

# Progress Towards Graph Optimization: Efficient Learning of Vector to Graph Space Mappings

Stefan Mautner<sup>1</sup> and Rolf Backofen<sup>1</sup>, Fabrizio Costa<sup>2</sup> \*

1- Albert-Ludwigs-University Freiburg - Department of Computer Science  
Georges-Koehler-Allee 106 - 79110 Freiburg - Germany

2- University of Exeter - Department of Computer Science  
Exeter EX4 4QF - United Kingdom

**Abstract.** Optimization in vector space domains is well understood. However, in high dimensional settings or when dealing with structured data such as sequences and graphs, optimization becomes difficult. A possible strategy is to map graphs to vector codes and use machine learning to learn a map from codes back to graphs. This in turn allows to employ standard optimization techniques over vectors to optimize graphs. Here we propose an approach to invert a vector mapping based on a combination of graph kernels and graph grammars. We evaluate the proposed approach in an artificial setup and on real molecular graphs.

## 1 Introduction

Graphs offer a natural and compact way to encode relational information as domain experts find it intuitive to describe problems in their field of expertise in terms of entities (nodes) and relations between entities (edges) and encode both qualitative and quantitative knowledge using categorical or numerical labels. In the traditional optimization setup, the desired objective is known and an efficient procedure exist to compute it. However there are several types of scenarios that significantly increase the difficulty of the optimization problem, ranging from a large number of variables, to many simultaneous objectives, to high computational costs (e.g. when the objective is measured using simulators or real experiments). Optimization problems over the domain of graphs incur in two additional types of difficulties: the first is the *mutable* state of nodes and edges, the second is the complexity of the graph domain. When problem instances have non mutable structure (i.e. only labels change, but not nodes and edges) one can resort to standard optimization techniques over the vectorized sequence of attributes; when both nodes and edges can change one must resort to non-standard combinatorial approaches. The complexity is then depending on the connectivity density and the size of the labels alphabet or if the labels are over the real field.

Recently a novel stream of machine learning concerned with generating complex artifacts (see for example the Constructive Machine Learning Workshop series) which falls mostly in the mutable, low degree, no real valued attributes case, of which drug design is a prototypical example. While bespoke computational

---

\*This project was funded by Deutsche Forschungs Gesellschaft (DFG grant BA 2168/3-3)

approaches to drug design date back to the 50s [1] and have a long tradition in chemoinformatics, the idea of treating the problem as an instance of a more general graph optimization problem has recently gained more attention, due in part to the application of deep learning techniques [2]. Some approaches [3] try to encode the molecules as strings (SMILES), other [4] use reinforcement learning to design molecules that optimize desired properties. Most approaches face the encoding-decoding problem, i.e. the need to encode graphs as vectors, perform changes in this space and then decode the result back to the graph space. Much effort is placed in trying to develop techniques that are end-to-end differentiable to fully leverage the potential offered by deep network approaches. Our idea is instead to map graphs to vector codes in very high dimensional spaces and use a multi-objective optimization scheme to learn a map from codes back to graphs. In this way we will be able to employ standard and well tested techniques, such as Bayesian Optimization, to perform graph optimization. Here we study the first step of this program and propose an approach to invert a vector mapping based on a combination of graph kernels and graph grammars.

## 2 Methods

Given a procedure to vectorize graphs, that is, a map from graphs to points in a vector space, possibly of a very high dimensionality, our objective is to invert such map and go from any point in the vector space, back to the corresponding graph. We make the following observations.

**High dimensionality.** Due to the *curse-of-dimensionality* it can be difficult to gauge distances, i.e. to ascertain if a given numerical value represents a large distance from the desired solution or if we are in fact reasonably close.

**Manifold linearity violation.** We cannot assume that the vectorization procedure will be capable to obtain a linear manifold, that is, that the geometric average of the codes assigned to two graphs  $G_A, G_B$  will correspond to the code of the “average graph”, where the meaning of average could be for example that it has the same *edit-distance* from both  $G_A$  and  $G_B$ .

**Domain bias.** Knowing that a graph has a desired distance from other graphs could not be enough in practice. Consider the following case: we are looking for the average graph and we select a subgraph  $G_C$  that is in common between two graphs  $G_A, G_B$  of size  $m$ , with, say, half the size  $|G_C| = m/2$ ; in this scenario  $G_C$  is the average graph between  $G_A$  and  $G_B$  but it is unlikely to be what we want (it does not seem to belong to the same *population* given its smaller size, it is just a fragment, not a whole instance). If we enforce a size constraint we still have to define how to *up-size* the graph. The answer would be a form of domain bias.

To address these issues we propose the following approach: given a code vector  $x$  we extract from a set of example graphs two sets of neighbours: a smaller set  $\ell$  that we call *landmarks*, and a larger set  $k$ . The neighbors are obtained by first mapping graphs to high dimensional vector spaces and then performing the neighbor search in this space. We then cast the inverse mapping

problem in a multi objective optimization problem with three objectives: 1) minimize the distance discrepancy with respect to  $\ell$  neighbors; 2) minimize the size discrepancy to the  $\ell$  neighbors and 3) maximize the ranking score of  $k$  neighbors. Objective 1) allows us to gauge distances w.r.t. the landmarks, objective 2) forces the solution to have roughly the same size as the average landmark and objective 3) accounts for the domain bias, since the score is a function of all the  $k$  neighbors and can capture dependencies on an extended set of features.

