pm})$, where \mathcal{P} is a set of patients and \mathcal{M} is a set of medicines. \mathcal{E}_{pm} is the set of edges. If a patient p_i takes a medicine m_j , there will be an edge e_{ij} between them, otherwise none. The weight w_{ij} of the edge between patient p_i and medicine m_j is defined as the total times of patient p_i takes the medicine m_j .

Definition 3 (Patient-Disease Bipartite Graph). The patient-disease bipartite graph is denoted as $G_{pd} = (\mathcal{P} \cup \mathcal{D}, \mathcal{E}_{pd})$, where \mathcal{P} is a set of patients and \mathcal{D} is a set of diseases. \mathcal{E}_{pd} is the set of edges. If a patient p_i is diagnosed with a disease d_j , there will be an edge e_{ij} between them, otherwise none. The weight w_{ij} is set to 1 when the edge e_{ij} exists.

Figure 2 illustrates a heterogeneous graph by constructing patient-disease, patient-medicine bipartite graphs from MIMIC-III, and linking them to medical knowledge graphs, ICD-9 ontology, and DrugBank. Finally, we formally define the safe medicine recommendation problem as follows.

PROBLEM 1 (Safe Medicine Recommendation). Given a patient p and his/her diagnoses dataset \mathcal{D}_p , recommending safe medicines for each $d \in \mathcal{D}_p$ is predicting edges from p to medicines dataset \mathcal{M} . The output is a set of medicines \mathcal{M}_p with minimum drug-drug interactions.

2.2. Model Description and Optimization

In this section, we propose embedding learning approaches to encode the heterogeneous graph in the latent space and its optimization method.

Medical Knowledge Graph Embedding

A medical knowledge graph $G = (\mathcal{N}, \mathcal{R})$ is a multi-relational graph, in which entities \mathcal{N} and relations \mathcal{R} can be different types. For a triple, $(h, r, t) \in G$, we use bold letters $\mathbf{h}, \mathbf{r}, \mathbf{t}$ to denote the corresponding embedding representations of h, r, t . Plenty of graph embedding methods has been proposed to encode a multi-relational graphs into a continuous vector space. Translation-based models [10, 11, 12] regard the relation r in each (h, r, t) as a translation from h to t within the low dimensional space, i.e., $\mathbf{h} + \mathbf{r} - \mathbf{t}$, and perform much more effectively and efficiently than conventional models. TransR [12] is a state-of-the-art translation-based embedding approach. It represents entities and relations in distinct vector space bridged by relation-specific matrices to get better graph representations.

Consider the above reason, we set entities embeddings $\mathbf{h}, \mathbf{t} \in \mathbb{R}^k$ and relations embeddings $\mathbf{r} \in \mathbb{R}^d$. And we set a projection matrix $\mathbf{H}_r \in \mathbb{R}^{k \times d}$, which projects entities from entity space to relation space. We define the translations between entities and get the energy function $z(h, r, t)$ as:

$$z(\mathbf{h}, \mathbf{r}, \mathbf{t}) = b - \|\mathbf{h}\mathbf{H}_r + \mathbf{r} - \mathbf{t}\mathbf{H}_r\|_{L1/L2} \quad (1)$$

where b is a bias constant.

Then, the conditional probability of a triple (h, r, t) is defined as follows:

$$P(h|r, t) = \frac{\exp\{z(\mathbf{h}, \mathbf{r}, \mathbf{t})\}}{\sum_{\hat{h} \in \mathcal{N}} \exp\{z(\hat{\mathbf{h}}, \mathbf{r}, \mathbf{t})\}} \quad (2)$$

and $P(t|h, r), P(r|h, t)$ can be defined in the analogous manner. We define the likelihood of observing a triple (h, r, t) as:

$$\mathcal{L}(h, r, t) = \log P(h|r, t) + \log P(t|h, r) + \log P(r|h, t) \quad (3)$$

We define an objective function by maximizing the conditional likelihoods of existing triples in G :

$$\mathcal{L}_G = \sum_{(h,r,t) \in G} \mathcal{L}(h, r, t) \quad (4)$$

Based on Eq.(4), the objective functions of medicine and disease knowledge graph $G_m = (\mathcal{N}_m, \mathcal{R}_m), G_d = (\mathcal{N}_d, \mathcal{R}_d)$ can be defined respectively:

$$\mathcal{L}_{G_m} = \sum_{(h_m, r_m, t_m) \in G_m} \mathcal{L}(h_m, r_m, t_m) \quad (5)$$

$$\mathcal{L}_{G_d} = \sum_{(h_d, r_d, t_d) \in G_d} \mathcal{L}(h_d, r_d, t_d) \quad (6)$$

