

CrossWalk: Learning Representations Across Multiple Graphs

Devin de Hueck
Red Hat Inc.
Boston University
ddehueck@redhat.com

ABSTRACT

Graphs provide a wonderful way to interpret the complex relationships found in data. However, this data may be spread across multiple domains, e.g. a python package can be characterized by a language graph from documentation, a python code graph, and a dependency graph. This same problem emerges across a wide variety of data-hungry industries where no single data domain is sufficient to represent the entity being sought after i.e. a python package, a user in social media network, or a web page. In this paper we help address this problem with a two-fold contribution: (1) We introduce the PyPI Graph Dataset that includes three graph data domains. (2) We benchmark this dataset with a novel unsupervised method that extends random walk models to cross multiple domains that we coin CrossWalk.

CCS CONCEPTS

• **Information systems** → **Document representation**; • **Computing methodologies** → *Unsupervised learning*; *Transfer learning*.

ACM Reference Format:

Devin de Hueck. 2018. CrossWalk: Learning Representations Across Multiple Graphs. In *Woodstock '18: ACM Symposium on Neural Gaze Detection, June 03–05, 2018, Woodstock, NY*. ACM, New York, NY, USA, 5 pages. <https://doi.org/10.1145/1122445.1122456>

1 INTRODUCTION

One of deep learning’s most fundamental advances is the ability to learn meaningful representations from raw data. Recently there has been greater emphasis on learning representations from graph data. These methods are concerned with embedding nodes, subgraphs, or entire graphs into some lower dimensional latent space.

Recent graph representation learning methods [5, 6, 12] interpret a graph G as a set of vertices or nodes V , a set of edges E and each node $v_i \in V$ has a feature matrix X_i . In the case that each node has graphical features one must first learn a vector representation to place in each X_i . For example, in a citation graph each node has textual data, i.e. the paper itself, which can be interpreted as a language graph [11], but in the X_i paradigm this is not acknowledged. Instead, popular datasets used in graph representation learning literature provide these embedding in place of learning another graph

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

Woodstock '18, June 03–05, 2018, Woodstock, NY

© 2018 Association for Computing Machinery.
ACM ISBN 978-1-4503-XXXX-X/18/06...\$15.00
<https://doi.org/10.1145/1122445.1122456>

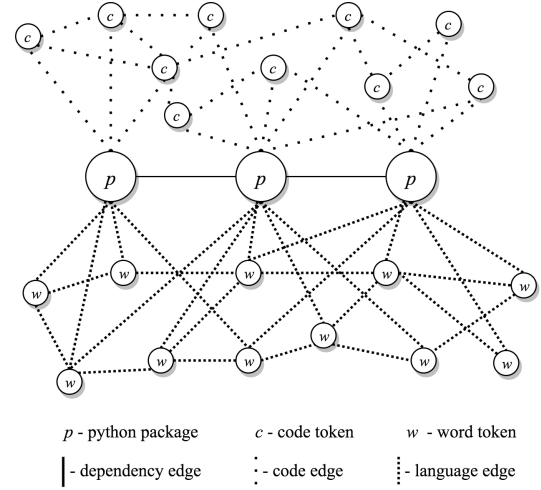


Figure 1: CrossWalk’s Interpretation of the PyPI Graph dataset. Notice the three distinct graphs, i.e domains, and how each relate to a python package.

representation. As a result, methods that follow the X_i cannot learn from the raw data directly.

Furthermore, these datasets poorly represent a class of problems where multiple graphs describe a single entity. A python package is such an entity as each package is characterized by a language graph from its documentation, a code graph, a dependency graph, and possibly more. In this example, no single graph can claim to be sufficient to fully represent a python package. Therefore, to represent this entity, we require a model that can learn representations which encode the relevant information across multiple graph domains with raw data.

