



**HAL**  
open science

## Link Prediction in Multi-layer Networks and Its Application to Drug Design

Maksim Koptelov, Albrecht Zimmermann, Bruno Crémilleux

► **To cite this version:**

Maksim Koptelov, Albrecht Zimmermann, Bruno Crémilleux. Link Prediction in Multi-layer Networks and Its Application to Drug Design. 17th International Symposium on Advances in Intelligent Data Analysis, Oct 2018, 's-Hertogenbosch,, Netherlands. pp.175-187. hal-01939475

**HAL Id: hal-01939475**

**<https://hal.science/hal-01939475>**

Submitted on 29 Nov 2018

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Link Prediction in Multi-Layer Networks and its Application to Drug Design

Maksim Koptelov<sup>[0000-0001-9065-2827]</sup>, Albrecht Zimmermann<sup>[0000-0002-8319-7456]</sup>, and Bruno Crémilleux

Normandie Univ, UNICAEN, ENSICAEN, CNRS - UMR GREYC  
{maksim.koptelov,albrecht.zimmermann,bruno.cremilleux}@unicaen.fr

**Abstract.** Search of valid drug candidates for a given target is a vital part of modern drug discovery. Since the problem was established, a number of approaches have been proposed that augment interaction networks with, typically, two compound/target similarity networks. In this work we propose a method capable of using an arbitrary number of similarity or interaction networks. We adapt an existing method for random walks on heterogeneous networks and show that adding additional networks improves prediction quality.

**Keywords:** Chemoinformatics · Link prediction · Multi-layer graphs.

## 1 Introduction

Predicting links between biological or chemical compounds, and targets, such as therapeutic targets, binding sites or disease phenotypes, is an integral part of research in biology and medicinal chemistry. While the main approach to reliably identifying such links still depends on *in vitro* testing, computational methods are employed more and more frequently to fine-tune the set of candidates to be tested *in silico*, cutting down on time and money invested in real-world testing.

A number of methods have been introduced since the problem was first formulated in this form. A straight-forward manner consists of formulating a classification problem: given a particular target, and a number of compounds that have been tested against it, one decides on a representation for the compounds and creates a binary prediction problem that can be solved using any number of existing machine learning techniques. The problem can also be turned around, treating a compound as the class, and targets as data instances, or learning on both entities' representation [13].

A problem such approaches have to overcome is sparsity: whether it is because a target has only recently been identified, because a disease is rare (hence commercially unattractive), or because the relation between certain compounds and targets has not been evaluated for plausible biological reasons, the total space of possible links remains largely under-explored. Concretely, this means that *negative examples* are often not available, ruling certain techniques out.

A semantically similar problem setting that also faces the sparsity problem is that of product recommendation, and recommender systems have therefore been

adapted for the problem setting [7]. A simple recommender system-like approach implements, for instance, the reasoning that if two compounds are both linked with several shared targets, and one of them is linked to an additional one, it is reasonable to assume that the other one should be as well. Such an approach has the advantage of exploiting information that is not directly linked to the chosen target yet is still faced with a sparsity problem since, as mentioned, most compounds do not have many links to start with. This makes some approaches to recommendation, such as matrix factorization, difficult or impossible to use.

One relatively recent proposal to address this problem consists of using network data: vertices represent entities, i.e. compounds and targets, and edges between them their relation. While this does not solve the sparsity problem per se, it allows to introduce additional information: chemical or genetic similarity, for instance, or drug-drug/target-target interactions reported in the literature. There are two obvious questions related to this idea: 1) Which information sources should one use? 2) How can different networks be integrated?

Current solutions often limit themselves to a single similarity network for both compounds and targets, choosing a single similarity measure such as the Tanimoto, Cosine, Simcomp [4] similarity, or Smith-Waterman scores [11] (typically based on empirical validation). In addition, networks are not integrated as such but typically treated separately, with the similarity networks inducing new edges in the interaction network.

In our work, we propose to use a multi-layer network to solve this problem. We illustrate our proposal in the context of ligand-protein interactions, ligands being organic compounds, and proteins biological targets identified as relevant for diseases. Instead of picking and choosing between different sources of information, we propose to use **all of them**, exploiting different similarity measures and interaction information available. Our main contribution is an improvement on the previously introduced NRWRH method [2] that allows us to exploit multi-layer networks assembled from an **arbitrary** number and type of layers. As we show in the experimental evaluation, the algorithm effectively exploits the combination of different types of incomplete data to perform drug-target prediction.

