

CS 11-747 Neural Networks for NLP

Model Interpretation

Danish Pruthi

April 28, 2020

Why interpretability?

- **Task:** predict probability of death for patients with pneumonia
- **Why:** so that high-risk patients can be admitted, low risk patients can be treated as outpatients
- $AUC_{\text{Neural networks}} > AUC_{\text{Logistic Regression}}$
- Rule based classifier

$\text{HasAsthma}(X) \rightarrow \text{LowerRisk}(X)$



Example from Caruana et al.

Why interpretability?

- Legal reasons: uninterpretable models are banned!
 - GDPR in EU necessitates "right to explanation"
- Distribution shift: deployed model might perform poorly *in the wild*
- User adoption: users happier with explanations
- Better Human-AI interaction and control
- Debugging machine learning models

Dictionary definition

interpret verb

in·ter·pret | \in-'tər-prət , -pət\

interpreted; interpreting; interprets

Definition of *interpret*

transitive verb

1 : to explain or tell the meaning of : present in understandable terms

// *interpret* dreams

// needed help *interpreting* the results

Only if we could
understand

model.ckpt

Two broad themes

- What is the model learning?
- Can we explain the outcome in "understandable terms"?

global interpretation

local interpretation



Comparing two directions

What is the model learning?

- Input: a model M, a **(linguistic) property P**
- Output: extent to which M captures P
- Techniques: classification, regression
- Evaluation: implicit

