Graph Neural Networks for NLP

Zhaofeng Wu

Outline

- Graphs in NLP and ML
- Graph neural networks (GNNs)
- Applying GNNs to linguistic graphs: Infusing Finetuning with Semantic Dependencies

Graphs in NLP and ML



Running example: WSJ #20209013



"A similar technique is almost impossible to apply to other crops, such as cotton, soybeans and rice."



Running example: WSJ #20209013 "A similar technique is almost impossible to apply to other crops, such as cotton, soybeans and rice."



Phrase structure tree; parsed by http://corenlp.run



Running example: WSJ #20209013 "A similar technique is almost impossible to apply to other crops, such as cotton, soybeans and rice."



Stanford Dependencies (UD) tree; from Ivanova (2012)



Running example: WSJ #20209013 "A similar technique is almost impossible to apply to other crops, such as cotton, soybeans and rice."



DELPH-IN MRS-derived dependencies (DM) graph; from Oepen et al. (2015)



7

Running example: WSJ #20209013 "A similar technique is almost impossible to apply to other crops, such as cotton, soybeans and rice."



Discourse Representation Graph (DRG) graph; from Oepen et al. (2020)





Coreference graph; from https://demo.allennlp.org/coreference-resolution









Citation Graph



https:// www.connectedpapers.c om/main/ 204e3073870fae3d05bc bc2f6a8e263d9b72e776/ Attention-is-All-you-Need/graph

PageRank



"The founders of Google computed the Perron-Frobenius eigenvector of the web graph and became billionaires." — Preface to Spectra of Graphs by Brouwer and Haemers



Graph Neural Networks





• Undirected G = (V, E) with adjacency matrix $A \in \mathbb{R}^{n \times n}$; n = |V|, m = |E|



- Input
 - Node feature $X \in \mathbb{R}^{n \times d_v}$
 - (Sometimes) Edge feature $X^e \in \mathbb{R}^{m \times d_e}$

• Undirected G = (V, E) with adjacency matrix $A \in \mathbb{R}^{n \times n}$; n = |V|, m = |E|



- Input
 - Node feature $X \in \mathbb{R}^{n \times d_v}$
 - (Sometimes) Edge feature $X^e \in \mathbb{R}^{m \times d_e}$
- Output: node representation $Z \in \mathbb{R}^{n \times d_o}$

• Undirected G = (V, E) with adjacency matrix $A \in \mathbb{R}^{n \times n}$; n = |V|, m = |E|



- Input
 - Node feature $X \in \mathbb{R}^{n \times d_v}$
 - (Sometimes) Edge feature $X^e \in \mathbb{R}^{m \times d_e}$
- Output: node representation $Z \in \mathbb{R}^{n \times d_o}$
- How do we perform this update while leveraging the graph structure?

• Undirected G = (V, E) with adjacency matrix $A \in \mathbb{R}^{n \times n}$; n = |V|, m = |E|









 $H^0 = X$







 $H^0 = X$ $\mathbf{m}_{i}^{l+1} = \sum M_{l}(\mathbf{h}_{i}^{l}, \mathbf{h}_{j}^{l}, \mathbf{x}_{ij}^{e;l})$ $j \in \mathcal{N}(i)$







 $H^0 = X$ $\mathbf{m}_{i}^{l+1} = \sum M_{l}(\mathbf{h}_{i}^{l}, \mathbf{h}_{j}^{l}, \mathbf{x}_{ij}^{e;l})$ $j \in \mathcal{N}(i)$ $\mathbf{h}_{i}^{l+1} = U_{t}(\mathbf{h}_{i}^{l}, \mathbf{m}_{i}^{l+1})$







 $H^0 = X$ $\mathbf{m}_{i}^{l+1} = \sum M_{l}(\mathbf{h}_{i}^{l}, \mathbf{h}_{j}^{l}, \mathbf{x}_{ij}^{e;l})$ $j \in \mathcal{N}(i)$ $\mathbf{h}_{i}^{l+1} = U_{t}(\mathbf{h}_{i}^{l}, \mathbf{m}_{i}^{l+1})$ $Z = H^L$







