

RESEARCH

Drug-drug interaction prediction based on co-medication patterns and graph matching

Wen-Hao Chiang¹, Li Shen², Lang Li³ and Xia Ning^{4*}

Abstract

Background: The problem of predicting whether a drug combination of arbitrary orders is likely to induce adverse drug reactions is considered in this manuscript.

Methods: Novel kernels over drug combinations of arbitrary orders are developed within support vector machines for the prediction. Graph matching methods are used in the novel kernels to measure the similarities among drug combinations, in which drug co-medication patterns are leveraged to measure single drug similarities.

Results: The experimental results on a real-world dataset demonstrated that the new kernels achieve an area under the curve (AUC) value 0.912 for the prediction problem.

Conclusions: The new methods with drug co-medication based single drug similarities can accurately predict whether a drug combination is likely to induce adverse drug reactions of interest.

Keywords: drug-drug interaction prediction; drug combination similarity; co-medication; graph matching

Introduction

Drug-Drug Interactions (DDIs) and the associated Adverse Drug Reactions (ADRs) represent a consistent detriment to the public health in the United States. DDIs have accounted for approximately 26% of the ADRs, occurred among 50% of the hospitalized patients [1], and caused nearly 74,000 emergency room visits and 195,000 hospitalizations annually in the US [2]. Apart from these, because of the common practice of co-medication among elderly Americans, particularly co-medication of more than two drugs, the high-order drug-drug interactions and their associated ADRs have imposed significant scientific and public health challenges. The National Health and Nutrition Examination Survey [3] reports that more than 76% of the elderly Americans take two or more drugs every day. Another study [4] estimates that about 29.4% of elderly American patients take six or more drugs every day. However, for most of such high-order DDIs, their mechanisms are unknown.

In this manuscript, novel approaches to predicting whether high-order drug combinations are likely to induce ADRs are presented. The prediction problems are

formulated as a binary classification problem and support vector machines (SVMs) are used for the prediction. Novel kernels over drug combinations of arbitrary orders are developed within the framework of SVMs. These kernels are constructed using drug co-medication information to measure single drug similarities and graph matching on drug combination graphs to measure drug combination similarities. A comparison on the new kernels with other convolutional kernels and probabilistic kernels on drug combinations is also conducted. The experimental results demonstrate that the new kernels outperform the others and can accurately predict whether a drug combination is likely to induce ADRs of interest with an AUC value 0.912. To the best of our knowledge, this manuscript represents the first effort in predicting DDIs for drug combinations of arbitrary orders.

Background

Drug-drug interactions

Significant research efforts have been dedicated to detect pairwise drug-drug interactions (DDIs) [5, 6] in recent years. Existing methods either extract DDI pairs mentioned in medical literature or Electronic Health Records (EHRs) [4], or predict/score DDI pairs from various drug/target information [7]. While most

*Correspondence: xning@iupui.edu

⁴Department of Computer & Information Science, Indiana University - Purdue University Indianapolis, 46202 Indianapolis, USA. Email: xning@iupui.edu

Full list of author information is available at the end of the article

of the existing DDI studies are focused on interactions between a pair of drugs (i.e., order-2 DDIs), understanding high-order DDIs and their associated ADRs has attracted increasing attention recently [2, 8]. These emerging methods on high-order DDI studies are largely focused on how to discover high-order DDIs through mining frequent itemsets (i.e., drug combinations) from EHRs efficiently. Most recent work also includes pattern discovery from directional high-order DDIs [9] and directional high-order DDI prediction [10].

Graph matching

Graph matching is to find the optimal vertex correspondence between two graphs [11, 12]. Graph matching problems can be broadly classified into two categories. The first category is exact graph matching, which is to find the graph and subgraph isomorphisms so that the mapping of vertices between two graphs is bijective and edge-preserving (i.e., vertices connected by an edge in one graph are mapped to vertices in the other graph that are also connected by an edge). The second category is inexact graph matching, which allows errors (e.g., different types of matched vertices in attributed graphs) during matching, and thus it is to minimize the total errors in finding optimal graph matching. Typical algorithms for graph matching include spectral methods [13], probabilistic methods [14], tree search [15], etc.

Definitions and notations

We use d_i to represent a drug, and $D^k = \{d_1, d_2, \dots, d_k\}$ to represent a combination of k drugs, where k is the number of unique drugs in D^k (i.e., $k = |D^k|$) and thus the order of D^k . A drug combination D^k is defined when the drugs and only the drugs in D^k are taken simultaneously. There are no orderings among the drugs in a drug combination. When no ambiguity is raised, we drop the superscript k in D^k and represent a drug combination as D . An event is referred to as a patient taking a drug combination. In addition, in this manuscript, all vectors (e.g., \mathbf{c}) are represented by bold lower-case letters and all matrices (e.g., X) are represented by upper-case letters. Row vectors are represented by having the transpose superscript \top , otherwise by default they are column vectors. Table 1 summarizes the important notations in the manuscript.

Methods

We formulate the problem of predicting whether high-order drug combinations induce a particular ADR as a binary classification problem, and solve the classification problem within the framework of kernel methods and support vector machines (SVMs). In this

manuscript, we consider myopathy as the ADR in particular. The central concept of SVM-based classification methods is that “similar” instances are likely to share similar labels, and thus the key is to capture and measure the “similarities” among instances (i.e., drug combinations in our ADR prediction problem) via kernels. In the case of drug combinations, we hypothesize that if two drug combinations share similar pharmaceutical, pharmacokinetic and/or pharmacodynamic properties, they may induce similar ADRs. Therefore, the question boils down to effectively representing and measuring the similarities in terms of such properties. To this end, we develop various kernels over drug combinations. A key property of such kernels as will be discussed later is that they are able to deal with drug combinations of arbitrary orders. These kernels are constructed using single drug similarities, which incorporate various drug information that could relate to DDIs. Here we decompose the discussion on such kernels from three aspects: 1). single drug similarities (SDS) as in Section , 2). our new kernel based on matching similar drugs in drug combination graphs in Section , and 3). other convolutional kernels [16] in Section . Given these kernels, we further employ the freely available *SVM-Light* software to build up the binary classifiers and conduct our experiments based on such classifiers [17].

Single drug similarities

We use two different approaches to measuring single drug similarities (SDS). The first approach measures single drug similarities based on their intrinsic properties that can be represented by their 2D structures [18]. The second approach measures the similarities in a more data-driven fashion based on the co-occurrence patterns among drugs.

SDS from drug 2d structures

A straightforward way to measure SDSs between two drugs is to look at their structures, which ultimately determine their physicochemical properties. We use Extended Connectivity Fingerprints (ECFP) [19] of length 2,048 to represent drug 2D structures. Each of the fingerprint dimensions corresponds to a substructure among the drugs of interest. The binary values in the fingerprints represent whether a drug has the corresponding substructure or not. We use a vector $\mathbf{x}_i \in \mathbb{R}^{2048}$ to represent the fingerprint for drug d_i . The SDS between two drugs from their 2D structures, denoted as SDS_{2d} , is calculated as the Tanimoto coefficient between their ECFP fingerprints [20]. Tanimoto coefficient between two sets is defined as follows,

$$\text{Tanimoto}(S_1, S_2) = \frac{|S_1 \cap S_2|}{|S_1| + |S_2| - |S_1 \cap S_2|}, \quad (1)$$

where $|S|$ is the cardinality of set S . Thus, SDS_{2d} is defined as

$$\text{SDS}_{2d}(\mathbf{d}_i, \mathbf{d}_j) = \text{Tanimoto}(\{\mathbf{x}_i\}, \{\mathbf{x}_j\}), \quad (2)$$

where $\{\mathbf{x}_i\}$ represents the set of substructures that \mathbf{d}_i has in its fingerprint \mathbf{x}_i .

SDS based on co-mediations

We develop a new approach to measuring the SDS between two drugs by looking at whether they are often involved in co-mediations with similar other drugs, respectively. The hypothesis is that drugs that are respectively taken together with other similar drugs may share similar therapeutic purposes and target similar therapeutic targets, and thus behave similarly in inducing ADRs. Such data-driven co-medication based SDSs have a potential advantage over SDS_{2d} in that they leverage the signals from ADRs information directly that may not be captured or explained by drug 2D structures or other features on individual drugs. Such co-medication based SDS is denoted as SDS_{cm} .

