Implicit Representations of Meaning in Neural Language Models

- Jacob Andreas
- LINGO.csail.mit.edu
- MIT CSAII



Implicit Representations of Meaning in Neural Language Models [ACL 2021]



Belinda Li



Max Nye



[Brown et al. 2020; example from Marcus & Davis 2020 / Charniak 1972]



[Brown et al. 2020; example from Marcus & Davis 2020 / Charniak 1972]

get a top." "I will get Jack a top," said Janet.



get a top." "I will get Jack a top," said Janet.

[Brown et al. 2020; example from Marcus & Davis 2020 / Charniak 1972]



get a top." "I will get Jack a top," said Janet.

[Brown et al. 2020; example from Marcus & Davis 2020 / Charniak 1972]

Language comprehension???



Neural sequence models



John has a book. Mary has an apple. He gave her his





















Janet went to the store to get Jack a top. But Jack already has a colorful top.

Janet went to the s She gave it to him.

[Heim 83, "File Change Semantics"; Kamp 81, "Discourse Representation Theory"; Groenendijk & Stokhof 91, "Dynamic Predicate Logic"]

Dynamic Semantics

Janet and Penny went to the store to get Jack a top. But Jack already has a colorful top. Penny possesses Janet Jack top

World models & language models

Janet and Penny went to the store to get Jack a top. But Jack already has a colorful top.

vector representations

Representations in language models

p("Jack will get a top" | ...)

♥
semantic probe ---->

Building the probe

Building the probe

Decode facts about an entity from the <u>encoding</u> of its first mention.

Building the probe

Evaluation

TextWorld

You are navigating through a house. You've just entered a serious study. There is a gross looking mantle in the room. It has nothing on it. You see a closed rusty toolbox. Now why would someone leave that there? Looks like there is a locked door. Find the key to unlock the door. You should try going east.

What fraction of entities are exactly reconstructed?

Evaluation: does it work?

TextWorld

What fraction of entities are exactly reconstructed?

Evaluation: does it work?

What kind of training matters?

Evaluation: does it work?

T5 T5, no fine-tuning random init random init, no f-t

TextWorld

What kind of training matters?

Evaluation: locality

Evaluation: locality

58.5% / 64.8% accuracy

Language models as world models

There's a locked wooden door leading east [...] you open the door.

Next, (6b) gets uttered, which prompts the listener to update card 1 by adding the entry "hit 2", and to update card 2 by adding "was hit by 1". He now has F2, still a two card file, but a different one:

Language models as file cards

Building states from scratch

Building states from scratch

What's still missing

Attribution of model errors:

p(generation is semantically acceptable)

p(probe is accurate)

Does grounded training improve accuracy?

of training examples

Would ground-truth states improve accuracy?

% of training examples with state labels

Quantification: There are twenty-three reindeer; most of them have red noses.

Implication and counterfactuals: If Pat goes to the party, so will Jan. If Pat had gone to the last one, Mo would have gone too. Pat will go to the party this time.

Summary

Language produce (rudimentary) representations of world states, and these states can be manipulated with predictable effects on model output.

But far from 100% reliable; lots of open capture and how to improve them.

- questions about what these representations

RESEARCHERS

Belinda Li

Max Nye

Thank you!

Sponsors

MachineLearningApplications@CSAIL

Summary

Language produce (rudimentary) representations of world states, and these states can be manipulated with predictable effects on model output.

But far from 100% reliable; lots of open capture and how to improve them.

- questions about what these representations