The rest of the paper is organized as follows. Section 2 provides basic notations and problem formulation. Section 3 discusses related work in the given field. Section 4 describes how we adapt existing algorithms to the multi-layer setting. Section 5 describes how we prepare and integrate different data sources, the experimental setup, and presents empirical results. Finally, Section 6 concludes and outlines future work.

## 2 Definitions and problem formulation

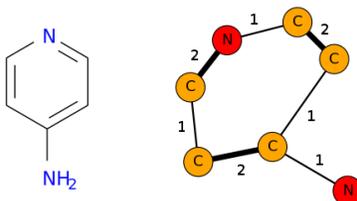
### 2.1 Basic notations

The most important concept in our work is that of a graph.

**Definition 1 (Labeled Graph).** A labeled graph is a tuple  $\langle V, E, \lambda_v, \lambda_e \rangle$ , with  $V$  a set of vertices,  $E \subseteq V \times V$  a set of edges,  $\lambda_v : V \mapsto \mathcal{A}_v$  a labeling function

mapping vertices to elements of an alphabet of possible vertex labels, and  $\lambda_e : E \mapsto \mathcal{A}_e$  a labeling function for edges. We call the degree of a vertex the number of edges in which it is involved:  $\text{deg}(v) = |\{(u, w) \in E \mid u = v \vee w = v\}|$ .

We exploit this representation in two ways: First, ligands are represented by their *molecular 2D-structure*, with  $\mathcal{A}_v$  a subset of atoms, and  $\mathcal{A}_e = \{\text{single covalent bond, double covalent bond, triple covalent bond}\}$ . For example, *Pyridin-4-amine* has chemical equation  $C_5H_6N_2$  and can be presented as in Fig. 1.



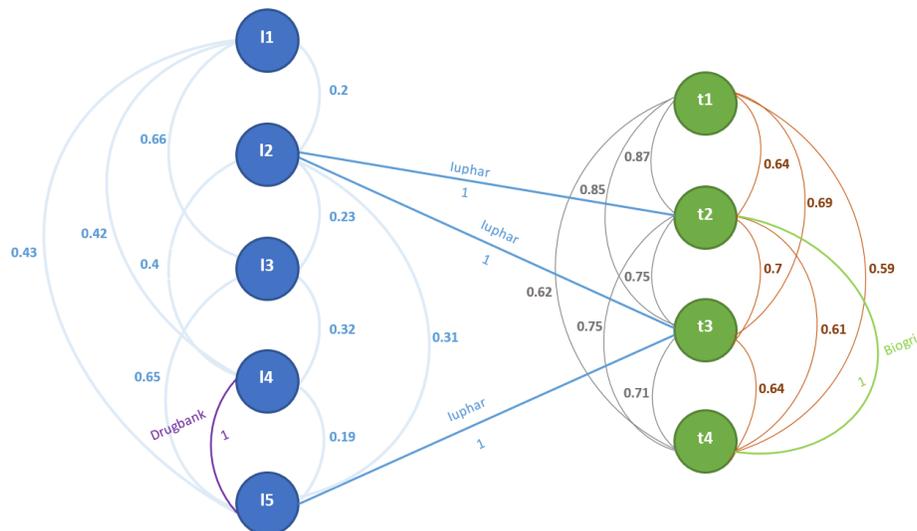
**Fig. 1.** Example of a molecule 2D representation (left) and its corresponding graph in hydrogen suppressed form (right)

Second, the relationships between ligands, proteins, or between ligands and proteins, are represented as *networks*. These include ligand-ligand (ll) and protein-protein (pp) similarity networks, in which  $\mathcal{A}_v$  is the set of ligand/protein identifiers, respectively, and  $\mathcal{A}_e = [0, 1]$ . The other type of network are interaction networks, both ligand-ligand/protein-protein interaction networks derived from the literature, with  $\mathcal{A}_e = \{0, 1\}$ , and ligand-protein (lp) interaction networks, which contain two set of vertices  $V_l, V_p$  and edges  $\forall (u, v) \in E : u \in V_l, v \in V_p$ , and  $\mathcal{A}_e = \{0, 1\}$  or  $\mathcal{A}_e = \mathbb{R}$ . The former labeling is usually derived from the latter by thresholding. An example of relationships between ligands and/or proteins as networks is presented in Fig. 2.