Bipartite Graph Embedding

Different from the medical knowledge graph, the patient-disease, patient-medicine are bipartite graphs. A bipartite graph has only one single type of relations. For a bipartite graph, LINE [13] model achieves the state-of-the-art performance of encoding the

entities into a continuous vector space while preserving co-relations information of the graph. Hence, we follow LINE and set patients, medicines, and diseases embeddings $\mathbf{p}, \mathbf{m}, \mathbf{d} \in \mathbb{R}^k$. We present the process of encoding patient-medicine bipartite graph as follow.

Given a patient-medicine bipartite graph $G_{pm} = (\mathcal{P} \cup \mathcal{M}, \mathcal{E}_{pm})$. We first define the conditional probability of that a patient p_i in set \mathcal{P} takes medicine m_j in set \mathcal{M} as follow:

$$P(m_j|p_i) = \frac{\exp\{z(\mathbf{p}_i, \mathbf{m}_j)\}}{\sum_{\hat{m}_j \in \mathcal{P}_i} \exp\{z(\mathbf{p}_i, \hat{\mathbf{m}}_j)\}} \quad (7)$$

where $z(\mathbf{p}_i, \mathbf{m}_j) = \mathbf{m}_j^T \cdot \mathbf{p}_i$, \mathbf{p}_i is the embedding vector of the patient p_i in \mathcal{P} , and \mathbf{m}_j is the embedding vector of medicine m_j in \mathcal{M} . Eq. (7) defines a conditional distribution $P(\cdot|p_i)$ over all medicines in \mathcal{M} . The empirical distribution $\hat{P}(\cdot|p_i)$ is defined as $\hat{P}(m_j|p_i) = \frac{w_{ij}}{sum_i}$, where w_{ij} is the weight of the edge e_{ij} and $sum_i = \sum_j w_{ij}$ is the total times that the patient p_i takes medicines. We maximize the following objective function:

$$\mathcal{L}_{G_{pm}} = - \sum_{p_i \in \mathcal{P}} \lambda_i d(\hat{P}(\cdot|p_i), P(\cdot|p_i)) \quad (8)$$

where $d(\cdot, \cdot)$ is the distance between two distributions. In this paper, we use KL-divergence to compute $d(\cdot, \cdot)$. As sum_i is different from patients, we use $\lambda_i = sum_i$ in the objective function to represent the personalization of the patient p_i in the graph. After omitting some constants, we have:

$$\mathcal{L}_{G_{pm}} = \sum_{e_{ij} \in \mathcal{E}_{pm}} w_{ij} \log(P(m_j|p_i)) \quad (9)$$

For the patient-disease bipartite graph, we can get the object function $\mathcal{L}_{G_{pd}}$ in the analogous manner:

$$\mathcal{L}_{G_{pd}} = \sum_{e_{ij} \in \mathcal{E}_{pd}} w_{ij} \log(P(d_j|p_i)) \quad (10)$$

Optimization and Training

To learn the medical knowledge graph and bipartite graphs embeddings simultaneously, an intuitive approach is to collectively embed the four graphs (patient-medicine bipartite graph, patient-disease bipartite graph, medicine knowledge graph, disease knowledge graph) by maximizing the sum of the four logarithm likelihood objective functions just as follow:

$$\mathcal{L}(X) = \mathcal{L}_{G_m} + \mathcal{L}_{G_d} + \mathcal{L}_{G_{pm}} + \mathcal{L}_{G_{pd}} + \gamma C(X) \quad (11)$$

where X stands for the embeddings $\mathbb{R}_k, \mathbb{R}_d$ of entities and relations in the heterogeneous graph we construct, γ is a hyper-parameter weighting the regularization factor

$C(X)$, which is defined as follows:

$$\begin{aligned}
C(X) = & \sum_{n_m \in \mathcal{N}_m} [||n_m|| - 1]_+ + \sum_{n_d \in \mathcal{N}_d} [||n_d|| - 1]_+ \\
& + \sum_{r_m \in \mathcal{R}_m} [||r_m|| - 1]_+ + \sum_{r_d \in \mathcal{R}_d} [||r_d|| - 1]_+ \\
& + \sum_{p \in \mathcal{P}} [||p|| - 1]_+ + \sum_{d \in \mathcal{D}} [||d|| - 1]_+ + \sum_{m \in \mathcal{M}} [||m|| - 1]_+
\end{aligned} \tag{12}$$

where $[x]_+ = \max(0, x)$ denotes the positive part of x . The regularization factor will normalize the embeddings during learning. And we adopt the asynchronous stochastic gradient algorithm (ASGD) [14] to maximize the transformed objective function.