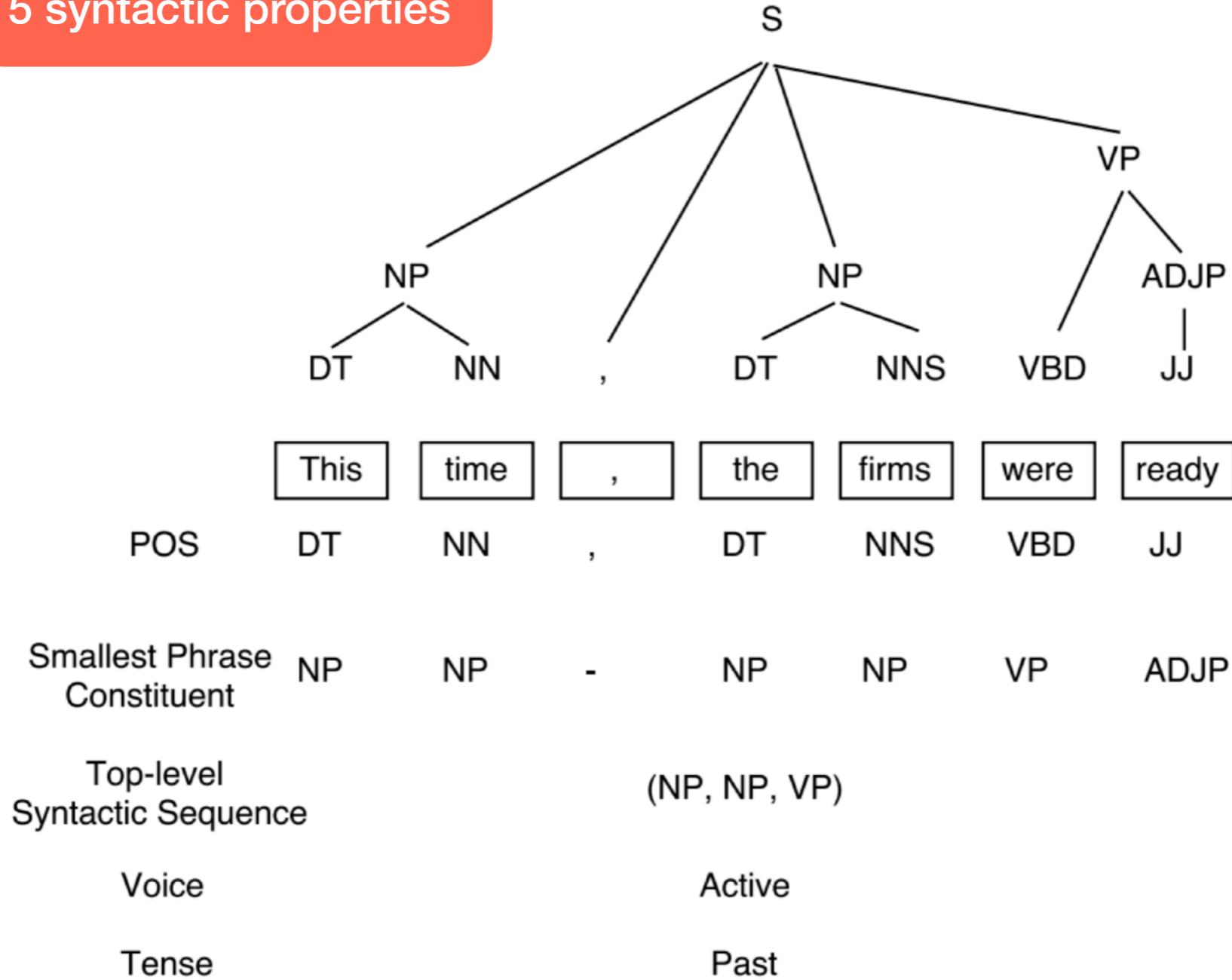
Explain the prediction

- Input: a model M, **a test example X**
- Output: an explanation E
- Techniques: varied ...
- Evaluation: complicated