Special Case: CNN as Message Passing

W 1	W 2	W3
W 4	W 5	W 6
W 7	W8	W9
L	L	a
		d
		f





Special Case: CNN as Message Passing

$\mathbf{h}_{i}^{l+1} = w_{1}^{l}\mathbf{h}_{a}^{l} + w_{2}^{l}\mathbf{h}_{b}^{l} + \dots + w_{5}^{l}\mathbf{h}_{i}^{l} + \dots + w_{9}^{l}\mathbf{h}_{b}^{l}$

W 1	W2	,	W3	
W 4	W 5	,	W6	
W 7	W8	,	W9	
		1	a	
			d	
			f	







Attempt #1:
$$\mathbf{h}_{i}^{l+1} = \sigma \left(\sum_{j \in \mathcal{N}(i)} W^{l} \mathbf{h}_{j}^{l} \right)$$



Attempt #1:
$$\mathbf{h}_{i}^{l+1} = \sigma \left(\sum_{j \in \mathcal{N}(i)} W^{l} \mathbf{h}_{j}^{l} \right)$$

Attempt #2: $\mathbf{h}_{i}^{l+1} = \sigma \left(\sum_{j \in \mathcal{N}(i) \cup \{i\}} W^{l} \mathbf{h}_{j}^{l} \right)$



Attempt #1:
$$\mathbf{h}_{i}^{l+1} = \sigma \left(\sum_{j \in \mathcal{N}(i)} W^{l} \mathbf{h}_{j}^{l} \right)$$

Attempt #2: $\mathbf{h}_{i}^{l+1} = \sigma \left(\sum_{j \in \mathcal{N}(i) \cup \{i\}} W^{l} \mathbf{h}_{j}^{l} \right)$
Attempt #3: $\mathbf{h}_{i}^{l+1} = \sigma \left(\sum_{j \in \mathcal{N}(i) \cup \{i\}} \frac{1}{\sqrt{|\mathcal{N}(i)|^{l}}} \right)$





- Localized filters/kernels
- Translation invariance
- Hierarchical/multi-scale





• Message passing (spatial): $\mathbf{h}_{i}^{l+1} = \sigma \left(\sum_{j \in \mathcal{N}(i) \cup \{i\}} \frac{1}{\sqrt{|\mathcal{N}(i)| + 1} \sqrt{|\mathcal{N}(j)| + 1}} W^{l} \mathbf{h}_{j}^{l} \right)$



- Message passing (spatial): $\mathbf{h}_{i}^{l+1} = \sigma \left(\int_{j \in \mathcal{I}_{i}} \mathbf{h}_{i}^{l+1} \right)$
- Matrix version: $H^{l+1} = \sigma \left(\tilde{D}^{-\frac{1}{2}} \tilde{A} \tilde{D}^{-\frac{1}{2}} H^l W \right)$

$$\sum_{i \in \mathcal{N}(i) \cup \{i\}} \frac{1}{\sqrt{|\mathcal{N}(i)| + 1}\sqrt{|\mathcal{N}(j)| + 1}} W^{l} \mathbf{h}_{j}^{l}$$

$$W^{l} \text{ where } \tilde{A} = A + I_{n}, \tilde{D} = \text{diag}\left(\sum_{i} \tilde{A}_{ii}\right)$$



- Message passing (spatial): $\mathbf{h}_i^{l+1} = \sigma \Big(\int_{j \in \mathcal{I}_i} \mathbf{h}_j^{l+1} \Big) = \sigma \Big(\int_{j \in \mathcal{I}_i} \mathbf{h}_j^{l+1} \Big) = \sigma \Big(\int_{j \in \mathcal{I}_i} \mathbf{h}_j^{l+1} \Big)$
- Matrix version: $H^{l+1} = \sigma \left(\tilde{D}^{-\frac{1}{2}} \tilde{A} \tilde{D}^{-\frac{1}{2}} H^l W \right)$
- Motivation in Spectral Graph Theory

$$\sum_{i \in \mathcal{N}(i) \cup \{i\}} \frac{1}{\sqrt{|\mathcal{N}(i)| + 1}\sqrt{|\mathcal{N}(j)| + 1}} W^{l} \mathbf{h}_{j}^{l}$$

$$W^{l} \text{ where } \tilde{A} = A + I_{n}, \tilde{D} = \text{diag}\left(\sum_{i} \tilde{A}_{ii}\right)$$