Optimizing objective functions Eq. (5), Eq. (6), Eq. (9) and Eq. (10) in Eq.(11) are computationally expensive, as calculating them need to sum over the entire set of entities and relations. To address this problem, we use the negative sampling method [15] to transform the objective functions.

For Eq.(5) and Eq.(6), we should transform $\log P(t|h, r)$, $\log P(r|h, t)$, $\log P(h|r, t)$ in Eq.(3). Taking $P(t|h, r)$ as an example, we maximize the following objective function instead of it:

$$\begin{aligned}
& \log \sigma(z(\mathbf{h}, \mathbf{r}, \mathbf{t})) \\
& + \sum_{n=1}^{C_1} E_{\tilde{h}_n \sim z_{neg}(\{(h, r, t)\})} [\sigma(z(\tilde{\mathbf{h}}_n, \mathbf{r}, \mathbf{t}))]
\end{aligned} \tag{13}$$

where C_1 is the number of negative examples, $\sigma(x) = 1/(1 + \exp(-x))$ is the sigmoid function. $\{(h, r, t)\}$ is the invalid triple set, and z_{neg} is a function randomly sampling instances from $\{(h, r, t)\}$. When a positive triple $(h, r, t) \in G$ is selected, to maximize Eq.(13), C_1 negative triples are constructed by sampling entities from an uniform distribution over \mathcal{N} and replacing the head of (h, r, t) . The transformed objective of $\log P(r|h, t)$, $\log P(t|h, r)$ are maximized in the same manner, but for $\log P(r|h, t)$, the negative relations are sampled from a uniform distribution over \mathcal{R} to corrupt the positive relation $r \in (h, r, t)$. We iteratively select random mini-batch from the training set to learn embeddings until converge.

For Eq. (9), we also use the negative sampling method to transform it to the following objective function:

$$\begin{aligned}
& \log \sigma(z(\mathbf{p}_i, \mathbf{m}_j)) \\
& + \sum_{n=1}^{C_2} E_{\tilde{m}_n \sim z_{neg}(\tilde{m})} [\log \sigma(z(\mathbf{p}_i, \tilde{\mathbf{m}}_n))]
\end{aligned} \tag{14}$$

where $\sigma(x) = 1/(1 + \exp(-x))$ is the sigmoid function, C_2 is the number of negative edges. $z_{neg}(\tilde{m}) \propto \text{sum}_{\tilde{m}}^{3/4}$ according to the empirical setting of [15], $\text{sum}_{\tilde{m}}$ is the total number of times that the medicine \tilde{m} is taken by patients. we can simplify Eq.(10) and maximize it in the same way.

Finally, we can efficiently learn the embeddings of different types of parts in the heterogeneous graph.

2.3. Safe Medicine Recommendation Process

In this section, we present how to recommend safe medicines based on the learned embeddings and diagnoses of a given patient. For an existing patient p , we use the learned embedding \mathbf{p} to predict new medicine recommendations. For a new given patient p , we first use diseases embeddings of p 's diagnoses to represent \mathbf{p} , and then recommend safe medicines for p , as shown in Figure 2.

New Patient Model

We aim to present a new patient p by his/her diagnoses embeddings. We should consider the time sequence of diseases that a patient is diagnosed, especially for the patient with multiple diseases. Assume a patient p in the hospital or on medication is associated with n ranked diseases according to their timestamps in an increasing order. Then, the patient embedding \mathbf{p} can be encoded as follow:

$$\mathbf{p} = \sum_{t=1}^n \exp^{-t} \cdot \mathbf{d}_t \quad (15)$$

where \mathbf{d}_t is the t -th embedding of disease d_t .