**Definition 2 (Connected component).** Given a graph  $G$ , we call a subgraph  $G' = \langle V', E', \lambda_v, \lambda_e \rangle, V' \subseteq V, E' \subseteq E$  a connected component (CC) iff for any two vertices  $u, v \in V$ , there exists a path  $\{(v_1, v_2), \dots, (v_{m-1}, v_m)\}, v_i \in V, (v_i, v_{i+1}) \in E$ , such that  $v_1 = u, v_m = v$  and there is no supergraph of  $G', G'' = \langle V'', E'', \lambda_v, \lambda_e \rangle, V'' \supset V', E'' \supset E'$  that is a CC.

**Definition 3 (Multi-layer graph).** A multi-layer graph is a tuple  $\langle V, E, \lambda_v, \lambda_e \rangle$ , with  $V$  a set of vertices,  $E$  a multi-set of edges, i.e. tuples  $(u, v), u, v \in V$ . In a multi-layer graph,  $E$  can be decomposed into disjunct sets  $E_l \subseteq V \times V$ , called layers,  $E = \bigcup_i E_{l_i}$ .

As becomes clear from this definition, an arbitrary number of networks can be aggregated into a multigraph, as long as there is overlap in their vertex sets. Trivially, even networks with disjunct vertex sets can be aggregated but since such vertices will not have any edges in the graphs from which they are missing, this will probably be of little use.



**Fig. 2.** An example of a multi-layer graph with 6 networks: ligand-protein network is in deep blue (IUPHAR), ligand-ligand networks are in light blue (ligand similarity network) and violet (DrugBank), protein-protein networks are in green (BioGrid), grey (protein similarity network based on substrings) and brown (protein similarity network based on motifs)

**Definition 4 (Motif).** *Protein motifs are patterns defined using biochemical background knowledge, often expressed in the form of regular expressions.*<sup>1</sup>

**Definition 5 (Tanimoto coefficient).** *The Tanimoto coefficient of two vectors  $\mathbf{x}, \mathbf{y} \in \{0, 1\}^d$  is calculated as:  $\text{coeff}_{\text{Tanimoto}}(\mathbf{x}, \mathbf{y}) = \frac{\mathbf{x} \cdot \mathbf{y}}{\|\mathbf{x}\|^2 + \|\mathbf{y}\|^2 - \mathbf{x} \cdot \mathbf{y}}$*

## 2.2 Problem formulation

The problem setting we address in this paper is one of link prediction between ligands (drug candidates) and proteins (biological targets).

**Definition 6 (Ligand-protein activity prediction).** *For a given number of ligand-protein activity networks  $G_{lp}^i = \langle V_l \cup V_p, E_i, \lambda_v, \lambda_{e_i} \rangle$ , with  $u \in V_l$  labeled with ligands identifiers,  $v \in V_p$  labeled with protein identifiers,  $\forall (u, v) \in E, u \in V_l, v \in V_p$ , and  $\mathcal{A}_e = \{0, 1\}$ , ligand-ligand networks  $G_l^i = \langle V_l, E_l^i, \lambda_v, \lambda_{e_l^i} \rangle$ , protein-protein networks  $G_p^i = \langle V_p, E_p^i, \lambda_v, \lambda_{e_p^i} \rangle$  and a given  $(u, v) \notin E, u \in V_l, v \in V_p$  predict, whether  $\lambda_e((u, v)) = 1$ .*

We limit ourselves to the relatively easier task of predicting whether there is activity or not, leaving the prediction of its *strength* as future work.

<sup>1</sup> An open-access database is available at <http://prosite.expasy.org>

### 3 Related work

The literature on compound-target activity prediction, even using networks, is too vast to discuss here. We therefore present a number of works illustrating the characteristics we discussed in the introduction. Ligand-protein activity, the use case we explore here, has been addressed in [13], which selects a ligand and target similarity measure each, and multiplies activity vectors of known ligands/targets with the similarity to new ligands/targets to derive predictions. In [14], the same group used ligand structural and pharmacological similarity, as well as genetic protein similarity, mapped ligands and targets into a shared feature space and predicted activity. The authors of [3] used three networks: ligand-ligand similarity, target-target similarity, ligand-target activity, evaluated four ligand similarity measures, settling on Tanimoto distance. The proposed method, NWNBI, exploits similarity weights and log-values of activity measurements to perform four-step network traversals. In [2], ligand similarity is calculated as weighted average of *two* similarity measures, and combined with a target similarity, and the interaction network into a three-layer network, which they refer to as “heterogeneous”. They simulate random walk with restart by matrix multiplication, and show that only using a single similarity measure or ignoring the interaction network deteriorates results. Three networks are also used in [7], the authors discuss different options for similarity measures, and perform low-rank matrix factorization on the adjacency/similarity matrices. They address sparsity by giving non-existing links a small non-negative weight. Ligand-protein activity is also the subject of [1], which exploits the three-layer network to perform weighted nearest-neighbor classification. Gene-disease interactions have been considered in [12], using three layers, simulating random walk by matrix multiplication, using different numbers of steps for the two similarity networks. Using a similar bi-random walk idea, [9] consider microRNA-disease interactions, exploiting a three-layer network. The random walk with restart in [8] is symmetric (and functionally the same as in [2]), with the similarity networks constructed by averaging two similarity measures. They evaluate different parameter settings.

