

GLoMo: Unsupervisedly Learned Relational Graphs as Transferable Representations

Zhilin Yang, Jake Zhao, Bhuwan Dhingra, Kaiming He, William W. Cohen, Ruslan Salakhutdinov, Yann LeCun

Presenter: Arshdeep Sekhon

<https://qdata.github.io/deep2Read>

Motivation

- CNNs and RNNs are tailored to grids or sequences
- To compensate for not being able to model graphical representations: use deep model expressiveness
- transfer learning standard approach: train on large dataset and then apply to extract features

Real Data has graph structures

- Real data has graph structures
 - parse trees, knowledge graphs, coreference resolution
 - same across all language tasks
- (-) expensive to acquire: human curated
- automatically-induced structures are mostly limited to one task (example NRI, Attention is all you need, Non-local neural networks)

GLOMO Objectives

- instead of unary feature learning : learn graphical features
- to learn latent relational graphs : input units are nodes (example words in a sentence)

- Graphs from Low-level unit MOdeling
- train a *neural* network that learns latent graphs from large scale unsupervised data
- transfer this *neural* network to downstream task: outputs new graphs

Why?

separate the features that represent the semantic meaning of each unit and the graphs that reflect how the units may interact. Ideally, the graphs capture task-independent structures underlying the data, and thus become applicable to different sets of features.

1 Method

- Model

- I. Unsupervised Relational Graph Learning
- II. Feature Predictor

2 Experiments

Given a sample/sequence:

$$(x_1, x_2, \dots, x_T) \quad (1)$$

- Learn G a $T \times T$, matrix where $G_{i,j}$ represents affinity between x_i and x_j

- Graph Predictor: $G = g(x)$, $G \in R^{L \times T \times T}$
- L the number of layers that produce graphs
- Feature Predictor $f(G, x)$
- During test time, in a downstream task, extract G using Graph predictor

1 Method

- Model
- I. Unsupervised Relational Graph Learning
- II. Feature Predictor

2 Experiments

Unsupervised Learning: Graph Predictor

- Made of two CNNs: key CNN and a query CNN
- input to both CNNs is x
- output of key CNN : (k_1, \dots, k_T)
- output of query CNN : q_1, \dots, q_T
- Compute Graph using:

$$G_{ij}^l = \frac{\text{ReLU}(k_i^{lT} q_j^l + b)^2}{\sum_i^l \text{ReLU}(k_i^{lT} q_j^l + b)^2} \quad (2)$$

- same as attention weights other than square and ReLU
- *we use ReLUs to enforce sparsity and the square operations to stabilize training. Moreover, we employ convolutional networks to let the graphs G be aware of the local order of the input and context, up to the size of each unit's receptive field.*

1 Method

- Model
- I. Unsupervised Relational Graph Learning
- **II. Feature Predictor**

2 Experiments

Feature Predictor

- $\mathbf{f}_t^l = v((\sum_j \mathbf{G}_{jt}^l \mathbf{f}_{jt}^{l-1}), \mathbf{f}_t^{l-1})$
- at $t = 0$, $\mathbf{f}_t^0 = \mathbf{x}_t$ input embeddings
- v can be MLP or GRU Cell

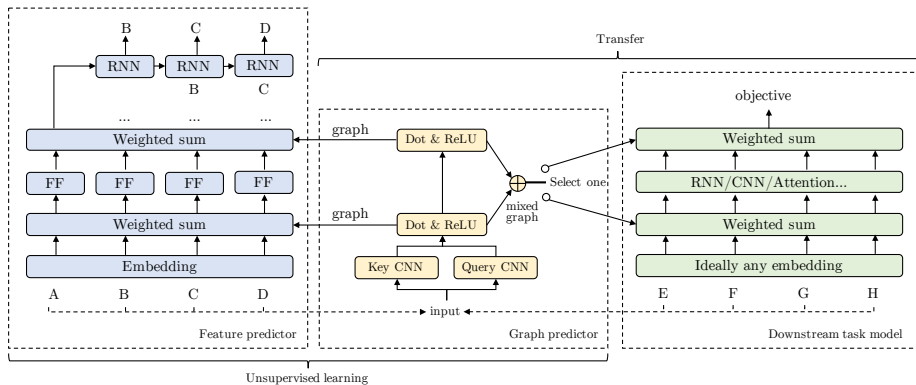
Feature Predictor

- Use \mathbf{f}_t^L to initialize hidden state
- predict units following x_t
- Objective function:

