Multi-Label Robust Factorization Autoencoder and its Application in Predicting Drug-Drug Interactions

Xu Chu, Yang Lin, Jingyue Gao, Jiangtao Wang, Yasha Wang, Leye Wang

{chu_xu, bdly, gaojingyue1997, jiangtaowang, wangyasha}@pku.edu.cn, wly@cse.ust.hk

Abstract

Drug-drug interactions (DDIs) are a major cause of preventable hospitalizations and deaths. Predicting the occurrence of DDIs helps drug safety professionals allocate investigative resources and take appropriate regulatory action promptly. Traditional DDI prediction methods predict DDIs based on the similarity between drugs. Recently, researchers revealed that the predictive performance can be improved by better modeling the interactions between drug pairs with bilinear forms. However, the shallow models leveraging bilinear forms suffer from limitations on capturing complicated nonlinear interactions between drug pairs. To this end, we propose Multi-Label Robust Factorization Autoencoder (abbreviated to MuLFA) for DDI prediction, which learns a representation of interactions between drug pairs and has the capability of characterizing complicated nonlinear interactions more precisely. Moreover, a novel loss called CuXCov is designed to effectively learn the parameters of MuLFA. Furthermore, the decoder is able to generate high-risk chemical structures of drug pairs for specific DDIs, assisting pharmacists to better understand the relationship between drug chemistry and DDI. Experimental results on real-world datasets demonstrate that MuLFA consistently outperforms state-ofthe-art methods; particularly, it increases 21.3% predictive performance compared to the best baseline for top 50 frequent DDIs. We also illustrate various case studies to demonstrate the efficacy of the chemical structures generated by MuLFA in DDI diagnosis.

Introduction

Drug-drug interactions (DDIs) are common situations in which a drug affects the efficacy and safety of another drug when both are administered together, resulting in many adverse drug reactions (ADRs) that may cause severe injuries or even be responsible for deaths (Oato et al. 2016). Most DDIs are discovered by accident once a drug is already on the market (Percha and Altman 2013). However, early detection of DDIs at preclinical stage based on data such as drug chemical structures helps drug safety professionals allocate investigative resources and take appropriate regulatory action (Zhang et al. 2015). In fact, predicting DDIs based on chemical structures is possible since the concept that similar chemical structures bring about similar biological properties has been employed over the years by medicinal chemists (Traphagen 2002; Gedeck and Lewis 2008). With the accumulation of massive adverse events data caused by DDI collected by systems such as FDA Adverse Event Reporting System (FAERS)¹, using computational methods to predict DDI becomes feasible, and DDI prediction is drawing increasing attention of the AI research community.

Since ADRs (e.g., Nausea, Emesis, High blood pressure, etc.) associated with DDIs are important for both clinical and pharmaceutical decisions (Vilar, Friedman, and Hripcsak 2017), we can classify DDIs into different types according to different ADRs in DDI prediction. The DDI prediction problem studied in this paper is defined as: Predicting the occurrence of different types of DDIs between a pair of drugs based on the drug features (e.g., chemical structures).

In literature, similarity-based methods have been widely applied for DDI prediction (Vilar et al. 2012; Zhang et al. 2015; Abdelaziz et al. 2017; Kastrin, Ferk, and Leskoek 2018). These methods first calculate the similarities between each pair of drugs based on independently extracted drug features, and then based on those similarities, they predict the type of DDIs between drug pairs. The idea of the prediction is that, if drug A is similar to drug B, then the drugs that have DDIs with drug A are likely to have the same type of DDIs with drug B. Recently, researchers prove that the predictive performance can be improved by better modeling the interactions between drug pairs by bilinear forms (Jin et al. 2017). Despite the impressive results achieved with this approach, the question remains as to whether there is a better approach that could be used to capture the complicated nonlinear interactions between drug pairs more precisely.

Nowadays, deep representation learning methods has been found advantageous in modeling the complicated non-linear relations (Krizhevsky, Sutskever, and Hinton 2012; Collobert et al. 2011). DDI prediction is actually a multi-label classification problem. A natural idea is to harness the power of representation learning to learn a representation of interactions between drug pairs that is efficient for classification. With adequate labeled data, supervised methods are encouraged to represent the classes of DDIs in a linearly separable way. However, of all possible combinations of two drugs, only a small proportion of drug pairs are labeled with DDIs. With insufficient labeled data, the supervised algorithms suffer severe overfitting and would achieve low predictive performance. On the other hand, unsupervised algo-

¹https://open.fda.gov/data/faers/

rithms such as autoencoders allow us to exploit information hidden in unlabeled data and therefore improve performance of DDI prediction. However, the representations learnt by unsupervised methods would in general entangle factors related to the types of DDIs with other class-irrelevant factors and therefore introduce undesired bias for DDI prediction. The aforementioned analysis inspires us that if we manage to disentangle the categorization factors across all factors, then we may use a supervised learning signal to train the representation of categorization factors. At the same time an unsupervised learning signal could be employed to exploit hidden information from large unlabeled data, regularizing the supervised learning process and thereby enhancing the generalization of model by restraining overfitting. The single-label learning method factorization autoencoder (FAE) (Cheung et al. 2015) introduced a dimerous representation. FAE considers the class label to be part of the representation and the remaining part encode the class-irrelevant factors. To disentangle the categorization factors across all latent factors, FAE introduced a mini-batch based crosscovariance loss termed XCov that penalizes the covariance matrix of each dimension in the class-relevant coding part and each dimension in the class-irrelevant coding part.

However, simply extending the vanilla FAE to highdimensional multi-label situations such as DDI prediction would fail to achieve the best result. The reason is that XCov estimates the cross-covariance in every mini-batch separately. When the batch size is small, the cross-covariance estimator employed by XCov would result in gradient descent directions with large variance and thus hurt the performance. On the other hand, a large batch size method tends to converge to sharp minimizers of the training function and result in a degradation in the quality of the model as measured by the ability to generalize (Keskar et al. 2017). We introduce a novel mini-batch based robust cumulative crosscovariance loss CuXCov that approximates the full-batch statistics, which guarantees a more accurate estimation and allows for better classification performance as well as robust representations of interactions between drug pairs.

