

Drug–drug interaction prediction with learnable size-adaptive molecular substructures

Arnold K. Nyamabo, Hui Yu, Zun Liu and Jian-Yu Shi

Corresponding authors: Hui Yu, School of Computer Science, Northwestern Polytechnical University, Xi'an 710072, China; E-mail: huiyu@nwpu.edu.cn, Jian-Yu Shi, School of Life Sciences, Northwestern Polytechnical University, Xi'an 710072, China; E-mail: jianyushi@nwpu.edu.cn

Abstract

Drug–drug interactions (DDIs) are interactions with adverse effects on the body, manifested when two or more incompatible drugs are taken together. They can be caused by the chemical compositions of the drugs involved. We introduce gated message passing neural network (GMPNN), a message passing neural network which learns chemical substructures with different sizes and shapes from the molecular graph representations of drugs for DDI prediction between a pair of drugs. In GMPNN, edges are considered as gates which control the flow of message passing, and therefore delimiting the substructures in a learnable way. The final DDI prediction between a drug pair is based on the interactions between pairs of their (learned) substructures, each pair weighted by a relevance score to the final DDI prediction output. Our proposed method GMPNN-CS (i.e. GMPNN + prediction module) is evaluated on two real-world datasets, with competitive results on one, and improved performance on the other compared with previous methods. Source code is freely available at <https://github.com/kanz76/GMPNN-CS>.

Key words: Drug–Drug Interaction; Substructure extraction; Substructure interaction; Multi-type interactions; Gated Message Passing Neural Network

Introduction

Drug–drug interactions (DDIs) refer to the phenomenon, with adverse effects on an organism, triggered when two or more drugs are taken together. The co-administration of more than one drug can be justified by the fact that some diseases are just too complex to be treated with a single drug [1], or that multiple diseases warrant multiple medications [2]. DDIs are therefore the risks, sometimes life-threatening [2], that come with the therapeutic benefits sought in multiple medications. The assessment of these risks have prompted many studies and research work aimed at identifying whether two or more given drugs are safe to be taken together.

The identification of DDIs is usually performed in pharmaceutical research/setting through extensive experimental testings (*in vitro*) and clinical trials. However, the huge number of combinations of drugs that should be considered for experimental testings makes this process highly expensive and quasi-impossible, even with high-throughput methods [3]. Computational methods (*in silico*) can thus be used as a cheap, yet effective and fast alternative to alleviate this problem by predicting potential DDIs based on the knowledge distilled from already known DDIs.

In the last few years, there has been significant progress with exciting results from the many learning based methods

Arnold K. Nyamabo is currently pursuing his master's degree in computer science at Northwestern Polytechnical University. He is interested in graph representation learning and applications.

Hui Yu received the master's and Ph.D. degrees from Northwestern Polytechnical University, Xi'an, China, where he is currently an Associate Professor. His research interest includes bioinformatics, machine learning and data mining.

Zun Liu received the master's and Ph.D. degrees from Northwestern Polytechnical University, Xi'an, China, where he is currently a lecturer. His research interest includes bioinformatics, machine learning and information security.

JianYu Shi received his Ph.D. degree from Northwestern Polytechnical University, Xi'an, China, where he is currently a Professor. His research interests include bioinformatics, cheminformatics and artificial intelligence.

Submitted: 27 July 2021; Received (in revised form): 12 September 2021

(i.e. both machine and deep learning methods) proposed for the identification of potential DDIs between combinations of two drugs, also known as DDI prediction task. However, most of these methods are limited in the way that they represent drugs as inputs, and perform the DDI prediction task.

Based on the medicinal chemistry knowledge [4], which states that a drug is simply an entity composed of different functional groups/chemical substructures which determine all of its pharmacokinetic (how it is handled by an organism) and pharmacodynamic (how it affects an organism) properties, and ultimately all of its interactions, we propose a simple, yet competitive computational method for DDI prediction (including specifying the type of the interaction) between two drugs equipped with this inductive bias. *First*, a drug is represented as a graph based on its molecular graph representation. In order to extract the substructures of a drug from its graph, we propose a message passing neural network where edges have learnable weights, constrained within the interval [0, 1]. These weights can be considered as gates to delimit the substructures, with the effect of producing flexible-sized and irregular-shaped substructures. *Second*, the DDI prediction between two drugs is based on the interaction scores between their learned substructures, each one weighted by a learnable weight using a co-attention mechanism [5] (or interaction map). As a byproduct, our method can give a hint at what substructures might be the cause of a DDI occurrence. *Third*, we evaluate our model on two real-world datasets: (1) Drugbank dataset which contains 1704 drugs, 191 808 DDI pairs and 86 interaction types. (2) Twosides dataset with 645 drugs, 4649 441 DDI pairs and 1317 interaction types. Experiments are conducted under two settings: transductive setting, where the training and test sets of DDI pairs share the same drugs; and inductive setting (also referred to as cold-start), where drugs in the two sets have to be different. The latter is more challenging than the former. Both datasets are used in the transductive setting where our method showed competitive results on DrugBank dataset, and outperformed previous methods on the Twosides dataset. In the inductive scenario, only the DrugBank dataset is considered, and our method shows better results.

Related work

This section discusses previous works in DDI prediction task from two perspectives: (1) How drugs are represented, and (2) how the actual prediction is performed.

Drug representation

The vast majority of existing DDI prediction methods represent drugs using molecular fingerprints [6–8], and/or other drug profiles such as side effects [6, 7], binding targets [8], transporters, enzymes, pathways or the combination of two or more of these features [7, 9, 10]. Molecular fingerprints [11, 12] are binary vectors whose elements indicate the presence (1) or absence (0) of a specific chemical substructure. The other profiles are similarly represented as binary vectors indicating the presence or absence of a particular profile, say, a specific side effect or binding target. Some methods [13–18] perform even further preprocessing by representing a drug as a similarity vector which indicates how similar it is to other drugs in the aforementioned representation spaces using similarity measures such as cosine similarity, Jaccard similarity. This is guided by the assumption that similar/dissimilar drugs are likely to have similar/dissimilar biological activities [17]. The downside with these representations is that they are hand-crafted, limited to the current

state of human knowledge. There is not enough flexibility to discover beyond what is encoded in them by a domain expert. For instance, with fingerprint representation, there is no way to discover beyond the chemical substructures that are already predefined, especially when dealing with new drugs. Additionally, some features such as side effects are not always available, especially in early drug development, and therefore impeding the methods that rely on them to be used.

In recent years, graph neural networks (GNNs) [19–22], deep learning models designed for graph-structured data, have been applied for learnable task-tailored representations of chemical molecules in general, and drugs in particular, with improved performance in molecule-related tasks [23–26]. However, this is usually optimized to learn the representation of a drug as a whole entity without giving too much consideration to the principal actors of DDIs, that is the functional groups/chemical substructures that constitute the drug molecule. Some of recent methods [27, 28] are proposed with the consideration of substructures involvement in DDIs. However, nodes' hidden representations (also referred to as patch representations) at each GNN layer are directly treated as substructure representations of the drugs. This approach produces substructures of regular shapes whose diameters/sizes are defined by the receptive field of the GNN layer. In this work, in order to extract substructures from a molecular graph, we propose a method that directly learns substructures of different sizes and shapes of the molecule.

DDI prediction

Approaches for DDI prediction can roughly be classified into two categories: one category where the drugs form a graph or network, and the other where they are considered independent from one another. In the latter, in order to perform the DDI prediction between two drugs, their representations are aggregated (e.g. summation, concatenation) and then fed into a linear or nonlinear classifier for prediction [10, 13, 29–31]. In the network-based category, drugs are assumed to form an interconnected system where the drugs are the nodes, and the edges can either represent the DDIs between the nodes [7, 14, 16, 17, 32, 33] or the similarities between the drugs, based on their representations [6, 18]. Different graph specific methods, including label propagation [6], matrix factorization [7, 8, 16], graph auto-encoders [18, 33] are applied to these derived networks to either conduct the prediction or first learn low-dimensional representations of the drugs and then perform the DDI prediction. The advantage of network-based methods is the addition of the drug interconnected system's topological information to the drug representations which can boost the performance. However, this approach does not work in inductive setting. In this work, our method considers drugs as independent entities, and therefore can be used both in inductive and transductive settings. Furthermore, we leverage the co-attention mechanism (same as in [27, 30]) between the learned substructures of a drug pair so that each drug can communicate to the other which substructures are really relevant. This has the effect that drugs are not completely handled independently.