- Message passing (spatial): $\mathbf{h}_i^{l+1} = \sigma \Big(\int_{i \in I} \mathbf{h}_i^{l+1} \Big) = \sigma \Big(\int_{i \in I} \mathbf{h}_i^{l+1} \Big)$
- Matrix version: $H^{l+1} = \sigma \left(\tilde{D}^{-\frac{1}{2}} \tilde{A} \tilde{D}^{-\frac{1}{2}} H^{l} W \right)$
- Motivation in Spectral Graph Theory
 - Spectral decomposition of the normalized graph Laplacian $L = I_n D^{-\frac{1}{2}}AD^{-\frac{1}{2}} = U\Lambda U^{\top}$

$$\sum_{e \in \mathcal{N}(i) \cup \{i\}} \frac{1}{\sqrt{|\mathcal{N}(i)| + 1}\sqrt{|\mathcal{N}(j)| + 1}} W^{l} \mathbf{h}_{j}^{l} \right)$$

$$\mathcal{N}^{l} \text{ where } \tilde{A} = A + I_{n}, \tilde{D} = \text{diag}\left(\sum_{i} \tilde{A}_{ii}\right)$$



• Message passing (spatial): $\mathbf{h}_{i}^{l+1} = \sigma \Big(\int_{i \in i} \mathbf{h}_{i}^{l+1} \Big) = \sigma \Big(\int_{i \in i} \mathbf{h}_{i}^{l+1} \Big)$

• Matrix version: $H^{l+1} = \sigma \left(\tilde{D}^{-\frac{1}{2}} \tilde{A} \tilde{D}^{-\frac{1}{2}} H^{l} W \right)$

- Motivation in Spectral Graph Theory
 - Spectral decomposition of the normalized graph Laplacian $L = I_n D^{-\frac{1}{2}}AD^{-\frac{1}{2}} = U\Lambda U^{\top}$
 - Generalizing convolution to "graph convolution": $\mathbf{h}_{G}^{*}\mathbf{g} = F^{-1}(F(\mathbf{h}) \odot F(\mathbf{g})), F(\mathbf{h}) = U\mathbf{h}$

$$\sum_{\substack{\in \mathcal{N}(i) \cup \{i\}}} \frac{1}{\sqrt{|\mathcal{N}(i)| + 1}\sqrt{|\mathcal{N}(j)| + 1}}} W^{l} \mathbf{h}_{j}^{l} \right)$$

$$\mathcal{N}^{l} \text{ where } \tilde{A} = A + I_{n}, \tilde{D} = \text{diag}\left(\sum_{i} \tilde{A}_{ii}\right)$$


Why "Convolutional"? The Long Answer

• Message passing (spatial): $\mathbf{h}_{i}^{l+1} = \sigma \Big(\int_{i \in i} \mathbf{h}_{i}^{l+1} \Big) = \sigma \Big(\int_{i \in i} \mathbf{h}_{i}^{l+1} \Big)$

• Matrix version: $H^{l+1} = \sigma \left(\tilde{D}^{-\frac{1}{2}} \tilde{A} \tilde{D}^{-\frac{1}{2}} H^{l} W \right)$

- Motivation in Spectral Graph Theory
 - Spectral decomposition of the normalized graph Laplacian $L = I_n D^{-\frac{1}{2}}AD^{-\frac{1}{2}} = U\Lambda U^{\top}$
 - Generalizing convolution to "graph convolution": $\mathbf{h} *_{G} \mathbf{g} = F^{-1}(F(\mathbf{h}) \odot F(\mathbf{g})), F(\mathbf{h}) = U\mathbf{h}$
 - Approximation leads to the matrix version above

$$\sum_{A \in \mathcal{N}(i) \cup \{i\}} \frac{1}{\sqrt{|\mathcal{N}(i)| + 1}} W^{l} \mathbf{h}_{j}^{l}} W^{l} \mathbf{h}_{j}^{l}}$$

$$V^{l} \text{ where } \tilde{A} = A + I_{n}, \tilde{D} = \text{diag}\left(\sum_{i} \tilde{A}_{ii}\right)$$





