Commonsense Reasoning with Implicit Knowledge in Natural Language

*Pratyay Banerjee, *Swaroop Mishra, *Kuntal Pal,

*Arindam Mitra, Chitta Baral





Motivation

Abductive NLI	Social IQA	Physical IQA
 Obs1: Jim was working on a project. Jim found he was missing an item. Jim needed a certain animal for it. Obs2: Luckily, he found it on a nearby shelf Knowledge: Peyton eventually found it before Peyton needed to determine that something is missing. Kendall never found it, as a result Kendall wants to lodge a missing complaint. 	 Context: Remy was an expert fisherman and was on the water with Kai. Remy baited Kai's hook. Question: What will Remy want to do next? cast the line put the boat in the water invite Kai out on the boat Knowledge: Alex baits Pat's hook as a result others want to cast their line. 	 Goal: When doing sit-ups: place your tongue in the roof of your mouth. It will stop you from straining your neck. place your elbow in the roof of your mouth. It will stop you from straining your neck. Knowledge: How to Do Superbrain Yoga. Place your tongue on the roof of your mouth.

- Can we answer commonsense questions regarding Abductive, Social Interactions and Physical Interactions?
- Can we improve our answers with implicit knowledge in natural language form ?
- How can we find relevant knowledge ?
- How to reason with retrieved implicit relevant knowledge ?
- What kind of training strategies can we use ?
- How do they compare with respect to large / complex models ?

Knowledge Retrieval Pipeline



• Knowledge Source

- SIQA ----- Atomic (Partially Derived)
- PIQA ----- Wikihow (Relevant)

Knowledge Infusion Modes

• The top 10 knowledge statements(K_{ij}) retrieved for each options (A_i) are infused with the transformer encoders in four different ways

Training Strategies

- **REVISION**
 - Pre-Trained (MLM+NSP [BERT] and MLM [RoBERTa]) on K_D
 - Fine-Tuned on D
- OPENBOOK
 - Fine-Tuned on D+S, S is a subset of K_D
- **REVISION+OPENBOOK**
 - Pre-Trained (MLM+NSP [BERT] and MLM [RoBERTa]) on K_n
 - Fine-Tuned on D+S, S is a subset of K_D

 $\rm K_{\rm D}$ - Respective knowledge sources for given 3 datasets D

Flow-Diagram of our Weighted-Sum Model



A sample from PIQA dataset

Results : Training Strategies with Knowledge infusion modes

		BERT				RoBERTa			
Dataset	Strategy	Concat	Max	Sim-Sum	Wtd-Sum	Concat	Max	Sim-Sum	Wtd-Sum
aNLI	OPENBOOK REVISION REVISION & OPENBOOK	$\begin{array}{c} 73.9 \pm \ 0.8 \\ 72.7 \pm \ 0.3 \\ 74.4 \pm \ 0.2 \end{array}$	73.7 ± 0.1 N/A 74.3 ± 0.1	73.5 ± 0.7 N/A 74.0 ± 0.9	$73.3 \pm 1.0 \\ { m N/A} \\ \underline{75.1} \pm 0.4$	83.9 ± 0.5 82.4 84.2 ± 0.7	$\begin{array}{c} 80.8 \pm \ 0.9 \\ {\rm N/A} \\ 81.4 \pm \ 0.8 \end{array}$	81.7 ± 0.6 N/A 82.6 ± 0.6	84.4 ± 0.4 N/A 86.7 ± 0.6
PIQA	OPENBOOK REVISION REVISION & OPENBOOK	67.8 ± 0.4 74.5 ± 0.3 67.7 ± 0.1	72.4 ± 0.6 N/A 73.8 ± 0.8	72.6 ± 1.2 N/A 76.8 ± 0.5	72.5 ± 0.1 N/A <u>76.8</u> ± 0.3	$\begin{array}{c} 74.8 \pm \ 0.5 \\ 75.2 \pm \ 0.8 \\ 75.4 \pm \ 0.7 \end{array}$	75.2 ± 0.9 N/A 76.2 ± 0.8	75.6 ± 0.7 N/A 76.8 ± 0.4	$77.1 \pm 0.2 \\ N/A \\ 80.2 \pm 0.6$
SIQA	OPENBOOK REVISION REVISION & OPENBOOK	$\begin{array}{c} 70.1 \pm \ 0.8 \\ 69.5 \pm \ 0.9 \\ 68.8 \pm \ 0.4 \end{array}$	67.8 ± 0.1 N/A 66.6 ± 0.4	70.0 ± 0.7 N/A 68.9 ± 0.1	$\frac{70.2 \pm 0.4}{\text{N/A}} \\ 69.3 \pm 0.6$	76.5 ± 0.7 76.8 ± 0.3 78.2 ± 0.3	77.2 ± 0.6 N/A 77.4 ± 0.9	77.4 ± 0.2 N/A 76.7 ± 0.5	$78.3 \pm 0.5 \\ N/A \\ \textbf{79.5} \pm 0.9$

