

Drug Target Interaction Prediction with Non-random Missing Labels

Sheng Ni, Chen Lin, Xiangxiang Zeng

*Department of Computer Science
Xiamen University
Xiamen, China
chenlin@xmu.edu.cn*

Yun Liang

*Department of Information
South China Agricultural University
Guangzhou, China*

Abstract—Drug-Target Interaction (DTI) prediction plays an important role in drug discovery and drug repurposing. DTI prediction is usually modeled as a binary classification problem. Unlike previous studies which label unknown DTIs as negative samples, we assume the unknown DTIs are labels that are missing not at random. For example, negative DTI labels are more likely to be missing because biomedical researchers prioritize to study DTIs that are more likely to be positive. We introduce a novel probabilistic model, Factorization with Non-random Missing Labels (FNML), for DTI prediction. FNML models the generative process for the DTI labels (i.e. the labels are positive or negative) and responses (i.e. the labels are observed or missing). In particular, the probability of observing or missing a label is associated with the sign of the label. We also conduct comprehensive experiments to validate the robust performance of the proposed models.

Index Terms—Missing Not At Random, Drug Target Interaction Prediction, Probabilistic Factor Models

I. INTRODUCTION

Drug-Target Interaction (DTI) is fundamental to drug discovery and design. As biochemical experimental methods for DTI identification are extremely costly and time-consuming, computational DTI prediction methods have received a growing popularity in literature. Traditional computational methods to predict DTIs mainly include ligand-based methods [1] and molecule docking methods [2]. Ligand-based methods are ineffective when target proteins have little binding ligands, while molecular docking methods are computationally costly and fail to offer accurate predictions when 3D structures of target proteins are not available [3]. To overcome these problems, many machine learning-based methods have been proposed for inferring DTI. The majority of existing machine learning-based methods treat DTI prediction as a binary classification task, where known DTIs are labeled as positive and unknown DTIs are labeled as negative [4]. To address the imbalanced problem arisen from the binary classification scheme, many research has attempted to extract a subset of reliable negative samples, e.g. by random sampling [5] or by Positive Unlabel Learning (PU Learning) [6].

Instead of labeling the unknown DTIs as negative, we argue that it is more natural to consider the unknown DTIs, i.e. DTIs that are neither identified in vivo to be positive nor

experimentally validated to be negative (non-interacting drug-target pairs), as missing labels. Furthermore, our assumption in this work is that labels are not missing at random. This is an intuitive and reasonable assumption, because researchers will use their domain expertise to filter DTIs with a high possibility to be positive and prioritize validations for these DTIs in vivo. For example, researchers find the efficacy target of a drug based on principles of biochemistry, biophysics, genetics and chemical biology. If ample evidences exist to support positive interactions with the target, then the possibility of a positive DTI is high, and the researchers are likely to conduct in vivo experiments. On the contrary, drug target interactions that are less likely to be positive are more likely to be ignored by researchers and their labels are likely to be missing.

A. Contribution

Our contribution in this work is a novel Factorization with Non-random Missing Labels model (FNML). To the best of our knowledge, this is the first time missing not at random theory is applied in DTI identification. The inputs of FNML are feature vectors of drugs and targets, the partially observed labels, and the fully observed responses (i.e. labels are given or missing). We allow the feature vectors to be learnt and/or integrated from heterogenous sources. The labels and responses are binary variables. The FNML model mimics the probabilistic procedures to generate labels from feature vectors and responses from labels. Specifically, the labels are related to feature vectors of both drugs and targets, and a hidden matrix mapping from the drug features to target features. The possibility of giving a response is associated with the sign of the label.

We conduct comprehensive experiments on the latest DTI database. Experimental results show that the FNML model outperforms state-of-the-art DTI prediction methods in terms of Area Under Receiver Operating Characteristic curve (AU-ROC) and Area Under Precision Recall curve (AUPR), which are the most commonly adopted metrics to evaluate DTI prediction performance. We also show that our models provide robust performance enhancement, despite of the input features.

Chen lin is supported by NSFC under grant No.61472335.