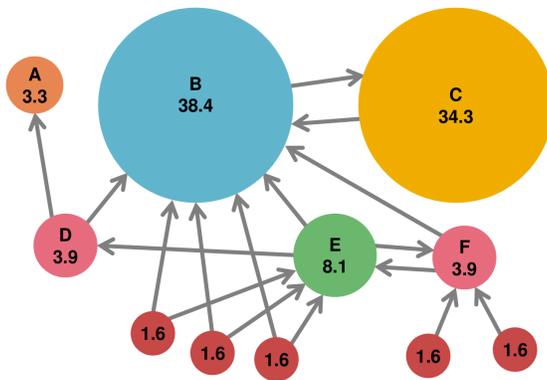
## 4 Exploring a multi-layer graph

As the preceding section shows, the standard setting employed consists of three networks, and to adhere to this setting, authors either choose a single similarity measure empirically, or combine similarity measures via user-specified weights. Instead, we propose to combine all available networks into a multi-graph having more than three layers. Once we have such a network, the question is how to exploit it, however, and here we hew close to the literature.

### 4.1 The random walk model

A long-established method for exploring a network is the random walk [10], which proceeds roughly as follows: starting from a randomly selected node, it

performs walks along edges of the graph at random. In every step, the edge to follow is chosen uniformly from all outgoing links (in the case of an unweighted graph) or proportional to link weights (in the case of a weighted graph). Node importance is based on how frequently the walker visits the node: a node with higher frequency is considered more important than a node with a low value (Fig. 3). This idea can be modified in a number of ways to improve network



**Fig. 3.** Nodes importance example in a graph, taken from [6]

exploration: the walker can be constrained to perform at most  $\overline{max\_steps}$  steps, to not visit any of the last  $c$  vertices it encountered, or with small probability  $1 - \beta$  the process can be restarted at any time to avoid getting trapped by those vertices it mustn't visit. The product of the probabilities of edges the walker traversed gives the cumulative probability of a path between two nodes and can be used to *predict* a link between a starting node and an end node: if the path probability is greater than a given threshold, a new edge is predicted.

To extend this approach to multi-layer graphs, one needs to add how to choose the layer to walk in. We propose to select a network uniformly at random from the set of networks, and multiply the path strength by  $\frac{1}{|\{G_i^i\}| + |\{G_p^j\}| + |\{G_{lp}^k\}|}$ . Repeat the process until a user-defined target vertex is reached or the maximum number of steps have been performed. Due to the randomized nature, random walks are usually repeated several times to derive more robust estimates.

## 4.2 Network-based Random Walk on Multi-layer network

Instead of explicitly random walking as described above, random walks are often simulated via matrix multiplication of transition matrices. This is notably the approach proposed in [2], abbreviated as NRWRH. They define a transition matrix  $M = \begin{bmatrix} M_{pp} & M_{pl} \\ M_{lp} & M_{ll} \end{bmatrix}$ , in the manner described above, i.e. uniform probabilities for lp/pl-transitions, proportional probabilities for similarities, with an additional user-specified parameter  $\lambda \in [0, 1]$  affecting moves from ligands to proteins and