We use two vectors $\mathbf{c}_i^+ \in \mathbb{R}^n$ and $\mathbf{c}_i^- \in \mathbb{R}^n$ (n is the total number of drugs) to represent the co-medication information for drug \mathbf{d}_i . The j -th dimension ($j = 1, \dots, n$) in $\mathbf{c}_i^+/\mathbf{c}_i^-$ corresponds to drug \mathbf{d}_j , and the value on the j -th dimension in $\mathbf{c}_i^+/\mathbf{c}_i^-$ is the co-medication frequency of \mathbf{d}_i and \mathbf{d}_j in all the events with/without ADRs. Both \mathbf{c}_i^+ and \mathbf{c}_i^- values are then normalized into probabilities. The normalized \mathbf{c}_i^+ and \mathbf{c}_i^- are further concatenated into one vector \mathbf{c}_i , that is, $\mathbf{c}_i = [\mathbf{c}_i^+; \mathbf{c}_i^-]$, for \mathbf{d}_i . The SDS_{cm} between drug \mathbf{d}_i and \mathbf{d}_j is calculated as the cosine similarity between \mathbf{c}_i and \mathbf{c}_j . The reason why we use \mathbf{c}_i^+ and \mathbf{c}_i^- to construct \mathbf{c}_i instead of co-medication frequencies from all events with and without ADRs together is that the co-medication patterns from the two types of events can be very different, and thus one unified co-medication vector for both of them could not necessarily capture discriminative information among drugs.

Drug combination kernels from graph matching

We formulate the problem of comparing drug combination similarities through matching drug combination graphs, and develop a graph-matching based kernel for drug combination similarities. Specifically, for a drug combination $D_p = \{\mathbf{d}_{p1}, \mathbf{d}_{p2}, \dots, \mathbf{d}_{pk_p}\}$, we construct a complete graph \mathcal{G}_p of k_p nodes, in which each node represents a drug in D_p , and all the nodes are connected to one another. Thus, the similarity between drug combination D_p and D_q can be measured based on how \mathcal{G}_p and \mathcal{G}_q match to each other. In matching such graphs, we consider SDSs so that drugs that are similar to each other should be matched, and the

graph matching procedure should maximize the overall SDSs from matched drugs. The underlying assumption is that if two drug combinations share similar drugs, they could have similar ADRs. Figure 1 illustrates the idea of complete graph matching for two drug combinations, in which the drugs connected by dash lines are matched between D_p and D_q . The similarity calculated from graph matching over two drug combinations, denoted as \mathcal{S}_{gm} , will be the sum of SDSs from matched drugs. \mathcal{S}_{gm} will be further converted to a valid kernel, denoted as \mathcal{K}_{gm} .

Graph matching algorithm for \mathcal{K}_{gm}

The drug combination graph matching problem can be solved as a well known linear sum assignment problem (LSAP) [21]. The objective is to minimize the total cost of matching vertices in two graphs, and thus to find the graph matching with minimal total cost. In the case of high-order drug combinations, we define the cost of matching two drugs \mathbf{d}_i and \mathbf{d}_j as the dissimilarity between the drugs, that is,

$$\text{cost}(\mathbf{d}_i, \mathbf{d}_j) = 1 - \text{SDS}(\mathbf{d}_i, \mathbf{d}_j), \quad (3)$$

where $\text{cost}(\mathbf{d}_i, \mathbf{d}_j)$ is the cost between \mathbf{d}_i and \mathbf{d}_j , SDS can be either SDS_{2d} or SDS_{cm} . Thus, if two drugs are very similar (i.e., large SDS), the cost of matching them will be small and therefore they are more likely to be matched.

Therefore, the graph matching can be solved by solving the following LSAP problem:

$$\begin{aligned} \min_X \quad & \text{trace}(C(\mathcal{G}_p, \mathcal{G}_q)X^T) \\ \text{subject to} \quad & X \in \mathcal{P}, \\ & \mathcal{P} := \{X \mid X \in \mathbb{R}^{k_p \times k_q}, X_{i,j} \in \{0, 1\}, \\ & \sum_{i=1}^{k_p} X_{i,j} \leq 1, \sum_{j=1}^{k_q} X_{i,j} \leq 1, \\ & \sum_{i=1}^{k_p} \sum_{j=1}^{k_q} X_{i,j} = \min(k_p, k_q)\}, \end{aligned} \quad (4)$$

where $\text{trace}()$ is the trace of a matrix; and k_p and k_q are the number of vertices in \mathcal{G}_p and \mathcal{G}_q (and thus the order of D_p and D_q), respectively; $C(\mathcal{G}_p, \mathcal{G}_q) \in \mathbb{R}^{k_p \times k_q}$ is the pairwise drug-matching cost matrix for two drug combinations D_p and D_q ($C(i, j) = \text{cost}(\mathbf{d}_{pi}, \mathbf{d}_{qj})$, $\mathbf{d}_{pi} \in D_p, \mathbf{d}_{qj} \in D_q$). In Problem 4, X is the assignment matrix to match \mathcal{G}_p and \mathcal{G}_q (i.e., to assign a vertex in \mathcal{G}_p to a vertex in \mathcal{G}_q), in which all the values are either 0 or 1, both the row sum and the column sum are either 0 or 1 (i.e., a vertex is either matched or not; if it is matched, it is matched to only one vertex in

the other graph), and thus the sum of all the values is exactly the minimal of k_p and k_q (i.e., the vertices in the small graph have to be all matched). Essentially, X assigns each of the vertices in the smaller graph of \mathcal{G}_p and \mathcal{G}_q to exactly one vertex in the larger graph. The optimization problem in 4 can be solved by the Hungarian algorithm [22]. The drug-combination similarity \mathcal{S}_{gm} is then calculated as

$$\mathcal{S}_{\text{gm}}(D_p, D_q) = \text{trace}(J - C(\mathcal{G}_p, \mathcal{G}_q)X^\top), \quad (5)$$

where $J \in \mathbb{R}^{k_p \times k_q}$ is a matrix of all 1's.

The drug-combination similarity matrix \mathcal{S}_{gm} is always symmetric but not necessarily positive semi-definite, and thus not always a valid kernel. To convert \mathcal{S}_{gm} to a valid kernel \mathcal{K}_{gm} , we follow the approach in Saigo *et al.*[23]. Specifically, we first conduct an eigenvalue decomposition on \mathcal{S}_{gm} , subtract from the diagonal of the eigenvalue matrix its smallest negative eigenvalue, and reconstruct the original matrix from the altered decomposition. The resulted matrix is positive, semi-definite, and is used as \mathcal{K}_{gm} .

Convolutional drug-combination kernels

Drug combination kernels from common drugs

We define a drug-combination kernel, denoted as \mathcal{K}_{cd} , based on common drugs among drug combinations. \mathcal{K}_{cd} is calculated as the Tanimoto coefficient over the sets of drugs in the drug combinations, that is,

$$\mathcal{K}_{\text{cd}}(D_p, D_q) = \text{Tanimoto}(D_p, D_q), \quad (6)$$

where $\text{Tanimoto}()$ is defined as in Equation 1. It has been proved that Tanimoto coefficient is a valid kernel function [24]. \mathcal{K}_{cd} essentially measures the proportion of shared common drugs among two drug combinations. The underlying assumption is that if two drug combinations share many common drugs, they are likely to have similar properties.