In response to these issues, we present a new interpretation that allows for learning across multiple graph domains with raw data. Discarding the X_i interpretation, let G_s be a supergraph with the set of vertices V_s being the set of entities we would like to learn a representation for. Instead of a feature matrix X_i for each node in V , V_s is now connected to a set of graphs: $D = \{G_0 \dots G_n\}$. Each $G_i \in D$ can be of two kinds - global or local. A global graph G_g has a set of edges E_g between the global supergraph nodes V_s (the entities we seek to represent) so there is no need for a new set of nodes. However, a local graph G_l has a new set of nodes V_l and - to relate the supergraph’s nodes to this local graph - E_l also includes edges between V_l and V_s in addition to edges between V_l .

In the case of our python package example, each node in V_s corresponds to a python package in our dataset and D contains three graphs: a language graph, a code graph, and a dependency

graph (See Figure 1). Of these three, the dependency graph is the only global graph as the python packages are connected between themselves. In the case of the local language graph we have a new set of nodes - the words. By connecting each package node in V_s to the word nodes we require the package vector to act as a document vector as shown in [7]. We follow the same idea for the code graph where each local node is a syntax token (perhaps taken from an abstract syntax tree [2]). With this super graph a good python package representation must encode the relevant information in all three graphs.

This new interpretation is the first contribution of the CrossWalk model. The second is the extension of random walk methods [1, 3, 7, 8, 11] to learn these cross-graph representations from raw graph data. Exactly how this is done is left to section 4. Lastly, To validate the CrossWalk model we require a dataset with all raw information available. Unfortunately, using the datasets consistently referenced in the graph representation learning literature is not possible as they follow the feature matrix paradigm outlined above. In response, we have created our own dataset of python packages - the PyPI graph dataset. The PyPI dataset currently has 200k python packages two graphs from our running example: (1) a text graph and (2) a dependency graph. We find that applying CrossWalk to the PyPI dataset validates our novel interpretation and random walk method so that we can learn from raw graph data with raw graphical features.

2 RELATED WORK

Graph Convolutional Networks: This class of networks aggregate local features generalizing the idea of Convolutional Neural Networks commonly used in the field of Computer Vision.

Random Walk Methods: DeepWalk [11], node2vec [3], sub2vec [1] all extend Mikolov et al.’s skip-gram language model [8] to arbitrary graphs with the use of random walks. A random walk begins at a randomly chosen node v , then another node is randomly chosen from v ’s 1-degree neighbors and so forth until a given length is met. As pointed out in [11], a sentence is analogous to a random walk through a language graph so the skip-gram language model [8] can be applied to a corpus of random walks to learn node representations. DeepWalk and word2vec both seek to find the parameters θ that maximize the following:

$$\arg \max_{\theta} \prod_{v_c, v_t}^W p(v_t | v_c; \theta) \quad (1)$$

Where W is all context node v_c and target node v_t pairs over all windows in the random walk corpus. The later work, node2vec [3] extends DeepWalk [11] by introducing hyperparameters that control a random walks tendency to conduct a Depth First Search (DFS) or a Breadth First Search (BFS) through the graph. They find that BFS helps to encode structural roles in a graph while DFS is most helpful for learning community structure.

While [3, 11] are concerned with learning only node representations sub2vec [1] can learn representations of subgraphs. If DeepWalk is to graphs as word2vec is to text then sub2vec is to graphs as doc2vec is to text documents. Just as the paragraph vector in [7] holds the "memory" of a document the subgraph vector in sub2vec maintains the "memory" of a subgraph. These two methods now

seek to maximize the updated version of equation (1) with the subgraph vector s :

$$\arg \max_{\theta} \prod_{v_c, v_t}^W p(v_t | v_c, s; \theta) \quad (2)$$

Other: To the best of our knowledge there is only one other paper that directly addresses the issue of cross-domain representations in the space of graph neural networks (GNNs). [10] introduces Deep Multi-Graph Embedding (DMGE) which relies on a multiple-gradient descent optimizer to balance learning across domains.

Not only do we learn cross-domain representations we also learn a representation for subgraph. Sub2Vec as detailed in [?] builds on the formalization in [11] which itself extends the skip-gram model [8] to arbitrary graphs.

