



UNIVERSITEIT VAN AMSTERDAM

MSC ARTIFICIAL INTELLIGENCE
MASTER THESIS

Knowledge Generation

by
FLORIAN WOLF
12393339

December 9, 2020

48 Credits
April 2020 - December 2020

Supervisor:
Dr Peter BLOEM
Chiara SPRUIJT

Aessor:
Dr Paul GROTH



Contents

1	Introduction	3
1.1	Motivation	3
1.2	Expected Contribution	3
1.3	Research Question	3
2	Related Work	3
2.1	Graph VAE	3
2.2	RGCN	3
2.3	Embedding models	4
3	Background	4
3.1	The Graph VAE – one shot method	4
3.1.1	VAE	4
3.1.2	MLP	5
3.1.3	Graph convolutions	5
3.1.4	RGCN	5
3.1.5	Graph VAE	6
3.2	Graph Matching	6
3.2.1	Permutation Invariance	6
3.2.2	Graph matching algorithms	6
3.2.3	Graph Matching Loss	7
3.3	Knowledge Graphs	7
3.4	Ranger Optimizer	7
4	Methodology	7
4.1	Knowledge graph representation	7
4.2	Graph VAE	8
4.3	Loss function	8
5	Experiments	8
5.1	Datasets and Training	8
5.2	Link Prediction	9
5.3	Sanity Checks	9
5.4	Interpolate Latent Space	9

5.5	Subgraph Generation	9
5.6	Syntax coherence	9
6	Results	9
7	Discussion & Future Work	10
7.1	Discuss Aspect 1	10
7.2	Future Work	10
7.3	Use Normalizing Flows	10
7.4	Recurrent VAE	10
7.5	Graph matching	11

Abstract

We generate Knowledge! [1]

1 Introduction

Here comes a beautiful introduction. Promise!

1.1 Motivation

Computer vision reached a point, where semantics, entities and relations can be inferred in an simple image. Would it not be fantastic to be able to apply this to text too? Could we train a model to learn the semantics of a KG?

1.2 Expected Contribution

This thesis is aimed to be at fundamental research and provide insight, into if further research in this direction would be meaningful.

1.3 Research Question

How well can a VAE learn the underlying semantics of a KG?

2 Related Work

This section presents previous work which inspired and layed the fundamentals for this thesis.

2.1 Graph VAE

Present different papers with graph VAEs

- Belli recurrent VAE
- GraphVAE paper
- some more

The model architecture presented in the GraphVAE paper is the starting point of our model.

2.2 RGCN

Present the relational graph convolution model paper by Kipf and maybe others

2.3 Embedding models

Present RASCAL and one or two more.

Embeddings

TransE represents entities in in low-dimensional embedding. The relationships between entities are represented by the vector between two entities [2]. (How are different relation between the same entities represented?)

OntoUSP

This method learns a hierarchical structure to better represent the relations between entities in embedding space.

3 Background

In this section we will go over related work and relevant background information for our model and experiments. The depth of the explanation is adopted to the expected prior knowledge of the reader. The reader is supposed to know the basics of machine learning and deep learning, including probability theory and basic knowledge on neural networks and their different architectures. Basic principles such as forward pass, backpropagation and convolutions are expected to be understood. Further the use and functionality of deep learning modules such as the model, the optimizer and the terms target and prediction should be known. This also includes being familiar with the training and testing pipeline of a model in deep learning. Should all these boxes be checked, then we can expect to get a deeper understanding of the magic behind the VAE and its differences to a normal autoencoder. After that we will present how convolutional layers can be used on graphs. Of course where there is a layer there is a model, thus we are presenting the graph convolutional network (GCN). Closing the circle we show how we can adopt the VAE to graph convolutions. Wrapping things up, we present the state of the art algorithms for graph matching, which will be util to allow permutation invariance when matching prediction and target.

3.1 The Graph VAE – one shot method

3.1.1 VAE

The VAE as first presented by [3] is an unsupervised generative model in form of an autoencoder, consisting of an encoder and a decoder. We will The architecture of the VAE differs from a common autoencoder by having a stochastic module between encoder and decoder. The encoder can be represented as recognition model with the probability $p_{\Theta}(\mathbf{z} | x)$ with x being the variable we want to inference and z being the latent representation given an observed value of x . The encoder parameters are represented by Θ . Similarly, We denote the decoder as $p_{\Theta}(\mathbf{x} | z)$, which given a latent representation z produces a probability distribution for the possible values, corresponding to the input of x . This will be the base architecture of all our models in this thesis.