Method

In this section, we mathematically formulate the problem we are trying to solve and present the building blocks of our method, including the input format and all involved computational steps. The overall framework is illustrated in Figure 1.

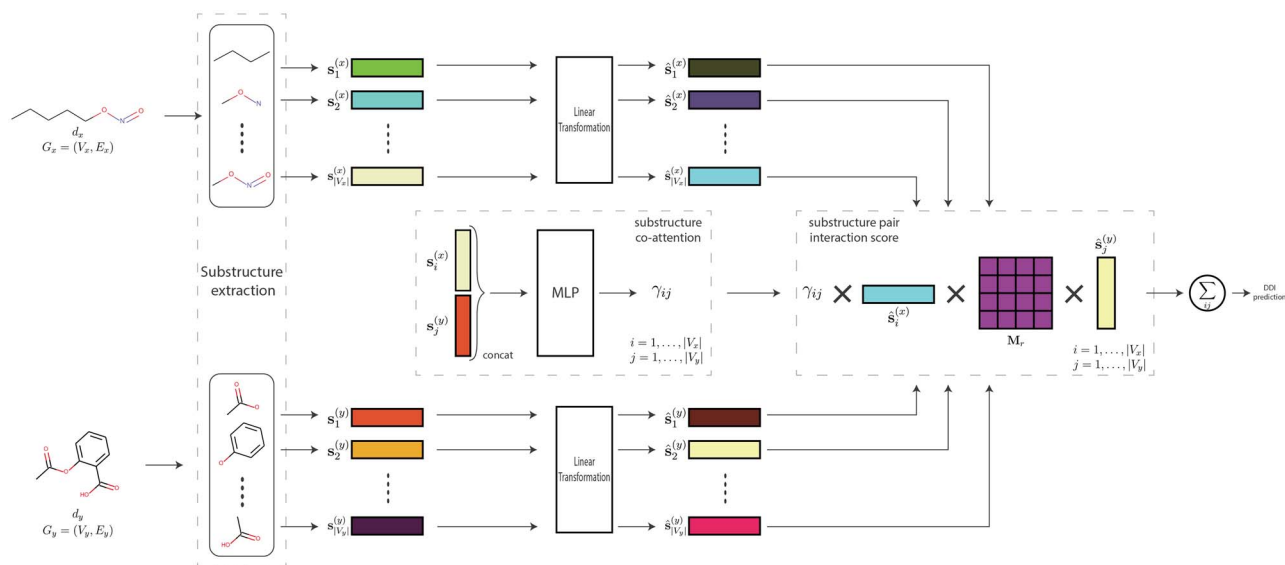


Figure 1. Overview of our method workflow. Given a DDI tuple (d_x, r, d_y) , d_x and d_y are represented in their molecular graph representation G_x and G_y , respectively, and r is represented as the learnable matrix $M_r \in \mathbb{R}^{b \times b}$. After substructure extraction (Section 3.3), G_x produces $|V_x|$ ($\{s_i^{(x)}\}_{i=1}^{|V_x|}$) substructures where each node is the center of a substructure. G_y undergoes the same process. A cross-substructure attention γ_{ij} is computed between these extracted substructures to learn how relevant they are to one another (Eq. 22). Each substructure undergoes a linear transformation, and the final DDI prediction is the sum of interaction scores between every pair of drugs d_x and d_y 's substructures, each one weighted by the co-attention weight γ_{ij} (Eq. 17).

Problem formulation

Given a set of drugs \mathcal{D} and a set of interaction types \mathcal{R} , the DDI prediction task can be regarded as a function $f: \mathcal{D} \times \mathcal{R} \times \mathcal{D} \rightarrow [0, 1]$. That is, given a triplet of two drugs and an interaction of a certain type, the task predicts the probability that this type of interaction will occur between the two drugs. The goal is to find an approximation of f , given a dataset of known DDIs $\mathcal{M} = \{(d_x, r, d_y)\}_{i=1}^N \subset \mathcal{D} \times \mathcal{R} \times \mathcal{D}$. The proposed method in this section is one such approximation.

Inputs

Drugs are represented as hydrogen-depleted undirected graphs $G = (V, E)$, where V is the set of nodes, representing atoms; $E \subset V \times V$ are the edges, representing the (covalent) bonds between the atoms. Each node v_i has a corresponding feature vector $\mathbf{x}_i \in \mathbb{R}^d$. Similarly, each edge $e_{ij} = (v_i, v_j)$ has a feature vector $\mathbf{x}_{ij} \in \mathbb{R}^d$. The features used for atoms and bonds are given in the Supplementary material Section 1. Note that at this initial stage, because the graph is undirected, the edges e_{ij} and e_{ji} are practically the same, and so are \mathbf{x}_{ij} and \mathbf{x}_{ji} .

Substructure extraction with gated message passing neural network

Given a general graph (i.e. not only a molecular graph), the representation of a substructure centered¹ at a node v_i can be regarded as the aggregation of all the nodes $\{v_j \mid v_j \in \mathcal{P}_{\rightarrow i}\}$ on all the paths (starting at end nodes²) in the graph that end at v_i . $\mathcal{P}_{\rightarrow i}$ means a path to node v_i . To refine the substructure even further, edges on the graphs can be considered as gates, that is,

with weight values constrained within $[0, 1]$. These gates control the flow of information along a path. A node v_j along a path to v_i is therefore weighted by the product of the weight values of the edges along the path that connect it to v_i . We assume that if there are many paths from v_j to v_i , which can happen if there is a cycle, each path is considered separately. Weighing a node by the product of edge values as proposed here has the end effect of producing substructures of different sizes and shapes. If along the path, an edge has weight value of ≈ 0 , the nodes on the rest of the path are cut off, therefore, delimiting the substructure in an irregular shape.

Each node in the graph can be treated as the center of a substructure in order to extract as many as possible substructures from the graph. However, this will require that each edge be transformed into a bidirectional edge (each direction with its own weight) because a node can be both a target node v_i (i.e. center) and a source node v_j . See Figure 2 for an illustrative example of the overall substructure extraction process.

Applying this concept to a drug for chemical substructure extraction, its molecular graphical representation $G = (V, E)$, which initially is undirected, is converted into a directed graph. Edges e_{ij} and e_{ji} become two separate edges. To highlight this difference, they are renamed $e_{i \rightarrow j}$ (edge from node v_i to node v_j) and $e_{j \rightarrow i}$, respectively. Each (directed) edge is assigned a learnable weight constrained within the range 0–1.

The generation of all the paths in the graph can be computationally inefficient; we therefore propose a message passing neural network (MPNN) [21] named gated message passing neural network (GMPNN) which can simulate this process. MPNN is a framework of multi-layer spatial convolutional GNNs. Each layer comprises three main components, namely, message passing (Eq. 1), aggregation (Eq. 2) and update (Eq. 3):

- 1 The node v_i is not literally at the center; it is simply the end node of all the paths.
- 2 These can be considered as peripheral nodes or nodes that are as far away as possible from the node being considered as the center.

$$\mathbf{m}_{j \rightarrow i}^{(l)} = M^{(l)}(\mathbf{h}_j^{(l-1)}, \mathbf{h}_i^{(l-1)}, \mathbf{q}_{j \rightarrow i}), \quad \forall j: v_j \in \mathcal{N}(v_i) \quad (1)$$

