

Learning beyond datasets: Knowledge Graph Augmented Neural Networks for Natural language Processing

Annervaz K M
Indian Institute of Science &
Accenture Technology Labs
annervaz@iisc.ac.in

Somnath Basu Roy Chowdhury
IIT Kharagpur
brcsomnath@ee.iitkgp.ernet.in

Ambedkar Dukkupati
Indian Institute of Science
ambedkar@iisc.ac.in

Abstract—Machine Learning has been the quintessential solution for many AI problems, but learning is still heavily dependent on the specific training data. Some learning models can be incorporated with a prior knowledge in the Bayesian set up, but these learning models do not have the ability to access any organised world knowledge on demand. In this work, we propose to enhance learning models with world knowledge in the form of Knowledge Graph (KG) fact triples for Natural Language Processing (NLP) tasks. Our aim is to develop a deep learning model that can extract relevant prior support facts from knowledge graphs depending on the task using attention mechanism. We introduce a convolution-based model for learning representations of knowledge graph entity and relation clusters in order to reduce the attention space. We show that the proposed method is highly scalable to the amount of prior information that has to be processed and can be applied to any generic NLP task. Using this method we show significant improvement in performance for text classification with News20, DBPedia datasets and natural language inference with Stanford Natural Language Inference (SNLI) dataset. We also demonstrate that a deep learning model can be trained well with substantially less amount of labeled training data, when it has access to organised world knowledge in the form of knowledge graph.

Index Terms—Deep Learning, Knowledge Graph, Natural Language Processing, Attention Mechanism,

I. INTRODUCTION

Today, machine learning is centered around models that can be trained on specific labeled and unlabeled training data available for a particular task. Although learning paradigms like Transfer Learning Pan and Yang [2010] attempt to incorporate learnings from one task for another, these are limited in scalability and very specific to tasks being dealt with. On the other hand, humans have intrinsic ability to elicit required past knowledge from the world on demand and combine with our new learned concepts to solve problems.

The question that we address in this paper is the following: Is it possible to develop learning models that can be trained in such a way that they can infuse such a general body of knowledge for prediction apart from learning from a training dataset.

By world knowledge, we mean a structured general purpose knowledge that need not be domain specific. One very well

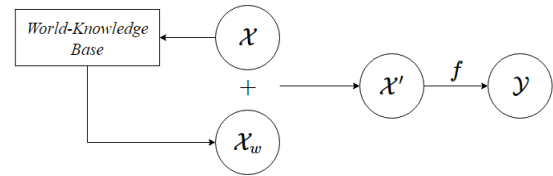


Fig. 1: The Basic Idea: \mathcal{X} is the feature input and \mathcal{Y} is the prediction. The relevant world knowledge for the task \mathcal{X}_w , is retrieved and augmented with the feature input before making the final prediction

known example for this is Knowledge Graphs¹. They represent information in the form of facts triplet, containing a subject entity, relation and object entity (example: $\langle \text{Italy, capital, Rome} \rangle$). The entities represent the nodes of the graph and their relation acts as edges. We represent this fact triple (subject entity, relation, object relation) by (h, r, t) . Practical knowledge bases congregate information from secondary databases or extract facts from unstructured text using various statistical learning mechanisms like NELL Mitchell et al. [2015] and DeepDive Niu et al. [2012]. There are human created knowledge bases as well, like Freebase (FB15k) Bollacker et al. [2008] and WordNet Miller et al. [1990]. The knowledge present in these knowledge bases includes common knowledge and partially covers common-sense knowledge and domain knowledge Song and Roth [2017]. Knowledge Graphs and Knowledge Bases are conceptually equivalent for our purpose and we will use the name interchangeably in this paper.

First, we illustrate the significance of world knowledge using a few examples. Consider the example of a Natural Language Inference (NLI) problem MacCartney [2009]. Consider the following two statements, A: The couple is walking on the sea shore and B: The man and woman are wide awake. Here, for a learning model to infer B from A, it should be able to access to the common

¹<https://googleblog.blogspot.in/2012/05/introducing-knowledge-graph-things-not.html>

knowledge that “*The man and woman and The couple means the same*” since this information may not be specific to a particular inference. Further, expecting a model to learn all such required information from just the labeled training data available for the task is next to impossible.

Consider the example of classifying the news snippet, Donald Trump offered his condolences towards the hurricane victims and their families in Texas. We cannot classify it as a political news unless we know the facts $\langle \text{Donald Trump, president, United States} \rangle$ and $\langle \text{Texas, state, United States} \rangle$. Our hypothesis is that machine learning models, apart from training them on the data with the ground-truth, if they can be trained to fetch relevant information from the structured knowledge bases the performance of the same can be improved.