The main contribution of the VAE is the so called reparameterization trick. By sampling from the latent prior distribution, we get stochastic module inside our model, which can not be backpropagates through and makes machine learning not possible. By placing the stochastic module outside the model FIGURE!!!, we can again backpropagate. We use the predicted latent space as mean and variance for a Gaussian normal distribution, from which we then sample ϵ , which acts as external parameter and does not need to be updated.

This makes the true posterior $p_{\theta}(\mathbf{z} | \mathbf{x})$ intractable. Thus, we assume that the prior to the decoder to be Gaussian with an approximately diagonal covariance, which gives us the approximated posterior.

$$\log q_{\phi}(\mathbf{z} | \mathbf{x}^{(i)}) = \log \mathcal{N}(\mathbf{z}; \mu^{(i)}, \sigma^{2(i)} \mathbf{I}) \quad (1)$$

This gives us computational freedom. meaning we can compute the posterior probability. Using the

Monte Carlo estimation of $q_\phi(\mathbf{z} | \mathbf{x})$ we get the so called estimated lower bound (ELBO):

$$y_k(\mathbf{x}, \mathbf{w}) = \sigma \left(\sum_{j=0}^M w_{kj}^{(2)} h \left(\sum_{i=0}^D w_{ji}^{(1)} x_i \right) \right) \quad (2)$$

We denote the first term the regularization term, as it forces the model into using a Gaussian normal prior. The second term represents the reconstruction loss, matching the prediction with the target.

While we use a discrete input space, the output space is a continuous probability. To generate final result, the prediction is used as binomial probability distribution, from which we then sample. Once training of a VAE is completed, the Decoder can be used on its own to generate new samples by using latent input signals [3].

3.1.2 MLP

The Multi-Layer Perceptron (MLP) was the beginning of machine learning models. Its properties as universal approximator has been discovered and widely studied since 1989. The innovation it brought to existing models was the hidden layer between the input and the output.

The mathematical definition of the MLP is rather simple. It takes linear input vector of the form x_1, \dots, x_D which is multiplied by the weight matrix $\mathbf{w}^{(1)}$ and then transformed using a non-linear activation function $h(\cdot)$. Due to its simple derivative, mostly the rectified linear unit (ReLU) function is used. This results in the hidden layer, consisting of hidden units. The hidden units get multiplied with the second weight matrix, denoted $\mathbf{w}^{(2)}$ and finally transformed by a sigmoid function $\sigma(\cdot)$, which produces the output. Grouping all weight and bias parameter together we get the following equation for the MLP:

$$y_k(\mathbf{x}, \mathbf{w}) = \sigma \left(\sum_{j=1}^M w_{kj}^{(2)} h \left(\sum_{i=1}^D w_{ji}^{(1)} x_i + w_{j0}^{(1)} \right) + w_{k0}^{(2)} \right) \quad (3)$$

for $j = 1, \dots, M$ and $k = 1, \dots, K$, with M being the total number of hidden units and K of the output.

Since the sigmoid function gives us a probability output, the main function of the MLP is as a classifier. Instead of the initial sigmoid function, it was found to also produce good results for multi label classification transforming the output with a softmax function instead. Images or higher dimensional tensors can be processed by flattening them to a one dimensional tensor. This makes the MLP a flexible and easy to implement model [4].

3.1.3 Graph convolutions

CNNs have shown great results in the field of images classification and object detection. This is due to the fact that a convolution layer takes into account the relation of one pixel to its neighbors. The same holds for graph CNNs where at each convolution the information of each node is passed on as messages to all its neighbors. Each convolution applies an activation function in between the steps. In this case the information is the edges each node has. To process node attributes and edge attributes we have to look at more complex models [5].

3.1.4 RGCN

Relational Graph Convolution Net (RGCN) was presented in [6] for edge prediction. This model takes into account features of nodes. Both the adjacency and the feature matrix are matrix-multiplied with the weight matrix and then with them-selves. The resulting vector is a classification of the nodes.

3.1.5 Graph VAE

Now we have all the building blocks for a Graph VAE. The encode can either be a MLP, a GCNN or an RGCN. The same holds for the decoder with the addition that model architecture needs to be inverted. An version of a Graph VAE presented in [7]. This model combines both the previous methods. The input graph undergoes relational graph convolutions before it is flattened and projected into latent space. After applying the reparametrization trick, a simple MLP decoder is used to regenerate the graph. In addition the model concatenates the input with a target vector y , which represents ????. The same vector is concatenated with the latent tensor. ***Elaborate why they do that***.