vice versa. Given a ligand  $l_i$ , a starting vector  $v_0 \in [0, 1]^{|V_l|+|V_p|}$  is initialized with 1 at the position for  $v \in V_l$ ,  $\lambda_v(v) = l_i$ ,  $\frac{1}{|\{(v,u) \in E_{lp}\}|}$  at the positions for the proteins linked to it, 0 otherwise. Protein entries in  $v_0$  are multiplied with  $1 - \eta$ , ligand entries with  $\eta$ , a user-defined parameter to bias the walk towards proteins ( $\eta < 0.5$ ), or ligands ( $\eta > 0.5$ ). The vector representing the probabilities that a walker starting with  $l_i$  finds itself in any of the nodes is calculated iteratively as  $p_{t+1} = (1-\beta)M^T p_t + \beta p_0$  until  $|p_{t+1} - p_t| < 10^{-10}$ . This can be understood as the random walker walking “in all directions at the same time”. The approach can be considered a simplified version of Personalized PageRank [5], simplified because edges are undirected and there is only a single starting vertex. Removing the starting vertices from the final state vector, and ranking entries gives predicted edges. We adapt this approach to a setting with  $|\{G_l^i\}| + |\{G_p^j\}| + |\{G_{lp}^k\}| \geq 3$ . While the algorithm stays essentially the same, we decompose the transition matrix into a matrix  $M$  for within-network/layer transitions, and a matrix  $N$  for between-network/layer transitions. We also do away with the user-dependent  $\lambda$ . Explicitly creating  $M$  in the manner shown above is easy for three layers but becomes much harder when different numbers can be involved. We hence con-

struct  $M = \begin{bmatrix} M_{G_p} & 0 & 0 \\ 0 & M_{G_l} & 0 \\ 0 & 0 & M_{G_{lp}} \end{bmatrix}$ , with  $M_{G_p} = \begin{bmatrix} M_{G_p}^1 & 0 & \dots & 0 \\ 0 & M_{G_p}^2 & \dots & 0 \\ 0 & 0 & \dots & M_{G_p}^{|\{G_p^i\}|} \end{bmatrix}$  derived

from protein-protein similarity networks ( $M_{G_l}$ ,  $M_{G_{lp}}$  accordingly). The tran-

sition matrix  $N = \begin{bmatrix} N_{G_p^1 \rightarrow G_p^1} & N_{G_p^2 \rightarrow G_p^1} & \dots & N_{G_{lp}^{G_p^i} \rightarrow G_p^1} \\ \dots & \dots & \dots & \dots \\ N_{G_p^1 \rightarrow G_{lp}^{G_p^i}} & N_{G_p^2 \rightarrow G_{lp}^{G_p^i}} & \dots & N_{G_{lp}^{G_p^i} \rightarrow G_{lp}^{G_p^i}} \end{bmatrix}$  explicitly models

possible layer transitions, with 1s on the main diagonal of a submatrix  $N_{G_j \rightarrow G_i}$  for all nodes present in both layers, 0s otherwise. Note that this means that transition matrixes from ligand to protein layers (and vice versa) have zeros everywhere including the main diagonal. The initial state vector  $v_0$  has dimensionality ( $|V_p| \cdot |\{G_{p_i}\}| + |V_l| \cdot |\{G_{l_i}\}| + |V_l \cup V_p| \cdot |\{G_{lp_i}\}|$ ) with entries for *all* vertices in *all* layers. It is initialized by setting the entry for the starting ligand and each linked protein to 1 in every network they are present. Matrices and state vectors are column-normalized – the entries of a column must sum to 1.

Our algorithm, NEtWork-basEd Random walk on MultI-layered NEtwork (NEWERMINE), is summarized in Algorithm 1.  $(M_{norm}N)_{norm}$  can be pre-computed, giving us a matrix that is functionally equivalent to  $M$  as defined in NRWRH, and used on every iteration of NEWERMINE to save computation time. At the end,  $v_{final}$  needs to be summarized by summing up for each vertex all corresponding entries, leading to a vector with dimensionality  $|V_l \cup V_p|$  from which the edge ranking can be derived.

## 5 Experimental Evaluation

In order to allow reproducibility of our work, we evaluated our approach on publicly available data. In this part we provide a description of the data used

---

**Algorithm 1:** The NEWERMINE algorithm

---

**Input** : adjacency matrix  $M$ , transition matrix  $N$ , *starting\_vertex*,  
 $max\_steps$ ,  $\eta$ ,  $\beta$ ,  $max\_diff$   
**Output:** Probability scores  $v_{final}$   
 $V_{0_l} \leftarrow$  initialize *starting\_vertex*  
 $V_{0_p} \leftarrow$  initialize targets for which an interaction with *starting\_vertex* is known  
 $V_0 \leftarrow (1 - \eta) \cdot V_{0_{l_{norm}}} + \eta \cdot V_{0_{p_{norm}}}$   
 $step \leftarrow 0$   
**repeat**  
   $step \leftarrow step + 1$   
   $V_{step} \leftarrow \beta \cdot (M_{norm}N)_{norm}V_{step-1} + (1 - \beta) \cdot V_0$   
**until**  $(|v_{step} - v_{step-1}| \leq max\_diff) \vee (step > max\_steps)$   
**return**  $v_{step}$

---

and the details of the experimental protocol. This is followed by the results and the discussion.