In this work, we propose a deep learning model that can extract relevant support facts on demand from a knowledge base Mitchell et al. [2015] and incorporate it in the feature space along with features learned from the training data. (shown in Figure 1). This is very challenging, as knowledge bases are typically have millions of fact triples. In this paper we use propose a deep learning mechanism to jointly model this look up mechanism along with the task specific training of the model. The look-up mechanism and model is generic enough so that it can be augmented to any task specific learning model to boost the performance. In this paper, we have established superior performance of the proposed KG-augmented models over vanilla model on text classification and natural language inference.

Although there are plethora of works for knowledge graph learning Nickel et al. [2016a] Mitchell et al. [2015] Niu et al. [2012] from natural language text, no attempt to augment learning models with knowledge graph information have been done. To the best of our knowledge this is the first attempt to incorporate such world knowledge for the learning models.

II. KNOWLEDGE GRAPH REPRESENTATIONS

The Knowledge Graph entities/relations have to be encoded to a numerical representation for processing. Before describing the model, we provide a brief overview of graph encoding techniques. Various KG embedding techniques can be classified at a high level into: *Structure-based embeddings* and *Semantically-enriched embeddings*.

Structure-based embeddings: TransE Bordes et al. [2013] is the introductory work on knowledge graph representation, which translated subject entity to object entity using one-dimensional relation vector ($h + r = t$). Variants of the TransE Bordes et al. [2013] model uses translation of the entity vectors over relation specific subspaces. TransH Wang et al. [2014b] introduced the relation-specific hyperplane to translate the entities. Similar work utilizing only the structure of the graph include ManifoldE Xiao et al. [2015b], TransG Xiao et al. [2015a], TransD Ji et al. [2015], TransM Fan et al. [2014], HolE Nickel et al. [2016b] and ProjE Shi and Weninger [2017].

Semantically-enriched embeddings: These embedding techniques learn to represent entities/relations of the KG along with its semantic information. Neural Tensor Network (NTN) Socher et al. [2011] was the pioneer work in this field which initialized entity vectors with the average word embeddings followed by tensor-based operations. Recent works involving this idea are “Joint Alignment” Zhong et al. [2015] and SSP Xiao et al. [2017]. DKRL Xie et al. [2016] is a KG representation technique which also takes into descriptive nature of text keeping the simple structure of TransE model. Pre-trained word2vec Mikolov et al. [2013] are used to form the entity representation by passing through a Convolutional Neural Network (CNN) Kim [2014] architecture constraining the relationships to hold.

In our experiments we have used the DKRL Xie et al. [2016] encoding scheme as it emphasizes on the semantic description of the text. Moreover, DKRL fundamentally uses TransE Bordes et al. [2013] method for encoding structural information. Therefore, we can retrieve relevant entities & relation and obtain the complete the fact using $t = h + r$. This reduces the complexity of fact retrieval as the number of entities/relations is much less compared to the number of facts, thus making the retrieval process faster.

III. THE PROPOSED MODEL

Conventional supervised learning models with parameters Θ , given training data x and label y , tries to maximize the following function

$$\max_{\Theta} P(y|x, \Theta)$$

The optimized parameters Θ are given as,

$$\Theta = \operatorname{argmax}_{\Theta} \log P(y|x, \Theta)$$

In this work, we propose to augment the supervised learning process by incorporation of world knowledge features x_w . The world knowledge features are retrieved using the data x , using a separate model where, $x_w = F(x, \Theta^{(2)})$. Thus, our modified objective function can be expressed as

$$\max_{\Theta} P(y|x, x_w, \Theta^{(1)})$$

where, $\Theta = \{\Theta^{(1)}, \Theta^{(2)}\}$. The optimized parameters can be obtained using the equation

$$\Theta = \operatorname{argmax}_{\Theta} \log P(y|x, F(x, \Theta^{(2)}), \Theta^{(1)})$$

The subsequent sections focus on the formulation of the function F which is responsible for fact triple retrieval using the data sample x . Here it is important to point out that, we are not assuming any structural form for P based on F . So the method is generically applicable to augment any supervised learning setting with any form for P , only constraint being P should be such that the error gradient can be computed with respect to F . In the experiments we have used softmax using the LSTM Greff et al. [2015] encodings of the input as the form for P . As for F , we use soft attention Bahdanau et al.