$$\max_t \sum_t \log P(x_{t+1}, \dots, x_{t+D} | x_t, \mathbf{f}_t^L) \quad (3)$$

- mask the convolutional filters and the graph G (see Eq. 1) in the network g to prevent the network from accessing the future

GLOMO



Differences with self-attention and predictive unsupervised learning

- Decoupling graphs and features: separate networks g and f
- Sparsity: a squared ReLU activation, most NLP graphs are sparse
We believe sparse structures reduce noise and are more transferable.
- Hierarchical graph representations: multiple layers of graphs, which allows us to model hierarchical structures in the data.

Latent Graph Transfer

- Say predictor takes in features $\mathbf{H} = (\mathbf{h}_1, \dots, \mathbf{h}_T)$
- produces $(\mathbf{h}'_1, \dots, \mathbf{h}'_T)$
- Get graphs from Graph predictor from all l layers \mathbf{G}
- $\Lambda^l = \prod_{i=1}^l \mathbf{G}^i$

-

$$\mathbf{M} = \sum_{l=1}^L m_g^l \mathbf{G}^l + \sum_{l=1}^L m_\Lambda^l \Lambda^l \quad (4)$$

- $\sum_l (m_\Lambda^l + m_g^l) = 1$
- use \mathbf{HM} as input

- Question Answering: SQuAD dataset
- Natural Language Inference: Multi-Genre NLI corpus: sentence pairs annotated with textual entailment information
- Sentiment Analysis: Movie Review from IMDB
- Transfer Setting: Wikipedia (700 million)

Results on Natural Language Tasks

Transfer method	SQuAD GloVe		SQuAD ELMo		IMDB GloVe	MNLI GloVe	
	<i>EM</i>	<i>F1</i>	<i>EM</i>	<i>F1</i>	<i>Accuracy</i>	<i>matched</i>	<i>mism.</i>
transfer feature only (baseline)	69.33	78.73	74.75	82.95	88.51	77.14	77.40
GLoMo on embeddings	70.84	79.90	76.00	84.13	89.16	78.32	78.00
GLoMo on RNN states	70.95	79.95	75.59	83.62	-	-	-

Ablation Study

Table 2: Ablation study.

Method	SQuAD GloVe		SQuAD ELMo		IMDB GloVe	MNLI GloVe	
	<i>EM</i>	<i>F1</i>	<i>EM</i>	<i>F1</i>	<i>Accuracy</i>	<i>matched</i>	<i>mism.</i>
GLoMo	70.84	79.90	76.00	84.13	89.16	78.32	78.00
- decouple	70.45	79.56	75.89	83.79	-	-	-
- sparse	70.13	79.34	75.61	83.89	88.96	78.07	77.75
- hierarchical	69.92	79.23	75.70	83.72	88.71	77.87	77.85
- unit-level	69.23	78.66	74.84	83.37	88.49	77.58	78.05
- sequence	69.92	79.29	75.50	83.70	88.96	78.11	77.76
uniform graph	69.48	78.82	75.14	83.28	88.57	77.26	77.50

Image Classification

Method / Base-model	ResNet-18	ResNet-34
baseline	90.93 ± 0.33	91.42 ± 0.17
GLoMo	91.55 ± 0.23	91.70 ± 0.09
ablation: uniform graph	91.07 ± 0.24	-

Conclusions

- transferable graph structures
- transfer graph learner module
- learn graph from large datasets and transfer to other tasks.