The decoder of the autoencoder can be utilized as a feature generator. With designed fabricared inputs, the generator could output feature vector associated with specific category. In the context of DDI prediction, the generator could output vectors describing high-risk chemical structures associated with specific types of DDIs and therefore providing hints for drug research and development process.

In summary, our contributions are summarized as follows:

- We proposed Multi-Label Robust Factorization Autoencoder (called MuLFA) for DDI prediction. MuLFA inherits the puissant expressing power of deep neural network to characterize the complicated nonlinear interactions between drug pairs and is capable of leveraging hidden information in unlabeled data.
- We proposed a robust cumulative cross-covariance loss CuXCov that approximate the full-batch statistics, which is designed to effectively learn the parameters of MuLFA by disentangling categorical factors across latent factors and thus improving the classification performance.

- We construct a dimerous representation, with which we could generate high-risk chemical structures for specific types of DDIs, assisting pharmacists to better understand the relationship between drug chemistry and DDI and providing hints in drug research and development process.
- Experimental results on real-world datasets demonstrate that MuLFA consistently outperforms state-of-the-art methods; particularly, it increases 21.3% predictive performance compared to the best baseline for top 50 frequent DDIs. We also illustrate various case studies to demonstrate the efficacy of the chemical structures generated by MuLFA in DDI diagnosis.

Related Work

In literature, there has been a long line of studies in DDI prediction based on preclinical data. From the methodological perspective, the most representative DDI prediction methods are two-stage similarity-based methods. Firstly, drug features are extracted for each drug independently, based on which similarites are calculated for all drug pairs. Secondly, based on the idea that similar drugs are also biologically similar, different strategies are employed to predict DDIs based on the similarites between drugs, e.g., nearest neighbor method (Vilar et al. 2012), label propagation method (Zhang et al. 2015), link prediction method (Abdelaziz et al. 2017; Kastrin, Ferk, and Leskoek 2018). Recently, researchers proved that the predictive performance can be improved by better modeling the interactions between drug pairs by bilinear forms (Jin et al. 2017). However, the shallow model being used suffers from limitations on capturing complicated nonlinear interactions between drugs pairs.

Inspired by the success of deep representation learning methods (Krizhevsky, Sutskever, and Hinton 2012; Collobert et al. 2011), we introduce a method to better capture the complicated interaction relationship between drug pairs as well as leveraging hidden information in unlabeled data.

Some previous works try to employ additional preclinical data, such as targets and enzymes, to enhance DDI prediction (Cheng and Zhao 2014; Takeda et al. 2017; Zhang et al. 2017). However, such additional data are not always available for all drugs of interest (Abdelaziz et al. 2017), limiting the usage scope of those methods. In future, we will consider how to incorporate more preclinical drug features, if such extra data can be obtained.

Preliminaries

We first define some notations to prepare our method.

Definition 1. Drug Chemical Structure Data

The drug chemical structure data contains substructure profiles of m drugs. Define set \mathcal{A} as $\mathcal{A} = \{d_1, d_2, ..., d_m\}$. The feature vector describing chemical structure of drug d_p is represented by a l-dimensional vector \mathbf{D}_p .

Definition 2. Chemical Structure Vector of a Pair of Drugs Let \mathcal{B} denote the set of all possible drug pairs in chemical structure data, i.e., $\mathcal{B} = \{(d_p, d_q) | d_p, d_q \in \mathcal{A}, 1 \leq p < q \leq m\}$. The chemical structure vector \mathbf{d}_{pq} of drug pair (d_p, d_q)

is the concatenation of substructure profile vectors \mathbf{D}_p and \mathbf{D}_q , namely, $\mathbf{d}_{pq}^T = (\mathbf{D}_p^T, \mathbf{D}_q^T)$.

Definition 3. *DDI Data*

The DDI data contains n drugs and v types of DDIs. We denote $\mathcal{C}=\{\tilde{d}_1,\tilde{d}_2,...,\tilde{d}_n\}$ as the set of n drugs and $\mathcal{R}=\{r_1,r_2,...,r_v\}$ as the set of v types of DDIs. Each type of DDI event corresponds to a specific type of adverse drug reaction². For a given type of DDI event $r_i\in\mathcal{R}, i=1,2,...,v$, the data only records credible drug pairs that could be the causing factor. We define set \mathcal{D} as $\mathcal{D}=\{(\tilde{d}_p,\tilde{d}_q)|\tilde{d}_p,\tilde{d}_q\in\mathcal{C},1\leq p< q\leq n,\tilde{d}_p \text{ and }\tilde{d}_q \text{ are reported to be associated with at least one }r_i,i=1,2,\cdots,v\}.$ We denote the set of drug pairs associate with r_i as \mathcal{E}_i , $\mathcal{E}_i\subseteq\mathcal{D}, i=1,2,...,v$.

Definition 4. Set of Labeled Drug Pairs and Set of Unlabeled Drug pairs

We let set \mathcal{D} be the set of labeled drug pairs for all v tasks. And set $\mathcal{F} = \mathcal{B} - \mathcal{D}$ be the set of unlabeled drug pairs³.

Definition 5. Set of Positive Samples and Set of Negative Samples for occurrence of the i-th Type of DDI r_i

Let set \mathcal{E}_i be the set of positive samples and set $\mathcal{G}_i = \mathcal{D} - \mathcal{E}_i$ be the set of negative samples for task r_i .

Problem Statement

The problem of DDI prediction is formulated as follows. **Input:**

- The set of chemical structure vectors $\{\mathbf{d}_k\}_{k=1}^{|\mathcal{B}|}$, where $\mathbf{d}_k = \mathbf{d}_{pq}$ for some p and q such that $(d_p, d_q) \in \mathcal{B}$.
- DDI training data: $\{r_i\}_{i=1}^v$ and the corresponding sets of positive samples and negative samples, $\{\mathcal{E}_i, \mathcal{G}_i\}_{i=1}^v$.

Output: Predicted occurrence of r_i of testing drug pairs for each type of DDI event $i = 1, 2, \dots, v$.