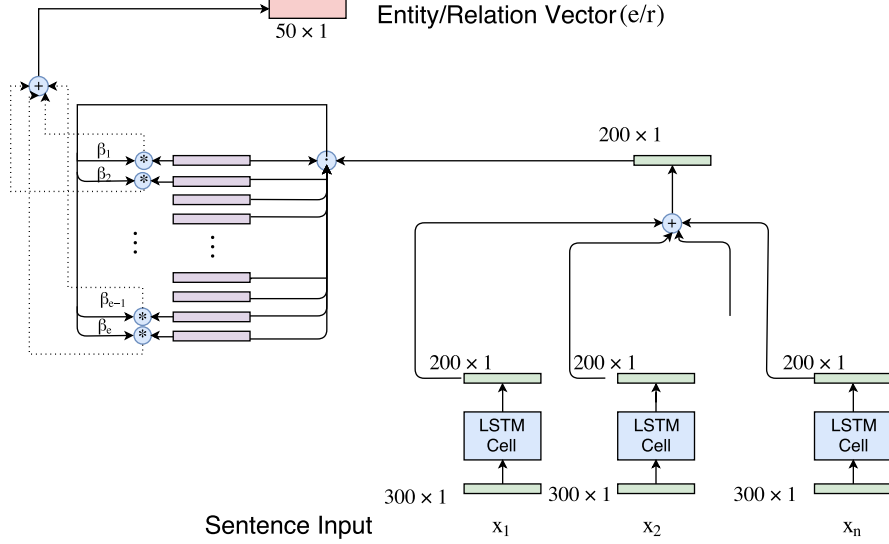


Fig. 2: Vanilla Entity/Relationship Retrieval Block Diagram

[2014]; Luong et al. [2015] using the LSTM encodings of the input and appropriate representations of the fact(s). Based on the representation used for the facts, we propose two models (a) Vanilla Model (b) Convolution-based entity/relation cluster representation, for fact retrieval in the subsequent sections.

A. Vanilla Model

The entities and relationships of the KG are encoded using DKRL, explained earlier. Let $e_i \in \mathbb{R}^m$ stand for the encoding of the entity i and $r_j \in \mathbb{R}^m$ stands for j^{th} relationship in the KG. The input text in the form of concatenated word vectors, $x = (x_1, x_2, \dots, x_T)$ is first encoded using an LSTM Greff et al. [2015] module as follows,

$$h_t = f(x_t, h_{t-1})$$

and

$$o = \frac{1}{T} \sum_{t=1}^T h_t,$$

$h_t \in \mathbb{R}^n$ is the hidden state of the LSTM at time t , f is a non-linear function and T is the sequence length. Then a context vector is formed from o as follows,

$$C = \text{ReLU}(o^T W),$$

where, $W \in \mathbb{R}^{n \times m}$ represent the weight parameters. The same procedure is duplicated with separate LSTMs to form two separate context vectors, one for entity retrieval (C_E) and one for relationship retrieval (C_R).

As the number of fact triples in a KG is in order of millions in the vanilla model we resort to generating attention over the entity and relation space separately. The fact is then formed

using the retrieved entity and relation. The attention for the entity, e_i using entity context vector is given by

$$\alpha_{e_i} = \frac{\exp(C_E^T e_i)}{\sum_{j=0}^{|E|} \exp(C_E^T e_j)}$$

where $|E|$ is the number of entities in the KG.

Similarly the attention for a relation vector r_i is computed as

$$\alpha_{r_i} = \frac{\exp(C_R^T r_i)}{\sum_{j=0}^{|R|} \exp(C_R^T r_j)}$$

where $|R|$ is the number of relations in the KG. The final entity & relation vector retrieval is computed by the weighted sum with the attention values of individual retrieved entity/relation vectors.

$$e = \sum_{i=0}^{|E|} \alpha_{e_i} e_i \quad r = \sum_{i=0}^{|R|} \alpha_{r_i} r_i$$

Figure 2 shows the schematic diagram for entity/relation retrieval. After the final entity and relation vectors are computed, we look forward to completion of the fact triple. The KG embedding technique used for the experiment is DKRL which inherently uses the TransE model assumption ($h + r \approx t$). Therefore, using the subject entity and relation we form the object entity as $t = e + r$. Thus the fact triplet retrieved is $\mathcal{F} = [e, r, e + r]$, where $\mathcal{F} \in \mathbb{R}^{3m}$. This retrieved fact information is concatenated along with the context vector (C) of input x obtained using LSTM module. The final classification label y is computed as follows,

$$\mathcal{F}' = \text{ReLU}(\mathcal{F}^T V)$$

$$y = \text{softmax}([\mathcal{F}' : C]^T U)$$

where, $V \in \mathbb{R}^{3m \times u}$ and $U \in \mathbb{R}^{2u \times u}$ are model parameters to be learned. \mathbf{y} is used to compute the cross entropy loss. We minimize this loss averaged across the training samples, to learn the various model parameters using stochastic gradient descent Bottou [2012]. The final prediction \mathbf{y} , now includes information from both dataset specific information and world knowledge to aid in enhanced performance. While jointly training the attention mechanism tunes itself to retrieve relevant facts that are required to do the final classification.