B. Related Work

One component of our work (i.e labels are generated by feature vectors learnt and fused from heterogenous information networks) is inspired by a recent work DTINet [5]. However, there are three key differences between our work and DTINet. (1) DTINet is based on deterministic matrix factorization, our work is based on probabilistic factor models. For example, the hidden feature space mapping matrix, labels, and responses are all random variables. This setting enables the FNML model to regulate the parameters (i.e. hidden feature space mapping matrix) by introducing appropriate priors. Therefore, performance on sparse dataset is improved. (2) DTINet is based on randomly missing responses, i.e. it samples uniformly a set of unknown DTIs as negative sample, while FNML is based on missing not at random theories. Statistical theory in [7] shows that applying a model based on missing at random assumptions can lead to biased parameter estimation on data sets with missing not at random entries. (3) DTINet adopts only a subset of unknown DTIs to preserve a balanced number of positive and negative samples, while our model uses all information in the data set.

We also want to distinguish our work with another line of research. Usually only positive DTIs are deposited in known databases. Due to the lack of negative samples, PU learning has been employed in DTI identification, e.g. to facilitate negative sample extraction [6]. PU learning does not explicitly associate the status of an instance (i.e. being labeled or unlabeled) with the value of its label. We also want to mention here that, although we experiment with datasets where only positive DTIs are deposited, FNML is extendable without difficulty to databases where positive and negative DTIs are available. Thus our model is applicable in more scenarios.

II. THE PROPOSED METHOD

We start with the problem definitions and notations in Sec. II-A. We then describe the proposed model FNML in Sec. II-B. Finally we present the inference algorithm in Sec. II-C.

A. Preliminaries

DTI identification is often modeled as a binary classification task. Formally, we are given $P \in \mathcal{R}^{N \times M}$ a set of DTI labels, where $p_{i,j} = 0$ indicates a negative interaction between drug i and target j , $p_{i,j} = 1$ indicates a positive DTI, the feature vectors on drug side $X \in \mathcal{R}^{N \times K}$, where $x_{i,k}$ represents drug i 's weight on drug feature k , the feature vectors on target side $Y \in \mathcal{R}^{M \times L}$, where $y_{j,l}$ represents target j 's weight on target feature l . The problem is to predict for a new drug-target pair $\langle i', j' \rangle$, the possibility of a positive DTI $p(p_{i',j'} = 1)$.

Similar to DTINet [5], we use a compact feature expression learnt from drug and protein networks. To extract features X, Y , we first create networks that involve drugs (for X) and proteins (for Y). We compute similarity score between each pair of nodes in the networks. Then, the diffusion component analysis (DCA) [8] is applied to learn a low-dimensional vector representation of each node of the drug network and

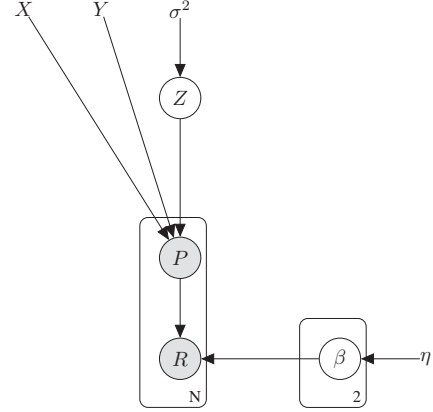


Fig. 1. Graphical Representation of the FNML model

protein network. Note here that X, Y can be extracted from a single network or an aggregation of several networks. The details of feature extraction are described in Sec. III.

In addition to the features X, Y and labels P , we make one essential modification to the problem definition. We assume that the inputs also contain responses $R \in \mathcal{R}^{N \times M}$, where $R_{i,j} = 0$ indicates an unknown DTI, $R_{i,j} = 1$ indicates a verified DTI (positive or negative). For positive responses $R_{i,j} = 1$, the labels $P_{i,j}$ are observed. For negative responses $R_{i,j} = 0$, the labels are hidden and unknown.

B. FNML Model

We use a factor model, which is depicted in Fig. 1. The features X, Y are in different dimensions. To associate the drug features with the target features, we introduce a hidden matrix $Z \in \mathcal{R}^{K \times L}$, where $Z_{k,l}$ is a projection that maps the drug feature k to the target feature l . We assume that Z is sampled from a Gaussian distribution,

$$\forall k, l, Z_{k,l} \sim \mathcal{N}(0, \sigma^2), \quad (1)$$

where σ^2 is the variance. We use zero mean to favor sparse feature mapping, i.e. a drug feature k is associated with a few target features.