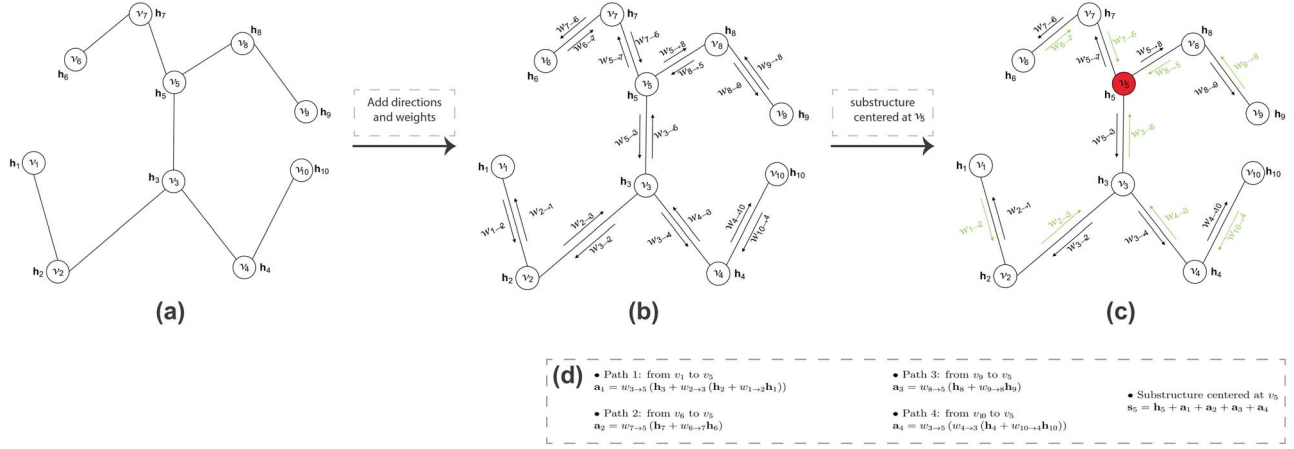


Figure 2. Example of substructure extraction on (a) a toy graph, v_i represents a node with associated feature vector h_i . (b) Every edge is transformed into a bidirectional one with weights w_{j-i} within $[0, 1]$ for each direction. (c) In order to extract the substructure centered at node v_5 (in red), all the paths (shown in green color) starting at end-nodes and ending at v_5 are generated. (d) The substructure centered at v_5 represented by vector s_5 is the aggregation of all the nodes along these paths, each node scaled by the product of the weights of the edges linking it to v_5 .

$$\mathbf{a}_i^{(l)} = A^{(l)} \left(\{\mathbf{m}_{j \rightarrow i}^{(l)}\}_{j: v_j \in \mathcal{N}(v_i)} \right) \quad (2)$$

$$\mathbf{q}_{j \rightarrow i}^{(l)} = U^{(l)} \left(\mathbf{q}_{j \rightarrow i}^{(l-1)}, \mathbf{a}_{j \rightarrow i}^{(l)} \right) \quad (6)$$

$$\mathbf{h}_i^{(l)} = U^{(l)} \left(\mathbf{h}_i^{(l-1)}, \mathbf{a}_i^{(l)} \right) \quad (3)$$

At the l^{th} iteration/layer ($l = 1, \dots, L$), in order to get the updated feature vector $\mathbf{h}_i^{(l)}$ of v_i , Eq. 1 flows information/message from its adjacent nodes ($\mathcal{N}(v_i)$) v_j 's feature vectors $\mathbf{h}_j^{(l-1)}$ of the previous iteration ($l-1$). This information can be constrained by or conditioned on node's v_i previous feature vector $\mathbf{h}_i^{(l-1)}$, and the features $\mathbf{q}_{j \rightarrow i}$, if any, of the edge from v_j and v_i , with the possibility of applying a nonlinear transformation $M^{(l)}$. In Eq. 2, all these messages are gathered using an aggregation function $A^{(l)}$ (e.g., sum, max, mean) to finally produce a newly updated feature vector $\mathbf{h}_i^{(l)}$ of v_i in Eq. 3. $U^{(l)}$ is an update function which can be as simple as the sum of its arguments, or a complex nonlinear function. In short, at each iteration, a node is passed messages (features) from its adjacent nodes. This has the effect that at iteration l , a node would be updated with the features of all the nodes that can be reached within a walk of length l . This is close to our idea of node aggregation along paths for substructure extraction. A walk in a graph has a redundancy in the nodes visited (a node can appear more than once in a walk, but at most once in a path); to get even closer to the path generation process, we will base our work on an MPNN variant, proposed in [26], named directed message passing neural network (D-MPNN). Here, in order to reduce the redundancy of nodes with standard MPNN, messages are passed between edges instead of nodes. The three components of MPNN in D-MPNN become:

$$\mathbf{m}_{k \rightarrow j}^{(l)} = M^{(l)} \left(\mathbf{h}_k, \mathbf{h}_j, \mathbf{q}_{k \rightarrow j}^{(l-1)} \right), \quad \forall k: v_k \in \mathcal{N}(v_j) \setminus \{v_i\} \quad (4)$$

$$\mathbf{a}_i^{(l)} = A^{(l)} \left(\{\mathbf{m}_{k \rightarrow i}^{(l)}\}_{k: v_k \in \mathcal{N}(v_i) \setminus \{v_i\}} \right) \quad (5)$$

The difference between MPNN and D-MPNN is that the former updates node features, while the latter updates edge features. In order to update edge $e_{j \rightarrow i}$, Eq. 4 passes messages from its neighboring edges $e_{k \rightarrow j}$, where the node v_j is the common vertex. Care is taken to remove the opposite direction to $e_{j \rightarrow i}$, that is, the edge $e_{i \rightarrow j}$ (hence, $\mathcal{N}(v_j) \setminus \{v_i\}$ in Eq. 4), so that information flows only in one direction, therefore reducing redundancy. After the final iteration L of edge features updates, nodes are represented as the aggregation of the features of all their incoming edges. For instance, node v_i 's final feature representation \mathbf{s}_i can be:

$$\mathbf{s}_i = \sum_{j: v_j \in \mathcal{N}(v_i)} \mathbf{q}_{j \rightarrow i}^{(L)} \quad (7)$$

Next, we present the computational steps of our message passing method GMPNN. Given a graph $G = (V, E)$ representing a drug:

- A nonlinear transformation is applied to the nodes for better feature representation:

$$\mathbf{h}_i = \text{MLP}_{\text{init}_n}(\mathbf{x}_i), \quad \forall v_i \in V \quad (8)$$

where $\text{MLP}_{\text{init}_n}(\cdot)$ is a multi-layer perceptron (MLP) for non-linear transformation. $\mathbf{h}_i \in \mathbb{R}^d$ is the updated feature vector of node v_i after transforming its original feature vector \mathbf{x}_i .

- The edge features are also transformed as follows:

$$\mathbf{h}_{j \rightarrow i} = \text{MLP}_{\text{init}_e}(\mathbf{x}_{ji}), \quad \forall e_{j \rightarrow i} \quad (9)$$

$$\mathbf{h}_{i \rightarrow j} = \text{MLP}_{\text{init}_e}(\mathbf{x}_{ij}), \quad \forall e_{i \rightarrow j} \quad (10)$$

$\mathbf{h}_{j \rightarrow i}, \mathbf{h}_{i \rightarrow j} \in \mathbb{R}^m$ are new feature vectors of edges $e_{j \rightarrow i}$ and $e_{i \rightarrow j}$, respectively. Note that even though the edges $e_{j \rightarrow i}$ and $e_{i \rightarrow j}$ are different, their features $\mathbf{h}_{j \rightarrow i}$ and $\mathbf{h}_{i \rightarrow j}$ are practically the same because $\mathbf{x}_{ji} = \mathbf{x}_{ij}$ (See Section 3.2).

- The weight (gate) $w_{j \rightarrow i} \in [0, 1]$ of edge $e_{j \rightarrow i}$ is initialized based on its incident nodes' features \mathbf{h}_j and \mathbf{h}_i , and its own features $\mathbf{h}_{j \rightarrow i}$:

$$\begin{aligned} o_{j \rightarrow i} &= \frac{1}{c} \left(\mathbf{h}_{j \rightarrow i}^T \text{MLP}_w(\mathbf{h}_j \parallel \mathbf{h}_i) \right) \\ w_{j \rightarrow i} &= \sigma(o_{j \rightarrow i}) \end{aligned} \quad (11)$$

where $\sigma(\cdot)$ is the sigmoid function used to constrain $w_{j \rightarrow i}$ within $[0, 1]$; T is the transposition operation; \parallel is the concatenation operation. $c(> 0)$ is a constant used to avoid saturation with gradient flow when using sigmoid function. We have set it to the degree of node v_j , the tail of the edge $e_{j \rightarrow i}$. Note that the edges $e_{j \rightarrow i}$ and $e_{i \rightarrow j}$ have different weights (i.e. $w_{j \rightarrow i} \neq w_{i \rightarrow j}$), because \parallel is not commutative.