## 5.1 Experimental Settings

**Datasets** In total we have used 4 datasets:

1. IUPHAR – an open-access database of ligands, biological targets and their interactions. We used version 2017.5 (released on 22/08/2017). The full dataset has 8978 ligands, 2987 proteins, and 17198 interactions (edges) between them<sup>2</sup>. In order to satisfy the designed setting conditions, we removed duplicate interactions (based on different affinity measures), leaving 12456 interactions in total. For existing interactions, we label an edge with 1 if the negative logarithm of the affinity measure is  $\geq 5$ , non-interacting otherwise.<sup>3</sup> We treat all affinity measures available in the data (pKi, pIC50, pEC50, pKd, pA2, pKB) as equivalent.
2. DrugBank (DB) – an open-access database of drug-drug interactions. We used version 5.0.11 (released 20-12-2017). It has 658079 interactions of 3138 distinct drugs. 242922 of these interactions involve 1254 distinct ligands that are present in IUPHAR. The database was also used as a source of 2D representations of ligands to compute ligand similarities.
3. BioGrid (BG) – an open-access database of protein-protein interactions mined from a corpus of biomedical literature. We used version 3.4.154 (25/10/2017). It has 1482649 interactions of 67372 distinct proteins. Only 15410 of these interactions involve proteins present in IUPHAR (1925 distinct proteins).
4. NCBI Protein database – The National Center for Biotechnology Information proteins database<sup>4</sup> was used to obtain amino acids sequences to represent targets. The data was parsed from the website of NCBI and mapped

<sup>2</sup> in ligands.csv, interactions.csv, and targets\_and\_families.csv, respectively

<sup>3</sup> Cutoff proposed by researchers from CERMN (<http://cermn.unicaen.fr>)

<sup>4</sup> <https://www.ncbi.nlm.nih.gov/protein/>

Data set	Entities	Relations	Sparsity	Network			
				Vertices	Edges	Sparsity	CC
IUPHAR	11965	12456	0.00017	11965	12456	0.00017	443
DrugBank	3138	658079	0.1337	1254	122808	0.15631	1
BioGrid	67372	1482649	0.00065	1898	8658	0.0048	11
Ligand similarity	6821	23259610	1	6821	23259610	1	1
NCBI	1818	1651653	1	1818	1651653	1	1

**Table 1.** Data set and network characteristics

to IUPHAR using the RefSeq attribute (human protein sequence identifier) available in IUPHAR. The database was accessed 20/12/2017.

Ligands were mapped between networks by numerical identifiers provided by IUPHAR as well as by INN (International Non-proprietary Name) and Common name attributes. Proteins were mapped by IUPHAR identifiers as well as by Human Entrez Gene attribute.<sup>5</sup> In total we have built 6 networks:

1. a drug interaction network based on DrugBank,
2. a drug similarity network based on similarities calculated using the Tanimoto coefficient on binary vectors constructed by frequent subgraphs,
3. the drug-target interaction network based on IUPHAR,
4. a target interaction network based on BioGrid, and
5. two target similarity networks calculated using the Tanimoto coefficient on feature vectors constructed by *frequent substrings* and *Prosite motifs*.

Similarity networks' edges were labeled with labels  $\in [0, 1]$ , interaction networks with labels  $\in \{0, 1\}$ . Table 1 shows the characteristics of the data sets, and of the networks we derived from them. It is noticeable how sparse the data is, and also how this sparsity translates into disconnected parts of the network. Sparsity might result in a low performance of the traditional recommender systems approaches, while disconnected networks are challenging for random walker approaches.

**Evaluation Protocol** To evaluate our approach, we used leave-one-out cross-validation: for each of the 12456 edges in the IUPHAR network, we remove it from the network, set the ligand as starting vertex, infer strengths for all possible ligand-target paths, remove ligand-target edges contained in the training data, and check whether the removed edge is found in the top-20 remaining paths<sup>6</sup> according to their strengths. If this is the case for an interacting edge, we consider it a *true positive*, otherwise a *false negative*. For negative examples, the relationship is inverse.