3 THE PYPI GRAPH DATASET

The Python Package Index (PyPI) is home to over two-hundred and twenty-thousand python packages¹. Finding similar packages in terms of use-cases, features, and performance can be very difficult without specific domain knowledge of a task or the python ecosystem. The PyPI graph dataset we provide contains a language graph, and a dependency graph domains which leans to a supergraph of python packages. With our goal being to jointly learn graph representations across multiple graph domains we have reinterpreted the graph representation learning paradigm. As a result no previous dataset gives us all multiple raw graph domains to learn representations across, so it is necessary we introduce our own. Furthermore, the PyPI graph dataset is applicable to many other problems. For example, a research paper can be described by its location within a citation graph, an authorship graph, and by the language graph of the paper itself. A web page can be described by a hyperlink graph and the language graph of the web page. Representing these types of entities with multiple graphs, i.e. multiple domains, requires learning representations that can cross domains and the PyPI Dataset gives a way to test methods that do just that.

3.1 The Language Graph

The language graph consists of either the package’s ReadMe file or its PyPI description (whichever is longer) for 190k python packages. If we were to learn a python package representation from just this graph we would end up with a reasonable measure of package similarity. However, we reiterate that this one domain is a part to the whole concept that is a python package so it can not account for the other domains in this dataset.

3.2 The Direct Dependency Graph

The direct dependency graph is presented as a list of edges between python packages. It is stored as a directed graph with 95k nodes with 1.5 million edges between them.

3.3 Data Collection

Data was collected using GitHub’s GraphQL API and PyPI’s API. We used PyPI as the entry point to every python package². From there

¹See <https://pypi.org/> for the most up-to-date statistics.

²See: <https://pypi.org/simple>

Table 1: Comparing PyPi Graph Dataset to other Graph Datasets with Text Features

| Dataset | # Nodes | # Edges | # Text Feature Size | Features Format |
|------------------------------|---------------|------------------|---------------------|---------------------|
| Cora | 2,708 | 5,429 | 1,433 | Binary Keyword |
| Citeseer | 3,312 | 4,732 | 3,703 | Binary Keyword |
| PubMed | 19,717 | 44,338 | 500 | TF-IDF |
| PyPi Dependency Graph | 95,951 | 1,509,378 | n/a | Raw Language |
| Reddit | 231,443 | 11,606,919 | 300 | GloVe Embeddings |

we can retrieve some textual data and links to GitHub repositories. Using GitHub’s GraphQL API we have easy access to the direct dependencies found in files like requirements.txt, Pipfile, setup.py, etc. The ReadMe files on GitHub tend to be much more detailed than the description on PyPi so we pull that along with topics set by the repository owner for potential downstream tasks.

3.4 Future Data to Collect

While we use just the language and direct dependency graph data we plan on adding more data domains in subsequent versions. The most obvious being the the code of the python packages. We do not yet have a quantitative measure of performance in this domain so we omit it from experiments in this paper. Additionally, another graph domain is a contributor graph where an edge exists between packages if they share a contributor. Lastly, we plan on adding non-graph data such as stars, PyPi downloads, number of versions, number of forks, etc. The abundance of potential data around python packages show that this is a worthwhile endeavor that we hope leads to more sophisticated methods that follow our supergraph interpretation.

4 CROSSWALK

In this section we introduce CrossWalk, an unsupervised model to learn representations across multiple graph domains. This fulfills our goal of being able to learn representations on graph data with graphical features. In the past, methods have relied on embedding the graphical features before any learning takes place ignoring any potential advantages of a jointly learned representation. In CrossWalk we use the raw data to learn these joint representations. Ultimately, there are two clear contributions of CrossWalk: (1) a new supergraph interpretation of graph data with graphical features and (2) a new random walk model that can learn representations across these multiple graph domains.