To further enhance the similarity between two drug combinations from their common drugs, we also define an order-2 \mathcal{K}_{cd} of drug combinations, denoted as $\mathcal{K}_{\text{cd}}^{(2)}$ (\mathcal{K}_{cd} in Equation 6 is correspondingly referred to as order-1 \mathcal{K}_{cd} and denoted as $\mathcal{K}_{\text{cd}}^{(1)}$). We first represent a drug combination $D = \{\mathbf{d}_1, \mathbf{d}_2, \dots, \mathbf{d}_k\}$ by all its single drugs and drug pairs, denoted as $D^{(2)} = \{\mathbf{d}_1, \mathbf{d}_2, \dots, \mathbf{d}_k, (\mathbf{d}_1, \mathbf{d}_2), (\mathbf{d}_1, \mathbf{d}_3), \dots, (\mathbf{d}_{k-1}, \mathbf{d}_k)\}$. Thus, $\mathcal{K}_{\text{cd}}^{(2)}$ on two drug combinations D_p and D_q can be calculated as the Tanimoto coefficient on $D_p^{(2)}$ and $D_q^{(2)}$, that is,

$$\mathcal{K}_{\text{cd}}^{(2)}(D_p, D_q) = \text{Tanimoto}(D_p^{(2)}, D_q^{(2)}). \quad (7)$$

Intuitively, $\mathcal{K}_{\text{cd}}^{(2)}$ better differentiates drug combinations with many shared drugs from those with fewer shared drugs than $\mathcal{K}_{\text{cd}}^{(1)}$. We only extend \mathcal{K}_{cd} to order 2 since higher-order extension does not lead to better performance according to our experimental results. According to Equation 6, when the order becomes much higher, $\text{Tanimoto}(D_p^{(n)}, D_q^{(n)})$ may become very small due to a rapid combinatorial growth in the denominator and the insufficient common drug n -tuples (i.e., the number in the nominator). Thus, \mathcal{K}_{cd} with extension to much higher order may lose the ability to differentiate drug combinations that contain more common drugs.

Drug combination kernels from drug similarities

The drug combination similarities can also be measured by the average drug similarities. The hypothesis is that if two drug combinations have drugs that are similar on average, they may share similar properties. If two drug combinations have drugs that are similar on average, they may share similar properties. Therefore, we define an average-drug-similarity based kernel for drug combinations, denoted as \mathcal{K}_{ds} , as follows,

$$\mathcal{K}_{\text{ds}}(D_p, D_q) = \frac{1}{k_p k_q} \sum_{\mathbf{d}_i \in D_p} \sum_{\mathbf{d}_j \in D_q} \text{SDS}(\mathbf{d}_i, \mathbf{d}_j), \quad (8)$$

where k_p and k_q are the order of D_p and D_q , respectively, and SDS can be $\text{SDS}_{2\text{d}}$ or SDS_{cm} . Intuitively, \mathcal{K}_{ds} tends to capture averaged and smoothed drug combination similarities. It has been proved that as long as the involved SDSs are valid kernels (i.e., positive semi-definite), \mathcal{K}_{ds} will also be a valid kernel [16].

Probabilistic drug combination kernels from drug sets

We apply an ensemble kernel for drug combinations based on the idea as in [25]. The key idea is to use a reproducing kernel to characterize sample similarities (i.e., SDS), and to use a probabilistic distance in the reproducing kernel Hilbert space (RKHS) to measure the ensemble similarity. The resulted ensemble similarity matrix is a valid kernel matrix, denoted as \mathcal{K}_{pb} . This ensemble involves an eigen value decomposition, during which, it is possible that some similarity matrices are deprecated numerically and it leads to defeats in \mathcal{K}_{pb} calculation. To deal with this issue, we increase the diagonals of involved square matrices by a small value to guarantee the positive semi-definite properties.

Materials

Mining drug combinations

We extract high-order drug combinations from FDA Adverse Event Reporting System (FAERS) [26]. We

use myopathy as the ADR of particular interest, and extract 64,892 case (myopathy) events, in which patients report myopathy after taking multiple drugs, and 1,475,840 control (non-myopathy) events, in which patients do not report myopathy after taking drugs. Each of these events involves a combination of more than one drug.

Among all the involved drug combinations, 10,250 unique drug combinations appear in both case and control events. For those 10,250 drug combinations, we use Odds Ratio (OR) to quantify their ADR risks. The OR for a drug combination D is defined based on the contingency table 2, that is, it is the ratio of the following two values: 1). the odds that the ADR occurs when D is taken (i.e., $\frac{n_1}{m_1}$ in Table 2); and 2). the odds that the ADR occurs when D is not taken (i.e., $\frac{n_2}{m_2}$ in Table 2). $OR < 1$ indicates the decreased risk of ADR after a patient takes the drug combination, $OR = 1$ indicates no risk change, and $OR > 1$ indicates the increased risk. In the 10,250 drug combinations, 8,986 combinations have $OR > 1$ and 1,264 combinations have $OR < 1$. These two sets of drug combinations are denoted as \mathcal{M}^0 and \mathcal{N}^0 , respectively. In addition to these combinations, there are 27,387 unique drug combinations that only appear in case events and 621,449 unique drug combinations that only appear in control events. These two sets are denoted as \mathcal{M}^+ and \mathcal{N}^- , respectively. The set of drug combinations in case events is denoted as \mathcal{M} (i.e., $\mathcal{M} = \mathcal{M}^+ \cup \mathcal{M}^0$), and the set of drug combinations in control events is denoted as \mathcal{N} (i.e., $\mathcal{N} = \mathcal{N}^- \cup \mathcal{N}^0$). All these four sets together define a high-order drug combination dataset from FAERS, denoted as \mathcal{D}_{FAERS} . Table 3 presents the statistics of \mathcal{D}_{FAERS} .

Training data generation

As shown in Table 3, \mathcal{M}^+ and \mathcal{M}^0 of \mathcal{D}_{FAERS} have fewer drug combinations than \mathcal{N}^- and \mathcal{N}^0 , and the drug combinations in \mathcal{M}^+ are very infrequent (average frequency 1.402). To use more frequent and more confident drug combinations from case events, we further pruned drug combinations from \mathcal{M}^+ and \mathcal{M}^0 as follows. From \mathcal{M}^+ , we retained the top 1,000 most frequent drug combinations. For \mathcal{M}^0 , we applied right-tailed Fisher’s exact test on the drug combinations to further test the significance of their ORs at 5% significance level. Then we retained drug combinations with statistically significant ORs. Thus, the pruned \mathcal{M}^+ and \mathcal{M}^0 contain statistically confident drug combinations, which are very likely to induce myopathy, and therefore, these drug combinations are labeled as positive instances for classification model learning.

We retained all 1,264 the drug combinations in \mathcal{N}^0 because this set is not large and contains informative

drug combinations that may or may not induce myopathy. We further prune \mathcal{N}^- and retain the top most frequent drug combinations. The drug combinations from \mathcal{N}^0 and the pruned \mathcal{N}^- are labeled as negative instances. To make the positive and negative training sets balance, we retained 2,200 drug combinations from \mathcal{N}^- . The pruned dataset from \mathcal{D}_{FAERS} is denoted as \mathcal{D}^* . Table 3 presents the description of \mathcal{D}^* . \mathcal{D}^* is the set of labeled drug combinations that are used for model learning. In \mathcal{D}^* , there are in total 1,210 drugs involved. 71 out of these 1,210 drugs induce myopathy on their own based on the Side Effect Resource (SIDER) [27]. This set of 71 drugs is denoted as \mathcal{D}_{Myo} .

Evaluation protocol and metrics

The performance of the different methods is evaluated through five-fold cross validation. The dataset is randomly split into five folds of equal size (i.e., same number of drug combinations). Four folds are used for model training and the rest fold is used for testing. This process is performed five times, with one fold for testing each time. The final result is the average out of the five experiments.

We use accuracy, precision, recall, F1 and AUC to evaluate the performance of the methods. Accuracy is defined as the fraction of all correctly classified instances (i.e., true positives and true negatives) over all the instances in the testing set. Precision is defined as the fraction of correctly classified positive instances (i.e., true positives) over all instances that are classified as positive instances (i.e., true positives and false positives). Recall is the fraction of correctly classified positive instances (i.e., true positives) over all positive instances in the testing set (i.e., true positives and false negatives). F1 is the harmonic mean of precision and recall. AUC score is the normalized area under the curve that plots the true positives against the false positives for different thresholds for classification [28]. Larger accuracy, precision, recall, F1 and AUC values indicate better classification performance.