**Quality Measures** To evaluate our methods we use several performance measures:

<sup>5</sup> Global Query Cross-Database Search System gene identifiers:  
<https://www.ncbi.nlm.nih.gov/gene>

<sup>6</sup> Precision at 20

- Accuracy: the ratio of true positives (TP) – drug-target links correctly classified as positives – and true negatives (TN) – drug-target links correctly classified as negatives – over all predictions:  $Acc = \frac{TP+TN}{TP+FP+TN+FN}$ .
- Area under receiver operating curve (AUC): evaluates whether true positives are usually ranked above or below false positives when sorting predictions by confidence.
- Precision: the ratio of TP over all drug-target links classified as positives:  $Prec = \frac{TP}{TP+FP}$ . Precision measures whether a model is specific enough to mainly classify links of the positive class as positive. This gives additional insight into the accuracy score.
- Recall: the ratio of TP over all positive links in the test data:  $Rec = \frac{TP}{TP+FN}$ . Recall measures whether a model is general enough to classify a large proportion of the positive class as positive.

In addition to this, we also report weighted versions of accuracy, precision, and recall that give us a more accurate assessment for unbalanced datasets. Due to the fact that the number of negative examples are smaller than that of the positives in our data, we assign a classification cost of 1 to positives and cost  $neg\_cost$  to negatives, derived by:  $neg\_cost = \frac{|D|}{2 \times |N|}$ , where  $|D|$  – number of examples,  $|N|$  – number of negative examples. We then perform evaluation based on the costs defined: FN and TN receive score  $neg\_cost$  for every negative example w.r.t. its real class, while FP and TP receives score 1 for positives.

**Implementation** We implemented NEWERMINE in Python<sup>7</sup>. We used the networkx library to model the multi-layer network, the NumPy library to perform all matrix computations and the sklearn library for cost-based evaluation.

## 5.2 Experimental results

**Using three-layer graphs** We first use NEWERMINE on a number of multi-graphs aggregated from three networks each, the ligand-target network, one ligand-ligand network, and one target-target network. This is the setting used in the papers discussed in Section 3.

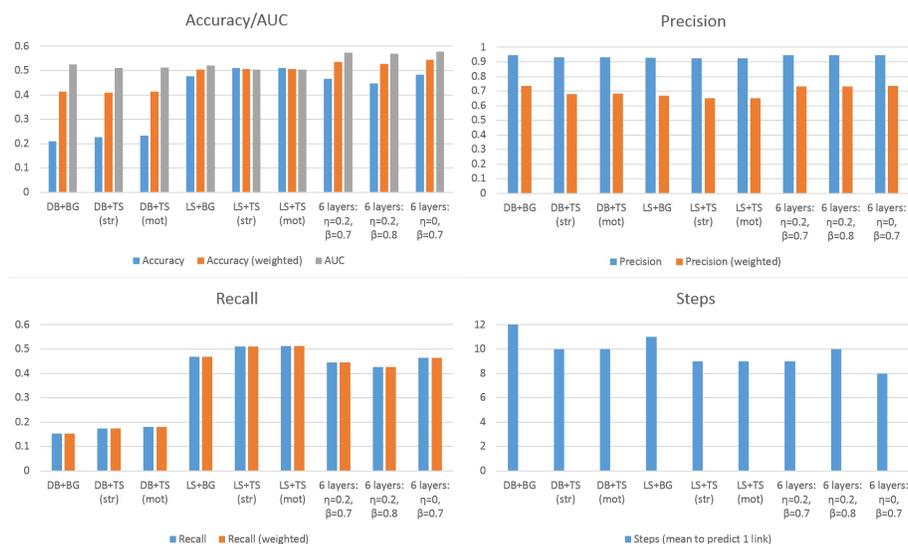
For the experiments we defined 6 possible combinations with IUPHAR, only ligand-target interaction network we have: (1) DrugBank + BioGrid, (2) DrugBank + Target similarity (TS) (substrings:str), (3) DrugBank + TS (motifs:mot), (4) Ligand similarity (LS) + BioGrid, (5) LS + TS (str), (6) LS + TS (mot). The basic properties of the combinations compared to the full graph are presented in Table 2. The results of the use of NEWERMINE with parameters  $\eta = 0.2$ ,  $\beta = 0.7$  (taken from [2]) are presented in Fig. 4. This is a rather conservative setting, equivalent to relatively few steps before the walker restarts.

The number of vertices in different networks depend on available IDs and structural information. In any case, the networks are sparse and they are not

<sup>7</sup> <https://zimmermanna.users.greyc.fr/supplementary-material.html>

Combination	Ligands	Targets	$ V $	$ E $	Sparsity	CC
DB + BG	7025	2307	9332	143922	0.003	87
DB + TS (str)	7025	2101	9126	1786917	0.042	103
DB + TS (mot)	7025	2101	9126	1786917	0.042	103
LS + BG	8056	2307	10363	23280724	0.434	21
LS + TS (str)	8056	2101	10157	24923719	0.4832	22
LS + TS (mot)	8056	2101	10157	24923719	0.4832	22
Six layers	8137	2502	10639	26706838	0.4719	1

**Table 2.** Basic properties of different combinations of networks



**Fig. 4.** Evaluation results of NEWERMINE for different combinations of three networks and the six-layer graph

fully connected. Using similarity networks alleviates this situation somewhat and combining *all* networks leads to a single connected component (bottom row).