Factorization Autoencoder and XCov Loss

Researchers have proposed a semi-supervised factorization autoencoder that could be used in single-label learning(Cheung et al. 2015). Specifically, given an input \mathbf{x} and its corresponding one-hot class label vector \mathbf{y} for a dataset \mathcal{D} , FAE learns the high-level representation (the last-layer of the encoder) in the form of concatenation of two vectors, i.e., $f_{\Theta}(\mathbf{x})^T = (\hat{\mathbf{y}}^T, \mathbf{z}^T)$. FAE considers the class label to be part of the high-level representation of its corresponding input. Using class labels, FAE incorporates supervised learning to a subset of high-level representation, transforming them into observed variable $\hat{\mathbf{y}}$. The remaining subset \mathbf{z} accounts for the remaining variation of dataset. To disentangle the categorization factors from other latent variables, FAE adds a mini-batch based cross-covariance loss (termed XCov). XCov loss prevents vector \mathbf{z} from encoding input

variations due to class label by penalizing the covariance matrix of each dimension in the class-relevant coding part and each dimension in the class-irrelevant coding part.

$$L_{XCov} = \frac{1}{2} \sum_{ij} \left[\frac{1}{N} \sum_{s} (\hat{y}_{i}^{s} - \bar{y}_{i})(z_{j}^{s} - \bar{z}_{j}) \right]^{2}.$$
 (1)

N is mini-batch size, and \bar{y}_i, \bar{z}_i denote means over examples. s is an index over examples and i,j index feature dimensions. In our problem, FAE cannot be directly used for two reasons: (1) FAE is a single-label learning method. (2) XCov estimates cross-covariance separately in each minibatch and could result in descent directions with large variance. To this end, we extend the FAE to address the two issues in next sections.

Methods

Framework Overview

The overall neural network within MuLFA is built with an autoencoder structure. Figure 1 shows the architecture of our proposed method. The network consists of H+1 layers where H is an even number. The first $\frac{H}{2}$ hidden layers are encoders to learn a representation of each input and the last $\frac{H}{2}$ hidden layers are decoders to reconstruct the input. For ease of illustration, we first fix some notations. Let

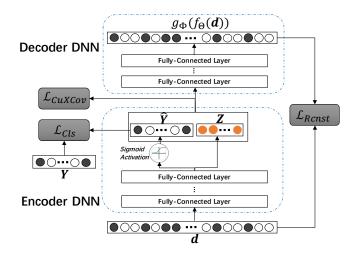


Figure 1: The architecture of MuLFA.

 $\mathbf{L}^{(0)} = \mathbf{d} \in \mathbb{R}^{2l}$ denote an input to the first layer and

$$\mathbf{L}^{(h)} = t^{(h)} ((\mathbf{W}^{(h)})^T \mathbf{L}^{(h-1)} + \mathbf{b}^{(h)}) \in \mathbb{R}^{d_h}$$
 (2)

be the output of the h-th layer, $h=1,2,\cdots,H$. d_h denotes the dimension of the output at the h-th layer and $t^{(h)}$ s are activation functions, which we take ReLU 5 for all hidden layers except the high-level representation layer, i.e., for h=1

²ADRs caused by 2 co-administered drugs rather than ADRs caused by a single drug.

³In general, databases recording chemical structure collect as much information as possible. It is reasonable to assume $\mathcal{D} \subseteq \mathcal{B}$.

 $^{^4}$ $|\mathcal{S}|$ denotes the cardinality of set \mathcal{S} . $\{e_i\}_{i=1}^L$ denotes a set and the index of element e_i in set ranges from 1 to L.

⁵We've tested other activations such as tanh and sigmoid by 10-fold cross-validation, among which ReLU performed best.

$$\frac{H}{2}$$

$$\mathbf{L}^{(\frac{H}{2})} = \begin{pmatrix} \hat{\mathbf{Y}} \\ \mathbf{Z} \end{pmatrix} = ((\hat{y}_1, \cdots, \hat{y}_v)^T, (z_1, \cdots, z_u)^T)^T \in \mathbb{R}^{v+u}.$$
(3)

$$\begin{cases} y_{i} = Sigmoid((\mathbf{W}_{i}^{(\frac{H}{2})})^{T} \mathbf{L}^{\frac{H}{2}-1} + b_{i}^{(\frac{H}{2})}), & i = 1, 2, \cdots, v; \\ z_{i} = (\mathbf{W}_{v+j}^{(\frac{H}{2})})^{T} \mathbf{L}^{\frac{H}{2}-1} + b_{v+j}^{(\frac{H}{2})}, & j = 1, 2, \cdots, u. \end{cases}$$

Where $\mathbf{W}^{\frac{H}{2}} = (\mathbf{W}_1^{(\frac{H}{2})}, \mathbf{W}_2^{(\frac{H}{2})}, \cdots, \mathbf{W}_{v+u}^{(\frac{H}{2})})$ and $\mathbf{b}^{\frac{H}{2}} = (b_1^{\frac{H}{2}}, b_2^{\frac{H}{2}}, \cdots, b_{v+u}^{\frac{H}{2}})^T$. The output of the top layer is \mathbf{L}^H . We define the function of encoder as f_{Θ} and the function of decoder as g_{Φ} . Θ and Φ denote the parameter space of encoder and decoder respectively.

$$f_{\Theta}(\mathbf{d}) = \begin{pmatrix} \hat{\mathbf{Y}} \\ \mathbf{Z} \end{pmatrix} \in \mathbb{R}^{v+u}, \quad g_{\Phi}(f_{\Theta}(\mathbf{d})) = \mathbf{L}^H \in \mathbb{R}^{2l}.$$
 (5)

Our goal of representation learning is to learn a representation $(\mathbf{L}^{(\frac{H}{2})})^T = (\hat{\mathbf{Y}}^T, \mathbf{Z}^T)$ that isolates the categorization factors from other latent factors. $\hat{\mathbf{Y}}$ codes the class labels of the input \mathbf{d} with \hat{y}_i denoting the probability of the occurrence of the i-th event r_i . \mathbf{Z} codes the class-irrelevant factors to serve for semi-supervised learning by preserving as many factors of variation in the data as possible for the sake of reconstruction of input \mathbf{d} .