Medicine Recommendation

Given a query patient p with the query disease d , i.e., $q = (p, d)$, we first project disease d and patient p into their latent space, and then select top- k safe medicines¹. More precisely, given a query $q = (p, d)$, for each medicine m which could be useful for p , we compute its ranking score as in Eq. (16), and then select the medicine m with the top- k highest ranking scores as the recommendation.

$$S(q, m_n) = \mathbf{p}^T \cdot \mathbf{m}_n - \sum_{o=1}^{n-1} \|\mathbf{m}_n + \mathbf{r}_{interaction} - \mathbf{m}_o\|_{L1/L2} \quad (16)$$

where \mathbf{p} is the representation of patient p and m_n is the n -th medicines to be considered from medicines \mathcal{M} based on the the already selected medicine m_1, \dots, m_{n-1} .

3. Experiments and Evaluation

We attempt to demonstrate the effectiveness of our recommendation method in this section, which is referred to as SMR in this paper. In particular, we expect to answer “**how well does our method compare with the competing techniques?**” in Section 3.2. The results show that our recommendation method significantly outperforms the three baselines. The detailed experimental settings of our evaluations are described in Section 3.1.

¹In MIMIC-III, patients in ICUs are sicker and usually need more medicines for a diagnosis. We set $k=3$ in this paper.

Table 2: Entities and relations in the heterogeneous graph.

Entities		Relations	
#Disease	6,985	#Medicine-related	71,460
#Medicine	8,054	#Disease-related	170,579
#Patient	46,520	#Patient-Disease	489,923
#Medicine-related	305,642	#Patient-Medicine	975,647

Table 3: Experiments on medicine group 1.

	Prediction accuracy	DDIs rate
Rule-based	0.3068	32.01%
K -Most frequent	0.3473	14.08%
LEAP	0.5582	1.66%
SMR	0.6113	0.17%

3.1. Experimental Settings

Data Sets

Our experiments are performed on the real EMRs datasets, MIMIC-III [16], and two medical knowledge graphs, ICD-9 ontology [9] and DrugBank [8]. These real datasets are publicly available in different forms.

- MIMIC-III (Medical Information Mart for Intensive Care III) collected bedside monitor trends, electronic medical notes, laboratory test results, and waveforms from the ICUs (Intensive Care Units) of Beth Israel Deaconess Medical Center between 2001 and 2012. It contains distinct 46,520 patients, 650,987 diagnoses and 1,517,702 prescription records that associated with 6,985 distinct diseases and 4,525 medicines.
- ICD-9 ontology² (International classification of diseases-version 9) contains 13,000 international standard codes of diagnoses and the relationships between them.
- DrugBank is a bioinformatics/cheminformatics resource which consists of medicine related entities. The medical knowledge graph version³ contains 8,054 medicines, 4,038 other related entities (e.g., protein or drug targets) and 21 relationships.

Heterogeneous Graph Construction We connect MIMIC-III, ICD-9 ontology, and DrugBank (medicine group 1) by constructing the patient-medicine bipartite graph and the patient-disease bipartite graph.

For the patient-disease bipartite graph, MIMIC-III provides ICD-9 codes for diagnoses, which implicitly the diagnoses of MIMIC-III can be linked to ICD-9 ontology by string matching. For the patient-medicine bipartite graph, the prescriptions in MIMIC-III consist of the drug information, e.g., the names, the duration, and the dosage. However, various names to a single type of medicine in MIMIC-III exist due to some noisy

²<http://biportal.bioontology.org/ontologies/ICD9CM>

³<http://wifo5-03.informatik.uni-mannheim.de/drugbank/>

Table 4: Experiments on medicine group 2.

	Prediction accuracy	DDIs rate
Rule-based	0.2736	27.01%
K-Most frequent	NA	NA
LEAP	NA	NA
SMR	0.5214	2.01%

words (20%, 50ml, glass bottle, etc.), which becomes an obstacle to link medicine names to DrugBank when directly applying to the string matching method. We use an entity linking method [17] instead to address this problem. Table 2 shows the statistic of the heterogeneous graph we construct. The heterogeneous graph will be used to learn low-dimension representations of entities and relations by the SMR framework. Afterward, we categorize the medicines in the heterogeneous graph into two groups: 1). The first group consists of all 4,525 medicines that are recorded in EMRs, and will be used as inputs of the baseline methods. 2). The second group contains 3,529 medicines that haven’t been observed in EMRs, and will be used as test data for cold start recommendation.

