

Poster: Knowledge Graph Completion to Predict Polypharmacy Side Effects

Brandon Malone¹[0000-0002-7027-3157], Alberto García-Durán¹, and Mathias Niepert¹

NEC Laboratories Europe, Kurfürsten-Anlage 36, 69115 Heidelberg, Germany
{brandon.malone,alberto.duran,mathias.niepert}@neclab.eu

Abstract. The polypharmacy side effect prediction problem considers cases in which two drugs taken individually do not result in a particular side effect; however, when the two drugs are taken in combination, the side effect manifests. In this work, we demonstrate that multi-relational knowledge graph completion achieves state-of-the-art results on the polypharmacy side effect prediction problem. Empirical results show that our approach is particularly effective when the protein targets of the drugs are well-characterized. In contrast to prior work, our approach provides more interpretable predictions and hypotheses for wet lab validation.

Keywords: Knowledge graph · embedding · side effect prediction.

1 Introduction

Disease and other health-related problems are often treated with medication. In many cases, though, multiple medications may be given to treat either a single condition or to account for co-morbidities. However, such combinations significantly increase the risk of unintended side effects due to unknown drug-drug interactions.

In this work, we show that multi-relational knowledge graph (KG) completion gives state-of-the-art performance in predicting these unknown drug-drug interactions. The KGs are multi-relational in the sense that they contain edges with different types. We formulate the problem as a multi-relational link prediction problem in a KG and adapt existing graph embedding strategies to predict the interactions. In contrast to prior approaches for the polypharmacy side effect problem, we incorporate interpretable features; thus, our approach naturally yields explainable predictions and suggests hypotheses for wet lab validation. Further, while we focus on the side effect prediction problem, our approach is general and can be applied to any multi-relational link prediction problem.

Much recent work has considered the problem of predicting drug-drug interactions (e.g. [2,13] and probabilistic approaches like [9]). However, these approaches only consider *whether* an interaction occurs; they do not consider the *type of interaction* as we do here. Thus, these methods are not directly comparable. The recently-proposed DECAGON approach [14] is most similar to ours; they

	Count
Proteins	19 089
Drugs	645
Protein-protein interactions	715 612
Drug-drug interactions	4 649 441
Drug-protein target relationships	11 501
Mono side effects	174 977
Distinct mono side effects	10 184
Distinct polypharmacy side effects	963

Table 1. Size statistics of the graph

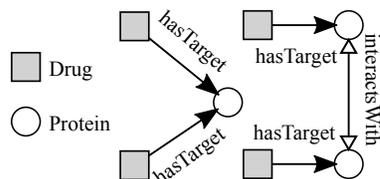


Fig. 1. Types of relational features.

also predict types of drug-drug interactions. However, they use a complicated combination of a graph convolutional network and a tensor factorization. In contrast, we use a neural KG embedding method in combination with a method to incorporate rule-based features. Hence, our method explicitly captures meaningful *relational features*. Empirically, we demonstrate that our method outperforms DECAGON in Section 4.

2 Datasets

We use the publicly-available, preprocessed version of the dataset used in [14].¹ It consists of a multi-relational knowledge graph with two main components: a protein-protein and a drug-drug interaction network. Known drug-protein target relationships connect these different components. The protein-protein interactions are derived from several existing sources; it is filtered to include only experimentally-validated physical interactions in human. The drug-drug interactions are extracted from the TWOSIDES database [11]. The drug-protein target relationships are experimentally-verified interactions from the STITCH [10] database. Finally, the SIDER [6] and OFFSIDES [11] databases were used to identify mono side effects of each drug. Please see Table 1 for detailed statistics of the size and density of each part of the graph. For more details, please see [14]. Each drug-drug link corresponds to a particular polypharmacy side effect. Our goal will be to predict missing drug-drug links.