**What is the model
learning?**

Source Syntax in NMT

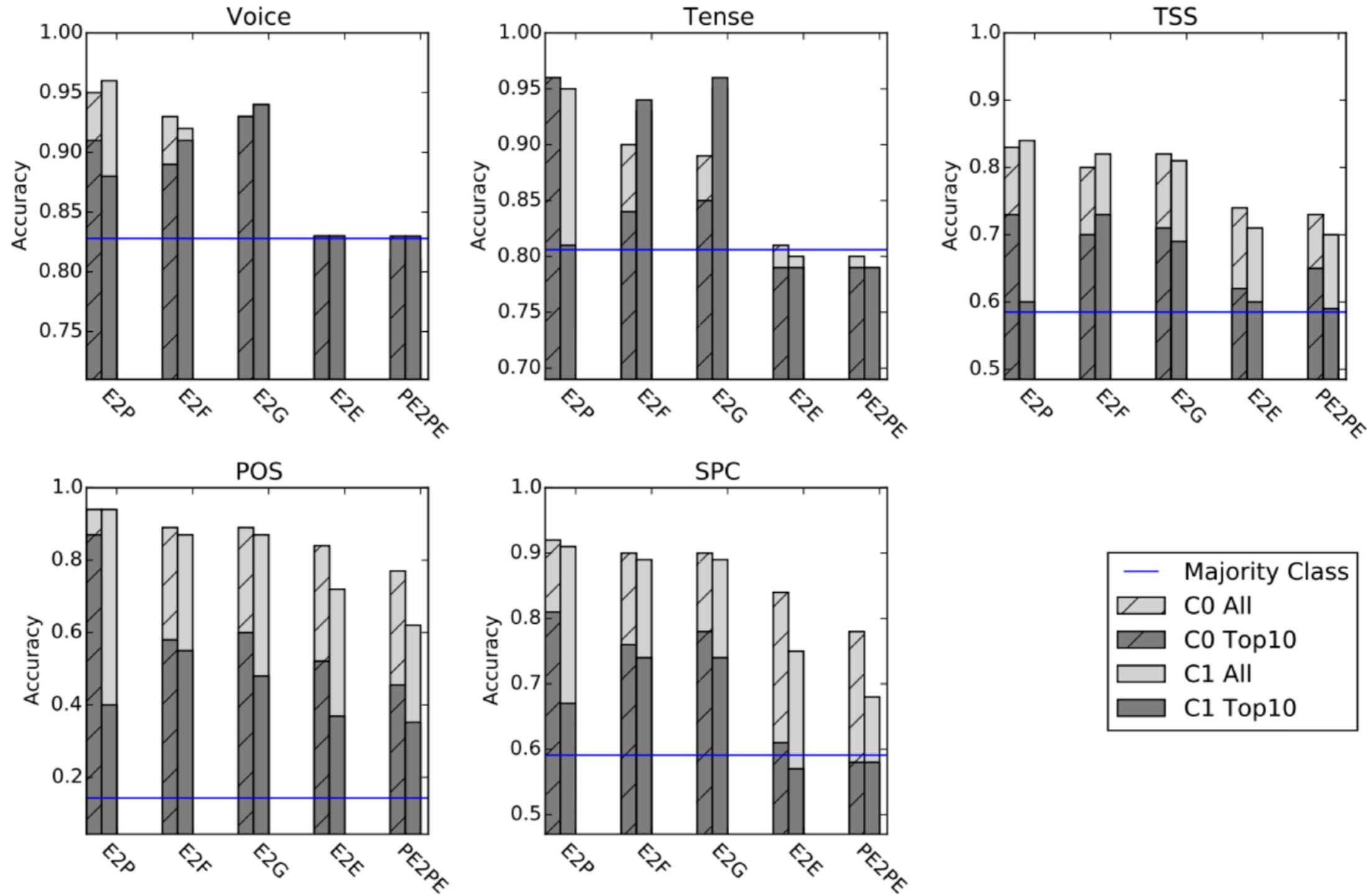