Baselines

We compare our SMR with the following state of the art methods:

- Rule-based method [2] recommends medicines based on mappings from existing medicine categories to diseases in the MEDI database [18]. For each disease, a drug is assigned to the patient according to the mappings.
- K -Most frequent method is a basic baseline which retrieves the top K medicines that most frequently co-occur with each disease as their recommendation. We set $K = 3$ in this paper.
- LEAP method [3] uses a Multi-Instance Multi-Label learning framework to train a predictive model taking disease conditions as input features and yielding multiple medicine labels as recommendations.

3.2. Evaluation Methods

To guarantee medicine recommendations generated by SMR work effectively, we evaluate four indice, the prediction accuracy, the ability to avoid adverse drug-drug interactions, the experienced clinical doctor assessments, and the capacity to process cold start problem. In all experiments, the ratio of training to validation to test sets is 0.7:0.1:0.2. The hyper-parameters was adjusted by a validation set.

Prediction Accuracy and DDIs Rate

We utilize Jaccard Coefficient to compare the similarity of the prescriptions generated by SMR and the corresponding prescriptions written by doctors. Given the recommendation medicines set M_i generated by SMR for a patient p_i , \tilde{M}_i is the medicines

set prescribed by doctors in the data. The mean of Jaccard coefficient is defined as follows,

$$Jaccard = \frac{1}{K} \sum_{i=1}^K \frac{|M_i \cap \check{M}_i|}{|M_i \cup \check{M}_i|} \quad (17)$$

where K is the number of samples in a test set. Table 3 shows the accuracy of the baselines and SMR on medicine group 1, the rule-based method performs the worst because it is the only one provides a general recommendation for a specific diagnosis and it is not able to endow personalized recommendations, especially for the patients with multiple diseases. The frequency of each medicine-disease pair remains high in ICUs. Hence, recommendations based on frequency, K -Most frequent method, also work deficiently. Our method SMR outperforms LEAP by 1.49% because more accurate medical facts are involved in medical knowledge graphs rather than the prescription information in EMRs.

We extract all adverse drug-drug interactions (DDIs) from DrugBank to evaluate whether medicine recommendations embrace unsafe DDIs. Table 3 shows the percentages of different medicine recommendations consisting of adverse DDIs. The result indicates that SMR can recommend most harmless medicines for patients as its drug interaction rate is the lowest. The rule-based method and K -Most frequent method select medicines by a greedy strategy only regardless of specific adverse DDIs. For the rarely used medicines and unknown DDIs in EMRs, SMR is more reliable than LEAP. The reason is that SMR can predict each patient-medicine link and compute potential hidden DDIs by the learned embeddings of medical knowledge graphs.

Cold Start

We evaluate the ability of baselines and SMR in addressing cold start medicine recommendations on the medicine group 2. The results are reported in Table 4. K -Most frequent method and LEAP are not applicable (NA) on recommending new medicines in the cold start scenario. Since our SMR process can present new medicines by the learned vector representations of used medicines, the potential patient-medicine links between cold start medicines and patients will be captured correspondingly. In other words, SMR can leverage not only the patient-medicine links in EMRs but also the medical knowledge graphs when recommending cold start medicines.

Clinical Assessment

We invited three experienced clinical experts to evaluate the effectiveness of the medicine recommendations by scoring on a 6-point scale: 5 corresponding to completely cover all diagnoses without DDIs; 4 to partially (at least 50%) diagnoses include without DDIs; 3 to completely cover all diagnoses with DDIs; 2 to less than 50% diagnoses without DDIs; 1 to partially (at least 50%) diagnoses covered with DDIs; 0 to less than 50% diagnoses with DDIs. The average score of three experts is used as the final clinical assessment score for each recommendation, as shown in Figure 3.

Case Study

In table 5 we illustrate two events of medicine recommendations on medicine group 1 for patients associated with multiple types of diseases. SMR is qualified to succeed in

Table 5: Examples of medicine recommendations generated by Baselines and SMR.