For $i=1,2,\cdots,v$, the learning of \hat{y}_i can be viewed as a task r_i . In DDI prediction, each \hat{y}_i is corresponding to a type of DDI event. Actually, different types of DDI events are related. For example, if a specific drug pair causes Nausea, then the specific drug pair is likely to cause Emesis. Thus the learning of $\hat{\mathbf{Y}}$ can benefit from Multi-task learning for better exploiting the relatedness among tasks. (Caruana 1997) characterized multi-task learning as an approach to inductive transfer that improves generalization and generalization error bounds (Baxter 1995) by using the domain information contained in the training signals of related tasks as an inductive bias. Regarding outputs of layers $\{\mathbf{L}_h\}_{h=1}^{H/2-1}$ as shared representation, hard under-sampled types of DDIs that could not be learnt in isolation are able to be learnt, and what is learnt for each task help other tasks be learnt better.

To factor the entangled source of variation relevant for categorization apart from other factors across the representation $\mathbf{L}^{(\frac{H}{2})}$, We introduce a mini-batch based robust cumulative cross-covariance CuXCov loss to approximate the full-batch statistics, aiming at minimizing the entries in cross-covariance matrix of $\hat{\mathbf{Y}}$ and \mathbf{Z} for all samples.

Loss Function

We now present details of how to train our model.

Goal The training objective of MuRFA is to minimize the weighted integration of 3 losses:

$$\min_{\Theta} \mathcal{L}_{Cls} + \beta \mathcal{L}_{CuXCov} + \gamma \mathcal{L}_{Rcnst}. \tag{6}$$

Where hyperparameters $\beta>0$ and $\gamma>0$ control relative weights of \mathcal{L}_{CuXCov} and \mathcal{L}_{Rcnst} over \mathcal{L}_{Cls} . \mathcal{L}_{Cls} penalizes the discrepancy between groundtruth labels and predicted occurrence probabilities of samples in labeled sets.

 \mathcal{L}_{CuXCov} penalizes the estimated values of entries in cross-variance matrix of $\hat{\mathbf{Y}}$ and \mathbf{Z} . \mathcal{L}_{Rcnst} is the general reconstruction error in autoencoders, penalizing the discrepancy between input \mathbf{d} and $g_{\Phi}(f_{\Theta}(\mathbf{d}))$ for all training samples.

CuXCov Loss When the number of different classes is large, (e.g., in DDI prediction, the known number of different classes of DDIs is more than 1,000.) the cross-covariance matrix between $\hat{\mathbf{Y}}$ and \mathbf{Z} would in general require a large sample size to achieve an accurate estimation. However, the XCov loss in (1) estimates cross-covariance separately in each mini-batch and could result in descent directions with large variance. To address this issue, inspired by (Chang, Xiang, and Hospedales 2018), we propose a cumulative loss CuXCov that approximates the full-batch cross-covariance matrix. This cumulative strategy can trace back to (Welford 1962), where the author proposed an accurate, one-pass, incremental approach to estimate the second central moment. Let $\Sigma_f^k, \Sigma_c^k, \bar{\Sigma}_m^k, \Sigma_a^k$ denote the full, cumulative, mini-batch, approximate cross-covariance estimator at the k-th training step respectively. The approximation works as follows:

$$\begin{cases} \Sigma_c^k = \alpha \Sigma_c^{k-1} + \Sigma_m^k, & with \ \Sigma_c^0 = \mathbf{0}, \\ p^k = \alpha p^{k-1} + 1, & with \ p^0 = 0. \end{cases}$$
 (7)

Where $\alpha \in [0,1]$ is the decay rate. Then let

$$\Sigma_a^k = \Sigma_c^k / p^k. \tag{8}$$

 Σ_a^k would start converging to Σ_f^k as k-th gets larger. Let

$$\begin{bmatrix} \hat{\tilde{\mathbf{Y}}} \\ \hat{\mathbf{Z}} \end{bmatrix} \in \mathbb{R}^{(v+u) \times N} \tag{9}$$

denote the high-level representation over a mini-batch with size N, where $\hat{\mathbf{Y}} = (\hat{\mathbf{Y}}^1, \hat{\mathbf{Y}}^2, \cdots, \hat{\mathbf{Y}}^N)$, and $\tilde{\mathbf{Z}} = (\mathbf{Z}^1, \mathbf{Z}^2, \cdots, \mathbf{Z}^N)$.

We write mini-batch cross-covariance estimator in matrix form

$$\Sigma_m^k = 1/N(\tilde{\hat{\mathbf{Y}}}(\mathbf{I} - \mathbf{e}\mathbf{e}^T))(\tilde{\mathbf{Z}}(\mathbf{I} - \mathbf{e}\mathbf{e}^T))^T = 1/N\tilde{\hat{\mathbf{Y}}}\mathbf{H}\tilde{\mathbf{Z}}^T \in \mathbb{R}^{v \times u}.$$
(10)

Where $\mathbf{e} \in \mathbb{R}^N$ is a column vector with all entries being 1, **I** is the identity matrix, and $\mathbf{H} = (\mathbf{I} - \mathbf{e}\mathbf{e}^T)(\mathbf{I} - \mathbf{e}\mathbf{e}^T)^T \in \mathbb{R}^{N \times N}$. From (7), (8) and (10), we have

$$\Sigma_a^k = 1/p^k (\alpha \Sigma_c^{k-1} + 1/N \tilde{\hat{\mathbf{Y}}} \mathbf{H} \tilde{\mathbf{Z}}^T) \in \mathbb{R}^{v \times u}.$$
 (11)