Graphs can be generated recursively or in an one-shot approach. This paper uses the second approach and generates the full graph in one go. ***Cite?***

3.2 Graph Matching

Intro to graph matching on sparse graphs.

3.2.1 Permutation Invariance

Permutation invariance refers to the invariance of a permutation of an object. An visual example is the the image generation of numbers. If the loss function of the model would not be permutation invariant, the generated image could show a perfect replica of the input number but due to positional permutation the loss function would penalize the model. OR: An example is in object detection in images. An object can have geometrical permutations such as translation, scale or rotation, none the less the model should be able to detect and classify it. In that case, the model is not limited by permutations and is there fore permutation invariant. In our case the object is a graph and the nodes can take different positions in the adjacency matrix. To detect similarities between graphs we apply graph matching.

3.2.2 Graph matching algorithms

These are three of the state of the art graph matching algorithms.

- Wasserstein
- Maxpooling
- one more

There are various graph matching algorithms. The one we will implement is the max-pooling (Finding Matches in a Haystack: A Max-Pooling Strategy for Graph Matching in the Presence of Outliers).

The max-pooling graph matching algorithm returns a symmetric affinity matrix for all nodes

The resulting similarity matrix gives us X^* which is continuous and therefore useless. To transform is to a discrete X we use the hungarian algorithm (GPU-accelerated Hungarian algorithms for the Linear Assignment Problem)

Hungarian algorithm

The hungarian algorithm is used to find the shortest path within a matrix. This could be the most efficient work-distribution in a cost matrix. Or ...

It consists of four steps of which the last two are repeated until convergence. This algorithm is not scalable. The Munks algorithm (reference) tackles this problem by?? and is scalable.

second option: Compare only graph structure. NX algorithms: Greedy, Shortest path ...

3.2.3 Graph Matching Loss

The loss is discribed in [7] as.

In contrast to his implementation we assume, that a node or edge can have none, one or multiple attributes. Therefore our attributes are also not sigmoided and do not sum up to one. This leads o the modification of term $\log F$ and $\log E$ where we do not matrix multiply over the attribute vector but take the BCE as over the rest of the points. KG can have multiple or no attributes vs molecular graphs can be one hot encoded.

Okay, further we need to treat the $\log_p E$ and $\log_p F$ just like $\log_p A$ and subtract the inverse. Otherwise the model learn to predict very high values only.

A note to the node features, these stay softmaxed and one-hot encoded since we will use them as node labels.

3.3 Knowledge Graphs

Knowledge Graphs are great! The best in the world. Knowledge graphs have very different formats. The datasets we will be sorting with are in rdf format. This format can include an defined onthology or not. This means the KG consists of triples subject, relation, object. when indexing these triples. we get a dense representation of the KG. About sparse KGs

3.4 Ranger Optimizer

A optimizer, which manged to improve results upon state of the art models. It combines rectified adam, lookahead and gradient centralization. HOW ABOUT A CITATION

4 Methodology

This section describes the methodology for this thesis. The first part includes the presentation of the model, the reprocessing of the input and the evaluation metrics. The second part describes the experimental setup and the different experimental runs. The work of this thesis has aimed to be fully reproducible, thus the code is opensourced and available on Github ¹.

4.1 Knowledge graph representation

The first step in our pipeline is the representation of the KG in tensor format. In order to represent the graph structure we use an adjacency matrix A of shape $n \times n$ with n being the number of nodes in our graph. The edge attribute or directed relations between the nodes are represented in the matrix E of shape $n \times n \times d_E$ with d_E being the number of edge attributes. Similarly for node attributes we have the matrix F of shape $n \times d_N$ with d_N number of node attributes. The input graph can have less nodes than the maximum n but not more. The diagonal of the adjacency matrix is filled with 1 if the indexed node exists, and with 0 otherwise. The number and encoding of the attributes must be predefined and cannot be changed after training. This way we can uniquely represent a KG.

Graph embeddings? unsupervised approach

¹***Thesis Repo Link***

4.2 Graph VAE

hi
Convolution part
RCGN relation Convolution neural net
MLP encoder
Latent space
reparametrization trick
MLP decoder
Graph matching
Discretization of prediction

4.3 Loss function

If loss function should be permutation invariant we need to do some kind of graph matching.
Different options for graph matching.
Maxpooling-algorithm:
Assumptions
Node to edge affinity equals 0

Self-loops are possible, adjacency matrix can be zero or one.