3 Methods

KG embedding methods learn vector representations for entities and relation types of a KG [1]. We investigate the performance of DISTMULT [12], a commonly-used KG embedding method whose symmetry assumption is well-suited to this problem due to the symmetric nature of the drug-drug (polypharmacy side effect) relation type. The advantage of KG embedding methods are their efficiency and their ability to learn fine-grained entity types suitable for downstream tasks without hand-crafted rules. These embedding methods, however, are less interpretable than rule-based approaches and cannot incorporate domain knowledge.

A *relational feature* is a logical rule which is evaluated in the KG to determine its truth value. For instance, the formula $(\text{drug}_1, \text{hasTarget}, \text{protein}_1) \wedge (\text{drug}_2, \text{hasTarget}, \text{protein}_1)$ corresponds to a binary feature which has value 1 if both drug_1 and drug_2 have protein_1 as a target, and 0 otherwise. In

¹ Available at <http://snap.stanford.edu/decagon>

this work, we leverage relational features modeling drug targets with the relation type `hasTarget` and protein-protein interactions with the relation type `interactsWith`. Figure 1 depicts the two features types we use in our polypharmacy model. For a pair of entities (\mathbf{h}, \mathbf{t}) , the relational feature vector is denoted by $\mathbf{r}_{(\mathbf{h}, \mathbf{t})}$. Relational features capture concrete relationships between entities; thus, as shown in Section 4, they offer explanations for our predictions.

KBLRN is a recently proposed framework for end-to-end learning of knowledge graph representations [4]. It learns a product of experts (PoE) [5] where each expert is responsible for one feature type. In the context of KG representation learning, the goal is to train a PoE that assigns high probability to true triples and low probabilities to triples assumed to be false. Let $\mathbf{d} = (\mathbf{h}, \mathbf{r}, \mathbf{t})$ be a triple. The specific experts we use are defined as

$$f_{(\mathbf{x}, \mathbf{L})}(\mathbf{d} | \theta_{(\mathbf{x}, \mathbf{L})}) = \begin{cases} \exp((\mathbf{e}_{\mathbf{h}} * \mathbf{e}_{\mathbf{t}}) \cdot \mathbf{w}^{\mathbf{x}}) \\ 1 \text{ for all } \mathbf{r}' \neq \mathbf{r} \end{cases} \text{ and } f_{(\mathbf{x}, \mathbf{R})}(\mathbf{d} | \theta_{(\mathbf{x}, \mathbf{R})}) = \begin{cases} \exp(\mathbf{r}_{(\mathbf{h}, \mathbf{t})} \cdot \mathbf{w}_{\mathbf{rel}}^{\mathbf{x}}) \\ 1 \text{ for all } \mathbf{r}' \neq \mathbf{r} \end{cases}$$

where $*$ is the element-wise product, \cdot is the dot product, $\mathbf{e}_{\mathbf{h}}$ and $\mathbf{e}_{\mathbf{t}}$ are the embedding of the head and tail entity, respectively, and $\mathbf{w}^{\mathbf{x}}, \mathbf{w}_{\mathbf{rel}}^{\mathbf{x}}$ are the parameter vectors for the embedding and relational features for relation type \mathbf{r} . The probability of triple $\mathbf{d} = (\mathbf{h}, \mathbf{r}, \mathbf{t})$ is now

$$p(\mathbf{d} | \boldsymbol{\theta}) = \frac{f_{(\mathbf{x}, \mathbf{L})}(\mathbf{d} | \theta_{(\mathbf{x}, \mathbf{L})}) f_{(\mathbf{x}, \mathbf{R})}(\mathbf{d} | \theta_{(\mathbf{x}, \mathbf{R})})}{\sum_{\mathbf{c}} f_{(\mathbf{x}, \mathbf{L})}(\mathbf{c} | \theta_{(\mathbf{x}, \mathbf{L})}) f_{(\mathbf{x}, \mathbf{R})}(\mathbf{c} | \theta_{(\mathbf{x}, \mathbf{R})})},$$

where \mathbf{c} indexes all possible triples. As proposed in previous work, we approximate the gradient of the log-likelihood by performing negative sampling [4].