We then assume that the binary label $P_{i,j}$ is generated from the following process:

$$\forall i, j, p(P_{i,j} = 1 | X, Y, Z) = \frac{1}{1 + \exp(-XZY)_{i,j}}. \quad (2)$$

The binary response is sampled from a Bernoulli distribution. The parameters of the Bernoulli distribution are related to the value of each $P_{i,j}$. Therefore we define $\beta_p \in \mathcal{R}^2, p \in \{0, 1\}$, $\forall p, \beta_{p,0} > 0, \beta_{p,1} > 0, \beta_{p,0} + \beta_{p,1} = 1$, we have:

$$\forall p \in \{0, 1\}, \beta_p \sim \text{Beta}(\eta), \quad (3)$$

$$\forall i, j, R_{i,j} \sim \text{Bern}(\beta_{P_{i,j},1}), \quad (4)$$

where $\eta \in \mathcal{R}^2$ is the hyperparameter for the Beta distribution.

C. Inference

The objective is to maximize the log-likelihood which consists of two terms. The first term is on partial observations, i.e. $R_{i,j} = 0$ and $P_{i,j}$ unknown. The second term is on full observations, i.e. $R_{i,j} = 1$ and known $P_{i,j}$.

$$\begin{aligned} \mathcal{L} &= \sum_{R_{i,j}=0} \log p(R_{i,j}|X, Y, \sigma^2, \eta) \\ &+ \sum_{R_{i,j}=1} \log p(R_{i,j}, P_{i,j}|X, Y, \sigma^2, \eta) \end{aligned} \quad (5)$$

Direct optimization for both terms in Equ. 5 is intractable, as they involve integration over continuous hidden variables. For example, $p(R|X, Y, \sigma^2, \eta) = \int_{P, Z, \beta} p(R|P, \beta, \eta)p(P|X, Y, Z)p(Z|\sigma^2)p(\beta|\eta)$. We employ variational inference [9] to infer the parameters. That is, we use the mean field assumption to factorize the posterior distribution:

$$q(Z, \beta, P|R, X, Y, \sigma^2, \eta) = q(P|\theta)q(Z|\mu, v)q(\beta|\rho), \quad (6)$$

It is convenient if $q(P|\theta)$, $q(Z|\mu, v)$, $q(\beta|\rho)$ are exponential distributions. We approximate the sigmoid function in Equ. 2 by an exponential distribution. We use the property that any sigmoid function $\sigma(\cdot)$ has a lower bound:

$$q(P|\theta) = \sigma(\theta) \geq \sigma(\zeta) \exp((\theta - \zeta)/2 - \lambda(\zeta)(\theta^2 - \zeta^2)), \quad (7)$$

where $\lambda(\zeta) = [\sigma(\zeta) - 1/2]/[2\zeta]$.

Maximizing the likelihood is equivalent to maximizing ELBO (Evidence Lower Bound):

$$\begin{aligned} \mathcal{L}(q(Z, \beta, P)) &= E_{q(Z, \beta, P)}[\ln P(R, P, Z, \beta)] \\ &- E_{q(Z, \beta, P)}[\ln q(Z, \beta, P)]. \end{aligned} \quad (8)$$

We divide the data objects into two disjoint sets, $s_1 = \{(i, j) \in \mathcal{R}^{N \times M} | R_{i,j} = 1\}$, and $s_2 = \{(i, j) \in \mathcal{R}^{N \times M} | R_{i,j} = 0\}$. First we derive the parameters μ, v of $\ln q(Z_{k,l}|\mu_{k,l}, v_{k,l})$:

$$\begin{aligned} \ln q(Z_{k,l}|\mu_{k,l}, v_{k,l}) &= \sum_{(i,j) \in s_1} E_{q(\beta)}[\ln p(R_{i,j}, P_{i,j}, Z, \beta)] \\ &+ \sum_{(i,j) \in s_2} E_{q(\beta, P_{i,j})}[\ln p(R_{i,j}, P_{i,j}, Z, \beta)] \\ &+ \text{const}, \end{aligned}$$

where *const* represents the irrelevant item. Removing irrelevant items we get:

$$\begin{aligned} \ln q(Z_{k,l}|\mu_{k,l}, v_{k,l}) &= [\sum_{(i,j) \in s_2} [(\theta_{i,j} - \frac{1}{2})X_{i,k} * (Y^T)_{l,j}] \\ &+ \sum_{(i,j) \in s_1} \frac{1}{2}X_{i,k} * (Y^T)_{l,j}]Z_{k,l} \\ &- \sum_{i,j} \lambda(\zeta_{ij}) * X_{i,k}^2 * (Y^T)_{l,j}^2 \\ &+ \frac{1}{\sigma^2} * Z_{k,l}^2 \end{aligned}$$