B. Pre-training KG Retrieval

The vanilla model attends over the entire entity/relation space which is not a good approach as we observed gradient for each attention value gets saturated easily. While training the classification and retrieval module together, the model tends to ignore the KG part and gradient propagates only through the classification module. This is expected to an extent as most pertinent information for the task at hand comes from the training samples, only background aiding information comes from KG. After few epochs of training, the KG retrieved fact always converged to a fixed vector. To overcome this problem, we attempted pre-training KG retrieval part separately. A pre-trained KG model is used to retrieve the facts and then concatenate with the classification module, while we allow error to be propagate through the pre-trained model, at the time of joint training. We infer that KG doesn't return noise and has essential information for the task as the separate KG part alone shows significant performance (59% for News20 & 66% for SNLI). Figure 3 depicts the entire training scheme. This procedure solved the issue of gradient saturation in the KG retrieval part at the time of joint training. The key problem that the attention mechanism has to cover a large space of entities/relation, remained.

C. Convolution-based Entity and Relation Cluster Representation

In this section, we propose a mechanism to reduce the large number of entities/relationships over which attention has to be generated in the knowledge graph. We propose to reduce the attention space by learning the representation of similar entity/relation vectors and attending over them.

In order to cluster similar entity/relation vectors, we used k -means clustering Bishop [2006] and formed l clusters with equal number of entity/relation vectors in each cluster. Each clusters were then encoded using convolutional filters. The output of the k -means clustering is a sequence of entity/relation vectors $\{e_1^T, e_2^T, \dots, e_q^T\}$, where $e_i \in \mathbb{R}^m$ and the number of elements in each cluster is given as $q = \lceil \frac{|E|}{l} \rceil$. For each cluster these vectors were stacked to form \mathcal{E} as the 2D input to the CNN encoder, where $\mathcal{E} \in \mathbb{R}^{m \times q}$. During experimentation for finding a suitable filter shape, it was observed that using 2-D filters the model failed to converge at all. Therefore, we inferred that the latent representation of two different indices in the vector e_i , should not be tampered using convolution. We then resorted to use 1-D convolution filters which slide only along the columns of \mathcal{E} , as shown Figure 4. The stride

length along y -axis is the window length k . The output of the convolution layer is expressed as,

$$\mathcal{E}'(i, j) = W^T [e_{i,j}, e_{i+1,j}, \dots, e_{i+k-1,j}]^T$$

where, $\mathcal{E}'(i, j)$ is the $(i, j)^{th}$ element of the output matrix \mathcal{E}' and $W \in \mathbb{R}^k$ is the convolution weight filter. A pooling layer followed the convolution layer in order to reduce the parameter space, we used 1-D window only along the y -axis similar to the convolutional kernel mentioned above. We used a two layered convolution network with the stride length k & max-pool windows n is adjusted to obtain output $E_i \in \mathbb{R}^m$, where i is the cluster index. Similar procedure of clustering followed by the encoding of the cluster entities is done for relations as well. Thus both the entity and relation space were reduced to contain fewer elements, one each for each cluster. After the compact entity space E and relation space R is formed, we followed the same steps as earlier for forming the attention, but now the training was more effective as the gradient was propagating effectively and was not choked by the large space. Also as convolution architecture is also simultaneously trained, attention mechanism was not burdened, as before to learn over the large space of entities and relationships.

One alternate point to mention here is regarding ranking/ordering items in the clusters, we have done experiments to verify the ordering does not affect the final result. We have verified this by randomly shuffling the entities/relationships in every clusters and the accuracy output remained within an error bound of $\pm 0.5\%$. In various permutations, the representations learned by the convolution operator for clusters varies, but it does not affect the overall results. Regarding the interpretation of what convolution operator learns, the operator is applied along each dimension of the entity/relationship vector, to learn a representation of the clusters. This representation includes information from relevant entities in the cluster, as the relevant entities varies across tasks, the representation learned using convolution also adapts accordingly. It is analogous to learning relevant features from an image, in our case the convolution layer learns the features focusing on relevant entities/relations in a cluster pertaining to the task.

IV. EXPERIMENTS AND EVALUATIONS

Our experiments were designed to analyze whether a deep learning model is being improved when it has access to KG facts from a relevant source. The selection of knowledge graph have to be pertinent to the task at hand, as currently there is no single knowledge base that contains multiple kinds of information and can cater to all tasks. We illustrate with results that the performance of a deep learning model improves when it has access to relevant facts. We also illustrate that as the model learns faster with access to knowledge bases, we can train deep learning models with substantially less training data, without compromising on the accuracy. In the subsequent section we briefly describe the datasets and associated Knowledge Bases used.