- Now that we have node feature vectors \mathbf{h}_i and edge weights $w_{j \rightarrow i}$, we are ready to apply the directed message passing mechanism for the simulation of aggregation of nodes along paths for substructure extraction. Since message is passed between edges instead of nodes, we propose to promote nodes' features to edge level. Edge $e_{j \rightarrow i}$'s new features become:

$$\mathbf{q}_{j \rightarrow i}^{(0)} = w_{j \rightarrow i} \mathbf{h}_j \quad (12)$$

where $\mathbf{q}_{j \rightarrow i}^{(0)} \in \mathbb{R}^f$. That is, $e_{j \rightarrow i}$ takes on the features \mathbf{h}_j of its tail node v_j scaled by its weight $w_{j \rightarrow i}$. Note that $e_{j \rightarrow i}$ previous features $\mathbf{h}_{j \rightarrow i}$ (Eq. 9) is different from $\mathbf{q}_{j \rightarrow i}^{(0)}$. The former was simply used to compute the weight $w_{j \rightarrow i}$ (Eq. 11), whereas the latter (containing node information) will be used in the actual message passing for node aggregation along paths as intended.

- The message passing, aggregation and update components of our proposed method are therefore defined as:

$$\mathbf{m}_{k \rightarrow j}^{(l)} = w_{j \rightarrow i} \mathbf{q}_{k \rightarrow j}^{(l-1)}, \quad \forall k : v_k \in \mathcal{N}(v_j) \setminus \{v_i\} \quad (13)$$

$$\mathbf{a}_{j \rightarrow i}^{(l)} = \sum_{k: v_k \in \mathcal{N}(v_j) \setminus \{v_i\}} \mathbf{m}_{k \rightarrow j}^{(l)} \quad (14)$$

$$\mathbf{q}_{j \rightarrow i}^{(l)} = \mathbf{q}_{j \rightarrow i}^{(0)} + \mathbf{a}_{j \rightarrow i}^{(l)} \quad (15)$$

Contrary to Eq. 4, Eq. 13 does not apply any transformations to feature vectors at each iteration l during message passing because we want to keep all the nodes that compose a substructure in the same feature space. However, at each iteration, features are scaled by the weight of the edge; this has the end effect of multiplying node features by the product of the weights of the edges linking it to the center node of the substructure. As an aside, this paradigm of not applying any transformations during message passing can also be regarded as an instance of the concept of simplifying graph convolutions [34].

- The final representation of a node v_i , after the last iteration L , which captures the substructure information of which it is the center is given by:

$$\mathbf{s}_i = f_{\text{sub}} \left(\mathbf{h}_i + \sum_{j: v_j \in \mathcal{N}(v_i)} \mathbf{q}_{j \rightarrow i}^{(L)} \right) \quad (16)$$

where $\mathcal{N}(v_i)$ is the set of nodes adjacent to v_i , and $f_{\text{sub}}(\cdot)$ is a nonlinear function implemented as an MLP. To obtain $\mathbf{s}_i \in \mathbb{R}^b$, we simply aggregate all the features $\mathbf{q}_{j \rightarrow i}^{(L)}$ from all of its incoming edges $e_{j \rightarrow i}, \forall j : v_j \in \mathcal{N}(v_i)$. \mathbf{s}_i is the vector representation of the learned/extracted substructure centered at v_i . Note that initially, the feature vector of node v_i (i.e. \mathbf{h}_i) contained only the information of a single atom, but now \mathbf{s}_i contains the features representing a substructure centered at that atom.

The application of our proposed method for substructure extraction in inductive setting is made possible by the fact that the values of edge weights (Eq. 11) depend only on node (atom) and edge (bond) features. If we encounter a new drug, the substructure extraction operation can still be performed because all the molecules share the same type of atoms and bonds.

DDI prediction

Given a DDI tuple (d_x, r, d_y) , the DDI prediction is determined as the join probability of the tuple:

$$P(d_x, r, d_y) = \sigma \left(\sum_i \sum_j |V_x| |V_y| \gamma_{ij} \hat{\mathbf{s}}_i^{(x)T} \mathbf{M}_r \hat{\mathbf{s}}_j^{(y)} \right) \quad (17)$$

- $\sigma(\cdot)$ is the sigmoid function.
- $\hat{\mathbf{s}}_i^{(x)}$ and $\hat{\mathbf{s}}_j^{(y)}$ are linear transformations of the substructure $\mathbf{s}_i^{(x)}$ and $\mathbf{s}_j^{(y)}$, respectively.

$$\hat{\mathbf{s}}_i^{(x)} = \mathbf{W}_{(x)} \mathbf{s}_i^{(x)}, \quad i = 1, \dots, |V_x| \quad (18)$$

$$\hat{\mathbf{s}}_j^{(y)} = \mathbf{W}_{(y)} \mathbf{s}_j^{(y)}, \quad j = 1, \dots, |V_y| \quad (19)$$

where $\mathbf{W}_{(x)}$, and $\mathbf{W}_{(y)} \in \mathbb{R}^{b \times b}$ are learnable transformation matrices.

- $\mathbf{M}_r \in \mathbb{R}^{b \times b}$ is the learnable representation matrix of the interaction type of r . To reduce the number of parameters, we constrain it to be a diagonal matrix.

$$\mathbf{M}_r = \text{diag}(\mathbf{m}_r) \quad (20)$$

$\text{diag}(\cdot)$ generates a diagonal matrix where \cdot is the diagonal and $\mathbf{m}_r \in \mathbb{R}^b$ is a learnable vector specific to the type of interaction r .

- Because the Twosides dataset defines multiple existing interactions between two given drugs (See Section 4.1 for more details), we redefine \mathbf{M}_r as follows for this dataset:

$$\mathbf{M}_r = \text{diag}(f_{\text{pred}}(\mathbf{m}_r)) \quad (21)$$

where f_{pred} is a nonlinear function implemented as an MLP to encourage similar or commonly co-occurring interaction types to have similar representations.

- $\gamma_{ij} \in [0, 1]$ is the cross-substructure interaction weight, also known as co-attention, between substructures (Eq. 16) $\mathbf{s}_i^{(x)}$ of drug d_x and $\mathbf{s}_j^{(y)}$ of d_y .

$$\gamma_{ij} = \text{softmax} \left(\text{MLP}_{\gamma}(\mathbf{s}_i^{(x)} \parallel \mathbf{s}_j^{(y)}) \right), \quad i = 1, \dots, |V_x| \quad (22)$$

$$j = 1, \dots, |V_y|$$

Here again to account for the multiplicity of interactions between two drugs in the Twosides dataset, for this dataset we propose:

$$\gamma_{ij} = \text{softmax} \left(\text{MLP}_{\gamma} \left(\mathbf{s}_i^{(x)} \parallel \mathbf{s}_j^{(y)} \parallel f_r(\mathbf{m}_r) \right) \right) \quad (23)$$

where f_r is an MLP and \mathbf{m}_r is the same as in Eq. 21. The goal is to have a co-attention score aware of the interaction type which is being considered.

The DDI prediction can therefore be considered as a binary prediction of a DDI tuple. Since only known DDIs are given in the dataset \mathcal{M} (See Section 3.1), they are considered as positive samples; negative samples are generated by corrupting either d_x or d_y . That is, given a known DDI tuple (d_x, r, d_y) , its derived negative sample is generated by either replacing d_x or d_y . We follow the strategy proposed in [35] for negative sample generation. The learning process of the whole model is done by minimizing the binary cross-entropy loss function given as:

$$\mathcal{L} = -\frac{1}{|\mathcal{M}|} \sum_{i:(d_x, r, d_y)_i \in \mathcal{M}} (\log(p_i) + \log(1 - p'_i)) \quad (24)$$

$|\mathcal{M}|$ is the number of DDI tuples in the dataset, p_i is the probability (Eq. 17) of a known DDI tuple and p'_i is the probability of its associated negative sample.

Experiments

Datasets

Two real-world datasets, downloaded using the `tdc` python library³ [36], were used for the evaluation of our method.