Since $Z_{k,l}$ follows a Gaussian distribution, the expectation and variance of the Gaussian distribution can be obtained by:

$$\begin{aligned} \mu_{k,l} &= \frac{\sum_{(i,j) \in s_2} (\theta_{i,j} - \frac{1}{2})X_{i,k} * Y_{j,l} + \sum_{(i,j) \in s_1} \frac{1}{2}X_{i,k} * Y_{j,l}}{2 * (\sum_{i,j} \lambda(\zeta_{i,j})X_{i,k}^2 Y_{j,l}^2 + \frac{1}{\sigma^2})}, \\ v_{k,l} &= \frac{1}{\sqrt{2 * (\sum_{i,j} \lambda(\zeta_{i,j})X_{i,k}^2 Y_{j,l}^2 + \frac{1}{\sigma^2})}}. \end{aligned}$$

Next, we derive parameter ρ for $\ln(\beta|\rho)$:

$$\begin{aligned} \ln q(\beta|\rho) &= \sum_{(i,j) \in s_1} E_{q(Z)} \ln p(R_{i,j}, P_{i,j}, Z, \beta) \\ &+ \sum_{(i,j) \in s_2} E_{q(Z, P_{i,j})} \ln p(R_{i,j}, P_{i,j}, Z, \beta) + \text{const} \end{aligned}$$

Expanding the two items $\sum_{(i,j) \in s_1} E_{q(Z)} \ln p(R_{i,j}, P_{i,j}, Z, \beta)$ and $\sum_{(i,j) \in s_2} E_{q(Z, P_{i,j})} \ln p(R_{i,j}, P_{i,j}, Z, \beta)$, then removing irrelevant items we have:

$$\begin{aligned} \ln q(\beta|\rho) &= (\sum_{(i,j) \in s_2} \theta_{i,j} R_{i,j} + \sum_{(i,j) \in s_1} P_{i,j} R_{i,j} + \eta_{10} - 1) \ln \beta_1 \\ &+ [\sum_{(i,j) \in s_2} \theta_{i,j} (1 - R_{i,j}) + \sum_{(i,j) \in s_1} P_{i,j} (1 - R_{i,j}) \\ &+ \eta_{11} - 1] \ln(1 - \beta_1) + [\sum_{(i,j) \in s_2} (1 - \theta_{i,j}) R_{i,j} \\ &+ \sum_{(i,j) \in s_1} (1 - P_{i,j}) R_{i,j} + \eta_{00} - 1] \ln \beta_0 \\ &+ [\sum_{(i,j) \in s_2} (1 - \theta_{i,j}) (1 - R_{i,j}) \\ &+ \sum_{(i,j) \in s_1} (1 - P_{i,j}) (1 - R_{i,j}) + \eta_{01} - 1] \ln(1 - \beta_0) \end{aligned}$$

Because β follows the Beta distribution, we have:

$$\begin{aligned} \rho_{0,0} &= \eta_{0,0}, \\ \rho_{0,1} &= \sum_{(i,j) \in s_2} (1 - \theta_{i,j}) + \eta_{0,1}, \\ \rho_{1,0} &= |s_1| + \eta_{1,0}, \\ \rho_{1,1} &= \sum_{(i,j) \in s_2} \theta_{i,j} + \eta_{1,1}, \end{aligned}$$

where $|s_1|$ is the number of elements in set s_1 . Next, we derive parameter θ for $\ln(P_{i,j}|\theta_{i,j})$:

$$\begin{aligned} \ln q(P_{i,j}|\theta_{i,j}) &= E_{q(Z, \beta)} [\ln p(R_{i,j}, P_{i,j}, Z, \beta)] \\ &= P_{i,j} * \ln[\exp(R_{i,j} * \psi(\rho_{1,0})) \\ &* \exp[(1 - R_{i,j}) * \psi(\rho_{1,1})] * \exp(-\psi(\rho_{1,0} \\ &+ \rho_{1,1})) * \exp(X\mu_2 Y)] + (1 - P_{i,j}) \\ &* \ln[\exp(R_{i,j} * \psi(\rho_{0,0})) * \exp[(1 - R_{i,j}) \\ &* \psi(\rho_{0,1})] * \exp(-\psi(\rho_{0,0} + \rho_{0,1}))] \end{aligned}$$

we define:

$$\begin{aligned} l_1 &= \exp(\psi(\rho_{1,1}) - \psi(\rho_{1,0} + \rho_{1,1}) + X_i \mu Y_j^T) \\ l_2 &= \exp(\psi(\rho_{0,1}) - \psi(\rho_{0,0} + \rho_{0,1})) \end{aligned}$$

Then we get the estimated value of $\theta_{i,j} = \frac{l_1}{l_1+l_2}$. Finally, for variational parameters ζ , we maximize Equ. 9:

$$\ln \sigma(\zeta_{i,j}) - \frac{\zeta_{i,j}}{2} - \lambda(\zeta_{i,j})[(X_i Z Y_j^T)^2 - (\zeta_{i,j})^2] \quad (9)$$

Making it equal to 0, we get the update formula:

$$\zeta_{i,j} = |X_i \mu Y_j^T|. \quad (10)$$

As shown in Alg. 1, in each iteration of the inference we alternatively optimize the variational parameters for $q(Z|\mu, v)$, $q(\beta|\rho)$, $q(P|\theta)$ and the parameters for the lower bound $\sigma(\zeta)$. In each iteration, we first obtain the optimal θ, μ, v, ρ and then we update ζ . The iteration is repeated until convergence is achieved.

```

input : P, R, X, Y
output:  $\mu, v, \rho, \theta, \zeta$ 
1 initialization;
2 repeat
3   for  $Z_{k,l} \in Z$  do
4      $\mu_{k,l} \leftarrow \frac{\sum_{(i,j) \in s_2} (\theta_{i,j} - \frac{1}{2}) X_{i,k} * Y_{j,l} + \sum_{(i,j) \in s_1} \frac{1}{2} X_{i,k} * Y_{j,l}}{2 * (\sum_{(i,j) \in s_2} \lambda(\zeta_{i,j}) X_{i,k}^2 Y_{j,l}^2 + \frac{1}{\sigma^2})}$ ;
5      $v_{k,l} \leftarrow \frac{1}{\sqrt{2 * (\sum_{(i,j) \in s_2} \lambda(\zeta_{i,j}) X_{i,k}^2 Y_{j,l}^2 + \frac{1}{\sigma^2})}}$ ;
6   end
7   for  $\beta$  do
8      $\rho_{0,0} \leftarrow \eta_{0,0}$ ;
9      $\rho_{0,1} \leftarrow \sum_{(i,j) \in s_2} (1 - \theta_{i,j}) + \eta_{0,1}$ ;
10     $\rho_{1,0} \leftarrow |s_1| + \eta_{1,0}$ ;
11     $\rho_{1,1} \leftarrow \sum_{(i,j) \in s_2} \theta_{i,j} + \eta_{1,1}$ ;
12  end
13  for  $(i, j) \in s_2$  do
14     $l_1 = \exp(\psi(\rho_{1,1}) - \psi(\rho_{1,0} + \rho_{1,1}) + X_i \mu Y_j^T)$ ;
15     $l_2 = \exp(\psi(\rho_{0,1}) - \psi(\rho_{0,0} + \rho_{0,1}))$ ;
16     $\theta_{i,j} \leftarrow \frac{l_1}{l_1 + l_2}$ ;
17  end
18  for  $(i, j) \in s_1 + s_2$  do
19     $\zeta_{i,j} \leftarrow |X_i \mu Y_j^T|$ ;
20  end
21 until convergence;

```

Algorithm 1: Inference for FNML

III. EXPERIMENT

A. Experimental Setup

Datasets. We use the same datasets as in [5]: i.e. the drug-target interaction labels are obtained from the latest version of DrugBank (version 3.0) [10]. This data set is referred to as the full data set. Only 0.18% of the drug-target interactions are labelled as positive, none is labelled as negative. As in [5], we also construct a sample dataset, where all the positive interactions are reserved and an equal number of unknown interactions are sampled to be negative. Statistics of the two data sets are shown in Tab. I.