4 Experimental results

We now empirically evaluate our proposed approach based on multi-relational knowledge graph completion to predict polypharmacy side effects.

Dataset construction We follow the common experimental design previously used [14] to construct our dataset. The knowledge graph only contains “positive” examples for which polypharmacy side effects exist. Thus, we create a set of negative examples by randomly selecting a pair of drugs and a polypharmacy side effect which does not exist in the knowledge graph. We ensure that the number of positive and negative examples of each polypharmacy side effect are equal. We then use stratified sampling to split the records in training, validation and testing sets.

We use an instance of the relational feature types depicted in Figure 1 if it occurs at least 10 times in the KG. We choose these relational feature types because they offer a biological explanation for polypharmacy side effects; namely, a polypharmacy side effect may manifest due to unexpected combinations or interactions on the drug targets.

Baselines We first compare our proposed approach to DECAGON [14]. Second, we consider each drug as a binary vector of indicators for each mono side effect and gene target. We construct training, validation and testing sets by concatenating the vectors of the pairs of drugs described above. We predict the likelihood of each polypharmacy side effect given the concatenated vectors.

Complete DECAGON dataset We first consider the same setting considered previously [14]. As shown in Table 2(top), our simple baseline, DISTMULT, and KBLRN all outperform DECAGON.

Drug-drug interactions only Next, we evaluate polypharmacy side effect prediction based solely on the pattern of other polypharmacy side effects. Specifically, we completely remove the drug-protein targets and protein-protein interactions from the KG; thus, we use only the drug-drug polypharmacy side effects in the training set for learning. We focus on DISTMULT and KBLRN since they outperformed the other methods in the first setting.

Surprisingly, the results in Table 2(middle) show that both DISTMULT and KBLRN perform roughly the same (or even *improve* slightly) in this setting, despite discarding presumably-valuable drug target information. However, as shown in Table 1, few drugs have annotated protein targets. Thus, we hypothesize that the learning algorithms ignore this information due to its sparsity.

Drugs with protein targets only To test this hypothesis, we remove all drugs which do not have any annotated protein targets from the KG (and the associated triples from the dataset). That is, the drug target information is no longer “sparse”, in that all drugs in the resulting KG have protein targets.

The results in Table 2(bottom) paint a very different picture than before; KBLRN significantly outperforms DISTMULT. These results show that the combination of learned (or embedding) features and relational features can significantly improve performance when the relational features are present in the KG.

Explanations and hypothesis generation The relational features allow us to explain predictions and generate new hypotheses for wet lab validation. We chose one of our high-likelihood predictions and “validated” it via literature evidence. In particular, the ranking of the drug combination CID115237 (paliperidone) and CID271 (calcium) for the side effect “pain” increased from 24 223 when using only the embedding features (of 58 029 pairs of drugs for which “pain” is not a known side effect) to a top-ranked pair when also using the relational features. Inspection of the relational features shows that the interaction between lysophosphatidic acid receptor 1 (LPAR1) and matrix metalloproteinase 2 (MMP2) is particularly important for this prediction. The MMP family is known to be associated with inflammation (pain) [7]. Independently, calcium already upregulates MMP2 [8]. Paliperidone upregulates LPAR1, which in turn has been shown to promote MMP activation [3]. Thus, paliperidone indirectly exacerbates the up-regulation of MMP2 already caused by calcium; this, then, leads to increased pain. Hence, the literature confirms our prediction discovered due to the relational features.

5 Discussion

We have shown that multi-relational knowledge graph completion can achieve state-of-the-art performance on the polypharmacy side effect prediction problem. Further, relational features offer explanations for our predictions; they can then be validated via the literature or wetlab. In the future, we plan to extend this work by considering additional features of nodes in the graph, such as Gene Ontology annotations for the proteins and chemical structure of the drugs.

