



Aggregating analyses in surveys and toolkits

Each analysis paper asks a very specific question.

How do we ask, *what does the field currently know about BERT?*

Answer: meta-studies compiling results

Analysis Methods in Neural Language Processing: A Survey

Yonatan Belinkov^{1,2} and James Glass¹

¹MIT Computer Science and Artificial Intelligence Laboratory

²Harvard School of Engineering and Applied Sciences

Cambridge, MA, USA

{belinkov, glass}@mit.edu

A Primer in BERTology: What we know about how BERT works

Anna Rogers, Olga Kovaleva, Anna Rumshisky

Department of Computer Science, University of Massachusetts Lowell
Lowell, MA 01854

{arogers, okovalev, arum}@cs.uml.edu



Aggregating analyses in surveys and toolkits

How do we ask, *what can I easily find out about my model?*

Answer: interpretability toolkits!

AllenNLP Interpret:
A Framework for Explaining Predictions of NLP Models

Eric Wallace¹ Jens Tuyls² Junlin Wang² Sanjay Subramanian¹
Matt Gardner¹ Sameer Singh²

¹Allen Institute for Artificial Intelligence ²University of California, Irvine
ericw@allenai.org, sameer@uci.edu

Input Reduction



Input Reduction removes as many words from the input as possible without changing the model's prediction.

Original Premise: Two women are wandering along the shore drinking iced tea.

Original Hypothesis: Two women are sitting on a blanket near some rocks talking about politics

Reduced Hypothesis: politics



Takeaways

Neural models are complex, fascinating objects that we don't currently understand, but we're making strides to understand them better!

A wide variety of analysis methods have been developed, for:

- Understanding a model's behavior on specific phenomena
- Understanding what a model learns about a topic or task
- Understanding what seemingly innocuous input changes make a model fail
- Many other things, with more coming every day!

These methods can be integrated into your future NLP projects!