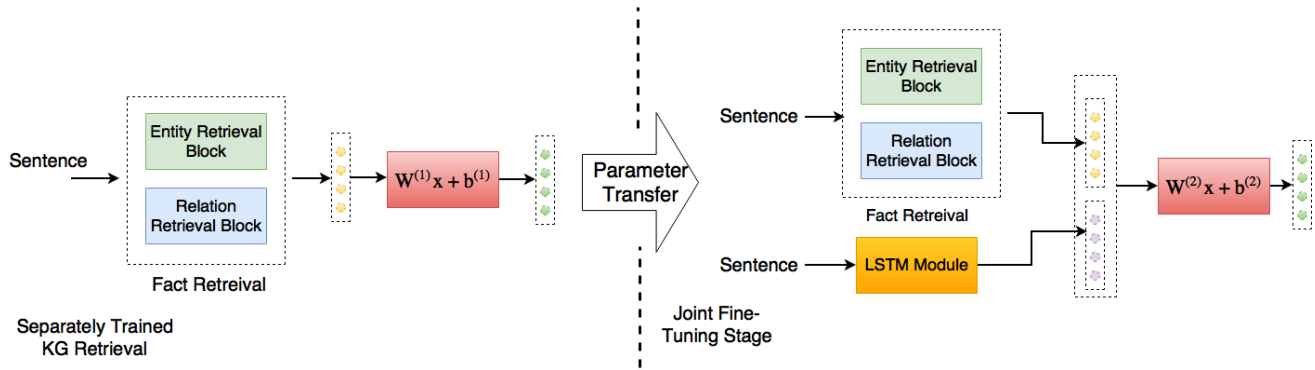


Fig. 3: Separately Training Knowledge Graph Retrieval and Jointly Training the Full Model

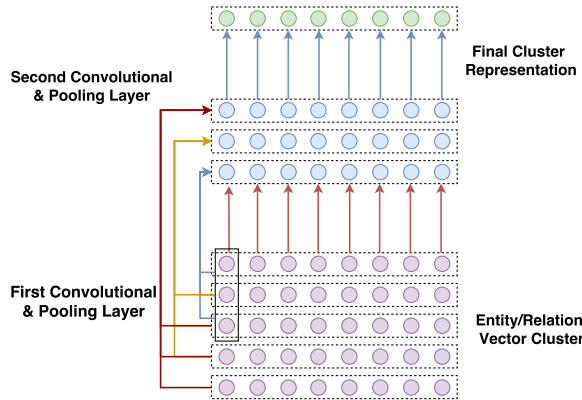


Fig. 4: Convolution model cluster representation

Datasets and Relevant Knowledge Graphs

In our experiments, we have used mainly the popular text classification dataset News20 Lichman [2013] and the Natural Language Inference dataset, Stanford Natural Language Inference (SNLI) corpus Bowman et al. [2015]. We have also done experiments on DBpedia ontology classification dataset², with a very strong baseline. These results are also discussed. These datasets are chosen as they share domain knowledge with two most popular knowledge bases, Freebase (FB15k) Bollacker et al. [2008] and WordNet (WN18) Bordes et al. [2013]. The training and test size of the datasets are mentioned in Table I.

Dataset	Train Size	Test Size	# Classes
News20	16000	2000	20
SNLI	549367	9824	3
DBpedia	553,000	70,000	14

TABLE I: Dataset Specifications

Freebase (FB15k) Bollacker et al. [2008] contains facts about people, places and things (contains 14904 entities, 1345

relations & 4.9M fact triples), which is useful for text classification in News20 Lichman [2013] dataset. On the other hand, WordNet (WN18) Bordes et al. [2013] (has 40943 entities, 18 relations & 1.5M fact triples) contains facts about common day-to-day things (example: furniture includes bed), which can help in inference tasks like SNLI. Both the knowledge bases are directed graphs, due to fewer number of relations WN18 is denser than FB15k. For experiments relating to both the datasets News20 & SNLI we used the standard LSTM as the classification module. As iterated earlier, our KG based fact retrieval is independent of the base model used. We show improvement in performance using the proposed models by KG fact retrieval. We use classification accuracy of test set as our evaluation metric.

A. Experimental Setup

All experiments were carried on a Dell Precision Tower 7910 server with Quadro M5000 GPU with 8 GB of memory. The models were trained using the Adam’s Optimizer Kingma and Ba [2014] in a stochastic gradient descent Bottou [2012] fashion. The models were implemented using TensorFlow Abadi et al. [2015]. The relevant hyper-parameters are listed in Table II. The word embeddings for the experiments were obtained using the pre-trained GloVe Pennington et al. [2014]³ vectors. For words missing in the pre-trained vectors, the locally trained GloVe vectors of the corresponding dataset was used.