Results

Overall performance

Table 4 presents the performance comparison among the four different kernels in combination with different single drug similarities on dataset \mathcal{D}^* . Kernel \mathcal{K}_{gm} with SDS_{cm} outperforms others in three (i.e., accuracy, F1 and AUC) out of five evaluation metrics. Specifically, in accuracy, \mathcal{K}_{gm} with SDS_{cm} outperforms the second best kernel \mathcal{K}_{gm} with SDS_{2d} at 0.84%. In F1, \mathcal{K}_{gm} with SDS_{cm} outperforms the second best kernel \mathcal{K}_{gm} with SDS_{2d} and \mathcal{K}_{ds} with SDS_{cm} at 0.98%. In AUC, \mathcal{K}_{gm} with SDS_{cm} outperforms the second best kernel order-2 \mathcal{K}_{cd} at 0.33%. In precision and recall,

\mathcal{K}_{gm} with SDS_{cm} is the second best kernel, whereas \mathcal{K}_{ds} with SDS_{cm} and \mathcal{K}_{ds} with $\text{SDS}_{2\text{d}}$, respectively, is the best one. Overall, \mathcal{K}_{gm} with SDS_{cm} has the best performance compared to other kernels. This indicates that it is effective to classify drug combinations by representing and comparing them as graphs (i.e., a set of drugs and their co-medication relation within the set), and measuring such graph similarities using their optimal matching (i.e., the optimal correspondence among drugs). In the following discussion, we use $\mathcal{K}_{\text{gm}}^{\text{cm}}$ to represent \mathcal{K}_{gm} with SDS_{cm} . More experimental results on other datasets are available in the supplementary materials (see Additional file 1).

SDS performance

Table 4 shows that SDS_{cm} on average outperforms $\text{SDS}_{2\text{d}}$ across different kernels (with a few exceptions on in precision for \mathcal{K}_{ds} and \mathcal{K}_{pb}). $\text{SDS}_{2\text{d}}$ considers drug intrinsic 2D structures. However, drug efficacy and side effects are the results of many complicated interactions and processes among drugs and various bioentities, which may not be sufficiently explained only by drug 2D structures. Compared to $\text{SDS}_{2\text{d}}$, SDS_{cm} measures drug similarity based on their co-medication patterns, which could be regarded as a high-level abstraction and representation of drug therapeutic properties that may or may not be explicitly explained by each drug and its intrinsic properties independently.

In \mathcal{K}_{cd} , order-2 representation (i.e., in $\mathcal{K}_{\text{cd}}^{(2)}$) for drug combinations outperforms order-1 representation (i.e., in $\mathcal{K}_{\text{cd}}^{(1)}$). In order-2 representation, in addition to single drugs, drug pairs are also used as a feature for a drug combination, which stresses the signals in drug combinations. This also conforms to common observations in other applications [29], in which higher-order features improve classification performance.

Classification

Figure 2 and 3 present the $\mathcal{K}_{\text{gm}}^{\text{cm}}$ prediction values with respect to drug combination orders. In Figure 2, \mathcal{M}^+ drug combinations have higher orders (on average 7.615 as in Table 3), and higher and mostly positive prediction values, while \mathcal{N}^- drug combinations have lower orders (on average 2.678), and lower and mostly negative prediction values. Meanwhile, the mis-classification typically happens on \mathcal{N}^- drug combinations of higher orders, and on \mathcal{M}^+ drug combinations of lower orders. Similar trends apply for \mathcal{M}^0 and \mathcal{N}^0 in Figure 3. This indicates that \mathcal{K}_{gm} and SDS_{cm} together are able to learn and make predictions that correspond to drug combination orders. In addition, drug combination order is correlated with their ADR labels.

Figure 4 presents the $\mathcal{K}_{\text{gm}}^{\text{cm}}$ prediction values with respect to drug combination frequencies for \mathcal{M}^+ and

\mathcal{N}^- . \mathcal{M}^+ drug combinations have lower frequencies (on average 5.520 as in Table 3), and higher and mostly positive prediction values, while \mathcal{N}^- drug combinations have higher frequencies (on average 42.082), and lower and mostly negative prediction values. For \mathcal{N}^- , the mis-classification typically happens on lower-frequency drug combinations (the mis-classification for \mathcal{M}^+ does not show strong patterns with respect to drug combination frequencies). As for \mathcal{M}^+ and \mathcal{N}^- , drug combination frequencies are used to define ADR labels. Figure 4 shows that $\mathcal{K}_{\text{gm}}^{\text{cm}}$ together are able to learn and make predictions that correspond to drug combination frequencies and thus ADR labels.

Figure 5 presents the $\mathcal{K}_{\text{gm}}^{\text{cm}}$ prediction values with respect to OR values for \mathcal{M}^0 and \mathcal{N}^0 . \mathcal{M}^0 drug combinations have higher OR values and also higher and mostly positive prediction values, while \mathcal{N}^0 drug combinations have lower OR values and also lower and mostly negative prediction values. For \mathcal{N}^0 , the mis-classification typically happens on drug combinations of higher OR values (close to 1 and thus more lean toward ADR; the mis-classification for \mathcal{M}^0 does not show strong patterns with respect to OR values). As we use OR to define ADR labels on \mathcal{M}^0 and \mathcal{N}^0 , Figure 5 shows $\mathcal{K}_{\text{gm}}^{\text{cm}}$ is able to make reasonably accurate prediction values on the drug combinations.

\mathcal{D}_{Myo} drug enrichment

Table 5 presents the average percentage of \mathcal{D}_{Myo} drugs among all the drug combinations. For each drug combination, the percentage is calculated as the number of its drugs that can cause myopathy on their own (i.e., drugs in \mathcal{D}_{Myo}) divided by the drug combination order. As Table 5 shows, top-10 mis-classified \mathcal{N} drug combinations (i.e., $\tilde{\mathcal{N}}^{10+}$) have almost twice as many \mathcal{D}_{Myo} drugs (30.7%) as those in \mathcal{N} (15.6%), and even more than those in \mathcal{M} drug combinations (24.3%). In addition, mis-classified \mathcal{N} drug combinations (i.e., $\tilde{\mathcal{N}}^+$) also have significantly more \mathcal{D}_{Myo} drugs (18.6%) than those in \mathcal{N} (15.6%). Since $\mathcal{K}_{\text{gm}}^{\text{cm}}$ matches similar drugs, high \mathcal{D}_{Myo} drug enrichment could be a primary reason for the mis-classification.

Top predictions

Top mis-classification on \mathcal{N}

Table 6 lists the top-10 (in terms of prediction values) drug combinations in \mathcal{N} (i.e., without myopathy) that are mis-classified as positive (i.e., with myopathy) by $\mathcal{K}_{\text{gm}}^{\text{cm}}$. For those drug combinations which appear in \mathcal{N}^0 , we present their OR values, otherwise only frequencies. Those top mis-classified \mathcal{N} drug combinations contain many single drugs, which on their own can induce myopathy (i.e., in \mathcal{D}_{Myo} , bold in Table 6). As a matter of fact, the percentage of \mathcal{D}_{Myo} drugs in top mis-classified

\mathcal{N} drug combinations is significantly higher than average. In Table 6, one special mis-classified \mathcal{N} drug combination is {lansoprazole omeprazole pantoprazole rabeprazole}, which does not contain any \mathcal{D}_{Myo} drugs. This set of drugs is commonly used as proton pump inhibitors (PPIs) to decrease the amount of acid produced in the stomach. Some case studies show evidence of causality between the PPI drug class and myopathy [30, 31].

Top prediction on \mathcal{M}

Table 7 presents the top-10 correctly predicted \mathcal{M} drug combinations by $\mathcal{K}_{\text{gm}}^{\text{cm}}$. These drug combinations are significantly enriched with \mathcal{D}_{Myo} drugs (i.e., drugs that can induce myopathy on their own). As Table 5 shows, $\tilde{\mathcal{M}}^{10+}$ has the most \mathcal{D}_{Myo} drugs (89.8%) compared to all the other sets and significantly more than \mathcal{M} . In particular, all of these combinations contain statin drugs (e.g., atorvastatin, simvastatin and rosuvastatin, etc.). These statin-related drugs have been studied in literature as a drug class that has high possibilities to induce myopathy [32, 33]. In addition, in Table 7, 4 out of the 6 \mathcal{M}^0 drug combinations among top 10 (i.e., the drug combinations that have OR values) have their OR values higher than average in \mathcal{M}^0 (31.998 as in Table 3), and 3 out of the 4 \mathcal{M}^+ drug combinations among top 10 (i.e., the drug combinations that do not have OR values) have their frequency higher than average in \mathcal{M}^+ (5.520 as in Table 3). In addition, among the top-20 drug combinations predicted by $\mathcal{K}_{\text{gm}}^{\text{cm}}$, 7 out of 12 \mathcal{M}^0 drug combinations have their OR values higher than average in \mathcal{M}^0 , and 5 out of 8 \mathcal{M}^+ drug combinations have their frequency higher than average in \mathcal{M}^+ . The average OR values of the top-10, top-20 and top-50 drug combinations from \mathcal{M}^0 predicted by $\mathcal{K}_{\text{gm}}^{\text{cm}}$ are 55.725, 42.956 and 42.114, respectively, and they are all higher than the average 31.998 for \mathcal{M}^0 . The average frequencies of the top-10, top-20 and top-50 drug combinations from \mathcal{M}^+ predicted by $\mathcal{K}_{\text{gm}}^{\text{cm}}$ are 8.250, 6.875 and 6.524, respectively, and they are also all higher than the average 5.520 on \mathcal{M}^+ . This indicates that $\mathcal{K}_{\text{gm}}^{\text{cm}}$ does learn signals from \mathcal{M} and correspondingly makes predictions.