Our goal is to minimize all entries in Σ_a^k . We define CuX-Cov loss as:

$$\mathcal{L}_{CuXCov} = trace((\Sigma_a^k)^T \Sigma_a^k)/2. \tag{12}$$

The gradient of CuXCov loss with respect to $\hat{\boldsymbol{Y}}$ and $\tilde{\boldsymbol{Z}}$ is:

$$\frac{\partial \mathcal{L}_{CuXCov}}{\partial \tilde{\mathbf{Z}}} = \frac{\alpha}{N(p^k)^2} (\Sigma_a^{k-1})^T \tilde{\hat{\mathbf{Y}}} \mathbf{H} + \frac{1}{N^2(p^k)^2} \tilde{\mathbf{Z}} \mathbf{H}^T \tilde{\hat{\mathbf{Y}}}^T \tilde{\hat{\mathbf{Y}}} \mathbf{H}.$$
(13)

$$\frac{\partial \mathcal{L}_{CuXCov}}{\partial \tilde{\mathbf{Y}}} = \frac{\alpha}{N(p^k)^2} \Sigma_a^{k-1} \tilde{\mathbf{Z}} \mathbf{H}^T + \frac{1}{N^2(p^k)^2} \tilde{\mathbf{Y}} \mathbf{H} \tilde{\mathbf{Z}}^T \tilde{\mathbf{Z}} \mathbf{H}^T.$$
(14)

Classification Loss Our model can achieve multi-label learning. The learnt $\hat{\mathbf{Y}} = (\hat{y}_1, \hat{y}_2, \cdots, \hat{y}_v)^T$ in representation for each input **d** predicts the probabilities of occurrence of r_i for $i = 1, 2, \cdots, v$. We use the labeled set for each task to supervise the training of \hat{y}_i s.

$$\mathcal{L}_{Cls} = -\sum_{\mathbf{d}_s \in \mathcal{E}_i \cup \mathcal{G}_i} \sum_{i=1}^v (\lambda y_i^s log(\hat{y}_i^s) + (1 - y_i^s) log(1 - \hat{y}_i^s))$$
(15)

Where i indexes tasks and s indexes examples. y_i^s is the label of s-th sample in labeled set with 1 coding the occurrence of event r_i and 0 otherwise. $\lambda>1$ denotes the relative confidence of positive samples over negative samples. In the DDI prediction case, we record a drug pair associated with a type of DDI r_i if the drug pair causes the corresponding adverse drug reaction. And no record of a drug pair associated with r_i does not induce the conclusion that the drug pair would never cause the corresponding adverse drug reaction. Thus, it is more appropriate to treat positive samples and negatives samples discriminatively, and put larger weights on positive samples.

Reconstruction Loss We often have a large amount of unlabled training data and relatively little labeled training data. To better explore information contained in unlabeled data, the architecture of autoencoder allows us to introduce the reconstruction loss over the whole training data.

$$\mathcal{L}_{Rcnst} = \sum_{\mathbf{d}_i \in \mathcal{B}} ||\mathbf{d}_i - g_{\Phi}(f_{\Theta}(\mathbf{d}_i))||^2.$$
 (16)

Where \mathcal{B} denotes the whole training set. Reconstruction loss also acts as a regularization term for the classification loss.

Canonical Samples Generator

After the auto-encoder is trained,

$$\hat{\Theta}, \hat{\Phi} = \arg\min_{\Theta, \Phi} \mathcal{L}_{Cls} + \beta \mathcal{L}_{CuXCov} + \gamma \mathcal{L}_{Rcnst}.$$
 (17)

decoding function $G_{\hat{\Phi}}(\mathbf{L}^H)$ can be used as sample generator given a high-level representation \mathbf{L}^{H_0} . Formally, we define **Definition 6.** The Canonical Sample with respect to the *i-th Category:* Let $\hat{\mathbf{Y}}_i$ be a vector with all entries being 0 except the *i-th* entry being 1. Let

$$\hat{\mathbf{Z}} = E_{\mathcal{B}}(\mathbf{Z}|\hat{\Theta}, \hat{\Phi}), \quad \mathbf{C}_i = G_{\hat{\Phi}}(\begin{pmatrix} \hat{\mathbf{Y}}_i \\ \hat{\mathbf{Z}} \end{pmatrix})$$
 (18)

We name C_i the feature vector of the canonical sample with respect to the *i*-th category, for $i=1,2,\cdots,v$. The corresponding sample is named as the canonical sample with respect to the *i*-th category.

Take the last u entries of $1/|\mathcal{B}|\sum_{\mathbf{d}_i\in\mathcal{B}}f_{\hat{\Theta}}(\mathbf{d}_i)$ as $\hat{\hat{\mathbf{Z}}}$ to esti-

mate
$$\hat{\mathbf{Z}}$$
. Then an approximation of \mathbf{C}_i is $\tilde{\mathbf{C}}_i = G_{\hat{\Phi}}(\begin{pmatrix} \hat{\mathbf{Y}}_i \\ \hat{\mathbf{Z}} \end{pmatrix})$.

In DDI prediction, when the inputs of MuLFA are chemical structure vectors of drug pairs, the canonical sample with respect to the i-th type of DDI event would be approximate canonical chemical structures of a pair of drugs that could jointly lead to the corresponding adverse drug reaction.

Experiments

Datasets

DDI Data The DDI data we use is from Twosides database⁶ (Tatonetti et al. 2012). It contains 645 drugs and 1318 types of DDIs, and in total 63473 drug pairs associated with DDI reports that makes the labeled set \mathcal{D} with $|\mathcal{D}|=63473$.

Drug Chemical Structure Data The chemical structure features we use are extracted from Pubchem⁷ substructure fingerprint, and are binary coded as an 881-bit feature vector, each bit representing a Boolean determination of the presence of a substructure. For fair comparison, we only extract structure features of all drugs appeared in Twosides database, namely, $\mathcal{A} = \mathcal{C}$ and $|\mathcal{B}| = \binom{645}{2}$.

Experiment I On Classification Performance

Methods for Comparison

Baselines We compared our model with the following methods.

- Nearest Neighbor (NN) method in (Vilar et al. 2012).
- Label Propagation (LP) method in (Zhang et al. 2015).
- Dyadic Prediction (DP) method in (Jin et al. 2017).

Variants of our methods We also studied the effect of different components proposed in our method. The networks were trained by back-propagation via Adam optimizer.

- MuLFA: Our proposed model.
- MuLFA-R: Our proposed model without considering the reconstruction penalty.
- MuLFA-X: Our proposed model without considering the cross-covariance penalty.
- MuLFA-X⁺ Our proposed model without **Z** in high-level representation, leading to no cross-covariance penalty.