5 syntactic properties



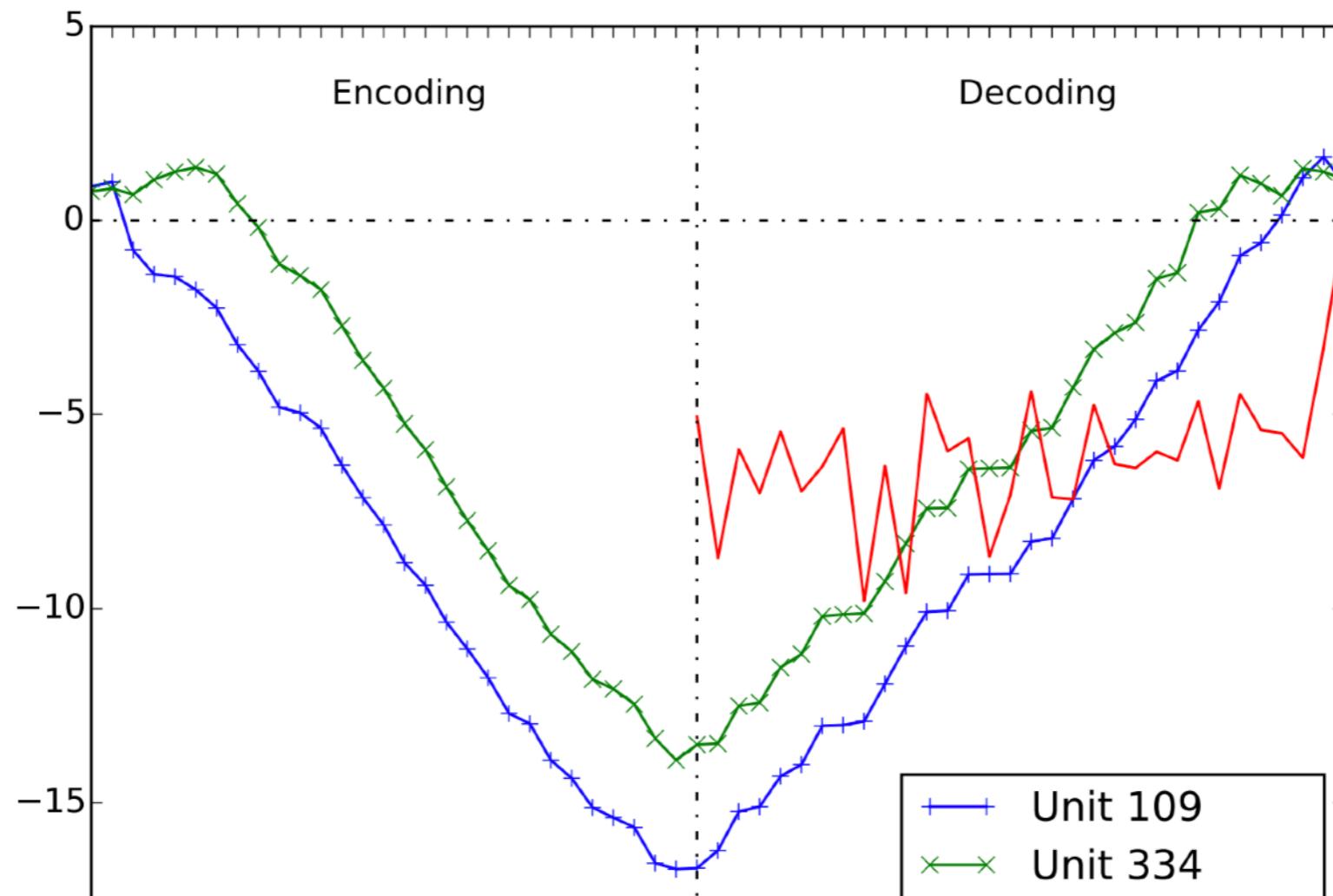
| Accuracy |
|----------|
| 82.8 |
| 92.8 |
| 82.7 |