TABLE I
STATISTICS OF THE DATASETS

Data	#Drugs	#Targets	#Positive	#Negative	#Unknown
Full	708	1,512	1,923	0	1,068,573
Sample	708	1,512	1,923	1,923	1,066,650

We use a variety of networks to extract features X, Y . The default feature vectors X are extracted from drug structure similarity network (denoted as ds), where the similarity score between two drugs is calculated using the Tanimoto coefficient [11] according to their chemical structures; The default feature vectors Y are extracted from protein sequence similarity network (denoted as ps), which is constructed by computing the Smith-Waterman score [12] of their primary sequences. In order to evaluate model performance with different features, we also use three extra drug networks: drug-drug interaction network (dd) [10], the drug-disease network (di) [13], the drug-side-effect network (de) [14] and two protein networks: the protein-disease association network (pd) [13], the protein-protein interaction network (pp) [15].

Evaluation. Throughout the experiment section, the major evaluation metric is Area Under Precision Recall curve (AUPR), which is commonly adopted in bioinformatic studies. An auxiliary evaluation metric is Area Under ROC curve (AUROC).

B. Results and Analysis

FNML Performance. We first evaluate the accuracy of DTI prediction of the proposed FNML model. The hyper-parameter settings are as follows. The number of dimensions for drug features are $K = 300$, for target features $L = 300$, hyper-parameters are $\sigma^2 = 1, \eta_0 = 1, \eta_1 = 1$. In this experiment, we use the default features X, Y . The code and data used in FNML are available at: <https://github.com/517515435/FNML>

We compare our FNML model with 5 state-of-the-art methods: (1) DeepWalk [16]: a similarity-based drug-target prediction method that enhances similarity computation by deep learning method within a linked tripartite network. (2) HNM [17]: a network model in which strength between a disease-drug pair is calculated through an iterative algorithm on the heterogeneous graph that also incorporates drug-target information. (3) NetLapRLS [18]: a manifold regularization semi-supervised learning method. (4) PUDTI [6]: an SVM-based optimization model that is trained on negative samples extracted based on positive-unlabeled learning. (5) DTINet [5]: a regression model that learns feature space mapping Z by the loss function $\min_Z \sum_{i,j} (P_{i,j} - (XZY)_{i,j})^2$. We do not change the default settings for all the above comparative methods.

In order to maintain the same experimental setup as [5], we perform the evaluation on two datasets. The first one is on the full dataset, i.e. we randomly segment the whole data set to 10 divisions and conduct 10-fold cross-validation. The second one is on the sample dataset, i.e. keeping the ratio of positive and negative samples to 1 : 1, we conduct random sampling for 10 times and the reported results are averaged over the 10 sets.

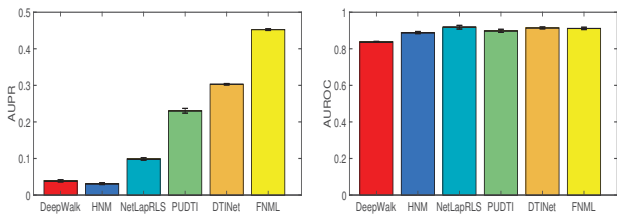


Fig. 2. On the full dataset, FNML significantly boosts AUPR while obtaining comparable AUROC.

The comparative performance on the full dataset is shown in Fig. 2. We can see that (1) FNML model significantly boosts the AUPR performance by 49.32%, compared with the best of state-of-the-art methods. The best comparative method is DTINet, which achieves a 30.29% AUPR. Our FNML model obtains a 45.23% AUPR. As AUPR is well regarded to be a more robust and accurate evaluation metric than AUROC [5], this observation demonstrates the potential of our model. (2) Most of the state-of-the-art methods yield very low AUPR results on the full dataset. This observation again reveals that obtaining a high AUPR performance is challenging on the full dataset. (3) In term of AUROC, the best result is obtained by NetLapRLS. However, the best comparative result is 91.78%, while FNML produces a comparable 91.12% AUROC.