- Node classification: softmax(\mathbf{z}_i)



- Node classification: softmax(\mathbf{z}_i)
- Graph classification: softmax $\left(\frac{1}{n}\sum_{i}\mathbf{z}_{i}\right)$







- Node classification: softmax(\mathbf{z}_i)
- Graph classification: softmax $\left(\frac{1}{n}\sum_{i}\mathbf{z}_{i}\right)$
- Link prediction: $p(A_{ij}) = \text{sigmoid}(\mathbf{z}_i^{\mathsf{T}}\mathbf{z}_j)$





Citation networks with partially labeled nodes





Method	Citeseer	Cora	Pubmed
ManiReg [3]	60.1	59.5	70.7
SemiEmb [28]	59.6	59.0	71.1
LP [32]	45.3	68.0	63.0
DeepWalk [22]	43.2	67.2	65.3
ICA [18]	69.1	75.1	73.9
Planetoid* [29]	64.7 (26s)	75.7 (13s)	77.2 (25)
GCN (this paper)	70.3 (7s)	81.5 (4s)	79.0 (38
GCN (rand. splits)	67.9 ± 0.5	80.1 ± 0.5	78.9 ± 0

	21.8
	26.7
	26.5
	58.1
	23.1
)	61.9 (185s)
)	66.0 (48s)
7	58.4 ± 1.7



t-SNE of hidden layer activation.



Description	Propagation model	Citeseer	Cora	Pubmed
Chebyshev filter (Eq. 5) $K = 3$	$\nabla^{K} T(\tilde{I}) V \Theta$	69.8	79.5	74.4
K = 2	$\sum_{k=0} I_k(L) \Lambda \Theta_k$	69.6	81.2	73.8
1 st -order model (Eq. 6)	$X\Theta_0 + D^{-\frac{1}{2}}AD^{-\frac{1}{2}}X\Theta_1$	68.3	80.0	77.5
Single parameter (Eq. 7)	$(I_N + D^{-\frac{1}{2}}AD^{-\frac{1}{2}})X\Theta$	69.3	79.2	77.4
Renormalization trick (Eq. 8)	$ ilde{D}^{-rac{1}{2}} ilde{A} ilde{D}^{-rac{1}{2}}X\Theta$	70.3	81.5	79.0
1 st -order term only	$D^{-\frac{1}{2}}AD^{-\frac{1}{2}}X\Theta$	68.7	80.5	77.8
Multi-layer perceptron	$X\Theta$	46.5	55.1	71.4









- Relations (edge types) ${\mathscr R}$



- Relations (edge types) ${\mathscr R}$

• GCN:
$$\mathbf{h}_{i}^{l+1} = \sigma \left(\sum_{j \in \mathcal{N}(i) \cup \{i\}} \frac{1}{c_{ij}} W^{l} \mathbf{h}_{j}^{l} \right)$$
 whe

ere $c_{ij} = \sqrt{|\mathcal{N}(i)| + 1}\sqrt{|\mathcal{N}(j)| + 1}$



• Relations (edge types) \mathscr{R}

• GCN:
$$\mathbf{h}_{i}^{l+1} = \sigma \left(\sum_{j \in \mathcal{N}(i) \cup \{i\}} \frac{1}{c_{ij}} W^{l} \mathbf{h}_{j}^{l} \right)$$
 where $c_{ij} = \sqrt{|\mathcal{N}(i)| + 1} \sqrt{|\mathcal{N}(j)| + 1}$
• RGCN: $\mathbf{h}_{i}^{l+1} = \sigma \left(\sum_{r \in \mathcal{R}} \sum_{j \in \mathcal{N}^{r}(i)} \frac{1}{c_{ijr}} W^{l}_{r} \mathbf{h}_{j}^{l} \right)$ where $c_{ijr} = |\mathcal{N}^{r}(i)|$