on accuracy using the majority class baseline

Source Syntax in NMT



Why neural translations are the right length?

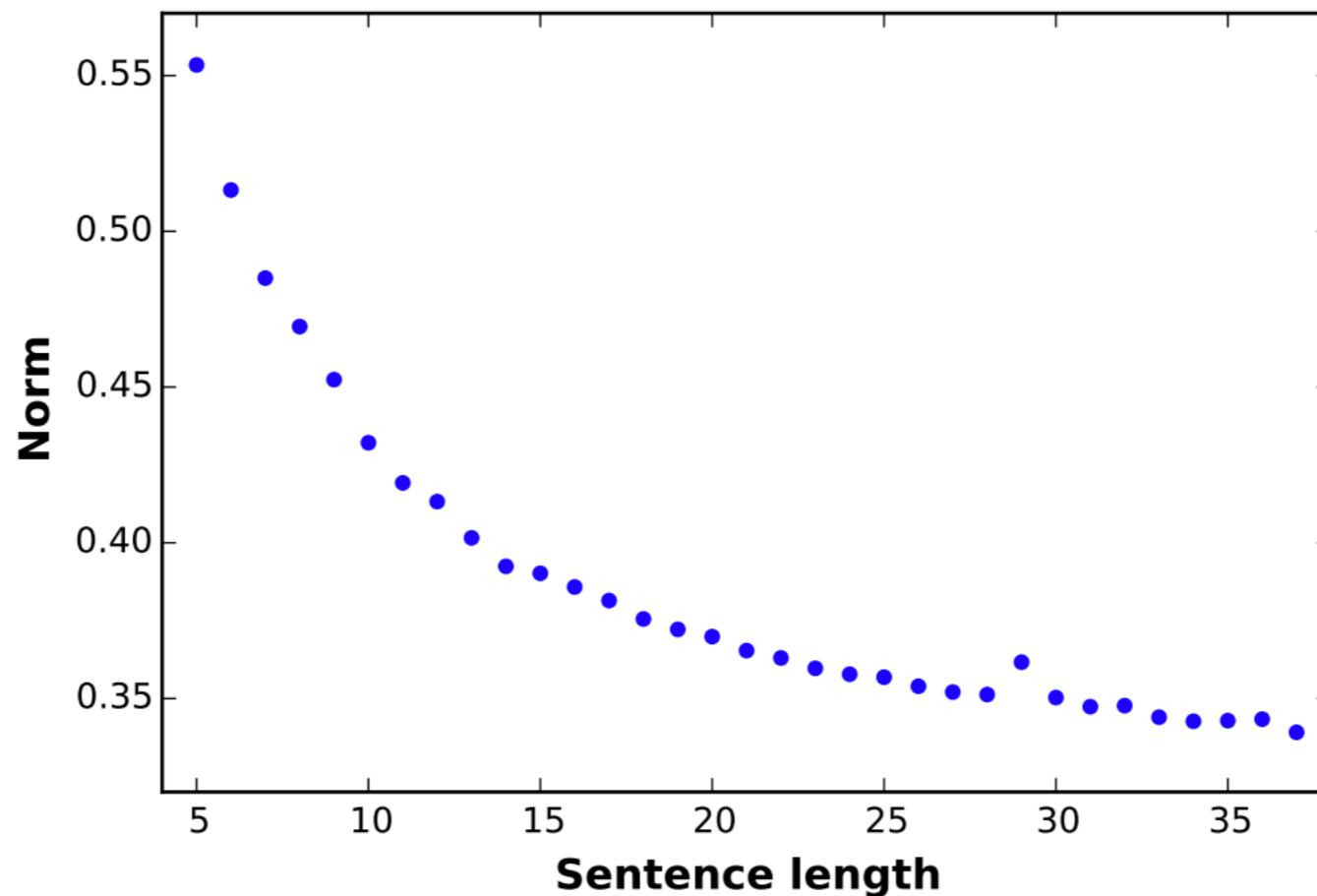


Note: LSTMs can learn to count, whereas GRUs can not do unbounded counting (Weiss et al. ACL 2018)

Fine grained analysis of sentence embeddings

- Sentence representations: word vector averaging, hidden states of the LSTM
- Auxiliary Tasks: predicting length, word order, content
- Findings:
 - hidden states of LSTM capture to a great deal length, word order and content
 - word vector averaging (CBOW) model captures content, length (!), word order (!!)

