

How Alignment Impact Language Learning Tasks? A Visual Analytics Tool for Interpreting, Debugging, and Tuning Text Entailment Model

Abstract— With the advances of deep learning, neural network based models have reset the state-of-the-art for nearly all linguistic tasks in natural language processing. The disruptive advance also brings enormous challenges. The opaque nature of neural network model leads to hard to debug system and difficult to interpret mechanism cite. In this work, we use textual entailment task as an example to illustrate how a flexible visualization system can help NLP researchers quickly identify the potential limitation of a models and probe the model inner states for interpreting key mechanism such as attention.

Index Terms—Natural Language Processing, Interpretable Machine Learning, Visual Analytics

1 INTRODUCTION

2 RELATED WORKS

- Depth learning vis
- Vis

How to interpret classification task

- What feature / dimension / part-of-the-input should change to alter the outcome
- What Examples can you provide that is similar to the current instance

Papers:

- Interacting with Predictions: Visual Inspection of Black-box Machine Learning Models, Use partial dependency (varying along one feature and observe how prediction changes)
- Why Should I Trust You? Explaining the Predictions of Any Classifier (LIME)

3 BACKGROUND

3.1 Textual Entailment Problem

- what is the textual entailment problem
- the importance of textual entailment problem
- how easily can the visualization method extends other NLP problem

Textual entailment is one of major natural language tasks, namely translation, summarization, question

4 ANALYSIS OF THE MODEL VIA PERTURBATION

4.1 Input Perturbation

4.2 Attention Perturbation

4.3 Prediction Perturbation

5 ATTENTION

5.1 Visualize Attention

5.2 Interpret Attention Via Model Perturbation

6 PUT IT ALL TOGETHER

6.1 Where Are the Mistakes?

6.2 What Does Attention Do?

6.3 Is Attention All You Need?

7 DISCUSSION

- how the visualization approaches generalizes to other NLP task
- limitation