Top non- \mathcal{D}_{Myo} prediction on \mathcal{M}

Table 8 presents the top-10 correctly predicted \mathcal{M} drug combinations by $\mathcal{K}_{\text{gm}}^{\text{cm}}$ that do not contain any drugs from \mathcal{D}_{Myo} (i.e., do not contain drugs that can induce myopathy on their own). 4 out of 7 drug combinations from \mathcal{M}^0 in this table have their OR values higher than the average in \mathcal{M}^0 (31.998 as in Table 3). The 3 drug combinations from \mathcal{M}^+ in this table have their frequency lower than the average in \mathcal{M}^+ (5.520 as in Table 3) but very close. In Table 8, 8 out of

top-10 drug combinations include alendronate. Case studies demonstrate that several events of severe muscle pain, which is the common symptom of myopathy, were reported after patients started therapy with alendronate [34], showing the association between the medical treatment with alendronate and myopathy.

Discussions

The experimental results show that the new methods with drug co-medication based single drug similarities outperform other kernels, such as convolutional kernels [16] and probabilistic kernels [25], and can accurately predict whether a drug combination is likely to induce ADRs of interest. The experimental results demonstrate the advance of such single drug similarities that leverage co-medication patterns among high-order drug-drug interactions, and also inspire further exploration that learns such similarities in a pure data-driven fashion without pre-defined kernels, for example, via manifold learning. Further research would also include learning drug representations in a data-driven fashion such that the representations better quantify drug similarities in terms of their co-medication patterns. Deep learning would be an optimistic option for such drug representation learning.

Conclusions

In this manuscript, SVM-based classification methods were developed to predict whether a drug combination of arbitrary orders is likely to induce adverse drug reactions. Novel kernels over drug combinations of arbitrary orders were developed for such classification. These kernels were constructed from various single-drug information including drug co-medication patterns, and compare drug combination similarities based on single drugs they have and the relations among the single drugs. Specifically, a novel kernel over drug combinations of arbitrary orders was developed based on graph matching over drug combination graphs. A dataset from FDA Adverse Event Reporting System (FAERS) was constructed to test the new methods. The experimental results demonstrated that the new methods with drug co-medication based single drug similarities and graph matching based kernels achieve the best AUC as 0.912. The prediction also revealed strong patterns among drug combinations (e.g., statin enriched) that may be highly correlated with their induced ADRs.

List of abbreviations

DDI: Drug-Drug Interactions; ADR: Adverse Drug Reaction; SDS: Single Drug Similarities; ECFP: Extended Connectivity Fingerprints; and OR: Odds Ratio.

Declarations

Ethics approval and consent to participate

Not applicable

Consent for publication

Not applicable

Availability of data and material

The data and materials will be made publicly available upon the acceptance of the manuscript.

Competing interests

The authors declare that they have no competing interests.

Funding

This material is based upon work supported by the National Science Foundation under Grant Number IIS-1566219 and IIS-1622526. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the National Science Foundation.

Author's contributions

Wen-Hao Chiang implemented the methods and conducted the experiments. Li Shen, Lang Li and Xia Ning developed the methods and designed the experiments. Xia Ning analyzed the experimental results. Wen-Hao Chiang and Xia Ning wrote the manuscript.

Author details

¹Department of Computer & Information Science, Indiana University - Purdue University Indianapolis, 46202 Indianapolis, USA. Email: chiangwe@iupui.edu. ²Department of Biostatistics, Epidemiology and Informatics, University of Pennsylvania, 19104 Philadelphia, USA. Email: Li.Shen@penncmedicine.upenn.edu. ³Department of Biomedical Informatics, Ohio State University, 43210 Columbus, USA. Email: Lang.Li@osumc.edu. ⁴Department of Computer & Information Science, Indiana University - Purdue University Indianapolis, 46202 Indianapolis, USA. Email: xning@iupui.edu.