Hyper-parameter	News20	SNLI
Batch size	256	1024
Learning rate	0.05	0.05
Word Vector Dimension	300	300
Sequence Length	300	85
LSTM hidden-state Dimension	200	200
KG Embedding Dimension	50	50
# Clusters	20	20
# Epochs	20	20

TABLE II: Hyper-parameters which were used in experiments for News20 & SNLI datasets

²<http://wiki.dbpedia.org/services-resources/dbpedia-data-set-2014>

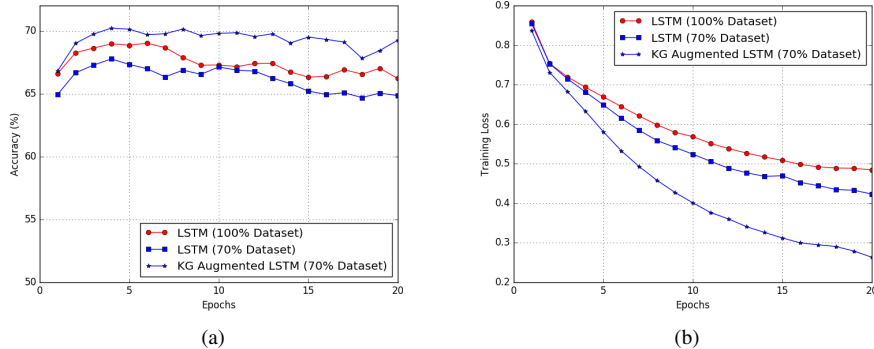


Fig. 5: (a) Accuracy Plot over training epochs for LSTM (using full & 70% dataset) and KG augmented LSTM (using 70% dataset) for SNLI task (b) Corresponding Training Loss plots for the aforementioned methods using SNLI dataset

B. Results & Discussion

Table III shows the results of test accuracy of the various methods proposed on the datasets News20 & SNLI. We observe that incorporation of KG facts using the basic vanilla model improves the performance slightly, as the retrieval model was not getting trained effectively. The convolution-based model shows significant improvement over the normal LSTM classification. While tuning the parameters of the convolution for clustered entities/relations it was observed that smaller stride length and longer max-pool window improved performance. For News20 dataset we show an improvement of almost 3% and for SNLI an improvement of almost 4%.

The work is motivated more from the perspective of whether incorporation of world knowledge will improve any deep learning model than from beating the state of the art performance. Although LSTM is used to encode the input for the model as well as the retrieval vector, as mentioned earlier, these two modules need not be same. For encoding the input any complex state-of-the-art model can be used. LSTM can be used to just generate the retrieval vector. On DBpedia ontology classification dataset, with a strong baseline of 98.6%, after augmenting it with KG (Freebase) using convolution based model we saw an improvement of 0.2%. Of course as the baseline is stronger, the improvement quantum has decreased. This is quite intuitive as the models are more complex and sufficient to learn from the data itself and the room available for further improvement is relatively less. Although the improvement is significant in weaker learning models, it is also capable of improving stronger baselines as is evident from the results on DBpedia dataset.

C. Reducing Dataset Size Requirements for Training Deep Learning Models

We hypothesized that as Knowledge Graph is feeding more information to the model, we can achieve better performance with less training data. To prove this we designed an experiment to compare the accuracy of the baseline model trained on

Model	Accuracy	
	News20	SNLI
Plain LSTM	66.75%	68.73%
Vanilla KG Retrieval	67.30%	69.20%
Convolution-based KG	69.34%	73.10%

TABLE III: Test accuracy of approaches in News20 using FB15K & SNLI datasets using WN18

full training data and compared it with the accuracy of the KG augmented model trained with just 70% of the training data for the SNLI dataset. The accuracy and training loss plots across training epochs is given in Figure 5. Even with just 70% of the data, KG augmented model is able to achieve better accuracy compared to the vanilla LSTM model trained on the full data. Figure 5(b) depicts the training loss for the 3 methods, the LSTM (with 70% data) shows less training loss than LSTM (with 100% data), wherever KG LSTM fits the training data significantly more. This clearly indicates that relevant information pertaining to the task is retrieved from the knowledge graph and the training loss reduction is not due to lesser data only. Also note that training loss is substantially less for KG LSTM compared to normal LSTM when the dataset size is reduced. This result is very promising, to reduce the humongous labeled training data requirement of large deep learning models, which is hard to come by.