- **DrugBank**: sourced from FDA/Health Canada drug labels, it contains 191 808 DDI tuples with 1706 drugs. Each drug is represented in SMILES from which molecular graphical representations are generated using the python library RDKit⁴. There are 86 interaction types describing how one drug affects the metabolism of another one. For example, *the excretion of Acamprosate can be decreased when combined with Acetylsalicylic acid (Aspirin)*. Each DDI pair is considered as a positive sample from which a negative sample was generated as mentioned in Section 3.4. In this dataset, there is only one interaction for each DDI tuple, that is, there are no two distinct tuples with the same drug pair but different interactions.

- **Twosides**: proposed by [37] after filtering the original TWO-SIDES side effects data [2]. It contains 4649 441 DDI triplets with 645 drugs, and 1317 interaction types. As opposed to the DrugBank dataset, these interactions are rather at the phenotypic level than metabolic. That is, here, interactions are simply side effects, such as *headache*, *pain in throat*. Furthermore, given two drugs, there can exist many such interactions between them, contrary to DDI tuples in DrugBank. As in [37], this dataset is further preprocessed by removing interaction types that occur in less than 500 DDI tuples in order to work only with commonly occurring types, thus, remaining with 963 interactions types, and 4576 287 DDI tuples.

The distributions of the DDI types in both datasets are given in the Supplementary materials Section 4.

Setup

Our model, named GMPNN-CS, was implemented in Pytorch⁵ [38] and Pytorch Geometric⁶ [39]. We used random search for hyper-parameters fine-tuning and decided on the best values based on the overall performance on validation set. We considered the following hyper-parameter settings. The number of message passing iterations L was searched from {3, 5, 7, 10, 15}; dimensions of \mathbf{h}_i (Eq. 8) and \mathbf{s}_i (Eq. 16) were searched from {64, 128}. The model was trained on mini-batches of 512 DDI tuples using the Adam optimizer [40] with a learning lr rate tuned from {1e-2, 1e-3, 1e-4}. Additionally, an exponentially decaying scheduler of 0.96^t (where t is the current epoch) was set on the learning rate. We found that the combination of $L = 10$, $\mathbf{h}_i \in \mathbb{R}^{64}$, $\mathbf{s}_i \in \mathbb{R}^{128}$, and $lr = 1e-3$ produced the best performance. The implementation details of MLP_{init-n} , MLP_{init-e} , MLP_w , f_{sub} , MLP_{γ} , f_r , and f_{pred} are given in the Supplementary materials Section 2.

Baselines

We compared our model with state-of-the-art methods, which similarly (1) work with molecular graphs as inputs; (2) integrate joint drug-drug information in some way during the learning process; (3) consider the involvement of substructures in DDI interaction prediction; and/or (4) work both in transductive and inductive settings.

- **MR-GNN** [28]: uses the representation at each graph convolution layer of nodes to capture substructures of different size for each drug. These representations are jointly fed into a recurrent neural network for a joint representation of a pair of drugs for DDI prediction.
- **MHCADDI** [30]: uses co-attention mechanism to integrate joint drug-drug information during the representation learning of individual drugs.
- **SSI-DDI** [27]: considers each node hidden features as substructures and then computes interactions between these substructures to determine the final DDI prediction.
- **GAT-DDI**: baseline we implemented using graph attention networks (GAT) [19] for drug representations which are directly used for DDI prediction.
- **GMPNN-U**: variant of our proposed method GMPNN-CS where the co-attention coefficient γ_{ij} (Eq. 22) is simply uniform, that is, $\gamma_{ij} = (|V_x||V_y|)^{-1}$.

3 <https://tdcommons.ai/>

4 <https://www.rdkit.org/>

5 <https://pytorch.org/>

6 <https://pytorch-geometric.readthedocs.io/>

We (re-)implemented these methods in Pytorch, some with little modifications from the original work for fair comparison and better performance. See Supplementary material Section 3 for details.

Results

Experimental results are presented in following metrics: the accuracy (ACC), the area under the receiver operating characteristic (AUC), the average precision (AP), the F1 score, the precision (P) and the recall (R).

Transductive setting

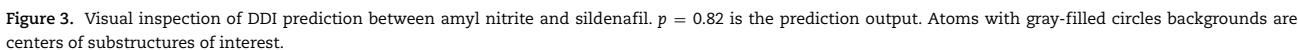
In transductive setting, the drugs used during training can also appear in the test set. In this setting, we split the datasets randomly based on DDI tuples. We performed a stratified split on both datasets based on the interaction types in order to keep the same proportions of interaction types in the training (60% of the data), validation (20%) and test (20%) sets. We did this three times, resulting in three stratified randomized folds. For each DDI tuple, a negative sample is generated as discussed in Section 3.4. They were generated before training to ensure that all the methods are trained on the same data. Each model, including our proposed method and all the baseline methods, were trained and tested on each one of these three stratified folds. The means and standard deviations of the results of each model from these three experiments are reported in Table 1. Improvement in performance in each metric score (without considering standard deviations) of our method GMPNN-CS with respect to the highest score of the baseline methods is shown in the bottom row of the table. On the DrugBank dataset, our method did not perform the best as we can see a decrease in scores. We can also see that the results of GMPNN-CS and its variant GMPNN-U are very close, giving the impression that co-attention did not work. An explanation for this behavior is given in Section 5. However, on the Twosides dataset there is an improvement with a significant margin over the other methods. Here, there is a big difference between GMPNN-U and GMPNN-CS because this is a dataset with multiple interactions between a pair of drugs. Due to computational limitation, we could not perform experiments with MHCADDI on this dataset; we only report the results (only AUC available) from the original paper. GAT-DDI does not seem to work well on this dataset, behaving just like a random classifier with scores of about 50%. This can be due to gradient vanishing/exploding or oversmoothing problem. It is not considered in the improvement computation (last row of Table 1) on this dataset. Furthermore, GMPNN-U is similar to GAT-DDI, except for the message passing component. The former uses our proposed GMPNN, whereas the latter uses GAT. Experimental results show that our message passing method performs better than GAT on both datasets. Furthermore, performance on each DDI type on both datasets is presented the Supplementary materials Section 5.

Inductive setting

Contrary to transductive setting, here the dataset is split based on the drugs. That is, the DDI tuples in training and test sets do not have overlapping drugs. This approximates a real-world scenario where there is a new drug for which there is no known prior associated drug interactions. It is also referred to as cold-start scenario in the literature. It is more challenging than the transductive setting. In the latter, the model only learns to generalize to unseen DDI tuples (with all the drugs already known

Table 1. Comparative evaluation (mean \pm std) in % in transductive setting. Best performance in each metric is shown in bold font. The last row shows the improvement in performance on the DDI prediction task in each metric by our method. It is the difference between our method's performance and the best performance among baseline methods. Negative values, therefore, mean our method did not perform the best, and positive ones indicate it performed the best. *GAT-DDI was not considered in the computation.

	DrugBank					Twosides				
	ACC	AUC	AP	P	R	ACC	AUC	AP	P	R
MR-GNN	96.04 \pm 0.05	98.87 \pm 0.04	98.57 \pm 0.06	94.48 \pm 0.08	97.78 \pm 0.03	76.23 \pm 0.23	85.00 \pm 0.22	84.32 \pm 0.35	72.82 \pm 0.44	83.70 \pm 0.39
MHCADDI	83.80 \pm 0.27	91.16 \pm 0.31	89.26 \pm 0.37	78.90 \pm 0.06	92.26 \pm 0.63	-	88.20	-	-	-
SSI-DDI	96.33 \pm 0.09	98.95 \pm 0.08	98.57 \pm 0.14	95.09 \pm 0.08	97.70 \pm 0.14	78.20 \pm 0.14	85.85 \pm 0.13	82.71 \pm 0.14	74.33 \pm 0.21	86.15 \pm 0.15
GAT-DDI	89.81 \pm 1.00	95.21 \pm 0.70	93.56 \pm 0.90	87.04 \pm 1.11	93.56 \pm 0.52	50.00	50.00	50.00	50.00	100
GMPNN-U	95.00 \pm 0.10	98.32 \pm 0.04	97.77 \pm 0.06	93.19 \pm 0.15	97.07 \pm 0.06	74.78 \pm 0.04	82.08 \pm 0.02	78.67 \pm 0.03	71.77 \pm 0.09	81.69 \pm 0.33
GMPNN-CS(Ours)	95.30 \pm 0.05	98.46 \pm 0.01	97.94 \pm 0.02	93.60 \pm 0.07	97.22 \pm 0.1	82.83 \pm 0.14	90.07 \pm 0.12	87.24 \pm 0.12	78.42 \pm 0.11	90.61 \pm 0.23
Improvement	-1.03	-0.49	-0.63	-1.49	-0.56	+4.63	+1.87	+2.92	+4.09	+4.46*