References

- Ramirez, E., Carcas, A.J., Borobia, A.M., Lei, S.H., Piñana, E., Fudio, S., Frias, J.: A pharmacovigilance program from laboratory signals for the detection and reporting of serious adverse drug reactions in hospitalized patients. *Clinical Pharmacology & Therapeutics* **87**(1), 74–86 (2010). doi:[10.1038/clpt.2009.185](https://doi.org/10.1038/clpt.2009.185)
- Percha, B., Altman, R.B.: Informatics confronts drug-drug interactions. *Trends in pharmacological sciences* **34**(3), 178–184 (2013). doi:[10.1016/j.tips.2013.01.006](https://doi.org/10.1016/j.tips.2013.01.006)
- The National Health and Nutrition Examination Survey. <http://www.cdc.gov/NCHS/NHANES.htm>
- Iyer, S.V., Harpaz, R., LePendu, P., Bauer-Mehren, A., Shah, N.H.: Mining clinical text for signals of adverse drug-drug interactions. *Journal of the American Medical Informatics Association* **21**(2), 353–362 (2014)
- Vilar, S., Uriarte, E., Santana, L., Lorberbaum, T., Hripcsak, G., Friedman, C., Tatonetti, N.P.: Similarity-based modeling in large-scale prediction of drug-drug interactions. *Nature protocols* **9**(9), 2147–2163 (2014)
- Hammann, F., Drewe, J.: Data mining for potential adverse drug–drug interactions. *Expert Opinion on Drug Metabolism & Toxicology* **10**(5), 665–671 (2014). doi:[10.1517/17425255.2014.894507](https://doi.org/10.1517/17425255.2014.894507). PMID: 24588496. <http://dx.doi.org/10.1517/17425255.2014.894507>
- Luo, H., Zhang, P., Huang, H., Huang, J., Kao, E., Shi, L., He, L., Yang, L.: Ddi-cpi, a server that predicts drug–drug interactions through implementing the chemical–protein interactome. *Nucleic Acids Research* (2014). doi:[10.1093/nar/gku433](https://doi.org/10.1093/nar/gku433)
- Harpaz, R., DuMouchel, W., Shah, N.H., Madigan, D., Ryan, P., Friedman, C.: Novel data-mining methodologies for adverse drug event discovery and analysis. *Clinical Pharmacology & Therapeutics* **91**(6), 1010–1021 (2012). doi:[10.1038/clpt.2012.50](https://doi.org/10.1038/clpt.2012.50)
- Ning, X., Schleyer, T., Shen, L., Li, L.: Pattern discovery from directional high-order drug-drug interaction relations. In: The 5th IEEE International Conference on Healthcare Informatics (2017). accepted
- Ning, X., Shen, L., Li, L.: Predicting high-order directional drug-drug interaction relations. In: The 5th IEEE International Conference on Healthcare Informatics (2017). accepted
- Conte, D., Foggia, P., Sansone, C., Vento, M.: Thirty years of graph matching in pattern recognition. *International journal of pattern recognition and artificial intelligence* **18**(03), 265–298 (2004)
- Foggia, P., Percannella, G., Vento, M.: Graph matching and learning in pattern recognition in the last 10 years. *International Journal of Pattern Recognition and Artificial Intelligence* **28**(01), 1450001 (2014)
- Caelli, T., Kosinov, S.: An eigenspace projection clustering method for inexact graph matching. *IEEE transactions on pattern analysis and machine intelligence* **26**(4), 515–519 (2004)
- Caetano, T.S., Caelli, T., Barone, D.A.C.: Graphical models for graph matching. In: Computer Vision and Pattern Recognition, 2004. CVPR 2004. Proceedings of the 2004 IEEE Computer Society Conference On, vol. 2, p. (2004). IEEE
- Messmer, B.T., Bunke, H.: A new algorithm for error-tolerant subgraph isomorphism detection. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **20**(5), 493–504 (1998)
- Haussler, D.: Convolution kernels on discrete structures. Technical Report UCS-CRL-99-10, University of California at Santa Cruz, Santa Cruz, CA, USA (1999)
- Joachims, T.: Making large-scale svm learning practical. Technical report, Technical Report, SFB 475: Komplexitätsreduktion in Multivariaten Datenstrukturen, Universität Dortmund (1998)
- Wale, N., Watson, I.A., Karypis, G.: Comparison of descriptor spaces for chemical compound retrieval and classification. *Knowl. Inf. Syst.* **14**(3), 347–375 (2008). doi:[10.1007/s10115-007-0103-5](https://doi.org/10.1007/s10115-007-0103-5)
- Scitegic Inc. <http://www.scitegic.com>
- Willett, P., Barnard, J.M., Downs, G.M.: Chemical similarity searching. *Journal of chemical information and computer sciences* **38**(6), 983–996 (1998)
- Burkard, R.E., Derigs, U.: The linear sum assignment problem, 1–15 (1980)
- Kuhn, H.W.: The hungarian method for the assignment problem. *Naval research logistics quarterly* **2**(1-2), 83–97 (1955)
- Saigo, H., Vert, J.-P., Ueda, N., Akutsu, T.: Protein homology detection using string alignment kernels. *Bioinformatics* **20**(11), 1682–1689 (2004). doi:[10.1093/bioinformatics/bth141](https://doi.org/10.1093/bioinformatics/bth141)
- Pizzuti, C., Ritchie, M.D., Giacobini, M.: Evolutionary computation, machine learning and data mining in bioinformatics: 7th european conference, evobio 2009 tübingen, germany, april 15-17, 2009 proceedings (2009)
- Zhou, S.K., Chellappa, R.: From sample similarity to ensemble similarity: probabilistic distance measures in reproducing kernel hilbert space. *IEEE Trans Pattern Anal Mach Intell* **28**(6), 917–929 (2006). doi:[10.1109/TPAMI.2006.120](https://doi.org/10.1109/TPAMI.2006.120)
- FDA Adverse Event Reporting System. <https://www.fda.gov/drugs/informationondrugs/ucm135151.htm>
- SIDER Side Effect Resource. <http://sideeffects.embl.de/se/C0026848/>
- Fawcett, T.: An introduction to roc analysis. *Pattern Recogn. Lett.* **27**(8), 861–874 (2006). doi:[10.1016/j.patrec.2005.10.010](https://doi.org/10.1016/j.patrec.2005.10.010)
- Min, M.R., Ning, X., Cheng, C., Gerstein, M.: Interpretable sparse high-order boltzmann machines. In: Artificial Intelligence and Statistics, pp. 614–622 (2014)
- Colmenares, E.W., Pappas, A.L.: Proton pump inhibitors: Risk for myopathy? *Annals of Pharmacotherapy* **51**(1), 66–71 (2017)
- Clark, D.W., Strandell, J.: Myopathy including polymyositis: a likely class adverse effect of proton pump inhibitors? *European journal of clinical pharmacology* **62**(6), 473–479 (2006)
- Joy, T.R., Hegele, R.A.: Narrative review: statin-related myopathy. *Annals of Internal Medicine* **150**(12), 858–868 (2009)
- Fernandez, G., Spatz, E.S., Jablecki, C., Phillips, P.S.: Statin myopathy: a common dilemma not reflected in clinical trials. *Cleve Clin J Med* **78**(6), 393–403 (2011)
- Wysowski, D.K., Chang, J.T.: Alendronate and risedronate: reports of severe bone, joint, and muscle pain. *Archives of internal medicine* **165**(3), 346–347 (2005)

Figures

Additional Files

Additional file 1 — Drug-Drug Interaction Prediction based on Co-Medication Patterns and Graph Matching (Supplementary Materials)
The additional file named as “supp.pdf” includes more experimental results on other datasets and it is provided in PDF format . Any PDF reader are recommended to view the file.

Tables

Table 1 Table of Notations

Notation	Description
d	Drug
D	Drug combination
\mathcal{G}	Complete graph for a drug combination
SDS_{2d}	Single drug similarity from drug 2d structures
SDS_{cm}	Single drug similarity based on co-medications
\mathcal{K}_{gm}	Kernel based on graph matching algorithm
\mathcal{K}_{cd}	Kernel from common drugs
\mathcal{K}_{ds}	Kernel from drug similarities
\mathcal{K}_{pb}	Probabilistic drug combination kernel

Table 2 Contingency Table

$\text{OR} = \frac{n_1/n_2}{m_1/m_2}$	ADR	no ADR	total
D	n_1	m_1	$n_1 + m_1$
$\setminus D$	n_2	m_2	$n_2 + m_2$
total	$n_1 + n_2$	$m_1 + m_2$	$n_1 + n_2 + m_1 + m_2$

In the table, n_1 is the number of events where D is taken with ADR occurring; m_1 is the number of events where D is taken without ADR occurring; n_2 is the number of events where D is not taken with ADR occurring; and m_2 is the number of events where D is not taken without ADR occurring, respectively.

Table 3 Data Statistics

dataset	stats	\mathcal{N}		\mathcal{M}	
		\mathcal{N}^-	\mathcal{N}^0	\mathcal{M}^0	\mathcal{M}^+
$\mathcal{D}_{\text{FAERS}}$	$\#\{D\}$	621,449	1,264	8,986	27,387
	$\#\{d\}$	1,209	417	881	1,201
	avgOrd	6.100	2.351	3.588	7.096
	avgFrq	1.761	225.317	13.730	1.402
	avgOR	-	0.546	16.343	-
\mathcal{D}^*	$\#\{D\}$	2,200	1,264	2,464	1,000
	$\#\{d\}$	562	417	692	679
	avgOrd	2.678	2.351	3.809	7.615
	avgFrq	42.082	225.317	20.565	5.520
	avgOR	-	0.546	31.998	-

In this table, “ $\#\{D\}$ ” and “ $\#\{d\}$ ” represent the number of drug combinations and the number of involved drugs, respectively. In each set of drug combinations, “avgOrd” is the the average order; “avgFrq” is the average frequency; and “avgOR” represents the average OR.

Table 4 Overall Performance Comparison

	\mathcal{K}		\mathcal{K}_{gm}		\mathcal{K}_{cd}		\mathcal{K}_{ds}		\mathcal{K}_{pb}	
	SDS	2d	cm	ord-1	ord-2	2d	cm	2d	cm	
acc	0.829	0.836	0.817	0.827	0.827	0.825	0.763	0.765		
pre	0.889	0.892	0.879	0.878	0.893	0.865	0.810	0.770		
rec	0.752	0.765	0.735	0.759	0.744	0.770	0.689	0.756		
F1	0.815	0.823	0.801	0.814	0.812	0.815	0.744	0.763		
AUC	0.898	0.912	0.907	0.909	0.900	0.900	0.843	0.853		

In this table, “acc”, “pre”, “rec”, “F1” and “AUC” represent accuracy, precision, recall, F1 and the area under a receiver operating characteristic curve, respectively. \mathcal{K}_{cd} with “ord-1” corresponds to $\mathcal{K}_{cd}^{(1)}$, and \mathcal{K}_{cd} with “ord-2” corresponds to $\mathcal{K}_{cd}^{(2)}$.

Table 5 Average Percentage (%) of \mathcal{D}_{Myo} Drugs ($\mathcal{K}_{gm}^{\text{cm}}$)

$\tilde{\mathcal{M}}^{10-}$	$\tilde{\mathcal{M}}^-$	\mathcal{M}	$\tilde{\mathcal{M}}^{10+}$	$\tilde{\mathcal{N}}^{10+}$	$\tilde{\mathcal{N}}^+$	\mathcal{N}	$\tilde{\mathcal{N}}^{10-}$
13.3	16.6	24.3	89.8	30.7	18.6	15.6	0.10

$\tilde{\mathcal{M}}^-$ and $\tilde{\mathcal{N}}^+$ are the sets of all mis-classified drug combinations in \mathcal{M} and \mathcal{N} by $\mathcal{K}_{gm}^{\text{cm}}$, respectively. $\tilde{\mathcal{M}}^{10-}$ is the set of the top-10 mis-classified drug combinations in $\tilde{\mathcal{M}}^-$ and $\tilde{\mathcal{N}}^{10+}$ is the set of the top-10 mis-classified drug combinations in $\tilde{\mathcal{N}}^+$ by $\mathcal{K}_{gm}^{\text{cm}}$. \mathcal{M}^{10+} and \mathcal{N}^{10-} are the sets of correctly classified drug combinations in \mathcal{M} and \mathcal{N} with the top-10 highest and lowest prediction values, respectively, by $\mathcal{K}_{gm}^{\text{cm}}$.