V. RELEVANT PREVIOUS WORK

The basic idea of infusing general world knowledge for learning tasks, especially for natural language processing, is not attempted before. For multi-label image classification, the use of KGs has been pursued recently by Marino et al. [2016]. In this work, they first obtain labels of the input data (using a different model), use these labels to populate features from the KG and in turn use these features back for the final classification. For NLP tasks the information needed may not necessarily depend on the final class, and we are directly using all the information available in the input for populating the relevant information from the knowledge graphs. Our attempt

³<http://nlp.stanford.edu/data/glove.840B.300d.zip>

is very different from Transfer Learning Pan and Yang [2010]. In Transfer Learning the focus is on training the model for one task and tuning the trained model to use it for another task. This is heavily dependent on the alignment between source task and destination task and transferred information is in the model. In our case, general world knowledge is being infused into the learning model for any given task. In the similar lines our work is different from Domain Adaptation Glorot et al. [2011] as well. There has been attempts to use world knowledge Song and Roth [2017] for creating more labeled training data and providing distant supervision etc. Incorporating Inductive Biases Ridgeway [2016] based on the known information about a domain onto the structure of the learned models, is an active area of research. However motivation and approach is obviously different from these works.

VI. CONCLUSION & FUTURE WORK

In this work we illustrated the need for incorporating world knowledge in training task specific models. We showed procedure of pre-training the retrieval model to overcome the training stagnation at the time of joint training. We also showed a convolution based architecture to reduce the attention space over entities and relationships that performed better. With significant improvements over the vanilla baselines for two well known datasets, we illustrated the effectiveness of the proposed methods in improving the performance of deep learning models. We showcased that the proposed method can be used to reduce the labeled training data requirements of the deep learning models. Although the work is presented mainly from natural language processing perspective and using LSTM as the vanilla model, the idea and the model is applicable for other domain task, with more complicated deep learning models as base. To the best of our knowledge this is the first attempt in infusing general world knowledge for task specific training of deep learning models.

Being the first work of its kind, that are lot of scope for future work. A more sophisticated model which is able to retrieve facts directly from millions of entries. Currently we have focused only on a flat attention structure, a hierarchical attention mechanism would be more suitable. The model uses soft attention to enable training by simple stochastic gradient descent. Hard attention over facts using reinforcement learning can be pursued further. This will further help in selection of multi-facts, that are not of similar type but relevant to the task. The convolution based model, helped to reduce the space over entities and relationships over which attention had to be generated. However more sophisticated techniques using similarity based search Mu and Liu [2017]; Wang et al. [2014a] can be pursued towards this purpose. The results from the initial experiments illustrates the effectiveness of the proposed models, advocating further investigations in these directions.

REFERENCES

- Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, and Craig Citro et al. TensorFlow: Large-scale machine learning on heterogeneous systems, 2015. URL <http://tensorflow.org/>. Software available from tensorflow.org.
- Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*, 2014.
- Christopher M Bishop. *Pattern recognition and machine learning*. springer, 2006.
- Kurt Bollacker, Colin Evans, Praveen Paritosh, Tim Sturge, and Jamie Taylor. Freebase: a collaboratively created graph database for structuring human knowledge. In *Proceedings of the 2008 ACM SIGMOD international conference on Management of data*, pages 1247–1250. AcM, 2008.
- Antoine Bordes, Nicolas Usunier, Alberto Garcia-Duran, Jason Weston, and Oksana Yakhnenko. Translating embeddings for modeling multi-relational data. In *Advances in neural information processing systems*, pages 2787–2795, 2013.
- Léon Bottou. *Stochastic Gradient Tricks*, volume 7700, page 430445. Springer, January 2012. URL <https://www.microsoft.com/en-us/research/publication/stochastic-gradient-tricks/>.
- Samuel R. Bowman, Gabor Angeli, Christopher Potts, and Christopher D. Manning. A large annotated corpus for learning natural language inference. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Association for Computational Linguistics, 2015.
- Miao Fan, Qiang Zhou, Emily Chang, and Thomas Fang Zheng. Transition-based knowledge graph embedding with relational mapping properties. In *PACLIC*, pages 328–337, 2014.
- Xavier Glorot, Antoine Bordes, and Yoshua Bengio. Domain adaptation for large-scale sentiment classification: A deep learning approach. In *Proceedings of the 28th international conference on machine learning (ICML-11)*, pages 513–520, 2011.
- Klaus Greff, Rupesh Kumar Srivastava, Jan Koutník, Bas R Steunebrink, and Jürgen Schmidhuber. Lstm: A search space odyssey. *arXiv preprint arXiv:1503.04069*, 2015.
- Guoliang Ji, Shizhu He, Liheng Xu, Kang Liu, and Jun Zhao. Knowledge graph embedding via dynamic mapping matrix. In *ACL (1)*, pages 687–696, 2015.
- Yoon Kim. Convolutional neural networks for sentence classification. *arXiv preprint arXiv:1408.5882*, 2014.
- Diederik Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- M. Lichman. UCI machine learning repository, 2013. URL <http://archive.ics.uci.edu/ml>.
- Minh-Thang Luong, Hieu Pham, and Christopher D Manning. Effective approaches to attention-based neural machine translation. *arXiv preprint arXiv:1508.04025*, 2015.