• Assume \mathscr{R} and \mathscr{N}^r capture self-loop



• Relations (edge types) \mathscr{R}

• GCN:
$$\mathbf{h}_{i}^{l+1} = \sigma \left(\sum_{j \in \mathcal{N}(i) \cup \{i\}} \frac{1}{c_{ij}} W^{l} \mathbf{h}_{j}^{l} \right)$$
 where $c_{ij} = \sqrt{|\mathcal{N}(i)| + 1} \sqrt{|\mathcal{N}(j)| + 1}$
• RGCN: $\mathbf{h}_{i}^{l+1} = \sigma \left(\sum_{r \in \mathcal{R}} \sum_{j \in \mathcal{N}^{r}(i)} \frac{1}{c_{ijr}} W_{r}^{l} \mathbf{h}_{j}^{l} \right)$ where $c_{ijr} = |\mathcal{N}^{r}(i)|$

- Assume \mathscr{R} and \mathscr{N}^r capture self-loop
- Can be used to model directed graphs



RGCN: Regularization Schlichtkrull et al. (2018)





RGCN: Regularization Schlichtkrull et al. (2018)





RGCN: Regularization Schlichtkrull et al. (2018)

• Basis-decomposition: $W_r^l = \sum_{rb}^{B} a_{rb}^l V_b^l$ b=1

• Block-diagonal-decomposition: $W_r^l = \bigoplus^B Q_{br}^l$

b=1





• GCN: $\mathbf{h}_{i}^{l+1} = \sigma \left(\sum_{j \in \mathcal{N}(i) \cup \{i\}} \frac{1}{c_{ij}} W^{l} \mathbf{h}_{j}^{l} \right)$ where $c_{ij} = \sqrt{|\mathcal{N}(i)| + 1} \sqrt{|\mathcal{N}(j)| + 1}$



• GCN:
$$\mathbf{h}_{i}^{l+1} = \sigma \left(\sum_{j \in \mathcal{N}(i) \cup \{i\}} \frac{1}{c_{ij}} W^{l} \mathbf{h}_{j}^{l} \right)$$
 where $c_{ij} = \sqrt{|\mathcal{N}(i)| + 1} \sqrt{|\mathcal{N}(j)| + 1} \sqrt{|$



• GCN:
$$\mathbf{h}_{i}^{l+1} = \sigma \left(\sum_{j \in \mathcal{N}(i) \cup \{i\}} \frac{1}{c_{ij}} W^{l} \mathbf{h}_{j}^{l} \right)$$
 where $c_{ij} = \sqrt{|\mathcal{N}(i)| + 1} \sqrt{|\mathcal{N}(j)| + 1} \sqrt{|$



GAT for Semi-Supervised Learning Velickovič et al. (2018)

Method

MLP

ManiReg (Belkin et al., 2006) SemiEmb (Weston et al., 2012) LP (Zhu et al., 2003) DeepWalk (Perozzi et al., 2014) ICA (Lu & Getoor, 2003) Planetoid (Yang et al., 2016) Chebyshev (Defferrard et al., 2016) GCN (Kipf & Welling, 2017) MoNet (Monti et al., 2016)

GCN-64* GAT (ours)

Cora	Citeseer	Pubmed
55.1%	46.5%	71.4%
59.5%	60.1%	70.7%
59.0%	59.6%	71.7%
68.0%	45.3%	63.0%
67.2%	43.2%	65.3%
75.1%	69.1%	73.9%
75.7%	64.7%	77.2%
81.2%	69.8%	74.4%
81.5%	70.3%	79.0%
$81.7\pm0.5\%$		$78.8\pm0.3\%$
$81.4\pm0.5\%$	$70.9\pm0.5\%$	$\textbf{79.0} \pm 0.3\%$
$\textbf{83.0}\pm0.7\%$	$\textbf{72.5}\pm0.7\%$	$\textbf{79.0} \pm 0.3\%$



GAT for Semi-Supervised Learning Velickovič et al. (2018)



GAT



Special Case: Transformer as Message Passing



Transformer. From Guo et al. (2018).





Special Case: Transformer as Message Passing



Transformer. From Guo et al. (2018).



Star-Transformer. From Guo et al. (2018).





Applying GNNs to Linguistic Graphs

Infusing Finetuning with Semantic Dependencies

Zhaofeng Wu, Hao Peng, and Noah Smith. TACL 2021.



Motivation



Motivation



Devlin et al. (2018); Liu et al. (2019b)

It was not bad.



Motivation



Hewitt and Manning (2019); Tenney et al. (2019); Liu et al. (2019a)

advmod	punct
ed Transfo	ormer
RB	JJ .
S Tagger	

It was not bad.