4.1 A New Interpretation

In the current graph representation learning paradigm a graph G has a set of vertices or nodes V , a set of edges E , and each node $v_i \in V$ has a feature matrix X_i . CrossWalk deals with the case when v_i has graphical properties. Instead of first embedding these graphical properties into a vector that can fit in X_i , CrossWalk aims to jointly learn a representation of the graphical properties of v_i and of G . To do this, we require a new interpretation.

Let G_s be a supergraph with the set of vertices V_s being the set of entities we would like to learn a representation for. Instead of a feature matrix X_i for each node in V , V_s is now connected to a set of graphs, i.e. domains: $D = \{G_0 \dots G_n\}$. Each $G_i \in D$ can

be of two kinds - global or local. A global graph G_g has a set of edges E_g between the global supergraph nodes V_s (the entities we seek to represent) so there is no need for a new set of nodes. However, a local graph G_l has a new set of nodes V_l and - to relate the supergraph’s nodes to this local graph - E_l also includes edges between V_l and V_s in addition to edges between V_l .

This interpretation successfully replaces X_i with the raw graph data that would’ve been represented separately. We can now focus on learning within this new paradigm.

4.2 Our Random Walk Model

With this new interpretation a simple and natural evolution of past random walk methods emerge. We extend the likelihood equation (2) across all of our graph domains D . Where $W^{(i)}$ is the set of all context, target vertex pairs in the i th domain created via windows over random walks (See Appendix for all hyperparameter details). Therefore, our maximum likelihood across all graph domains becomes:

$$\arg \max_{\theta} \prod_i^D \prod_{v_c, v_t}^{W^{(i)}} p(v_t | v_c; \theta) \quad (3)$$

In our formulation, we optimize θ , the graph representations through negative sampling a form of noise contrastive estimation [4, 8, 9]. Therefore our loss function across graph domains becomes:

$$\sum_{i=1}^D [\log \sigma(v_c^{(i)} \cdot v_t^{(i)}) + \sum_k \log \sigma(v_{neg}^{(i)} \cdot v_t^{(i)})] \quad (4)$$

Where k is the number of negative samples and $v_{neg}^{(i)}$ is a negatively sampled vertex vector from domain i . As discussed in section 4.1 a domain can be local or global, therefore θ consists of a single set of global parameters, θ_g and a set of local parameters for each local domain. As a result, the number of parameters in CrossWalk scales linearly with the number of local domains, but stays constant w.r.t to the global domains. Alternatively one could chose to introduce parameters that learn from a subset of graph domains in D . For example, if $D = \{G_0, G_1, G_2\}$ one could instantiate a set of parameters, $\theta_{0 \cup 2}$ that will only cross G_0 and G_2 , ignoring data from G_1 . As a result, CrossWalk can learn up to 2^{D-1} sets of parameters, i.e. representations:

$$N_{\theta} = \sum_{k=1}^D \binom{D}{k} = 2^{D-1} \quad (5)$$

However, even though you can doesn’t mean you should. Clearly, this growth rate w.r.t D is massively expensive in terms of memory. Instead, the benefit of this CrossWalk property is that researchers

Table 2: Node Classification into GitHub Topics

| Domain(s): | Text | Dependency Graph | Joint | X_i Joint |
|------------|---------|------------------|-----------|-------------|
| Method: | Doc2Vec | DeepWalk | CrossWalk | GraphSAGE |
| Micro-F1: | 44.96% | 34.74% | 46.89% | 17.80% |
| Macro-F1: | 34.61% | 21.20% | 36.22% | 3.83% |

Table 3: Link prediction results on PyPI Graph Dataset (Dependency Graph Domain)

| Domain(s): | Text | Dependency Graph | Raw Joint | X_i Joint |
|------------|---------|------------------|-----------|-------------|
| Method: | Doc2Vec | DeepWalk | CrossWalk | GraphSAGE |
| Micro-F1: | 65.94% | 88.71% | 86.52% | 76.56% |
| Macro-F1: | 65.93% | 88.70% | 86.52% | 75.50% |

can get multiple representations over various subsets of D which can be used in downstream tasks to evaluate the necessity of certain graph domains. For example, if $\theta_{0 \cup 2}$ does just as well on some task as θ_g then there is evidence G_1 is not necessary for such a task.