	ACC	AUC	AP	F1
\mathcal{M}_{S1} Partition (new drug \leftrightarrow new drug)				
MR-GNN	62.63 \pm 0.77	70.92 \pm 0.84	73.01 \pm 1.23	45.81 \pm 2.51
MHCADDI	66.50 \pm 0.62	72.53 \pm 0.92	71.06 \pm 1.61	67.21 \pm 0.59
SSI-DDI	65.40 \pm 1.30	73.43 \pm 1.81	75.03 \pm 1.42	54.12 \pm 3.46
GAT-DDI	66.31 \pm 0.61	72.75 \pm 0.78	71.61 \pm 1.00	68.68\pm0.60
GMPNN-U	67.90 \pm 0.50	73.76 \pm 0.90	74.58 \pm 1.06	63.73 \pm 0.80
GMPNN-CS(Ours)	68.57\pm0.30	74.96\pm0.40	75.44\pm0.50	65.32 \pm 0.23
\mathcal{M}_{S2} Partition (new drug \leftrightarrow old drug)				
MR-GNN	74.67 \pm 0.33	83.15 \pm 0.60	83.81 \pm 0.69	69.88 \pm 0.86
MHCADDI	70.58 \pm 0.94	77.84 \pm 1.08	76.16 \pm 1.45	72.74 \pm 0.65
SSI-DDI	76.38 \pm 0.92	84.23 \pm 1.05	84.94\pm0.76	73.54 \pm 1.50
GAT-DDI	69.83 \pm 1.41	77.29 \pm 1.63	75.79 \pm 1.95	73.01 \pm 0.85
GMPNN-U	77.00 \pm 0.60	83.92 \pm 0.50	84.14 \pm 0.74	77.29 \pm 0.50
GMPNN-CS(Ours)	77.72\pm0.30	84.84\pm0.15	84.87 \pm 0.40	78.29\pm0.16

chemical structure) are concerned [27]. Thus, drugs in \mathcal{D}_{new} and \mathcal{D}_{old} are not only different, but also share very little structurally.

Here are some visual examples of DDI prediction on DrugBank, which can give hints at what might be the cause of a DDI. First, pairs of substructures with top values of γ_{ij} (Eq. 22) are retrieved. Secondly, the weights of the edges of a substructure in the rendered figures are redefined as follows: for instance, if v_1 is the center of the substructure, and $v_1 \xleftarrow{w_{2,1}} v_2 \xleftarrow{w_{3,2}} v_3 \xleftarrow{w_{4,3}} v_4$ is one of the paths constituting the substructure, where $w_{j,i} (= w_{i \rightarrow j})$ (Eq. 11) is the learned weight of edge $e_{j \rightarrow i}$, in the visual examples (Figure 3-4), edges have values of the products of all the edge values before them in the path and their own values, that is, $v_1 \xleftarrow{w_{2,1}} v_2 \xleftarrow{w_{2,1}w_{3,2}} v_3 \xleftarrow{w_{2,1}w_{3,2}w_{4,3}} v_4$. If $w_{2,1} = 0.9, w_{3,2} = 0.7$ and $w_{4,3} = 0.6$, then we have $v_1 \xleftarrow{0.9} v_2 \xleftarrow{0.63=0.9 \times 0.7} v_3 \xleftarrow{0.378=0.9 \times 0.7 \times 0.6} v_4$. For simplicity, if an edge has more than one neighboring edges (i.e. appears in more than one path), we take the maximum value. Thirdly, these values are shown on corresponding edges and used as intensities to highlight (in green) edges in the figures, and the center nodes of substructures are highlighted with gray-filled circles. Figure 3 shows the drug sildenafil, a phosphodiesterase-5 (PDE5) inhibitor used for erectile

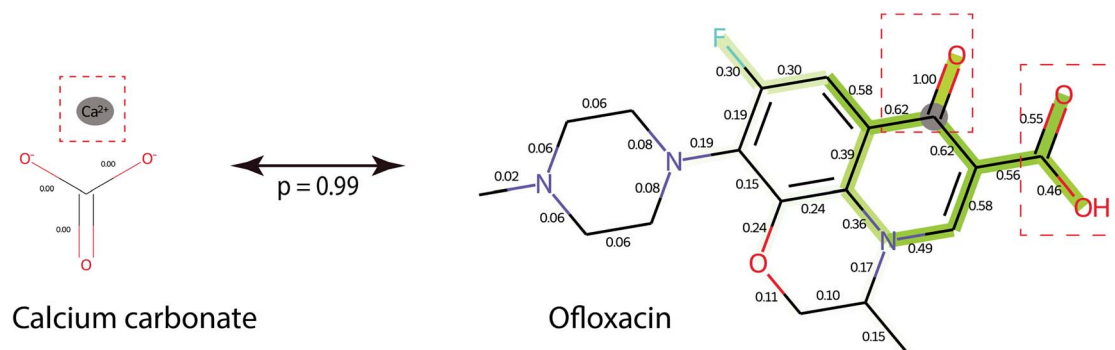


Figure 4. Visual inspection of DDI prediction between calcium carbonate and ofloxacin.

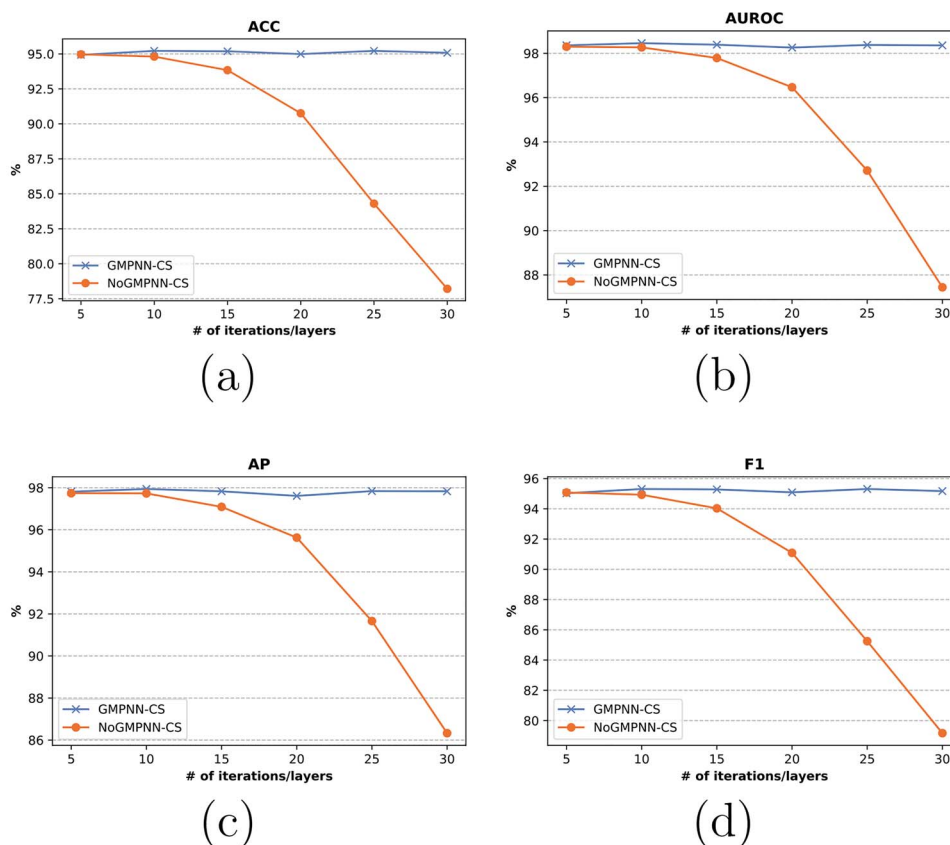


Figure 5. Comparison of performance between our method GMPNN-CS, and its variant without edge gate values, NoGMPNN-CS, or equivalently with constant gate values all set to 1. Contrary to NoGMPNN-CS, GMPNN-CS can be extended to many iterations (layers) with constant performance.

dysfunction, and amyl nitrite, a nitrate drug. It is contraindicated to take a PDE5 inhibitor with a nitrate drug simultaneously because it can cause a reduction in the blood pressure [29, 42]. We can see that nitrate group (in the dashed red box) of Amyl nitrite is very much involved in the DDI prediction outcome. Figure 4 is the example of ofloxacin (a fluoroquinolone) and calcium carbonate (an antacid). The carbonyl oxygen and the carboxylic acid of ofloxacin (both shown within red dashed boxes) can form a chelate with a metal ion (in this case the calcium of the calcium carbonate). Chelates have very poor water solubility and can therefore cause a significant reduction in the absorption of ofloxacin in the body [4].