Fine grained analysis of sentence embeddings

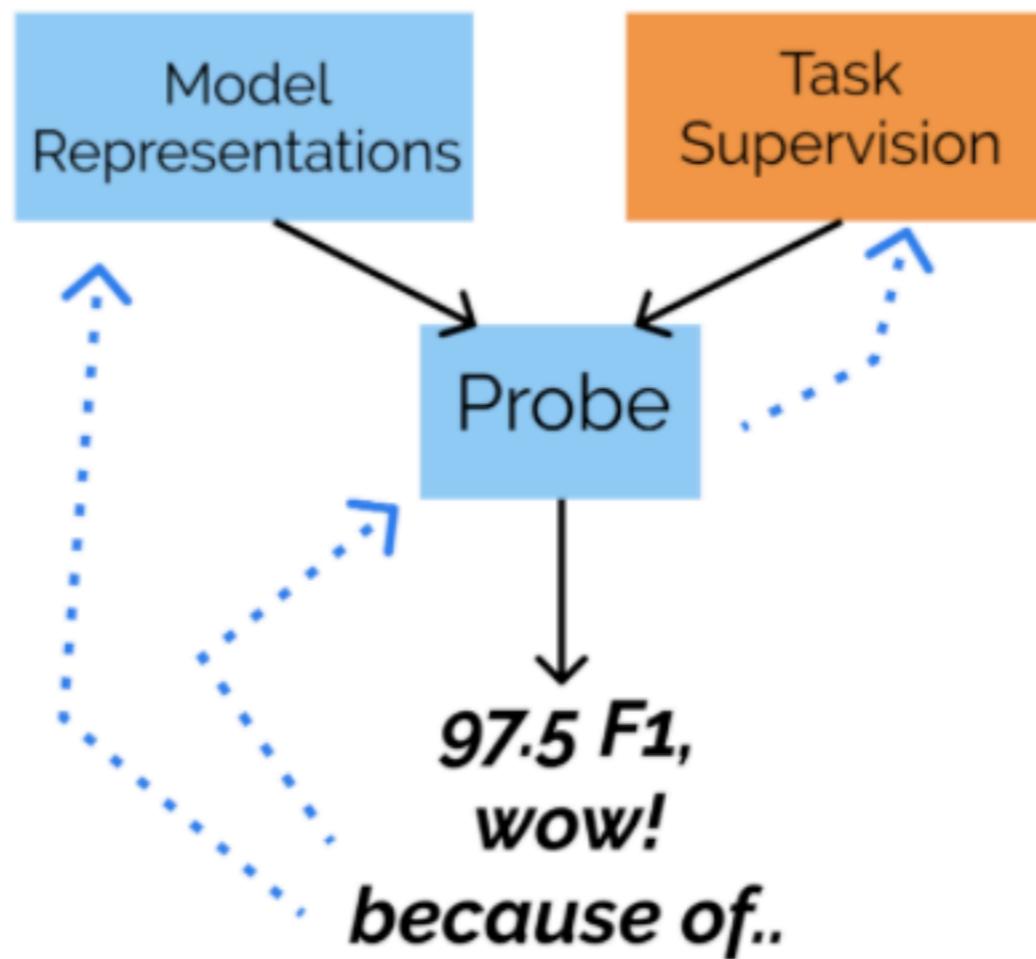


(b) Average embedding norm vs. sentence length for CBOW with an embedding size of 300.

What you can cram into a single vector: Probing sentence embeddings for linguistic properties

- "you cannot cram the meaning of a whole %&!\$# sentence into a single \$&!#* vector" – Ray Mooney
- Design 10 probing tasks: len, word content, bigram shift, tree depth, top constituency, tense, subject number, object number, semantically odd man out, coordination inversion
- Test BiLSTM last, BiLSTM max, Gated ConvNet encoder

