Tracing and Manipulating Intermediate Values in Neural Math Problem Solvers

<u>Yuta Matsumoto¹</u>, Benjamin Heinzerling², Masashi Yoshikawa¹,

Kentaro Inui^{1,2}

^{1.}Tohoku Universiy ^{2.}RIKEN

Summary

- How do neural networks process complex inputs?
- We focus on the intermediate values of simple math problems



- We show that...
 - intermediate values are stored in particular directions of model representations
 - model representations can be manipulated as if inputs change

Neural math problem solvers perform well

• Solving math problems requires storing intermediate values

$$54 - ((258 + 314) - (143 - 96))$$
$$= 54 - (572 - 47)$$
$$= 54 - 525$$
$$= -471$$

- We can guess that models capture those intermediate values
- Goal: Find out how intermediate values are stored

Contributions: Tracing and Manipulation

- Idea: focus on simple math problems and their intermediate values
- 0. Pre-training Transformer on simple math problems
 - e.g., (115-28)-(32-56), 25-(79+104),...

1. Tracing intermediate values by PCA

- find components where information on intermediate values is encoded

2. Manipulating intermediate values

- test these components influence model predictions

0. Pre-training on simple math problems

Method

- Prepare multiple patterns of equations
- Train a Transformer model on the equations

Result

- High performance
 - R^2 score = 0.999





2022/11/16





2022/11/16





Correlation with intermediate values

Q. Are there principal components that are highly correlated with intermediate values?

A. Yes

- Some components highly correlate with intermediate values.













Manipulation results

- Move weights of PC2 in layer 3
 - highly correlated (r=0.97) with
 b in a-(b-c)
- Model predictions change consistently
 - the intermediate values may be manipulated



Are we moving intermediate values as if inputs change?

- Prepare inputs with different intermediate values
 - e.g., 617 (100 602), 617 (101 602), ..., 617 (999 602)
- Compare principal components after changing inputs and after manipulating activations



Are we moving intermediate values as if inputs change? \rightarrow Yes

- High agreement
 - particularly around the original data
- We can change intermediate values by manipulating activations



Predicted: change activation

Actual: change the input

Summary

- We analyzed the representation of intermediate values in neural math problem solvers
- By Tracing, we found the directions which correlate with intermediate values
- By Manipulation, we showed that the directions are causally related with intermediate values