Evaluations We randomly selected 10% of drugs and masked all DDIs associated with these drugs for testing in alignment with (Zhang et al. 2015). DDIs associated with drugs not in testing set are used for training all models and we use 10-fold cross-validation to tune all hyperparameters of different methods. For testing data, we evaluate all methods on different collections of DDIs. For a given collection of DDIs, we randomly selected 50% of the testing set for evaluation and repeated the selection-evaluation process for 50 times. We report the mean and standard deviation of the Area Under Precision-Recall Curve(AUPR) over 50 repetitions. DDI data is highly unbalanced with small positive sample sets and much larger negative sample sets for different types of DDIs. It was shown in (Davis and Goadrich 2006) that the area under Receiver Operating Characteristic curve (AUROC) is not appropriate for unbalanced data and metrics such as AUPR should be used instead.

We adopt a different strategy in constructing the sets of negative samples compared with DR. DR took the complement of Twosides DDI interactions $\{\mathcal{G}_i \cup \mathcal{F}\}_{i=1}^v$ as negative samples, while we took $\{\mathcal{G}_i\}_{i=1}^v$ as negative samples and \mathcal{F} as the unlabeled set in our method. We examined the

⁶http://tatonettilab.org/resources/tatonetti-stm.html

⁷https://pubchem.ncbi.nlm.nih.gov/

drug pairs in the \mathcal{F} and found some drug pairs should not be co-administered, e.g., Carbamazepine (ID=2554) and Isoniazid (ID=3767) are a pair of drugs in \mathcal{F} . Concurrent use of Carbamazepine and Isoniazid may result in increased carbamazepine exposure and increased risk of isoniazid-induced hepatotoxicity (Wright, Stokes, and Sweeney 1982). For drug pairs in \mathcal{F} constructed from Twosides, no credible data can be observed and we are not able to judge whether the other drug pairs in \mathcal{F} interact or not, thus we take \mathcal{F} as the unlabeled set. For all methods, we utilized DDI interactions from Twosides $\{\mathcal{E}_i\}_{i=1}^v$ as positive samples.

Table 1: AUPR of MuLFA against Baselines.

Methods	Top 50 DDIs	Top 51-100 DDIs	Top 101-150 DDIs
NN	0.367(0.0030)	0.265(0.0023)	0.224(0.0028)
LP	0.360(0.0031)	0.254(0.0025)	0.211(0.0029)
DR	0.375(0.0025)	0.272(0.0019)	0.233(0.0026)
MuLFA	0.455(0.0029)	0.313(0.0027)	0.270(0.0033)

Table 2: AUPR of MuLFA against its Variants.

Methods	Top 50 DDIs	Top 51-100 DDIs	All DDIs
MuLFA-R	0.424(0.0035)	0.293(0.0028)	0.236(0.0023)
MuLFA-X	0.429(0.0028)	0.294(0.0026)	0.243(0.0025)
$MuLFA-X^+$	0.423(0.0028)	0.288(0.0029)	0.241(0.0020)
MuLFA	0.455(0.0029)	0.313(0.0027)	0.278(0.0019)

Results and Discussion Table 1 and 2 compare the performance of the proposed method against competing methods and the variants of MuLFA on different of collections of DDIs. The Top $X_1 - X_2$ DDIs in the table denotes the collection of X_1 -th to X_2 -th most frequent DDIs. The tables show that MuLFA consistently achieves higher AUPRs as compared to all competing methods at different settings. More concretely, from Table 1 we can see that DR and MuLFA outperform NN and LP consistently by better modeling the interactions between drug pairs. As the tasks become harder, namely, when the labeled sets become more and more unbalanced, the performances of similarities-based methods decay faster for failing to exploit the relatedness among tasks. Furthermore, our proposed method achieves improvement consistently over DR by up to 21.3% because the deep neural network framework is able to capture the complicated nonlinear interaction relationship between drug pairs more precisely. Table 2 shows the performances of MuLFA and its variants. MuLFA-R is a supervised model which fails to leverage the information contained in the large unlabeled set. MuLFA-X outperforms MULFA-X⁺ consistently demonstrating the effectiveness of including **Z** for encoding the class-irrelevant factors. MuLFA outperforms MuLFA-X significantly demonstrates the necessity of disentangle the categorization factors across all latent factors.

Sensitivity Analysis We study the sensitivity of two important parameters in our approach. The decay rate α in CuXCov loss and the relative confidence of positive samples over negative samples λ . We evaluate our model by grid

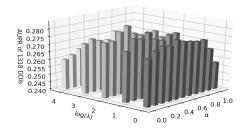


Figure 2: AUPRs of different combinations of (α, λ) .

search as shown in figure 2, computing AUPRs of predicting 1318 DDIs simultaneously for different combinations of (α,λ) . The best combination is $(\alpha,\lambda)=(0.3,4).$ Small values of α lead to estimations that lose too much information about early steps, while large values lead to too much emphasis on early steps and gain deficient information of a new batch. Small values of λ make the model fail to discriminate the positive samples from less reliable negative samples, while large values of λ prevent the model predicting an example being positive because of the high penalty.

CuXCov Loss vs XCov Loss We study the performances of models leveraging CuXCov loss against models leveraging XCov loss. Note that CuXCov loss degenerates to XCov loss when $\alpha = 0$. For $\alpha = \{0, 0.3\}$, we evaluate the predictive performances of 2 losses as shown in Figure 3. In the left panel of Figure 3, we study the influence of number of tasks on performance for mini-batch size being 200. We evaluate the methods under different dimensions of $\hat{\mathbf{Y}}$, i.e., the number of considered DDIs, from 400 to 1200 with step length being 200. The result demonstrates that CuXCov model outperforms XCov model consistently and decays slower than XCov model as dimension of $\hat{\mathbf{Y}}$ is getting higher. In the right panel of Figure 3, we study the influence of batch size on performance for all 1318 DDIs. We evaluate the methods under different batch sizes ranging from 50 to 3200. The result demonstrates that CuXCov model outperforms XCov model consistently. The performances of two models are increasing as batch size getting larger at first for better estimation of cross-covariance with smaller estimation variance. As the batch size continually getting larger, the performances of two models decay because of the degradation of generalization caused by convergence to sharp minimizers of the training function (Keskar et al. 2017).