Fig. 4 shows that using different three-layer graphs leads to rather different results. The arguably most notable result is that using ligand structural similarity instead of DrugBank network significantly improves accuracy and recall.

**Using the full, six-layer graph** The results for NEWERMINE on the full multi-layer graph are also presented in the Fig. 4. We show additional values for  $\eta$  and  $\beta$ :  $\eta = 0$  strongly biases the walk towards targets, we also consider  $\beta = 0.8$  for  $\eta = 0.2$ . Using more layers decreases recall somewhat, but improves weighted accuracy (taking the lower proportion of negative examples into account), AUC score and precision. Different parameter values do not have a large effect on the results but change running times: increasing  $\beta$  also increases the number of steps necessary for convergence, and decreasing  $\eta$  decreases this number.

## 6 Conclusion and perspectives

We have presented an approach for exploiting an arbitrary number of networks combined into a multi-layer network, proposing general matrix formulations to form intra- and inter-network transitions.

As we have demonstrated experimentally, combining different networks improves vertex reachability and therefore interaction prediction. So far, we have only exploited more than one protein similarity network, already achieving very good results. In future work, we intend to also integrate different ligand similarity semantics, and different databases indicating ligand-protein activity. Additionally, we intend to employ our approach for different target settings, e.g. for miRNG-disease links. Finally, we aim to move from the “active”/“inactive” setting to one where we predict the strength of the activity.

## References

1. Buza, K., Peska, L.: Aladin: A new approach for drug–target interaction prediction. In: ECML/PKDD. pp. 322–337. Springer (2017)
2. Chen, X., Liu, M.X., Yan, G.Y.: Drug–target interaction prediction by random walk on the heterogeneous network. *Molecular BioSystems* **8**(7), 1970–1978 (2012)
3. Cheng, F., Zhou, Y., Li, W., Liu, G., Tang, Y.: Prediction of chemical-protein interactions network with weighted network-based inference method. *PLoS one* **7**(7), e41064 (2012)
4. Hattori, M., Okuno, Y., Goto, S., Kanehisa, M.: Development of a chemical structure comparison method for integrated analysis of chemical and genomic information in the metabolic pathways. *JACS* **125**(39), 11853–11865 (2003)
5. Haveliwala, T.H.: Topic-sensitive pagerank: A context-sensitive ranking algorithm for web search. *TKDE* **15**(4), 784–796 (2003)
6. Leskovec, J., Rajaraman, A., Ullman, J.D.: Mining of massive datasets. Cambridge university press (2014)
7. Lim, H., Gray, P., Xie, L., Poleksic, A.: Improved genome-scale multi-target virtual screening via a novel collaborative filtering approach to cold-start problem. *Scientific reports* **6**, 38860 (2016)
8. Liu, Y., Zeng, X., He, Z., Zou, Q.: Inferring microRNA-disease associations by random walk on a heterogeneous network with multiple data sources. *TCBB* **14**(4), 905–915 (2017)
9. Luo, J., Xiao, Q.: A novel approach for predicting microRNA-disease associations by unbalanced bi-random walk on heterogeneous network. *Journal of biomedical informatics* **66**, 194–203 (2017)
10. Pearson, K.: The problem of the random walk. *Nature* **72**(1867), 342 (1905)
11. Smith, T., Waterman, M.: Identification of common molecular subsequences. *Molecular Biology* **147**, 195–197 (1981)
12. Xie, M., Hwang, T., Kuang, R.: Prioritizing disease genes by bi-random walk. In: PAKDD. pp. 292–303. Springer (2012)
13. Yamanishi, Y., Araki, M., Gutteridge, A., Honda, W., Kanehisa, M.: Prediction of drug–target interaction networks from the integration of chemical and genomic spaces. *Bioinformatics* **24**(13), i232–i240 (2008)
14. Yamanishi, Y., Kotera, M., Kanehisa, M., Goto, S.: Drug-target interaction prediction from chemical, genomic and pharmacological data in an integrated framework. *Bioinformatics* **26**(12), i246–i254 (2010)