Discussion and limitations

Every node/atom is considered the center of a substructure, and therefore there are as many substructures as there are nodes. Nodes that are adjacent to each other end up being centers to substructures that are similar, causing redundancy. This has the negative effect that γ_{ij} (Eq. 22) is overused by only a group of very similar substructures, with the tendency of the former to be uniform within groups of (adjacent) substructures. This can explain why, especially in the case of Drugbank, there is no much difference between GMPNN-U and GMPNN-CS (See Section 4.3.1). In future work, we think that, to solve this issue,

a clustering or pooling algorithm might be used to pool similar substructures together and retain only one representative substructure.

The importance of edge weights as gates for substructure extraction only becomes apparent when we increase the number of iterations L . Figure 5 shows the comparison between GMPNN-CS and a variant (NoGMPNN-CS) where we remove the weights altogether (or equivalently, the weights are all set to 1). We can see that with 5–10 iterations of message passing the difference is quite small, but becomes significant as the number of iterations increases. As an aside, these results also demonstrate the ability of our proposed method to be extended into deeper GNNs without degradation in performance. Additionally, drugs are mainly organic molecules with the majority of atoms being carbon, producing graphs where the majority of nodes are alike. This can affect the computation of edge weights (Eq. 11). For future work, we think additional information such as spatial location of atoms might be useful to make the difference more prominent.

The poor performance of our method on the DrugBank dataset (Table 1) can be linked to the imbalance state of this dataset as shown in the Supplementary materials Section 4 (Figure 2). Our method has difficulty generalizing on DDI types with very low frequency. In future work, we aim to investigate this issue further in order to come up with an effective approach for handling DDI types with very low frequencies.

Conclusion

We proposed GMPNN-CS, a computational method for DDI prediction. GMPNN-CS learns substructures of different sizes and shapes of drugs in order to infer whether a pair of drugs can cause a DDI based on their chemical substructures. We demonstrated empirically the effectiveness of GMPNN-CS using two real-world datasets, with significant performance improvement in one of them. Experiments are conducted in both transductive and inductive settings. A visual inspection (as conducted in the experiments) of substructure extraction and their involvement in a DDI prediction can be used as hint by both expert and non-expert users to interpret the results of a prediction.

Key Points

- Using chemical substructures of drugs to predict drug–drug interaction between a pair of drugs by also giving the specific type of the interaction (multi-type prediction).
- Proposing a message passing neural network named gated message passing neural network (GMPNN) for learning substructures of various (or adaptive) sizes and shapes from molecular graphs of drugs.
- Using co-attention mechanism to learn the relevance of each pairwise substructure interaction involvement in a drug–drug interaction.
- Proposed method able to be applied both in transductive and inductive (or cold-start) settings.

Funding

This work was supported by National Nature Science Foundation of China (Grant No. 61872297) and Shaanxi Provincial

Key Research & Development Program, China (Grand No. 2020KW-063).

References

1. Ryall KA, Tan AC. Systems biology approaches for advancing the discovery of effective drug combinations. *J Chem* 2015. ISSN 17582946;7(1):7–7. <https://doi.org/10.1186/s13321-015-0055-9>. <https://pubmed.ncbi.nlm.nih.gov/25741385/>.
2. Tatonetti NP, Ye PP, Daneshjou R, et al. Data-driven prediction of drug effects and interactions. *Sci Transl Med* mar 2012. ISSN 1946-6242;4(125):125ra31–125ra31. <https://doi.org/10.1126/scitranslmed.3003377>. <https://pubmed.ncbi.nlm.nih.gov/22422992/>, <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3382018/>.
3. Sun X, Vilar S, Tatonetti NP. High-throughput methods for combinatorial drug discovery. *Sci Transl Med* oct 2013. ISSN 19466234;5(205):205rv1–205rv1. [10.1126/scitranslmed.3006667](https://doi.org/10.1126/scitranslmed.3006667). <https://pubmed.ncbi.nlm.nih.gov/24089409/>.
4. Harrold MW, Zavod RM. Basic Concepts in Medicinal Chemistry. *Drug Dev Ind Pharm* jul 2014. ISSN 0363-9045;40(7):988–8. <https://doi.org/10.3109/03639045.2013.789908>. <https://www.tandfonline.com/doi/abs/10.3109/03639045.2013.789908>.
5. Jiasen Lu, Jianwei Yang, Dhruv Batra, et al. Hierarchical Question-Image Co-Attention for Visual Question Answering. In *NeurIPS*, Barcelona, Spain: Curran Associates Inc. volume 29, pages 289–97, may 2016. URL <http://arxiv.org/abs/1606.00061>.
6. Zhang P, Wang F, Hu J, et al. Label Propagation Prediction of Drug-Drug Interactions Based on Clinical Side Effects. *Sci Rep* jul 2015. ISSN 20452322;5(1):12339. <https://doi.org/10.1038/srep12339>. www.nature.com/scientificreports.
7. Yu H, Mao KT, Shi JY, et al. Predicting and understanding comprehensive drug-drug interactions via semi-nonnegative matrix factorization. *BMC Syst Biol* apr 2018. ISSN 17520509;12(S1):14. <https://doi.org/10.1186/s12918-018-0532-7>. <https://bmcsystbiol.biomedcentral.com/articles/10.1186/s12918-018-0532-7>.
8. Shi J-Y, Mao K-T, Yu H, et al. Detecting drug communities and predicting comprehensive drug-drug interactions via balance regularized semi-nonnegative matrix factorization. *J Chem* 2019. ISSN 1758-2946;11(1):28. <https://doi.org/10.1186/s13321-019-0352-9>.
9. Masumshah R, Aghdam R, Eslahchi C. A neural network-based method for polypharmacy side effects prediction. *BMC Bioinform* jul 2021. ISSN 1471-2105 (Electronic);22(1):1–7. <https://doi.org/10.1186/s12859-021-04298-y>.
10. Deng Y, Xu X, Yang Q, et al. A multimodal deep learning framework for predicting drug-drug interaction events. *Bioinformatics* 2020. ISSN 1367-4803;36(15):4316–22. <https://doi.org/10.1093/bioinformatics/btaa501>.
11. Durant JL, Leland BA, Henry DR, et al. Nourse. Reoptimization of MDL keys for use in drug discovery. *J Chem Inf Comput Sci* nov 2002. ISSN 00952338;42(6):1273–80. <https://doi.org/10.1021/ci010132r>. <https://pubs.acs.org/doi/abs/10.1021/ci010132r>.
12. Bolton EE, Wang Y, Thiessen PA, et al. PubChem: Integrated Platform of Small Molecules and Biological Activities. *Annu Rep Comput Chem* jan 2008. ISSN 15741400;4:217–41. [https://doi.org/10.1016/S1574-1400\(08\)00012-1](https://doi.org/10.1016/S1574-1400(08)00012-1).
13. Assaf Gottlieb, Gideon Y. Stein, Yoram Oron, et al. INDI: A computational framework for inferring drug interactions