Table 6 Top mis-Classified \mathcal{N} Drug Combinations by $\mathcal{K}_{gm}^{\text{cm}}$

N	prd	frq	OR	combinations
1	2.696	26	-	atorvastatin fenofibrate rosiglitazone simvastatin
2	2.507	26	-	allopurinol amlodipine atorvastatin levothyroxine naproxen omeprazole simvastatin
3	1.878	22	-	acetylsalicylicacid atorvastatin bisoprolol clopidogrel ramipril simvastatin
4	1.855	27	-	acetylsalicylicacid atenolol atorvastatin furosemide lansoprazole lisinopril nitroglycerin
5	1.785	21	-	citalopram clozapine isosorbidedemononitrate prochlorperazine simvastatin zopiclone
6	1.750	- 0.842	-	amlodipine bisoprolol pravastatin ramipril simvastatin spironolactone warfarin
7	1.696	22	-	amlodipine clopidogrel ibuprofen omeprazole ramipril simvastatin
8	1.669	29	-	bisoprolol flecainide ramipril simvastatin
9	1.613	35	-	aripiprazole atorvastatin bendroflumethiazide clozapine diazepam folicacid furosemide iron lactulose lansoprazole perindopril ramipril trimethoprim zopiclone
10	1.549	- 0.875	-	lansoprazole omeprazole pantoprazole rabeprazole

In this table, “prd”, “frq” and “OR” represent prediction values, frequency and odds ratio, respectively. Drugs in \mathcal{D}_{Myo} are marked in **bold**.

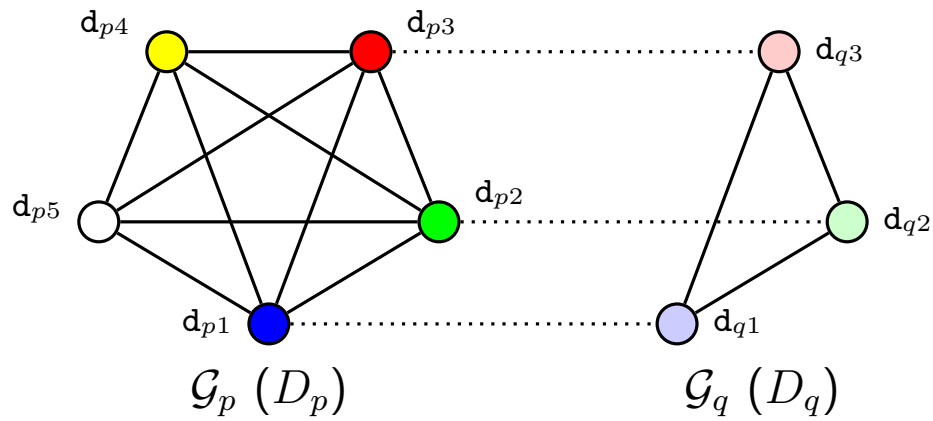


Figure 1 Graph matching for drug combinations

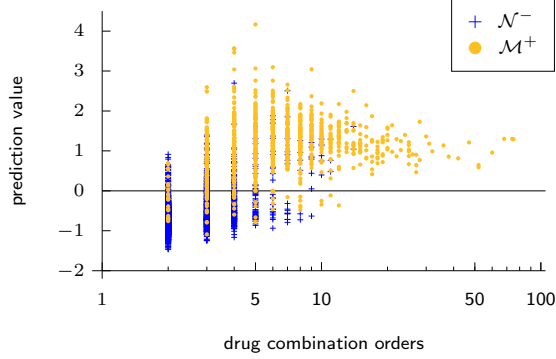


Figure 2 Orders vs. predictions in \mathcal{K}_{gm}^{cm}

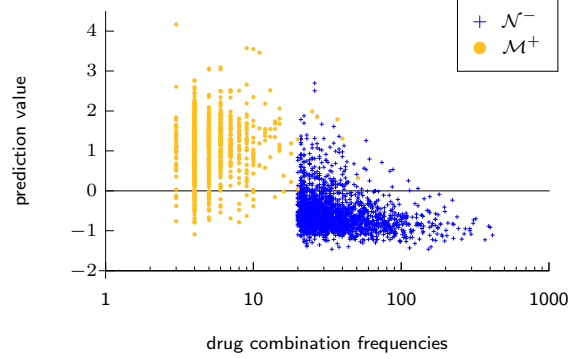


Figure 4 Frequencies vs. predictions in \mathcal{K}_{gm}^{cm}

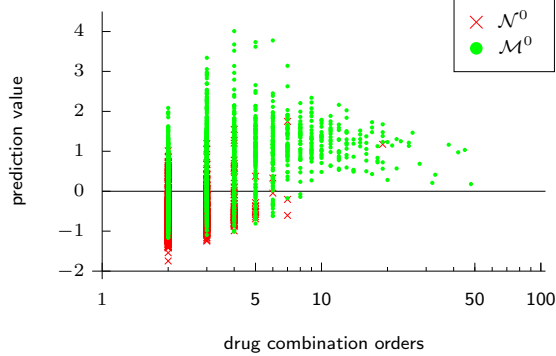


Figure 3 Orders vs. predictions in \mathcal{K}_{gm}^{cm}

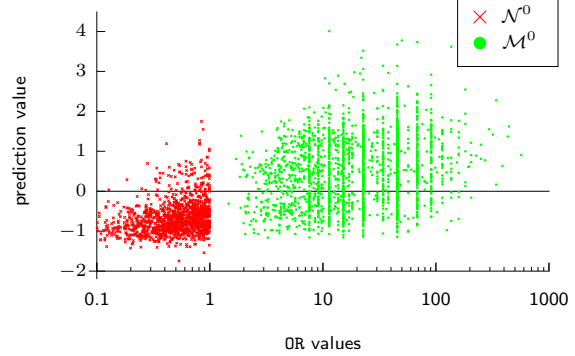


Figure 5 OR values vs. predictions in \mathcal{K}_{gm}^{cm}

Table 7 Top Predictions on \mathcal{M} by \mathcal{K}_{gm}^{cm}

N	prd	frq	OR	Combinations
1	4.167	3	-	atorvastatin lansoprazole pravastatin rosuvastatin simvastatin
2	4.009	-	11.372	atorvastatin pravastatin rosuvastatin simvastatin
3	3.776	-	50.043	atorvastatin fenofibrate metformin pravastatin rosuvastatin simvastatin
4	3.734	-	68.232	atorvastatin metformin pravastatin rosuvastatin simvastatin
5	3.676	-	45.487	atorvastatin lovastatin rosuvastatin simvastatin
6	3.618	-	136.470	atorvastatin pravastatin rosuvastatin simvastatin tadalafil
7	3.573	9	-	atorvastatin fenofibrate pravastatin simvastatin
8	3.552	10	-	atorvastatin ezetimibe fenofibrate rosuvastatin
9	3.519	-	22.746	atorvastatin ezetimibe rosuvastatin simvastatin
10	3.461	11	-	atorvastatin lansoprazole pravastatin simvastatin

In this table, “prd”, “frq” and “OR” represent prediction values, frequency and odds ratio, respectively. Drugs in \mathcal{D}_{Myo} are marked in **bold**.