Summing over the neighbors means summing over the whole column Normalize matrix with Frobenius Norm

Batch version: Only matmul and dot. keep dimension of S with shape (bs,n,n,k,k) When maxpooling, flatten Xs (n,k) for batch dot multiplication. This way (i think) we sum over all j nad b neighbors instead of taking the max.

Hungarian algorithm for discrimination of X

The hungarian algorithm as presented in section 2 return the shortest path in a matrix. We use this shortest path as bast match between the two graphs. The node paris identified as optimal are masked as 1 and the rest of the matrix as 0. This way we discretizice $X \star$ to X .

The equations

5 Experiments

This section presents the experiments we ran. We covered link and node prediction and compared those to SOTA scores. Further we ran experiments on investigating the coherence of the reproduced graph structure. Lastly we measured the adherence of our model to the KG's underlying syntax.

5.1 Datasets and Training

For this sake of meaningful results, we chose to use the earlier presented, two most popular dataset used in this field, FB15k237 and WN18rr.

Training models on each dataset for 333 epochs, without early stopping.

5.2 Link Prediction

We used link prediction as evaluation protocol and for comparison state of the art models. For each triple in the dataset, we remove the tail and combine it with all possible entities in our dataset. Even though this is called 'triple-corruption', also correct triples can be generated, which could appear in the trainset. These have to be filtered out before evaluation. Link prediction on unfiltered test data is termed 'raw'. Our model then computes the ELBO loss for all corrupted triples, which are sorted and stored in ascending order. The same procedure is repeated the triple's head in place of the tail.

The metrics used to evaluate the predictions, is the mean reciprocal rank (MRR) and hits at 1,3 and 10. The MRR ??? explain what the MRR is. Hits at 1 indicates what percentage of the triples in the test set have been ranked the highest compared to their corruptions. Similar, hits at 3 and 10, give a percentage of triples ranked in the top 3 and top 10. These metrics allow a fair comparison on a dataset between models, regardless of the ranking method of each model.

5.3 Sanity Checks

Check if adj matrix adheres to edge attribute matrix.

5.4 Interpolate Latent Space

We take two random triples and interpolate the latent space of these two triples. The interpolations result in: HERE AN EXAMPLE.

Further we go ahead and test what happens if we modify one latent dimension of a triple. HERE AN EXAMPLE. Can the model assign logical features to latent dimensions?

5.5 Subgraph Generation

Until now our model trained on only one triple per sparse graph. What will happen if we train it on more than one triples?

5.6 Syntax coherence

Check if generated triples follow basic logic.

- Generate triples by random signals
- Filter these triples on a certain relation
- Check if the entities are part of the linked class

Since there is no preset for how to check the semantics of a KG, we will use simple basic logical criteria. The generated triples are filtered for the relation 'is capital of', thus the subject entity should be a city and the object entity member of the class 'country', Hope this gives good results.

6 Results

Fabulous results only!

7 Discussion & Future Work

We will discuss certain aspects of our results and give advice on further research.

7.1 Discuss Aspect 1

Well well well

7.2 Future Work

Basing on the believe that the experiments were successful, we recommend:

- Improve the model (deeper)
- link the latent space to word signals e,g text
- prior not normal, e.g NF
- try a GAN

Good luck amigos!

Ideas

The bigger picture of this thesis is to efficiently generate a representation of the information hold in plain text. This representation has the form of a knowledge graph and consists of subject-relation-object triples.

Challenges:

- What knowledge are we looking for? I would like the model to focus on the most important information. This could for example be topic specific.
- Another option would be to extract based on a query or point of interest. Especially if we build the KG incrementally we can use this input as starting point and recursively build from there.

7.3 Use Normalizing Flows

- Is it possible to generate a graph from text using Normalizing flow (NF)?
- Can we train such a NF unsupervised, using the output distributions as loss function?
- NFs can build KG at once or incrementally. No one is using it, could have a reason.

7.4 Recurrent VAE

- Recursively generate graph using a Variational Auto Encoder.
- Use query as start node.
- expand graph on that node. Thiviyan and others work on VAEs as well. Seems an easier approach. Can analyze the latent space.