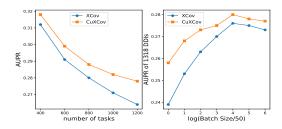


Figure 3: AUPRs of model leveraging CuXCov against model leveraging XCov loss in different settings.

Experiment II

Case Studies On Generated Chemical Structures

On Generated Canonical Chemical Structures of DDIs with High Co-occurrence Different types of DDIs are related. Usually, the ingenerate biological properties of DDIs with high co-occurrence are similar or even identical. For example, Sinus Tachycardia frequently co-occurs with Nausea and both of them can be activated by increased catecholamine release (Koch et al. 1990). Moreover, similar chemical structures may bring about similar medication effect mechanisms (Traphagen 2002), which may cause DDIs. Thus it is reasonable to conjecture that the chemical structures of drug pairs that cause highly frequent co-occurred DDIs should be similar to each other. The left panel of fig-

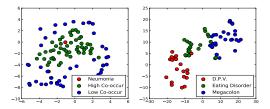


Figure 4: Visualization of generated canonical chemical structure vectors of some chosen DDIs by t-SNE.

ure 4 is the 2-dimensional t-SNE (Maaten and Hinton 2008) visualization of generated canonical chemical structure vectors of Neumonia, 50 most frequent DDI types co-occur with Neumonia and 50 least frequent DDI types co-occur with Neumonia. From the figure we can see that DDIs co-occur with Neumonia more frequently result in more similar generated canonical chemical structures. Megacolon, Eating Disorder and Disorder Perpheral Vascular are 3 randomly selected types of DDIs with low co-occurrence with each other, and the right panel of figure 4 is the visualization of canonical chemical structure vectors of these 3 DDIs accompanied with 20 frequent co-occurred DDIs for each respectively. We can see 3 clusters in the embedding space in alignment with our conjecture.

On Explanation of Generated Canonical Structures For better interpretability we first constructed a reference setby randomly selecting 10% of all drugs from Twosides database and masked all DDIs associated with these drugs. We used the hyperparameters tuned in experiment I and trained MuLFA with the remain data. After the network is trained, we use the decoder to generate the canonical samples for top 20 DDIs with highest AUPRs. We consult pharmacy experts on the generated canonical structures. They confirmed that two classes of frequent generated high-risk structures are actually structurally similar to an inhibitor and a substrate of Cytochromes P450. Inhibitors and substrates are of high-risk to interact with other drugs. The pharmacy experts also confirm and explain the efficacy of the generated chemical structures of drug pairs with respect to DDI type Hypoventilation.

On Understanding the Frequent Generated High-risk Structures We compare the feature vectors of the generated canonical samples with the vectors of drugs in the reference set, and we record the top 3 nearest neighbors for each generated drug structure. There are 120 counts in total, among of which structures resembling Valproic Acid (30/120) and structures resembling Fentanyl (25/120) are most frequent.

Cytochromes P450 (CYPs) are proteins of the superfamily that functions as an important enzyme system for drug metabolism. CYPs catalyze a wide range of oxidative reactions and are the most important pathway for drug metabolism. A drug can act as a substrate, an inducer or an inhibitor of CYPs. Inducers can increase the activity of the enzyme, and accelerate the metabolism of itself or other drugs. Inhibitors can attenuate the activity of the enzyme and slow down the metabolism of itself or other drugs. Inhibitors may also increase drug concentration that poses toxicity risks. Valproic acid is thought as an inhibitor of CYP450 2C9 (CYP2C9). Co-administration of valproic acid with drugs that are primarily metabolized by CYP2C9 may result in increased drug concentration and adverse reactions would be observed. Fentanyl is the substrate of CYP3A4. Adverse events may occur if Fentanyl is co-administered with drugs that are the inducers or inhibitors of CYP3A4.

On Explanation of Generated Canonical Sample We compare the generated canonical feature vector of DDI type Hypoventilation with the drugs in reference set and find Ilopost and Venlafaxine are two most structurally similar drugs. From the chemical structure perspective, Iloprost $(C_{22}H_{32}O_4)$ is an eicosanoid, derived from the cyclooxygenase pathway of arachidonic acid metabolism, functioning as a stable analog of prostacyclin (PGI2). Iloprost is thought to promote benefit in pulmonary arterial hypertension (PAH) through vasodilation, antiproliferative effects, and inhibition of platelet aggregation(Baker and Hockman 2005). Venlafaxine $(C_{17}H_{27}NO_2)$ is a cyclohexanol and phenylethylamine derivative that functions as a serotonin-norepinephrine reuptake inhibitor (SNRI). In vitro studies (Sarma 2010) suggest that Venlafaxine would impact platelet aggregation. When Venlafaxine and Iloprost are administered together, Venlafaxine would decrease the drug efficacy of Iloprost by affecting the inhibition of platelet aggregation. Thus hypoventilation would be observed as a consequence of inadequate treatment of PAH.

Conclusion

In this paper, we propose a novel semi-supervised representation learning approach MuLFA for DDI prediction. We construct a dimerous representation for drug pairs, with which we can not only predict different types of DDIs simultaneously but also generate high-risk chemical structures for specific types of DDIs. We conduct extensive experiments on large-scale real-world data. The results demonstrate better classification performance of MuLFA than state-of-theart prediction methods based on chemical structures. We also illustrate various case studies to demonstrate the efficacy of the chemical structures generated by MuLFA.