Table 8 Top Predictions without \mathcal{D}_{Myo} Drugs by \mathcal{K}_{gm}^{cm}

N	prd	frq	OR	Combinations
1	2.083	4	-	calcium clonazepam colestipol prednisone teriparatide
2	1.992	-	45.487	alendronate anastrozole desloratadine hydrochlorothiazide lisinopril triamterene valdecoxib vitaminc
3	1.968	-	17.058	alendronate raloxifene risedronate teriparatide
4	1.960	-	90.978	alendronate amlodipine atenolol clonazepam raloxifene teriparatide
5	1.901	-	45.489	alendronate fexofenadine hydrochlorothiazide omeprazole prednisone risedronate triamterene
6	1.850	5	-	alendronate fexofenadine levothyroxine nabumetone oxybutynin
7	1.849	-	113.720	alendronate calcium esomeprazole ibandronate levothyroxine rabeprazole
8	1.843	-	7.581	alendronate calciumgluconate teriparatide
9	1.838	-	22.744	alendronate calcium levothyroxine raloxifene teriparatide
10	1.834	4	-	calcium escitalopram iron ketorolac raloxifene teriparatide

In this table, “prd”, “frq” and “OR” represent prediction values, frequency and odds ratio, respectively.

RESEARCH

Drug-drug interaction prediction based on co-medication patterns and graph matching (supplementary materials)

Wen-Hao Chiang¹, Li Shen², Lang Li³ and Xia Ning^{4*}

Performance on Other Datasets

Table 1 Data Statistics of \mathcal{D}^{2000} and \mathcal{D}^{4000}

dataset	stats	\mathcal{N}		\mathcal{M}	
		\mathcal{N}^-	\mathcal{N}^0	\mathcal{M}^0	\mathcal{M}^+
\mathcal{D}^{4000}	$\#\{D\}$	1,000	1,000	1,000	1,000
	$\#\{d\}$	1,000	369	596	679
	avgOrd	2.585	2.340	4.770	7.615
	avgFrq	62.545	264.898	5.957	5.520
	avgOR	-	0.451	60.994	-
	$\#\{D\}$	1,000	-	-	1,000
\mathcal{D}^{2000}	$\#\{d\}$	426	-	-	679
	avgOrd	2.585	-	-	7.615
	avgFrq	62.545	-	-	5.520
	avgOR	-	-	-	-

In this table, " $\#\{D\}$ " and " $\#\{d\}$ " represent the number of drug combinations and the number of involved drugs, respectively. In each set of drug combinations, "avgOrd" is the average order; "avgFrq" is the average frequency; and "avgOR" represents the average OR.

In \mathcal{D}^* , we retain the 1,000 drug combinations from \mathcal{M}^+ and the 1,000 drug combinations with the largest OR values from \mathcal{M}^0 . The labels of these drug combinations from \mathcal{M}^+ and \mathcal{M}^0 remain positive. From \mathcal{N}^- in \mathcal{D}^* , we retain the top 1,000 most frequent drug combinations. From \mathcal{N}^0 in \mathcal{D}^* , we retain the 1,000 drug combinations with the smallest OR values. The labels of these drug combinations from \mathcal{N}^- and \mathcal{N}^0 remain negative. The pruned dataset is denoted as \mathcal{D}^{4000} .

In \mathcal{D}^{4000} , we retain the 1,000 drug combinations from \mathcal{M}^+ and 1,000 drug combinations from \mathcal{N}^- and the labels remain positive and negative, respectively. This pruned dataset based on \mathcal{D}^{4000} is denoted as \mathcal{D}^{2000} . Table 1 presents the statistics of \mathcal{D}^{4000} and \mathcal{D}^{2000} .

Table 2 presents the overall performance comparison on dataset \mathcal{D}^{2000} and \mathcal{D}^{4000} . In both dataset, kernel \mathcal{K}_{gm} with SDS_{cm} outperforms others in accuracy, recall, F1 and AUC. Specifically, in accuracy, \mathcal{K}_{gm} with

Table 2 Overall Performance Comparison of \mathcal{D}^{2000}

\mathcal{K}	SDS	\mathcal{K}_{gm}		\mathcal{K}_{cd}		\mathcal{K}_{ds}		\mathcal{K}_{pb}	
		2d	cm	ord-1	ord-2	2d	cm	2d	cm
\mathcal{D}^{2000}	acc	0.929	0.938	0.897	0.910	0.924	0.912	0.845	0.860
	pre	0.954	0.959	0.941	0.942	0.957	0.917	0.839	0.860
	rec	0.902	0.916	0.847	0.873	0.888	0.907	0.854	0.869
	F1	0.927	0.937	0.891	0.906	0.921	0.912	0.846	0.861
	AUC	0.973	0.976	0.964	0.968	0.970	0.960	0.920	0.926
\mathcal{D}^{4000}	acc	0.896	0.899	0.867	0.882	0.887	0.880	0.817	0.821
	pre	0.933	0.939	0.928	0.930	0.944	0.903	0.838	0.816
	rec	0.853	0.854	0.795	0.826	0.822	0.851	0.787	0.829
	F1	0.891	0.894	0.856	0.875	0.879	0.876	0.812	0.822
	AUC	0.953	0.962	0.954	0.958	0.951	0.945	0.898	0.908

In this table, "acc", "pre", "rec" and "AUC" represent accuracy, precision, recall, and the area under a receiver operating characteristic curve, respectively. In \mathcal{K}_{cd} , "ord-1" is the kernel, in which the similarity between two drug combinations is calculated by the feature vectors, whose dimensions correspond to only one drug. As for "ord-2", the similarity between two drug combinations is calculated by the feature vectors, which have dimensions that correspond to different pairs of drugs.

SDS_{cm} outperforms the second best kernel \mathcal{K}_{gm} with SDS_{2d} at 0.97% and 0.33% in \mathcal{D}^{2000} and \mathcal{D}^{4000} , respectively. In recall, \mathcal{K}_{gm} with SDS_{cm} outperforms the second best kernel \mathcal{K}_{gm} with SDS_{2d} at 1.55% and 0.12% in \mathcal{D}^{2000} and \mathcal{D}^{4000} , respectively. In F1, \mathcal{K}_{gm} with SDS_{cm} outperforms the second best kernel \mathcal{K}_{gm} with SDS_{2d} at 1.08% and 0.34% in \mathcal{D}^{2000} and \mathcal{D}^{4000} , respectively. In AUC, \mathcal{K}_{gm} with SDS_{cm} outperforms the second best kernel \mathcal{K}_{gm} with SDS_{2d} at 0.31% in \mathcal{D}^{2000} and outperforms the second best kernel ord-2 \mathcal{K}_{cd} at 0.42% in \mathcal{D}^{4000} . In precision, the best kernel \mathcal{K}_{gm} with SDS_{cm} outperforms the second best kernel \mathcal{K}_{ds} with SDS_{2d} at 0.21% in \mathcal{D}^{2000} , whereas the best kernel \mathcal{K}_{ds} with SDS_{2d} outperforms the second best kernel \mathcal{K}_{gm} with SDS_{cm} at 0.53% in \mathcal{D}^{4000} . In three datasets, \mathcal{D}^{2000} , \mathcal{D}^{4000} and \mathcal{D}^* , in general, \mathcal{K}_{gm} with SDS_{cm} has the best performance compared to other kernels with a few exceptions. This may show the effectiveness to classify drug combinations by measuring the similarities between graphs, which represent drug combinations.

In kernel \mathcal{K}_{gm} and \mathcal{K}_{pb} , SDS_{cm} outperforms SDS_{2d} on average. This indicates the consistency in three

*Correspondence: xning@iupui.edu

⁴Department of Computer & Information Science, Indiana University - Purdue University Indianapolis, 46202 Indianapolis, USA. Email: xning@iupui.edu

Full list of author information is available at the end of the article

datasets of better capability to measure the similarities based on co-medication patterns than drug 2D structures. For kernel \mathcal{K}_{cd} , ord-2 outperforms ord-1 in both datasets as in \mathcal{D}^* . This consists with our observation before. That is, the ord-2 representation, which contains drug pairs as a feature, can emphasize the co-occurrence patterns in drug combinations.

Author details

¹Department of Computer & Information Science, Indiana University - Purdue University Indianapolis, 46202 Indianapolis, USA. Email: chiangwe@iupui.edu. ²Department of Biostatistics, Epidemiology and Informatics, University of Pennsylvania, 19104 Philadelphia, USA. Email: Li.Shen@pennmedicine.upenn.edu. ³Department of Biomedical Informatics, Ohio State University, 43210 Columbus, USA. Email: Lang.Li@osumc.edu. ⁴Department of Computer & Information Science, Indiana University - Purdue University Indianapolis, 46202 Indianapolis, USA. Email: xning@iupui.edu.