The encoder decoder strategy has been applied successfully in many cases. Bellis thesis about modeling a graph from images can be applied to text as well by changing the input to a text vector [8]. Note that this is a supervised method and would require a dataset of text labeled with resulting KGs. To the best of my knowledge this does not exist yet.

Same holds for the contrastive world model by Kipf. If we input text vectors instead of an image the model could recognize different objects in the text as it does with pixels [1]. Here the model is trained unsupervised using a loss over the energy function of the graph embedding space TransE [2]. An open question is if this would work for our approach.

A bit more abstract is the idea of the feedback recurrent VAE [9] where sound signals are encoded and decoded. This could also be adopted to text for instance with one sentence at a time, or a fixed number of tokens. Here the text vector would be induced as the latent input to the decoder. This would mean finding an translation of the models latent space to the word embedding. The dataset would consist of positive and negative examples of resulting graphs over timesteps.

While text can be easily vectorized by word embeddings like word2vec, the graph representation seems more tricky. A reasonable approach following Bellis example would be to output a coordinates vector for the nodes and an adjantency matrix for their relationships. Here one node at a time is outputted and the relation is conditioned on the number of previous nodes. Alternatively we could make use of the graph embeddings RASCAL or TransE. Lastly the question remains how to model the graph when the nodes need to be predefined. A subgraph of DBpedia could be a good starting point.

Open Issues

A section where note will be made on open issues. These issues are ment to be discussed with Peter or Thivyian and ultimatly solved. No open issues should remain at the end of November.

7.5 Graph matching

Hungarian algorithm: Should we assume the affinity matrix is a cost or a profit matrix. If we assume it is a cost matrix and $n \neq k$ the algorithm returns an error while padding. If we assume it is a profit matrix and convert it to a profit matrix it can handle any dimension. The shortest path varies with both approaches. Further I assume the shortest path is what we are looking for and mask these entries with 1 and the rest of the matrix with 0.

Diagonal of A is zeros. As soon as I do this, code crashes. Think it is because of divided zero, Martin adds $1e-7$ to the X norm

Do we backpropagate through the graph matching? then all torch functions, otherwise np.

Martin takes 300 iterations for hungarian. also not looping over pairs.

References

1. Kipf, T., van der Pol, E. & Welling, M. Contrastive Learning of Structured World Models. *arXiv:1911.12247 [cs, stat]*. arXiv: 1911.12247. <http://arxiv.org/abs/1911.12247> (2020) (Jan. 2020).
2. Bordes, A., Usunier, N., Garcia-Duran, A., Weston, J. & Yakhnenko, O. in *Advances in Neural Information Processing Systems 26* (eds Burges, C. J. C., Bottou, L., Welling, M., Ghahramani, Z. & Weinberger, K. Q.) 2787–2795 (Curran Associates, Inc., 2013). <http://papers.nips.cc/paper/5071-translating-embeddings-for-modeling-multi-relational-data.pdf> (2020).
3. Kingma, D. P. & Welling, M. Auto-Encoding Variational Bayes. *arXiv:1312.6114 [cs, stat]*. arXiv: 1312.6114. <http://arxiv.org/abs/1312.6114> (2020) (May 2014).
4. Bishop, C. M. *Pattern recognition and machine learning* (springer, 2006).
5. Tiao, L., Elinas, P., Nguyen, H. & Bonilla, E. V. Variational Graph Convolutional Networks. en, 11.
6. Kipf, T. N. & Welling, M. Semi-Supervised Classification with Graph Convolutional Networks. en. *arXiv:1609.02907 [cs, stat]*. arXiv: 1609.02907. <http://arxiv.org/abs/1609.02907> (2020) (Feb. 2017).
7. Simonovsky, M. & Komodakis, N. GraphVAE: Towards Generation of Small Graphs Using Variational Autoencoders. *arXiv:1802.03480 [cs]*. arXiv: 1802.03480. <http://arxiv.org/abs/1802.03480> (2020) (Feb. 2018).

8. Belli, D. & Kipf, T. Image-Conditioned Graph Generation for Road Network Extraction. *arXiv:1910.14388 [cs, stat]*. arXiv: 1910.14388. <http://arxiv.org/abs/1910.14388> (2020) (Oct. 2019).
9. Yang, Y., Sautière, G., Ryu, J. J. & Cohen, T. S. Feedback Recurrent AutoEncoder. *arXiv:1911.04018 [cs, eess, stat]*. arXiv: 1911.04018. <http://arxiv.org/abs/1911.04018> (2020) (Feb. 2020).