Would semantics help?





• We show BERT/RoBERTa less prominently surface semantics...



- We show BERT/RoBERTa less prominently surface semantics...
- ... and the explicit incorporation of semantic information:



- We show BERT/RoBERTa less prominently surface semantics...
- ... and the explicit incorporation of semantic information:
 - 1. Improves downstream task performance



- We show BERT/RoBERTa less prominently surface semantics...
- ... and the explicit incorporation of semantic information:
 - 1. Improves downstream task performance
 - 2. Helps guard against frequent yet invalid heuristics


Introduction

- We show BERT/RoBERTa less prominently surface semantics...
- ... and the explicit incorporation of semantic information:
 - 1. Improves downstream task performance
 - 2. Helps guard against frequent yet invalid heuristics
 - 3. Better captures nuanced linguistic phenomena



Introduction

- We show BERT/RoBERTa less prominently surface semantics...
- ... and the explicit incorporation of semantic information:
 - 1. Improves downstream task performance
 - 2. Helps guard against frequent yet invalid heuristics
 - 3. Better captures nuanced linguistic phenomena
 - 4. Increases training sample efficiency



Operationalizing "Meaning"

is technique This

BV

DELPH-IN MRS-Derived Dependencies (**DM**; Ivanova et al., 2012)

Adapted from WSJ #20209013



impossible adopt to



Operationalizing "Meaning"

This technique is

BV



Stanford Dependencies (SD; de Marneffe et al., 2006)

Adapted from WSJ #20209013



impossible adopt to

DELPH-IN MRS-Derived Dependencies (**DM**; Ivanova et al., 2012)



7



This technique is impossible to adopt.

Probing model (Shi et al., 2016; Adi et al., 2017)





Probing model (Shi et al., 2016; Adi et al., 2017)

Ceiling model (Dozat and Manning, 2017, 2018)

8

Probing RoBERTa with Semantics

Probing - Ceiling; RoBERTa-base





Can we use semantics to augment pretrained transformers?



10



I so love to make slides.

11



I so love to make slides.







Che et al. (2019)

I so love to make slides.







I so love to make slides.



14



Schlichtkrull et al. (2017)





I so love to make slides.





I so love to make slides.



Experiments

- Dataset: GLUE (Wang et al., 2018)
- Backbone: RoBERTa (Liu et al., 2019b)
- Parser: SOTA DM parser with 92.5 labeled F1 (Che et al., 2019)
- Graph Encoder: RGCN (Schlichtkrull et al., 2017)
 - 2 layers
 - Hidden dimension $\in \{256, 512, 768\}$
- Epochs $\in \{3, 10, 20\}$, learning rate $\in \{1 \times 10^{-4}, 2 \times 10^{-5}\}$







Results





Results





Analysis: When Do Semantic Structures Help?

- Two datasets
 - 2019)
 - GLUE diagnostics tests the model capability in various **linguistic** phenomena (Wang et al., 2018)
- Examine a model trained on existing NLI datasets with synthetic NLI examples

HANS tests if a model uses invalid reasoning heuristics (McCoy et al.,





Analysis: HANS Lexical Overlap

The actor stopped the banker. do

RoBER

68.1

does not entail The banker stopped the actor.

Та	SIFT	
	71.0 (+2.9)	



Analysis: HANS Subsequence

The judges heard the actor resigned. The judges heard the actor. does not entail

RoBER

25.8



Та	SIFT	
	29.5 (+3.7)	



Analysis: HANS Constituent

If the actor slept, the senator ran. doe

RoBER

37.9

does not entail **The actor slept.**

Та	SIFT	
	37.6 (-0.3)	



Analysis: HANS Constituent

If the actor slept, the senator ran. doe **Before** the actor slept, the senator ran.