To address the manifold linearity problem we perform the search operations in the graph space, that is, we start from the landmark graphs and we iteratively modify them to obtain intermediate solutions that are always *feasible*. The feasibility is ensured by using a grammar to generate the candidates for the search that is induced on the local neighborhood of size  $k$ .

In summary, to build a graph corresponding to a given target point in a vector space we 1) start from a set of graphs whose code is close to the target, 2) use a graph grammar to perform a local search and we 3) optimize the result according to a multi objective score.

## 2.1 Graph Vectorization

To map graphs to vector spaces we make use of the graph kernel developed in [5], called Neighborhood Subgraph Pairwise Distance Kernel (NSPDK). The idea is to decompose a graph into small parts (i.e. subgraphs of a pre-defined type) and use them as features, that is create a high dimensional vector representation where each feature is associated to a subgraph and the associated value represents the number of times that the subgraph occurs in the graph instance. In NSPDK the subgraphs are built from small neighborhoods of increasing radii  $r < r_{max}$  (i.e. graphs induced by all nodes that are connected to a root vertex by a shortest path of at most  $r$  edges). Finally, all pairs of such subgraphs whose roots are at a distance not greater than  $d < d_{max}$  are considered as individual features. In a graph with  $n$  nodes, the number of features is proportional (with small multiplicative factor that depends on  $r_{max}$  and  $d_{max}$ ) to  $n$ . This allows us to use a hashing scheme to map these features to integers and use these directly as feature indicators (see [6] for details). While the domain of the possible feature indicators is large ( $10^5$ ), every single graph will have only a few hundreds non zero features. Using a sparse vector encoding allows us to obtain efficient storage and low run times.

## 2.2 Graph Grammar

A graph grammar is a generative device that can model a family of graphs via a finite set of production rules, i.e. instructions that prescribe how to build a derived graph starting from an original graph. Using a graph grammar we will generate a number of perturbed versions of any given graph, we then evaluate them to find the modified version that better fits our objective functions. The rules can be automatically induced from a dataset and ensure that the gener-

ated graphs belong to the same “type”, i.e. that they respect the feasibility constraints implicit in the defining dataset.

We use the Locally Substitutable Graph Grammar (LSGG) defined in [6]. The grammar is characterized by *core graphs* and *interface graphs*. Given a node  $v$ , a core graph is simply a neighborhood graph of radius  $R$  (also called the core *size*) rooted in  $v$ . An interface graph is the difference graph of two neighborhood graphs with the same root and different radii. The value of the difference between the radii is the interface *size*. The generation rules allow to swap cores provided that they are associated to the same interface or context.

### 2.3 Objective Functions

We consider three objectives.

The first objective function minimizes the distance discrepancy to the  $\ell$  landmark graphs  $g_i$  with  $i \in 1, \dots, \ell$ . Given a target code  $x$  we compute the distances  $d_i$  from each landmark. Given a candidate graph  $g$  to score, we vectorize  $g$  and calculate the distances  $d'_i$  to the landmarks. The first objective is then  $f^1(g) = \frac{1}{\ell} \sum_{i \in 1, \dots, \ell} |d_i - d'_i|$ .

The second objective function minimizes the size discrepancy with respect to the  $n$  neighbors, i.e.  $f^2(g) = ||g| - \frac{1}{n} \sum_{i \in 1 \dots n} |g_i||$ .

For the third objective we induce a Ranking SVM [7] on the  $k$  neighbors.

### 2.4 Optimization

We employ a simple greedy strategy: given the multiple objectives, we maintain the Pareto set of the non dominated solutions (i.e. the set of all solutions that are not outperformed on all objectives by another solution). We proceed in an iterative fashion and sample instances (graphs) from the Pareto set and from a separate set of instances that optimize each objective individually. For each of the sampled graphs we generate the set of all possible *perturbed graphs* employing the graph grammar. Each candidate graph is then scored by all objective functions and the Pareto set is updated, removing all graphs that are dominated on all objectives by other graphs in the set. The procedure is repeated until a perfect reconstruction is obtained (i.e.  $f^1(g) = 0$ ) or when a maximal number of iteration is reached.

