



A Benchmark for Rationalizing Natural Language Analogical Reasoning

Jiangjie Chen, Rui Xu, Ziquan Fu, Wei Shi, Zhongqiao Li, Xinbo Zhang, Changzhi Sun, Lei Li, Yanghua Xiao, Hao Zhou



Word Analogy Recognition

Query	${f Q}$) newton:english
Candidate	A) marx:german
answers	B) confucius:russian
	C) caesar: american
	D) plato:canadian

An analogical reasoning problem from The Bigger Analogy Test Set (BATS).

Word Analogy Recognition



An analogical reasoning problem from The Bigger Analogy Test Set (BATS).





Linear Analogy (Ethayarajh et al. 2019)

$$\overrightarrow{king} - \overrightarrow{man} + \overrightarrow{woman} = \overrightarrow{queen}$$

Previous work

- Methods: Hold a connectionist assumption
- Benchmarks: Evaluate pre-trained word representations for linear analogy

e.g. Word2Vec

Simple Binary Relations

Lexical, morphological, simple semantic relations.



Linear Analogy (Ethayarajh et al. 2019)

$$\overrightarrow{king} - \overrightarrow{man} + \overrightarrow{woman} = \overrightarrow{queen}$$

Previous work

- Methods: Hold a connectionist assumption
- Benchmarks: Evaluate pre-trained word representations for linear analogy

e.g. *Word2Vec*

Simple Binary Relations

Lexical, morphological, simple semantic relations.

Not Explainable

Unable to reveal human-like analogical reasoning process.

term2

Nationality

term1

Complex Analogy

Q) tea¹:teapot²:teacup³

A) passengers¹:bus²:taxi³

B) magazine¹:bookshelf²:reading room³

C) talents¹:school²:enterprise³

D) textbooks¹:bookstore²:printing factory³

An analogical reasoning problem from Civil Service Exams of China. (Translated)



Sourced from Civil Service Exams of China

Explainable

Free-text Explanations



Chinese & English

Image: Chinese Image: Chinese #Problems=1665 English #Expl. =5×1665 #Expl. =5×1251

Complex Analogy

Dataset	Lang.	Data Size (train / val / test)	# of Terms in Cand.	Has Expl.
SAT	En	0 / 37 / 337	2	×
Google	En	0 / 50 / 500	2	X
BATS	En	0 / 199 / 1,799	2	×
E-KAR	Zh	1,155 /165 / 335	$-\frac{1}{2}_{(64.5\%)}^{}, 3_{(35.5\%)}$	✓
	En	870 / 119 / 262	$2_{(60.5\%)}, 3_{(39.5\%)}$	1



Challenging

Sourced from Civil Service Exams of China



Free-text Explanations

How to Rationalize Analogical Reasoning?

Bilingual

Chinese & English

Structuremapping theory

(Minnameier et al, 2010)

Q) tea¹:teapot²:teacup³

A) passengers¹:bus²:taxi³

B) magazine¹:bookshelf²:reading room³

C) talents¹:school²:enterprise³

D) textbooks¹:bookstore²:printing factory³

Structuremapping theory

(Minnameier et al, 2010)

Abduction

Draw a *source structure* that may work for target.

Q) tea ¹ :teapot ² :teacup ³								
Source Structures	Container for is_a $teapot^2$	r holding <i>tea</i> ¹ \scale_is_a <i>teacup</i> ³	transpo teapot ²	teacup ³				

A) passengers¹:bus²:taxi³

B) magazine¹:bookshelf²:reading room³

C) talents¹:school²:enterprise³

D) textbooks¹:bookstore²:printing factory³





How to Rationalize Analogical Reasoning?











Lessons from Preliminary Exploration of E-KAR

Lesson 1: W2Vs and LMs both struggle at complex analogical reasoning.



Humans outperform SOTA models by large margins.

(Please check the paper for details.)

Lessons from Preliminary Exploration of E-KAR

Lesson 2: LMs struggle at rationalizing analogical reasoning.



 Poor quality of generated explanations, improvement over baseline but fall far behind gold.
 Gold explanations can be exploited by Analogical QA models to achieve nearly perfect results (97.7%).

Lessons from Preliminary Exploration of E-KAR

Lesson 2: LMs struggle at rationalizing analogical reasoning.

Error Analysis

- *1.Unable to generate <u>negated</u>* <u>facts</u> for refutation.
- 2. Generating <u>factually incorrect</u> <u>statements</u>.
- 3.Biasing towards <u>common</u> patterns.

Ex1. "term1" and "term2" has the same meaning. Ex2. "term1" is a "term2".

What is Next?

- What we have: A novel benchmark for rationalizing analogical reasoning, which is challenging, explainable and bilingual.
- Analogical reasoning by effectively interacting with various kinds of knowledge.
 - e.g. commonsense, factual and *cultural* knowledge.
- Generate reasons with evidence to rationalize reasoning.
 - Particularly, enable models to generate *negated* statements/facts.

Have Fun with E-KAR!







https://eval.ai/web/challenges/ challenge-page/1671/overview



jjchen19@fudan.edu.cn



https://jiangjiechen.github.io