- and their associated recommendations. *Mol Syst Biol*, 8: 592, 2012. ISSN 17444292. <https://doi.org/10.1038/msb.2012.26>. URL <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3421442/?report=abstract> <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3421442/>.
14. Zhang W, Chen Y, Liu F, et al. Predicting potential drug-drug interactions by integrating chemical, biological, phenotypic and network data. *BMC Bioinform* jan 2017. ISSN 1471-2105;18(1):18. <https://doi.org/10.1186/s12859-016-1415-9>. <https://pubmed.ncbi.nlm.nih.gov/28056782>, <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC5217341/>.
 15. Ferdousi R, Safdari R, Omid Y. Computational prediction of drug-drug interactions based on drugs functional similarities. *J Biomed Inform* jun 2017. ISSN 15320464;70:54–64. <https://doi.org/10.1016/j.jbi.2017.04.021>. <https://pubmed.ncbi.nlm.nih.gov/28465082/>.
 16. Zhang W, Chen Y, Li D, et al. Manifold regularized matrix factorization for drug-drug interaction prediction. *J Biomed Inform* dec 2018. ISSN 15320464;88:90–7. <https://doi.org/10.1016/j.jbi.2018.11.005>.
 17. Vilar S, Harpaz R, Uriarte E, et al. Drug-drug interaction through molecular structure similarity analysis. *J Am Med Inform Assoc*. Stockholm, Sweden: International Joint Conferences on Artificial Intelligence Organization nov 2012. ISSN 10675027;19(6):1066–74. <https://doi.org/10.1136/amiajnl-2012-000935>. <https://pubmed.ncbi.nlm.nih.gov/22647690/>.
 18. Ma T, Xiao C, Zhou J, et al. Drug Similarity Integration Through Attentive Multi-view Graph Auto-Encoders. In: *IJCAI*, Vol. 7, 2018, 3477–83.
 19. Petar Veličković, Arantxa Casanova, Pietro Liò, et al. Graph attention networks. In *ICLR*. Vancouver, Canada: International Conference on Learning Representations (ICLR). oct 2018. URL <https://arxiv.org/abs/1710.10903v3>.
 20. Michaël Defferrard, Xavier Bresson, and Pierre Vandergheynst. Convolutional Neural Networks on Graphs with Fast Localized Spectral Filtering. In *NeurIPS*. Barcelona, Spain: Curran Associates, Inc. volume 29, pages 3844–52, 2016. URL <http://papers.nips.cc/paper/6081-convolutional-neural-networks-on-graphs-with-fast-localized-spectral-filtering>.
 21. Gilmer J, Schoenholz SS, Riley PF, et al. Neural Message Passing for Quantum Chemistry. In: *ICML*. Sydney, NSW, Australia: PMLR 70 of ICML'17, 2017, 1263–72.
 22. Thomas N Kipf and Max Welling. Semi-Supervised Classification with Graph Convolutional Networks. In *ICLR*. Toulon, France: ICLR 2017. URL <https://openreview.net/forum?id=SJU4ayYgl>.
 23. Duvenaud D, Maclaurin D, Aguilera-Iparraguirre J, et al. Convolutional Networks on Graphs for Learning Molecular Fingerprints. In: *NeurIPS*. Montreal, Canada: Curran Associates, Inc. Vol. 28, 2015, 2224–32.
 24. Kearnes S, McCloskey K, Berndl M, et al. Molecular graph convolutions: moving beyond fingerprints. *J Comput Aided Mol Des* aug 2016. ISSN 15734951;30(8):595–608. <https://doi.org/10.1007/s10822-016-9938-8>.
 25. Wu Z, Ramsundar B, Feinberg EN, et al. MoleculeNet: A benchmark for molecular machine learning. *Chem Sci* jan 2018. ISSN 20416539;9(2):513–30. <https://doi.org/10.1039/c7sc02664a>. <https://pubs.rsc.org/en/content/articlehtml/2018/sc/c7sc02664a>, <https://pubs.rsc.org/en/content/articlelanding/2018/sc/c7sc02664a>.
 26. Yang K, Swanson K, Jin W, et al. Analyzing Learned Molecular Representations for Property Prediction. *J Chem Inf Model* aug 2019. ISSN 15205142;59(8):3370–88. <https://doi.org/10.1021/acs.jcim.9b00237>.
 27. Nyamabo AK, Yu H, Shi J-Y. SSI-DDI: substructure-substructure interactions for drug-drug interaction prediction. *Brief Bioinform* 2021 ISSN 1477-4054;bbab133. <https://doi.org/10.1093/bib/bbab133>.
 28. Nuo Xu, Pinghui Wang, Long Chen, et al. MR-GNN: Multi-resolution and dual graph neural network for predicting structured entity interactions. In *IJCAI*, Macao, China: International Joint Conferences on Artificial Intelligence Organization. volume 2019-Augus, pages 3968–74, 2019. ISBN 9780999241141. <https://doi.org/10.24963/ijcai.2019/551>.
 29. Huang K, Xiao C, Hoang TN, et al. CASTER: Predicting Drug Interactions with Chemical Substructure Representation. *AAAI* nov 2019;34(01):702–9. <http://arxiv.org/abs/1911.06446>.
 30. Andreea Deac, Yu-Hsiang Huang, Petar Veličković, et al. Drug-Drug Adverse Effect Prediction with Graph Co-Attention. In *ICML Workshop on Computational Biology*, New Orleans, Louisiana, United States: ICLR. may 2019. URL <http://arxiv.org/abs/1905.00534>.
 31. Ryu JY, Kim HU, Lee SY. Deep learning improves prediction of drug-drug and drug-food interactions. *Proc Natl Acad Sci U S A* may 2018. ISSN 10916490;115(18):E4304–11. <https://doi.org/10.1073/pnas.1803294115>. <https://www.pnas.org/content/115/18/E4304>, <https://www.pnas.org/content/115/18/E4304.abstract>.
 32. Hanchen Wang, Defu Lian, Ying Zhang, et al. Gognn: Graph of graphs neural network for predicting structured entity interactions. In *IJCAI*, Yokohama, Japan: International Joint Conferences on Artificial Intelligence Organization. volume 2, pages 1317–23, may 2020. <https://doi.org/10.24963/ijcai.2020/183>. URL <https://github.com/Hanchen-Wang/GoGNN>.
 33. Feng Y-H, Zhang S-W, Shi J-Y. DPDDI: a deep predictor for drug-drug interactions. *BMC Bioinform*. : . 2020 ISSN 1471-2105;21(1):419. <https://doi.org/10.1186/s12859-020-03724-x>.
 34. Felix Wu, Amauri Souza, Tianyi Zhang, et al. Simplifying Graph Convolutional Networks. In *ICML*, Long Beach, California, USA: PMLR. volume 97, pages 6861–71, 2019. URL <http://proceedings.mlr.press/v97/wu19e.html>.
 35. Zhen Wang, Jianwen Zhang, Jianlin Feng, et al. Knowledge Graph Embedding by Translating on Hyperplanes. In *AAI*. Québec, Canada: AAAI Press. pages 1112–9, 2014. URL <http://www.aaai.org/ocs/index.php/AAAI/AAAI14/paper/view/8531>.
 36. Huang K, Tianfan F, Gao W, et al. Therapeutics Data Commons: Machine Learning Datasets and Tasks for Therapeutics. *NeurIPS 2021, NeurIPS Datasets and Benchmarks*. virtual conference: Curran Associates, Inc. ArXiv 2021;abs/2102.0.
 37. Zitnik M, Agrawal M, Leskovec J. Modeling polypharmacy side effects with graph convolutional networks. *Bioinformatics* jul 2018. ISSN 14602059;34(13):i457–66. <https://doi.org/10.1093/bioinformatics/bty294>. <https://academic.oup.com/bioinformatics/article/34/13/i457/5045770>.
 38. Paszke A, Gross S, Massa F, et al. PyTorch: An Imperative Style, High-Performance Deep Learning Library. In: *NeurIPS*, Vol. 32, 2019, 8024–35.
 39. Fey M, Lenssen JE. Fast Graph Representation Learning with PyTorch Geometric. *ICLR Workshop*. Vancouver, Canada: Curran Associates, Inc. 2019; URL <http://arxiv.org/abs/1903.02428>.

40. Kingma DP, Ba J. Adam: A Method for Stochastic Optimization. In: ICLR. San Diego, United States: ICLR. 2015, URL <http://arxiv.org/abs/1412.6980>.
41. Hinton GE, Srivastava N, Krizhevsky A, et al. Improving neural networks by preventing co-adaptation of feature detectors. CoRR Arxiv paper. abs/1207.0, 2012. URL <http://arxiv.org/abs/1207.0580>.
42. Huang SA, Lie JD. Phosphodiesterase-5 (PDE5) Inhibitors In the Management of Erectile Dysfunction. *Pharm Ther* 2013;**38**(7):407–19. ISSN 1052-1372 (Print).