Diagnosis	Methods	Medicine Recommendations
Sepsis Acute respiratory failure Hypertension	Rule-based	Teicoplanin, Metoprolol
	<i>K</i> -Most frequent	Vancomycin, Furosemide, Metoprolol, Insulin
	LEAP	Vancomycin, Furosemide, Metoprolol Tartrate
	SMR	Vancomycin, Furosemide, Amlodipine, Norepinephrine, Acetaminophen
Type 2 diabetes Rheumatoid arthritis Hypertension Hyperlipidemia	Rule-based	Gliclazide, Phenylbutazone, Sulfasalazine, Fenofibrate
	<i>K</i> -Most frequent	Furosemide, Tolbutamide, Phenylbutazone, Metoprolol, Insulin, Acetaminophen
	LEAP	Metformin, Amethopterin, Amiloride/HCTZ, Fenofibrate
	SMR	Metformin, Insulin, Acetaminophen, Nifedipine, Fenofibrate

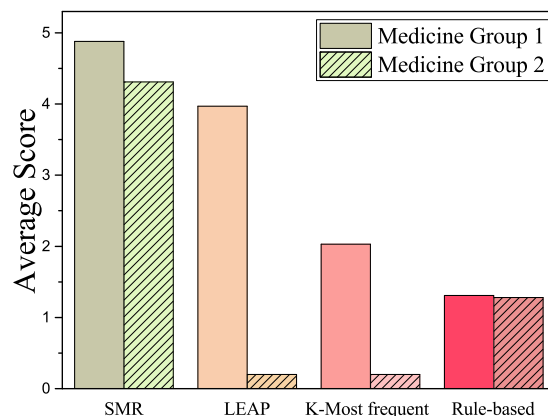


Figure 3: Clinical Assessment.

all these two cases when comparing it against other baselines. For the first patient, SMR recommended a set of medicines with 100% coverage, with Vancomycin for Sepsis, Norepinephrine, Acetaminophen for respiratory failure, Furosemide and Amlodipine for Hypertension. In contrast, other baselines are not capable of make an adequate consideration. The rule-based method adopted Teicoplanin, targeting Sepsis only and not appropriate. The *K*-Most frequent method and LEAP only selected Vancomycin for Sepsis and other medicines for Hypertension. For the second patient, SMR recommends more suitable medicines than LEAP and Rule-based method, i.e., Metformin and Insulin for Type 2 diabetes, Acetaminophen to release Rheumatoid arthritis, Nifedipine for Hypertension, and Fenofibrate for Hyperlipidemia. There is an adverse DDI among the medicines recommended by the *K*-Most frequent method. Tolbutamide and

Phenylbutazone can lead to harmful, potentially fatal effects when taken together. This case also indicates SMR can avoid the adverse DDIs when recommending medicines.

4. Related Work

In this section, we discuss related work, including medicine recommendation, medical knowledge graphs and embeddings.

Medicine Recommendation

As introduced in Section 1, two types of methods, rule-based protocols [1, 19, 2], and supervised-learning-based methods [3, 20], are currently utilizing EMRs to recommending medicines. Ideally, medicine recommendation systems aim to tailor treatment to the individual characteristics of each patient [21]. Hence, medicine recommendation has also received attention recently in genetics/genomics research fields. There are already existing medicine recommendation systems [22, 23] by leveraging genetics/genomics information of patients in current practice, such information is not yet widely available in everyday clinical practice, and is insufficient since it only addresses one of many factors affecting response to medication. A previous pre-print version of our work [24] mainly introduced our idea and recommendation model, while in this paper we provide a more in-depth analysis of the experiment and a summary of the relevant work.

Medical Knowledge Graphs and Embeddings

Recent evaluation efforts on knowledge graphs have focused on automatic knowledge base population and completion. Some knowledge graphs have also been constructed from huge volumes of medical databases over the last years, such as [25], Bio2RDF[26], and Chem2Bio2RDF[27]. Medical knowledge graphs contain an abundance of basic medical facts of medicines and diseases and provide a pathway for medical discovery and applications, such as effective medicine recommendations. Unfortunately, such medical knowledge graphs suffer from serious data incomplete problem, which impedes its application in the field of clinical medicine. Recently, Celebi et al. [28] proposed a knowledge graph embedding approach for drug-drug interaction prediction in realistic scenario. MedGraph [29] is a new EMR embedding framework which introduces a graph-based data structure to capture both structural visit-code collocation information structurally and temporal visit sequencing information. However, they only focused on the real-world biomedical embeddings learning and do not directly harness the medicine recommendation task, thus the specific clinical medicine information is not sufficiently tailored. In contrast, our model jointly embeds diseases, medicines, patients, and uses the learned embeddings to decompose the medicine recommendation into a linked prediction process while considering the patient’s diagnoses and adverse drug reactions.