5 REPRESENTATION EVALUATION

5.1 Node Classification

In order to evaluate the learned python package embeddings we use an extrinsic evaluation task. We define this task as the classification of python packages into different categories. GitHub allows for users to label their projects under various tags - as a vast majority of PyPI packages link to a GitHub repository we were able to extract a subset of python packages that are labeled by the creator into these topics. GitHub encourages users to classify their repositories "to help other people find and contribute to your project, you can add topics to your repository related to your project's intended purpose, subject area, affinity groups, or other important qualities."³ This further justifies our interpretation of topics as classification classes.

Following the example of [11] (maybe add others?) we use a one vs. rest logistic regression model and a shallow neural network to classify the learned representations. These model choices reflect our desire to intentionally cripple the transformation ability of the classifier so that the natural discriminatory ability of the learned embeddings is what is being evaluated.

This task empirically shows that no single domain is sufficient to capture a complete representation of a python package. [Perform Ablation Study!]

5.2 Link Prediction

The node classification task has a close association with the language domain in the PyPI graph dataset. To have a more complete benchmark of the dataset we use a link prediction task within the dependency graph domain. This task consists of predicting whether or not a connection should exist between node v and node u with inputs being the concatenation of the respective node embeddings.

To learn the mapping between embeddings and the existence of an edge we use a shallow NN implemented through sklearn's MLP classifier. Using 100k randomly selected edges we create another 100k false edges by randomly selecting source and destination nodes from the true edges. The values reported are the macro and micro F1 scores. We compare single domain representations created from only textual data and only graphical data with a jointly learned CrossWalk embedding.

The results in Table 3. show that while the text-only embeddings do 16% better than chance the graph embedding does the best which makes sense as this task comes from the dependency graph domain of the PyPI dataset. Most notably, the jointly learned embedding performs at just about the same score as the graph only domain. This joint representation is the same one used in the previous topic prediction task where the graph embedding did not do nearly as well as the text embedding. This shows we have achieved our goal of learning a single representation from raw graph data that captures information across multiple domains.

6 CONCLUSION

CrossWalk introduces a new interpretation of graphs with graphical features that allow us to learn joint representations across graph domains from the raw data. Due to this new interpretation current datasets used in the graph representation learning literature do not suffice. As a result, we introduce the PyPI Graph Dataset with raw text and graph data (with more to follow) that is representative of a class of problems where one graph data domain is not sufficient to represent the whole, i.e. research papers, webpages, social media users, etc. CrossWalk also consists of a simple extension to past data domains so that we can learn under this new interpretation. We then apply CrossWalk to the PYPI dataset and find that CrossWalk successfully learns representations that encode the information of the multiple graph domains it was trained on.

REFERENCES

- [1] Bijaya Adhikari, Yao Zhang, Naren Ramakrishnan, and B Aditya Prakash. 2017. Distributed representations of subgraphs. In *2017 IEEE International Conference on Data Mining Workshops (ICDMW)*. IEEE, 111–117.
- [2] Uri Alon, Meital Zilberstein, Omer Levy, and Eran Yahav. 2019. Code2Vec: Learning Distributed Representations of Code. *Proc. ACM Program. Lang.* 3, POPL, Article 40 (Jan. 2019), 29 pages. <https://doi.org/10.1145/3290353>

³<https://help.github.com/en/github/administering-a-repository/classifying-your-repository-with-topics>