- Revision with Openbook strategy works best across 3 datasets for each of 4 knowledge infusion modes
- Weighted sum model shows the best performance across all datasets

Results : Weighted Sum model Performance Comparison

Models/ Accuracy	aNLI		PIQA		SIQA	
	Val	Test	Val	Test	Val	Test
BERT	67.36	66.75	68.08	69.23	64.88	64.50
GPT-2 XL	N/A	N/A	70.20	69.50	47.50	45.30
RoBERTa	85.05	83.91	76.28	76.80	77.85	76.74
RoBERTa 5 Ensemble	N/A	83.22	N/A	79.66	N/A	78.68
$L2R^2$ [2020]	N/A	86.81	N/A	N/A	N/A	N/A
KagNet [2019]	N/A	N/A	N/A	N/A	65.05	64.59
GBR [2020]	N/A	N/A	N/A	N/A	75.64	76.25
UnifiedQA T5 11B [2020]	N/A	80.04	N/A	89.50	N/A	79.75
Ours: BERT + WS	74.60	74.96	76.82	72.28	70.21	67.22
Ours: $RoBERTa + WS$	85.90	84.18	80.20	78.24	79.53	78.00

- WS with BERT and RoBERTa is better than baseline BERT and RoBERTa
- Better than huge parameterized model (T5 in aNLI)
- WS shows better performance than complex Graph based models like KagNet or GBR

Analysis

Types of Error	aNLI	SIQA	PIQA
Annotation	41%	38%	10%
Model Prediction	48%	27%	29%
Distracting Knowledge	11%	35%	61%

Error Percentage

Knowledge	aNLI	SIQA	PIQA
Explicitly Present Implicitly Present Fully Irrelevant	$14\% \\ 55\% \\ 31\%$	$11\% \\ 59\% \\ 30\%$	$10\% \\ 51\% \\ 39\%$

Knowledge Classification - Correct Predictions

- For PIQA we find most of the errors are due to distracting knowledge
- Over 50% of the correct predictions for each of the 3 datasets have implicit knowledge

Analysis of the Weighted Sum model





- We show the impact of increasing knowledge sentences on the accuracy with Revision training strategy for both BERT and RoBERTa
 - It increases for SIQA and PIQA
- We also show the weights learned by the weight layer of the WS model
 - We find that the model learns lower weights where there is very less lexical overlap between knowledge retrieved and Question and Answer options

Conclusion

In this paper through 3 commonsense reasoning dataset,

- We perform a thorough analysis of transformers' (BERT and RoBERTa) ability to reason.
- We present four modes of knowledge infusion in transformer encoders which improves 2-9% of accuracy across the datasets
- We carry out an extensive investigation to study the impact of different knowledge sources and pre-training on such knowledge sources.





Thank You !!!

Contacts : Pratyay Banerjee, Swaroop Mishra, Kuntal Pal, Arindam Mitra, Chitta Baral

- **Paper** : <u>https://openreview.net/pdf?id=a4-fFL7aCi0</u>
- Email: {pbanerj6, srmishr1, kkpal, chitta}@asu.edu, arindam.mitra@microsoft.com