References

- [Abdelaziz et al. 2017] Abdelaziz, I.; Fokoue, A.; Hassanzadeh, O.; Zhang, P.; and Sadoghi, M. 2017. Large-scale structural and textual similarity-based mining of knowledge graph to predict drugdrug interactions. *Web Semant*. 44(C):104–117.
- [Baker and Hockman 2005] Baker, S. E., and Hockman, R. H. 2005. Inhaled iloprost in pulmonary arterial hypertension. *Annals of Pharmacotherapy* 39(7-8):1265–1274.
- [Baxter 1995] Baxter, J. 1995. Learning internal representations. In *Proceedings of the eighth annual conference on Computational learning theory*, 311–320. ACM.
- [Caruana 1997] Caruana, R. 1997. Multitask learning. *Machine learning* 28(1):41–75.
- [Chang, Xiang, and Hospedales 2018] Chang, X.; Xiang, T.; and Hospedales, T. M. 2018. Scalable and effective deep cca via soft decorrelation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 1488–1497
- [Cheng and Zhao 2014] Cheng, F., and Zhao, Z. 2014. Machine learning-based prediction of drugdrug interactions by integrating drug phenotypic, therapeutic, chemical, and genomic properties. *Journal of the American Medical Informatics Association Jamia* 21(2):278–86.
- [Cheung et al. 2015] Cheung, B.; Livezey, J. A.; Bansal, A. K.; and Olshausen, B. A. 2015. Discovering hidden factors of variation in deep networks. *International Conference on Learning Representations workshop*.
- [Collobert et al. 2011] Collobert, R.; Weston, J.; Bottou, L.; Karlen, M.; Kavukcuoglu, K.; and Kuksa, P. 2011. Natural language processing (almost) from scratch. *Journal of Machine Learning Research* 12(Aug):2493–2537.
- [Davis and Goadrich 2006] Davis, J., and Goadrich, M. 2006. The relationship between precision-recall and roc curves. In *Proceedings of the 23rd international conference on Machine learning*, 233–240. ACM.
- [Gedeck and Lewis 2008] Gedeck, P., and Lewis, R. A. 2008. Exploiting qsar models in lead optimization. *Current opinion in drug discovery & development* 11(4):569–575.
- [Jin et al. 2017] Jin, B.; Yang, H.; Xiao, C.; Zhang, P.; Wei, X.; and Wang, F. 2017. Multitask dyadic prediction and its application in prediction of adverse drug-drug interaction. In *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence*.
- [Kastrin, Ferk, and Leskoek 2018] Kastrin, A.; Ferk, P.; and Leskoek, B. 2018. Predicting potential drug-drug interactions on topological and semantic similarity features using statistical learning. *Plos One* 13(5):e0196865.
- [Keskar et al. 2017] Keskar, N. S.; Mudigere, D.; Nocedal, J.; Smelyanskiy, M.; and Tang, P. T. P. 2017. On large-batch training for deep learning: Generalization gap and sharp minima. In 5th International Conference on Learning Representations.
- [Koch et al. 1990] Koch, K. L.; Stern, R. M.; Vasey, M. W.; Seaton, J. F.; Demers, L. M.; and Harrison, T. S. 1990. Neu-

- roendocrine and gastric myoelectrical responses to illusory self-motion in humans. *Am J Physiol* 258(1):304–10.
- [Krizhevsky, Sutskever, and Hinton 2012] Krizhevsky, A.; Sutskever, I.; and Hinton, G. E. 2012. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, 1097–1105.
- [Maaten and Hinton 2008] Maaten, L. v. d., and Hinton, G. 2008. Visualizing data using t-sne. *Journal of machine learning research* 9(Nov):2579–2605.
- [Percha and Altman 2013] Percha, B., and Altman, R. B. 2013. Informatics confronts drugdrug interactions. *Trends in Pharmacological Sciences* 34(3):178–184.
- [Qato et al. 2016] Qato, D. M.; Wilder, J.; Schumm, L. P.; Gillet, V.; and Alexander, G. C. 2016. Changes in prescription and over-the-counter medication and dietary supplement use among older adults in the united states, 2005 vs 2011. *Jama Internal Medicine* 176(4):473.
- [Sarma 2010] Sarma, A. 2010. Venlafaxine-induced ecchymoses and impaired platelet aggregation. *European Journal of Haematology* 77(6):533–537.
- [Takeda et al. 2017] Takeda, T.; Hao, M.; Cheng, T.; Bryant, S. H.; and Wang, Y. 2017. Predicting drug-drug interactions through drug structural similarities and interaction networks incorporating pharmacokinetics and pharmacodynamics knowledge. *Journal of Cheminformatics* 9(1):16.
- [Tatonetti et al. 2012] Tatonetti, N. P.; Patrick, P. Y.; Daneshjou, R.; and Altman, R. B. 2012. Data-driven prediction of drug effects and interactions. *Science translational medicine* 4(125):125ra31–125ra31.
- [Traphagen 2002] Traphagen, L. M. 2002. Do structurally similar molecules have similar biological activity? *Journal of Medicinal Chemistry* 45(19):4350–8.
- [Vilar et al. 2012] Vilar, S.; Harpaz, R.; Uriarte, E.; Santana, L.; Rabadan, R.; and Friedman, C. 2012. Drugdrug interaction through molecular structure similarity analysis. *Journal of the American Medical Informatics Association Jamia* 19(6):1066.
- [Vilar, Friedman, and Hripcsak 2017] Vilar, S.; Friedman, C.; and Hripcsak, G. 2017. Detection of drug-drug interactions through data mining studies using clinical sources, scientific literature and social media. *Briefings in Bioinformatics*.
- [Welford 1962] Welford, B. P. 1962. Note on a method for calculating corrected sums of squares and products. *Technometrics* 4(3):419–420.
- [Wright, Stokes, and Sweeney 1982] Wright, J. M.; Stokes, E. F.; and Sweeney, V. P. 1982. Isoniazid-induced carbamazepine toxicity and vice versa: a double drug interaction. *New England Journal of Medicine* 307(21):1325–1327.
- [Zhang et al. 2015] Zhang, P.; Wang, F.; Hu, J.; and Sorrentino, R. 2015. Label propagation prediction of drug-drug interactions based on clinical side effects. *Scientific Reports* 5:12339.
- [Zhang et al. 2017] Zhang, W.; Chen, Y.; Liu, F.; Luo, F.; Tian, G.; and Li, X. 2017. Predicting potential drug-drug

interactions by integrating chemical, biological, phenotypic and network data. *Bmc Bioinformatics* 18(1):18.