- [3] Aditya Grover and Jure Leskovec. 2016. node2vec: Scalable feature learning for networks. In *Proceedings of the 22nd ACM SIGKDD international conference on Knowledge discovery and data mining*. 855–864.
- [4] Michael U Gutmann and Aapo Hyvärinen. 2012. Noise-contrastive estimation of unnormalized statistical models, with applications to natural image statistics. *Journal of Machine Learning Research* 13, Feb (2012), 307–361.
- [5] William L. Hamilton, Rex Ying, and Jure Leskovec. 2017. Inductive Representation Learning on Large Graphs. In *NIPS*.
- [6] Thomas N. Kipf and Max Welling. 2016. Semi-Supervised Classification with Graph Convolutional Networks. *arXiv:cs.LG/1609.02907*
- [7] Quoc V. Le and Tomas Mikolov. 2014. Distributed Representations of Sentences and Documents. *CoRR* abs/1405.4053 (2014). [arXiv:1405.4053](http://arxiv.org/abs/1405.4053) <http://arxiv.org/abs/1405.4053>
- [8] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. 2013. Distributed representations of words and phrases and their compositionality. In *Advances in neural information processing systems*. 3111–3119.
- [9] Andriy Mnih and Yee Whye Teh. 2012. A fast and simple algorithm for training neural probabilistic language models. *arXiv preprint arXiv:1206.6426* (2012).
- [10] Yi Ouyang, Bin Guo, Xing Tang, Xiuqiang He, Jian Xiong, and Zhiwen Yu. 2019. Learning Cross-Domain Representation with Multi-Graph Neural Network. *arXiv preprint arXiv:1905.10095* (2019).
- [11] Bryan Perozzi, Rami Al-Rfou, and Steven Skiena. 2014. DeepWalk: Online Learning of Social Representations. In *Proceedings of the 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* (New York, New York, USA) (KDD '14). ACM, New York, NY, USA, 701–710. <https://doi.org/10.1145/2623330.2623732>
- [12] Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio. 2018. Graph Attention Networks. *International Conference on Learning Representations* (2018). <https://openreview.net/forum?id=rjXmpikCZ> accepted as poster.
- [13] R. Zafarani and H. Liu. 2009. Social Computing Data Repository at ASU. <http://socialcomputing.asu.edu>

A REPRODUCIBILITY

There is a much needed and growing movement in the field of machine learning that encourages the researchers to ensure all results are easily reproduced. In the context of this paper and our results we make this priority.

We want to help create a standard that moves beyond providing a script or notebook that simply outputs the same numbers found in the paper. Instead, we want to build a tool that will allow other researchers to not just match our numbers but build on any aspect of our codebase that may benefit future research.

To accomplish this goal we first note that all of our code will be published on GitHub after the anonymous review process has been completed. Second, we will do the same for the PyPI graph dataset we introduce.

A.1 The Code

The code we plan to publish will include all hyperparameters and decisions made that may affect any of the results shown. This includes jupyter notebooks that run the two evaluation tasks on the PyPI dataset and a larger repository of the CrossWalk codebase. The CrossWalk codebase is written in python using the PyTorch framework and has data preprocessing scripts built in along with the ability to read from HDF5 if all training examples are too large to fit into RAM. Our choice to provide such features reflects our goal to allow users with varying levels of compute power be able to replicate our results and use their own data with ease.

A.2 The Dataset

In addition to releasing the raw data, we also plan on releasing the python scripts used to collect the data. As we plan are planning subsequent versions of the PyPI dataset we believe it is important

Table 4: Hyperparameters used

| Domains: | Text Only | Graph Only | Joint |
|---------------------|------------|------------|------------|
| Epochs: | 20 | 20 | 15 |
| # Examples: | 26,876,633 | 34,529,040 | 61,405,673 |
| Optimizer: | Adam | Adam | Adam |
| Learning Rate: | 1e-3 | 1e-3 | 1e-3 |
| Batch Size: | 4096 | 4096 | 4096 |
| # Negative Samples: | 2 | 2 | 2 |

to release these scripts to account for any differences between dataset versions. We also believe that allowing users to improve these scripts through an open source project will result in a stronger dataset and stronger models.

A.3 Experiment Hyperparameters

For the two experiments detailed in this paper we built one set of representations: text only, dependency graph only, and joint. We used the hyperparameters reported in Table 4. The random walk and tokenization procedure used to create training examples in the graph and text domain respectively are detailed in the codebase mentioned above.