5. Conclusion and Future Work

In this paper, we propose a novel framework SMR to recommend safe medicines for patients, especially for the patients with multiple diseases. SMR first constructs a

high-quality heterogeneous graph by bridging EMRs (MIMIC-III) and medical knowledge graphs (ICD-9 ontology and DrugBank). Then, SMR jointly embeds diseases, medicines, patients, and their corresponding relations into a shared lower dimensional space. Finally, SMR uses the embeddings to decompose the medicine recommendation into a linked prediction process considering the clinical diagnoses and adverse DDI reactions. Extensive experiments on real world datasets are conducted and demonstrate the effectiveness of SMR. In future work, we will improve the linking accuracy by considering more information of patients, such as the clinical outcomes and demographics.

6. Acknowledgment

This work was supported by National Science Foundation of China with Grant Nos. 61906037 and U1736204; the Fundamental Research Funds for the Central Universities with Grant No.4009009106 and 22120200184; the CCF-Baidu Open Fund.

References

- [1] Z. Chen, K. Marple, E. Salazar, G. Gupta, L. Tamil, A physician advisory system for chronic heart failure management based on knowledge patterns, *Theory and Practice of Logic Programming* 16 (5-6) (2016) 604–618.
- [2] D. Almirall, S. N. Compton, M. Gunlicks-Stoessel, N. Duan, S. A. Murphy, Designing a pilot sequential multiple assignment randomized trial for developing an adaptive treatment strategy, *Statistics in Medicine* 31 (17) (2012) 1887–1902.
- [3] Y. Zhang, R. Chen, J. Tang, W. F. Stewart, J. Sun, Leap: Learning to prescribe effective and safe treatment combinations for multimorbidity, in: *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, ACM, 2017, pp. 1315–1324.
- [4] M. Panagioti, J. Stokes, A. Esmail, P. Coventry, S. Cheraghi-Sohi, R. Alam, P. Bower, Multimorbidity and patient safety incidents in primary care: a systematic review and meta-analysis, *PloS One* 10 (8) (2015) e0135947.
- [5] D. N. Juurlink, M. Mamdani, A. Kopp, A. Laupacis, D. A. Redelmeier, Drug-drug interactions among elderly patients hospitalized for drug toxicity, *Journal of The American Medical Association* 289 (13) (2003) 1652–1658.
- [6] I. R. Edwards, J. K. Aronson, Adverse drug reactions: definitions, diagnosis, and management, *The Lancet* 356 (9237) (2000) 1255–1259.
- [7] L. L. Leape, D. W. Bates, D. J. Cullen, J. Cooper, H. J. Demonaco, T. Gallivan, R. Hallisey, J. Ives, N. Laird, G. Laffel, et al., Systems analysis of adverse drug events, *Journal of The American Medical Association* 274 (1) (1995) 35–43.
- [8] V. Law, C. Knox, Y. Djoumbou, T. Jewison, A. C. Guo, Y. Liu, A. Maciejewski, D. Arndt, M. Wilson, V. Neveu, et al., Drugbank 4.0: shedding new light on drug metabolism, *Nucleic Acids Research* 42 (D1) (2014) D1091–D1097.