Method	AuROC	AuPR	AP@50
Baseline	0.896	0.859	0.812
DECAGON (values reported in [14])	0.872	0.832	0.803
DISTMULT	0.923	0.898	0.899
KBLRN	0.899	0.878	0.857
DISTMULT (drug-drug interactions only)	0.931	0.909	0.919
KBLRN (drug-drug interactions only)	0.894	0.886	0.892
DISTMULT (drugs with protein targets only)	0.534	0.545	0.394
KBLRN (drugs with protein targets only)	0.829	0.797	0.774

Table 2. The performance of each approach on the pre-defined test set. The measures are: area under the receiver operating characteristic curve (AuROC), area under the precision-recall curve (AuPR), and the average precision for the top 50 predictions for each polypharmacy side effect (AP@50). The best result within each group is in bold.

References

- Bordes, A., Usunier, N., García-Durán, A., Weston, J., Yakhnenko, O.: Translating embeddings for modeling multi-relational data. In: *Advances in Neural Information Processing Systems 26* (2013)
- Cheng, F., Zhao, Z.: Machine learning-based prediction of drug-drug interactions by integrating drug phenotypic, therapeutic, chemical, and genomic properties. *Journal of the American Medical Informatics Association* **21**(e2), e278–e286 (2014)
- Fishman, D.A., Liu, Y., Ellerbroek, S.M., Stack, M.S.: Lysophosphatidic acid promotes matrix metalloproteinase (MMP) activation and MMP-dependent invasion in ovarian cancer cells. *Cancer Research* **61**(7), 3194–3199 (2001)
- García-Durán, A., Niepert, M.: KBLrn: End-to-end learning of knowledge base representations with latent, relational, and numerical features. In: *Proceedings of the 34th Conference on Uncertainty in Artificial Intelligence* (2018)
- Hinton, G.E.: Training products of experts by minimizing contrastive divergence. *Neural computation* **14**(8), 1771–1800 (2002)
- Kuhn, M., Letunic, I., Jensen, L.J., Bork, P.: The SIDER database of drugs and side effects. *Nucleic Acids Research* **44**(D1), D1075–D1079 (2016)
- Manicone, A.M., McGuire, J.K.: Matrix metalloproteinases as modulators of inflammation. *Seminars in Cell & Developmental Biology* **19**(1), 34 – 41 (2008)
- Munshi, H.G., Wu, Y.I., Ariztia, E.V., Stack, M.S.: Calcium regulation of matrix metalloproteinase-mediated migration in oral squamous cell carcinoma cells. *Journal of Biological Chemistry* **277**(44), 41480–41488 (2002)
- Sridhar, D., Fakhraei, S., Getoor, L.: A probabilistic approach for collective similarity-based drug–drug interaction prediction. *Bioinformatics* **32**(20), 3175–3182 (2016)
- Szklarczyk, D., Santos, A., von Mering, C., Jensen, L.J., Bork, P., Kuhn, M.: STITCH 5: augmenting protein-chemical interaction networks with tissue and affinity data. *Nucleic Acids Research* **44**, D380–D384 (2016)
- Tatonetti, N.P., Ye, P.P., Daneshjou, R., Altman, R.B.: Data-driven prediction of drug effects and interactions. *Science Translational Medicine* **4**(125), 125ra31 (2012)
- Yang, B., tau Yih, S.W., He, X., Gao, J., Deng, L.: Embedding entities and relations for learning and inference in knowledge bases. In: *Proceedings of the 3rd International Conference on Learning Representations* (2015)

13. Zhang, W., Chen, Y., Liu, F., Luo, F., Tian, G., Li, X.: Predicting potential drug-drug interactions by integrating chemical, biological, phenotypic and network data. *BMC Bioinformatics* **18**(18) (2017)
14. Zitnik, M., Agrawal, M., Leskovec, J.: Modeling polypharmacy side effects with graph convolutional networks. *Bioinformatics* **34**(13), 457–466 (2018)