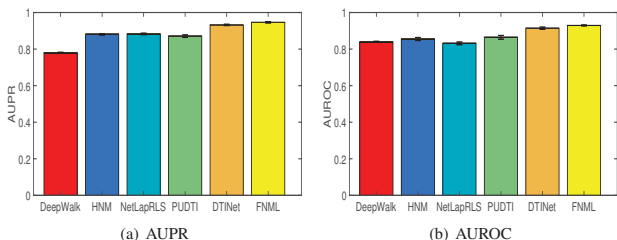


Fig. 3. On the sample dataset, FNML outperforms state-of-the-art methods in terms of both AUPR and AUROC.

The comparative performance on the sample dataset is shown in Fig. 3. We can see that (1) FNML model achieves a better AUPR than all state-of-the-art methods. The best comparative method is DTINet, which achieves 93.20%. Our FNML model obtains a 94.66% AUPR. (2) Most of the state-of-the-art methods have a higher AUPR result on the sample dataset than the full dataset, due to the balanced ratio of positive and negative samples. (3) FNML model outperforms all state-of-the-art models in AUROC performance. The best comparative method is again DTINet, which achieves 91.41%. Our FNML model obtains a 92.93% AUROC. (4) Surprisingly, DeepWalk has a lowest AUPR performance on the sample set. A possible reason is that the network representation extracted by deepwalk is based on homogeneous network structure, and thus is not accurate.

FNML Performance with Different Features. We next study how FNML model performs with different features. We use various combination of X and Y as inputs. That is, we extract X from the four networks on the drug side (i.e. dd, di, de, ds) respectively, extract Y from the three networks

on the protein side (i.e. pp, pd, ps) respectively, and use the 12 combinations as inputs to train the model. The predictions are tested on the full dataset.

We compare the AUPR and AUROC performance of FNML and DTINet. As shown in Fig. 4, FNML outperforms DTINet in most cases. FNML generates better AUPR results for 10 feature combinations out of 12. In term of AUROC, FNML is better for 7 feature combinations. The result shows that the performance improvement is stable. Change of feature representations does not affect FNML's ability to learn a better feature mapping space.

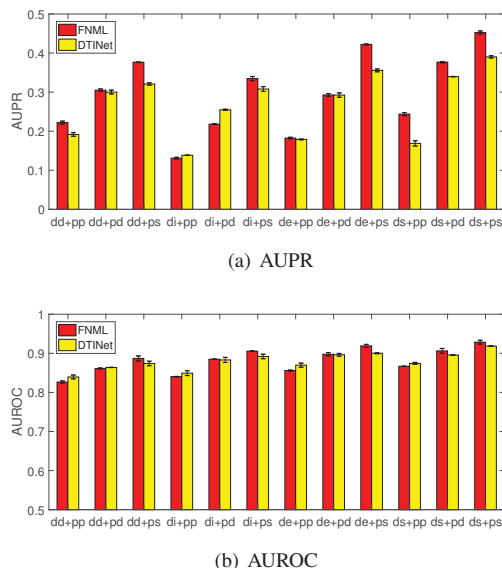
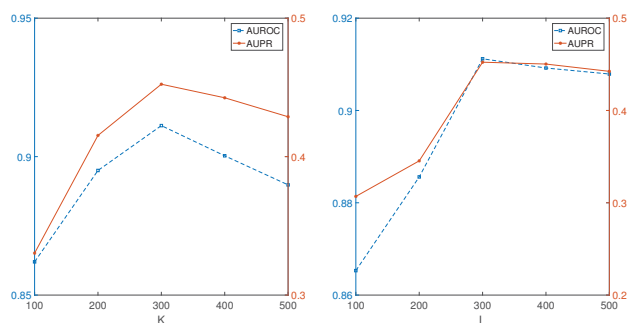


Fig. 4. FNML model consistently outperforms DTINet with different feature inputs.

Number of dimensions. We next study the effects of number of dimensions K, L . We first fix $L = 300$ and tune from $K = 100$ to $K = 500$. We can see from Fig. 5(a) that the best number of drug features is around 300. Then, we fix $K = 300$ and tune from $L = 100$ to $L = 500$. As shown in Fig. 5(b), the best number of target features is 300. An appropriate number of drug features is important. When the number of drug features is too large or too small i.e. $K \geq 400$ or $K \leq 200$, we observe a descent fall in both AUPR and AUROC. However, the model performance is less sensitive to the number of target features. For $L > 300$, AUPR and AUROC remain the same.