- [9] L. M. Schriml, C. Arze, S. Nadendla, Y.-W. W. Chang, M. Mazaitis, V. Felix, G. Feng, W. A. Kibbe, Disease ontology: a backbone for disease semantic integration, *Nucleic Acids Research* 40 (D1) (2011) D940–D946.
- [10] A. Bordes, N. Usunier, A. Garcia-Duran, J. Weston, O. Yakhnenko, Translating embeddings for modeling multi-relational data, in: *Advances in neural information processing systems*, 2013, pp. 2787–2795.
- [11] Z. Wang, J. Zhang, J. Feng, Z. Chen, Knowledge graph embedding by translating on hyperplanes., in: *Proceedings of the 28th AAAI Conference on Artificial Intelligence*, 2014, pp. 1112–1119.
- [12] Y. Lin, Z. Liu, M. Sun, Y. Liu, X. Zhu, Learning entity and relation embeddings for knowledge graph completion., in: *Proceedings of the 29th AAAI Conference on Artificial Intelligence*, 2015, pp. 2181–2187.
- [13] J. Tang, M. Qu, M. Wang, M. Zhang, J. Yan, Q. Mei, Line: Large-scale information network embedding, in: *Proceedings of the 24th International Conference on World Wide Web, International World Wide Web Conferences Steering Committee*, 2015, pp. 1067–1077.
- [14] B. Recht, C. Re, S. Wright, F. Niu, Hogwild: A lock-free approach to parallelizing stochastic gradient descent, in: *Advances in Neural Information Processing Systems*, 2011, pp. 693–701.
- [15] T. Mikolov, I. Sutskever, K. Chen, G. S. Corrado, J. Dean, Distributed representations of words and phrases and their compositionality, in: *Advances in Neural Information Processing Systems*, 2013, pp. 3111–3119.
- [16] A. E. Johnson, T. J. Pollard, L. Shen, L.-w. H. Lehman, M. Feng, M. Ghassemi, B. Moody, P. Szolovits, L. A. Celi, R. G. Mark, Mimic-iii, a freely accessible critical care database, *Scientific Data* 3.
- [17] M. Wang, J. Zhang, J. Liu, W. Hu, S. Wang, X. Li, W. Liu, Pdd graph: Bridging electronic medical records and biomedical knowledge graphs via entity linking, in: *International Semantic Web Conference*, Springer, 2017.
- [18] W.-Q. Wei, R. M. Cronin, H. Xu, T. A. Lasko, L. Bastarache, J. C. Denny, Development and evaluation of an ensemble resource linking medications to their indications, *Journal of the American Medical Informatics Association* 20 (5) (2013) 954–961.
- [19] M. Gunlicks-Stoessel, L. Mufson, A. Westervelt, D. Almirall, A pilot smart for developing an adaptive treatment strategy for adolescent depression, *Journal of Clinical Child & Adolescent Psychology* 45 (4) (2016) 480–494.
- [20] P. Zhang, F. Wang, J. Hu, R. Sorrentino, Towards personalized medicine: leveraging patient similarity and drug similarity analytics, *Proceedings of AMIA Summits on Translational Science* 2014 (2014) 132.

- [21] G. H. Fernald, E. Capriotti, R. Daneshjou, K. J. Karczewski, R. B. Altman, Bioinformatics challenges for personalized medicine, *Bioinformatics* 27 (13) (2011) 1741–1748.
- [22] M. Rosen-Zvi, A. Altmann, M. Prosperi, E. Aharoni, H. Neuvirth, A. Sönnnerborg, E. Schülter, D. Struck, Y. Peres, F. Incardona, et al., Selecting anti-hiv therapies based on a variety of genomic and clinical factors, *Bioinformatics* 24 (13) (2008) i399–i406.
- [23] C. C. Bennett, K. Hauser, Artificial intelligence framework for simulating clinical decision-making: A markov decision process approach, *Artificial Intelligence in Medicine* 57 (1) (2013) 9–19.
- [24] M. Wang, M. Liu, J. Liu, S. Wang, G. Long, B. Qian, Safe medicine recommendation via medical knowledge graph embedding, arXiv preprint arXiv:1710.05980.
- [25] P. Ernst, A. Siu, G. Weikum, Knowlife: a versatile approach for constructing a large knowledge graph for biomedical sciences, *BMC Bioinformatics* 16 (1) (2015) 157.
- [26] M. Dumontier, A. Callahan, J. Cruz-Toledo, P. Ansell, V. Emonet, F. Belleau, A. Droit, Bio2rdf release 3: a larger connected network of linked data for the life sciences, in: *Proceedings of the 2014 International Conference on Posters & Demonstrations Track-Volume 1272*, CEUR-WS. org, 2014, pp. 401–404.
- [27] B. Chen, X. Dong, D. Jiao, H. Wang, Q. Zhu, Y. Ding, D. J. Wild, Chem2bio2rdf: a semantic framework for linking and data mining chemogenomic and systems chemical biology data, *BMC Bioinformatics* 11 (1) (2010) 255.
- [28] R. Celebi, H. Uyar, E. Yasar, O. Gumus, O. Dikenelli, M. Dumontier, Evaluation of knowledge graph embedding approaches for drug-drug interaction prediction in realistic settings, *BMC bioinformatics* 20 (1) (2019) 1–14.
- [29] B. Hettige, Y.-F. Li, W. Wang, S. Le, W. Buntine, Medgraph: Structural and temporal representation learning of electronic medical records, in: *Proceedings of the 24th European Conference on Artificial Intelligence*, 2020.