RoBER

37.9

does not entail The actor slept.

entails

Та	SIFT	
	37.6 (-0.3)	



Analysis: GLUE Diagnostics

Pred-Arg Structure

I opened the door.

enta

does no

Logic I have no pet puppy. entagonal does not

		RoBERTa	SIF
ails ot entail	The door opened. I opened.	43.5	44.6 (-
ails ot entail	I have no corgi pet puppy. I have no pet.	36.2	38.3 (-







Analysis: GLUE Diagnostics

Lexical Semantics

I have a dog.

enta

does no

I live in Seattle.

enta

Knowledge

does no

		RoBERTa	SIF
ails ot entail	I have an animal. I have a cat.	45.6	44.8 (
ails ot entail	I live in the U.S. I live in Antarctica.	28.0	26.3 (







Analysis: Sample Efficiency

Use the same downsampled MNLI training set to train RoBERTa & SIFT





Analysis: Sample Efficiency

- Use the same downsampled MNLI training set to train RoBERTa & SIFT





Absolute Δ (SIFT - RoBERTa) on MNLI



Summary





Summary





I so love to make slides .



Further Readings

- A Comprehensive Survey on Graph Neural Networks (Wu et al., 2019)
- Deep Learning on Graphs: A Survey (Zhang et al., 2020)
- Graph Neural Networks: A Review of Methods and Applications (Zhou et al., 2021)
- Frameworks
 - DGL: <u>https://github.com/dmlc/dgl</u> lacksquare
 - PyTorch Geometric: <u>https://github.com/rusty1s/pytorch_geometric</u>





Questions?



References

Wanxiang Che, Longxu Dou, Yang Xu, Yuxuan Wang, Yijia Liu, and Ting Liu. 2019. HIT-SCIR at MRP 2019: A unified pipeline for meaning representation parsing via efficient training and effective encoding. In Proc. of MRP.

Justin Gilmer, Samuel S. Schoenholz, Patrick F. Riley, Oriol Vinyals, and George E. Dahl. 2018. Neural Message Passing for Quantum Chemistry. In Proc. of ICML.

Qipeng Guo, Xipeng Qiu, Pengfei Liu, Yunfan Shao, Xiangyang Xue, and Zheng Zhang. 2019. Star-Transformer. In Proc. of NAACL.

Angelina Ivanova, Stephan Oepen, Lilja Øvrelid, and Dan Flickinger. 2012. Who did what to whom?: A contrastive study of syntactosemantic dependencies. In Proc. of LAW.

Thomas N. Kipf and Max Welling. 2017. Semi-supervised classification with graph convolutional networks. In Proc. of ICLR.

Stephan Oepen, Omri Abend, Lasha Abzianidze, Johan Bos, Jan Hajic, Daniel Hershcovich, Bin Li, Tim O'Gorman, Nianwen Xue, and Daniel Zeman. 2020. MRP 2020: The Second Shared Task on Cross-Framework and Cross-Lingual Meaning Representation Parsing. In Proc. of CoNLL.



References

Stephan Oepen, Marco Kuhlmann, Yusuke Miyao, Daniel Zeman, Silvie Cinkova, Dan Flickinger, Jan Hajic, and Zdenka Uresova. 2015. Semeval 2015 task 18: Broad-coverage semantic dependency parsing. In Proc. of SemEval.

Michael Schlichtkrull, Thomas N. Kipf, Peter Bloem, Rianne Van Den Berg, Ivan Titov, and Max Welling. 2018. Modeling relational data with graph convolutional networks. In Proc. of European Semantic Web Conference.

Petar Velickovič, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, and Yoshua Bengio. 2018. Graph attention networks. In Proc. of ICLR.

Zonghan Wu, Shirui Pan, Fengwen Chen, Guodong Long, Chengqi Zhang, and Philip S. Yu. 2019. A Comprehensive Survey on Graph Neural Networks. In IEEE Transactions on Neural Networks and Learning Systems.

Zhaofeng Wu, Hao Peng, and Noah A. Smith. 2021. Infusing Finetuning with Semantic Dependencies. In TACL.

Ziwei Zhang, Peng Cui, and Wenwu Zhu. 2020. Deep Learning on Graphs: A Survey. In IEEE Transactions on Knowledge and Data Engineering.

Jie Zhou, Ganqu Cui, Shengding Hu, Zhengyan Zhang, Cheng Yang, Zhiyuan Liu, Lifeng Wang, Changcheng Li, and Maosong Sun. 2021. Graph Neural Networks: A Review of Methods and Applications. In Al Open.

ı. al

51