## 3 Empirical evaluation

**Artificial Grammar:** We evaluate the efficiency of the proposed approach on an artificial dataset, i.e. a set generated by sampling a pre-specified graph grammar. We generate 30 random graphs with 8 vertices, with a maximum vertex degree of 3, and a node (edge) label dictionary size of 4 (2). We induce the LSGG (i.e. we collect all pairs of cores and corresponding interfaces). Starting from each initial graph we apply all possible grammar production rules, remove duplicates, select at random 550 instances. We repeat this procedure three times obtaining graphs that are progressively more diverse. We use 500 instances for

the training set and the remaining 50 as targets. For each target graph we extract 10 landmarks and 50 neighbors from the training set. We induce the graph grammar on the 50 neighbors and we run our optimization procedure. We measure the fraction of times that we can perfectly reconstruct the target (which is quite a stringent requirement). Table 1 shows the success rate for different parameters of the grammar. Note that limiting the cores to size 0 (i.e. allowing to change only a single node label at each move) degrades the performances. Larger cores allow a richer set of moves, and the capacity to change the size of the graph (e.g. replacing a small core with a larger one expands the size of the graph). Larger interface size create very specific substitution rules that model well the domain constraints but can hinder the search. Note that limiting the grammar to the closest 50 neighbors allows to sample more than 70% of the true grammar and it is therefore not a severely limiting factor in the overall design.

C	I	R
2	.5	.62 (31)
2	1	.60 (30)
2	2	.20 (10)
1	.5	.62 (31)
1	1	.42 (21)
1	2	.10 (5)
0	.5	.04 (2)
0	1	.04 (2)
0	2	.00 (0)

Table 1: Fraction of successful reconstructions out of 50 graph generated at random from a grammar. C is the core size, I is the interface size. The number in parenthesis is the number of perfectly reconstructed graphs.

PubChem AID	I=1	I=2
AID119	.85 (17)	.75 (15)
AID1345082	.45 (9)	.85 (17)
AID624202	.80 (16)	.75 (15)
AID977611	.70 (14)	.55 (11)

Table 2: Fraction of successful reconstructions out of 20 chemical compounds selected at random as targets. I is the interface size. The core size was 3. The number in parenthesis is the number of perfectly reconstructed graphs.

The hyper parameters of the proposed approach include the graph vectorizer radius  $r$  and distance  $d$ , the graph grammar core  $C$  and interface  $I$ , size and the optimization procedure number of landmarks  $\ell$  and number of neighbors  $n$ . The optimized values, obtained under cross-validation, are  $r = 3$ ,  $d = 3$ ,  $\ell \in [10, 20]$ ,  $n \in [50, 100]$ . The number of neighbors  $n$  affects the total number of grammar rules extracted and hence the complexity of the searched space and the run time, going from 18 minutes for  $n = 75$  to 39 minutes for  $n = 400$ . A Python implementation is available from [github.com/smautner/reconstruct](https://github.com/smautner/reconstruct).

**Chemical compounds:** We queried the PubChem repository to extract 4 bio assays containing more than 5000 active compounds with the keyword “cancer”. Molecular graphs are on average 3 times larger than the examples generated in the artificial setup (12.5 vs 30 nodes on average) and thus computationally more challenging. To better capture entire cycles (a feature known to be important when predicting chemical properties) we increased the core radius.

To limit the number of possible rule applications we selected rules that occurred at least twice in the 50 neighbors. We investigate two settings in table 2. The interface size 2 version is only slightly worse but completes all runs on 32Gb of memory while the interface size 1 version performs better as one would expect, but even on 128Gb of memory, 25 construction attempts out of 80 failed running out of memory.

## 4 Conclusions

We have exhibited a procedure that can learn how to *perfectly* reconstruct more than 70% of the test graphs once given their corresponding vector code and information on relatively few instances in their neighborhood.

As described in the Introduction, this routine is a useful step in the formulation of a more ambitious generalized graph Bayesian optimization procedure and hence it represent significant progress towards its development.

In the future we will analyze the dependency between the quality of the reconstruction and the graph grammar parameters and we will extended the approach to tackle more challenging graph types, i.e. graphs with larger label dictionaries and with real valued attributes.

## References

- [1] Markus Hartenfeller and Gisbert Schneider. Enabling future drug discovery by de novo design. *Wiley Interdisciplinary Reviews: Computational Molecular Science*, 1(5):742–759, 2011.
- [2] Peter Ertl, Richard Lewis, Eric Martin, and Valery Polyakov. In silico generation of novel, drug-like chemical matter using the lstm neural network. *arXiv preprint arXiv:1712.07449*, 2017.
- [3] Matt J Kusner, Brooks Paige, and José Miguel Hernández-Lobato. Grammar variational autoencoder. *arXiv preprint arXiv:1703.01925*, 2017.
- [4] J. You, B. Liu, R. Ying, V. Pande, and J. Leskovec. Graph Convolutional Policy Network for Goal-Directed Molecular Graph Generation. *ArXiv e-prints*, June 2018.
- [5] Fabrizio Costa and Kurt De Grave. Fast neighborhood subgraph pairwise distance kernel. In *Proceedings of the 27th International Conference on International Conference on Machine Learning*, pages 255–262. Omnipress, 2010.
- [6] Fabrizio Costa. Learning an efficient constructive sampler for graphs. *Artif. Intell.*, 2017.
- [7] Thorsten Joachims. Optimizing search engines using clickthrough data. In *Proceedings of the eighth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 133–142. ACM, 2002.