IV. CONCLUSION

We propose a novel DTI prediction model based on the assumption that unknown DTI labels are missing not at random. By associating the status of a DTI being labelled or unknown to the sign of the DTI label, our proposed FNML model can learn a better feature mapping from drug feature space to target feature space. We experimentally demonstrate that FNML outperforms state-of-the-art computational DTI identification methods. This work sheds some insights into fully exploiting the information in unknown DTIs. Our future



(a) K number of drug features (b) L number of protein features
 Fig. 5. AUPR and AUROC performance of our model with different number of drug and protein features.

directions include analyzing the missing mechanisms and enhancing the DTI prediction performance by an ensemble scheme.

REFERENCES

- [1] Keiser, M. J. et al. Relating protein pharmacology by ligand chemistry. *Nature biotechnology* 25, 197206 (2007).
- [2] Cheng, A. C. et al. Structure-based maximal affinity model predicts small-molecule druggability. *Nat. Biotechnol.* 25, 7175 (2007).
- [3] Chen, X. et al. Drug-target interaction prediction: databases, web servers and computational models. *Brief. Bioinform.* 17, 696712 (2016).
- [4] Ding, H., Takigawa, I., Mamitsuka, H., and Zhu, S. Similarity-based machine learning methods for predicting drug-target interactions: a brief review. In *Briefings in bioinformatics*, 15(5), 734-747.
- [5] Luo Y, Zhao X, Zhou J, et al. A network integration approach for drug-target interaction prediction and computational drug repositioning from heterogeneous information. In *Nature Communications*, 2017, 8(1).
- [6] Peng L, Zhu W, Liao B, et al. Screening drug-target interactions with positive-unlabeled learning. In *Scientific Reports*, 2017, 7(1): 8087.
- [7] R. J. A. Little and D. B. Rubin. *Statistical Analysis with Missing Data*, 1987.
- [8] Cho H, Berger B, Peng J. Diffusion Component Analysis: Unraveling Functional Topology in Biological Networks In *Research in Computational Molecular Biology*. Springer International Publishing, 2015:62-64.
- [9] Bishop C M. *Pattern Recognition and Machine Learning*. In *Information Science and Statistics*. Springer-Verlag New York, Inc. 2006.
- [10] Knox C, Law V, Jewison T, et al. DrugBank 3.0: a comprehensive resource for Omics research on drugs. In *Nucleic Acids Research*, 2011, 39(Database issue):D1035.
- [11] 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. In *Journal of the American Chemical Society* 125, 1185311865 (2003).
- [12] Smith, T. F. and Waterman, M. S. Identification of common molecular subsequences. In *Journal of molecular biology* 147, 195197 (1981).
- [13] Davis, A. P., Murphy, C. G., Johnson, R., Lay, J. M., Lennon-Hopkins, K., Saraceni-Richards, C., Sciaky, D., King, B. L., Rosenstein, M. C., Wiegiers, T. C., et al. The comparative toxicogenomics database: update 2013. *Nucleic acids research*, 41(D1), D1104D1114.
- [14] Kuhn, M., Campillos, M., Letunic, I., Jensen, L. J., and Bork, P. A side effect resource to capture phenotypic effects of drugs. *Molecular systems biology*, 6(1), 343.
- [15] Keshava Prasad T S, Goel R, Kandasamy K, et al. Human Protein Reference Database—2009 update. *Nucleic Acids Research*, 2009, 37(Database issue):767-72.
- [16] Zong N, Kim H, Ngo V, et al. Deep Mining Heterogeneous Networks of Biomedical Linked Data to Predict Novel Drug-Target Associations. In *Bioinformatics*, 2017, 33(15).
- [17] Wang W, Yang S, Zhang X, et al. Drug repositioning by integrating target information through a heterogeneous network model. In *Bioinformatics*, 2014, 30(20):2923-2930.
- [18] Xia Z, Wu L Y, Zhou X, et al. Semi-supervised drug-protein interaction prediction from heterogeneous biological spaces. In *Bmc Systems Biology*, 2010, 4(S2):1-16.