- Bill MacCartney. *Natural language inference*. PhD thesis, Citeseer, 2009.
- Kenneth Marino, Ruslan Salakhutdinov, and Abhinav Gupta. The more you know: Using knowledge graphs for image classification. *arXiv preprint arXiv:1612.04844*, 2016.
- Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*, 2013.
- George A Miller, Richard Beckwith, Christiane Fellbaum, Derek Gross, and Katherine J Miller. Introduction to wordnet: An on-line lexical database. *International journal of lexicography*, 3(4):235–244, 1990.
- Tom M Mitchell, William W Cohen, Estevam R Hruschka Jr, Partha Pratim Talukdar, Justin Betteridge, Andrew Carlson, Bhavana Dalvi Mishra, Matthew Gardner, Bryan Kisiel, Jayant Krishnamurthy, et al. Never ending learning. In *AAAI*, pages 2302–2310, 2015.
- Yadong Mu and Zhu Liu. Deep hashing: A joint approach for image signature learning. In *AAAI*, pages 2380–2386, 2017.
- Maximilian Nickel, Kevin Murphy, Volker Tresp, and Evgeniy Gabrilovich. A review of relational machine learning for knowledge graphs. *Proceedings of the IEEE*, 104(1):11–33, 2016a.
- Maximilian Nickel, Lorenzo Rosasco, Tomaso A Poggio, et al. Holographic embeddings of knowledge graphs. In *AAAI*, pages 1955–1961, 2016b.
- Feng Niu, Ce Zhang, Christopher Ré, and Jude W Shavlik. Deepdive: Web-scale knowledge-base construction using statistical learning and inference. *VLDS*, 12:25–28, 2012.
- Sinno Jialin Pan and Qiang Yang. A survey on transfer learning. *IEEE Transactions on knowledge and data engineering*, 22(10):1345–1359, 2010.
- Jeffrey Pennington, Richard Socher, and Christopher D Manning. Glove: Global vectors for word representation. In *EMNLP*, volume 14, pages 1532–1543, 2014.
- Karl Ridgeway. A survey of inductive biases for factorial representation-learning. *arXiv preprint arXiv:1612.05299*, 2016.
- Baoxu Shi and Tim Weninger. Proje: Embedding projection for knowledge graph completion. In *AAAI*, pages 1236–1242, 2017.
- Richard Socher, Cliff C Lin, Chris Manning, and Andrew Y Ng. Parsing natural scenes and natural language with recursive neural networks. In *Proceedings of the 28th international conference on machine learning (ICML-11)*, pages 129–136, 2011.
- Yangqiu Song and Dan Roth. Machine learning with world knowledge: The position and survey. *arXiv preprint arXiv:1705.02908*, 2017.
- Jingdong Wang, Heng Tao Shen, Jingkuan Song, and Jianqiu Ji. Hashing for similarity search: A survey. *arXiv preprint arXiv:1408.2927*, 2014a.
- Zhen Wang, Jianwen Zhang, Jianlin Feng, and Zheng Chen. Knowledge graph embedding by translating on hyperplanes. In *AAAI*, pages 1112–1119, 2014b.
- Han Xiao, Minlie Huang, Yu Hao, and Xiaoyan Zhu. Transg: A generative mixture model for knowledge graph embedding. *arXiv preprint arXiv:1509.05488*, 2015a.
- Han Xiao, Minlie Huang, and Xiaoyan Zhu. From one point to a manifold: Knowledge graph embedding for precise link prediction. *arXiv preprint arXiv:1512.04792*, 2015b.
- Han Xiao, Minlie Huang, Lian Meng, and Xiaoyan Zhu. Ssp: Semantic space projection for knowledge graph embedding with text descriptions. In *AAAI*, pages 3104–3110, 2017.
- Ruobing Xie, Zhiyuan Liu, Jia Jia, Huanbo Luan, and Maosong Sun. Representation learning of knowledge graphs with entity descriptions. In *AAAI*, pages 2659–2665, 2016.
- Huaping Zhong, Jianwen Zhang, Zhen Wang, Hai Wan, and Zheng Chen. Aligning knowledge and text embeddings by entity descriptions. In *EMNLP*, pages 267–272, 2015.