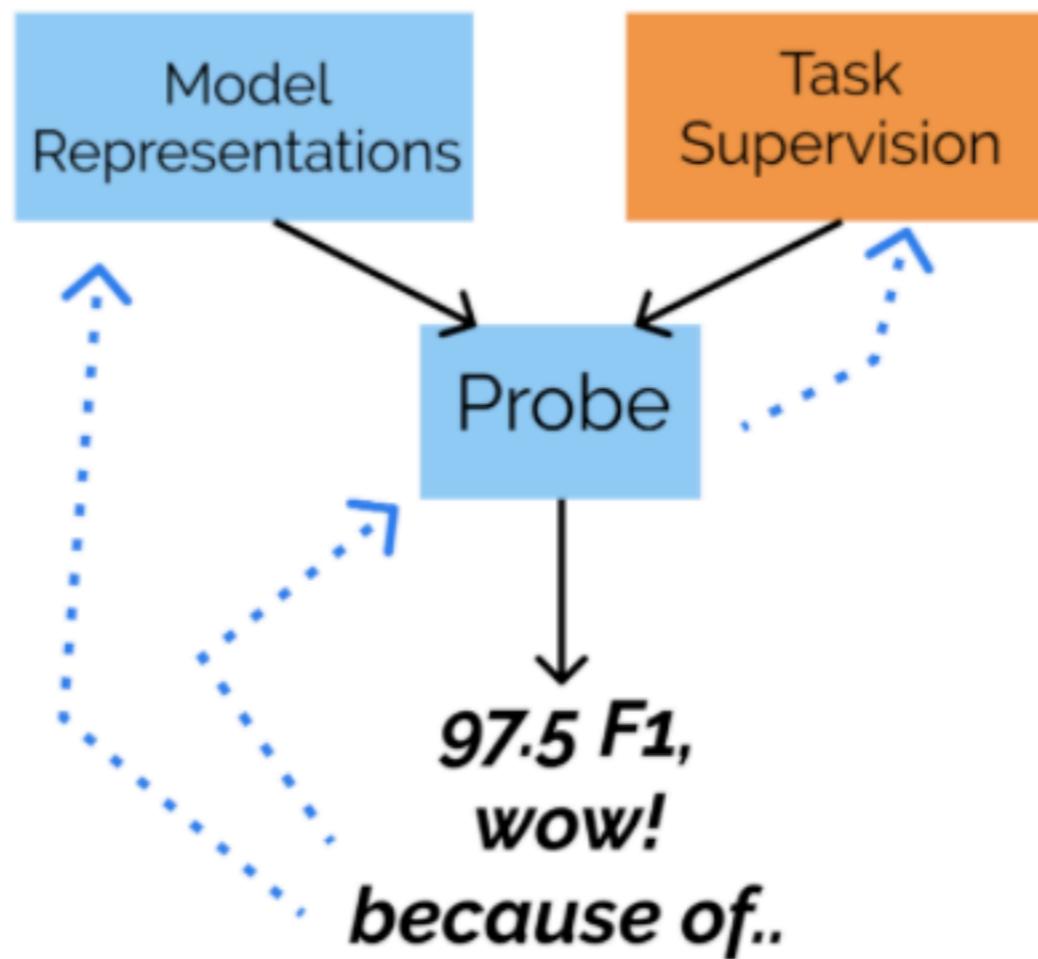
Issues with probing



Probing turns supervised tasks into tools for interpreting representations. But the use of supervision leads to the question, did I interpret the representation?

Or did my probe just learn the task itself?

Issues with probing



Probing turns supervised tasks into tools for interpreting representations. But the use of supervision leads to the question, did I interpret the representation?

Or did my probe just learn the task itself?

Minimum Description Length (MDL) Probes

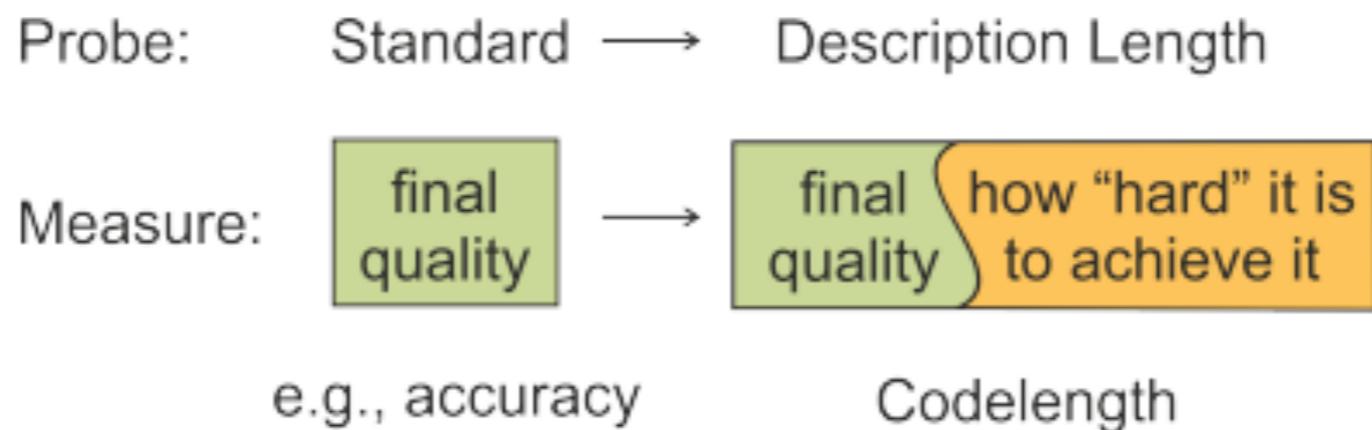


Figure 1: Illustration of the idea behind MDL probes.

- Characterizes both **probe quality** and **the amount of effort** needed to achieve it
 - More informative and stable

Summary: What is the model learning?

<https://boknilev.github.io/nlp-analysis-methods/table1.html>

Explain the prediction

How to evaluate?

Training Phase

Some $x, f(x)$ pairs



Test Phase

Input x
Predict $f(x)$



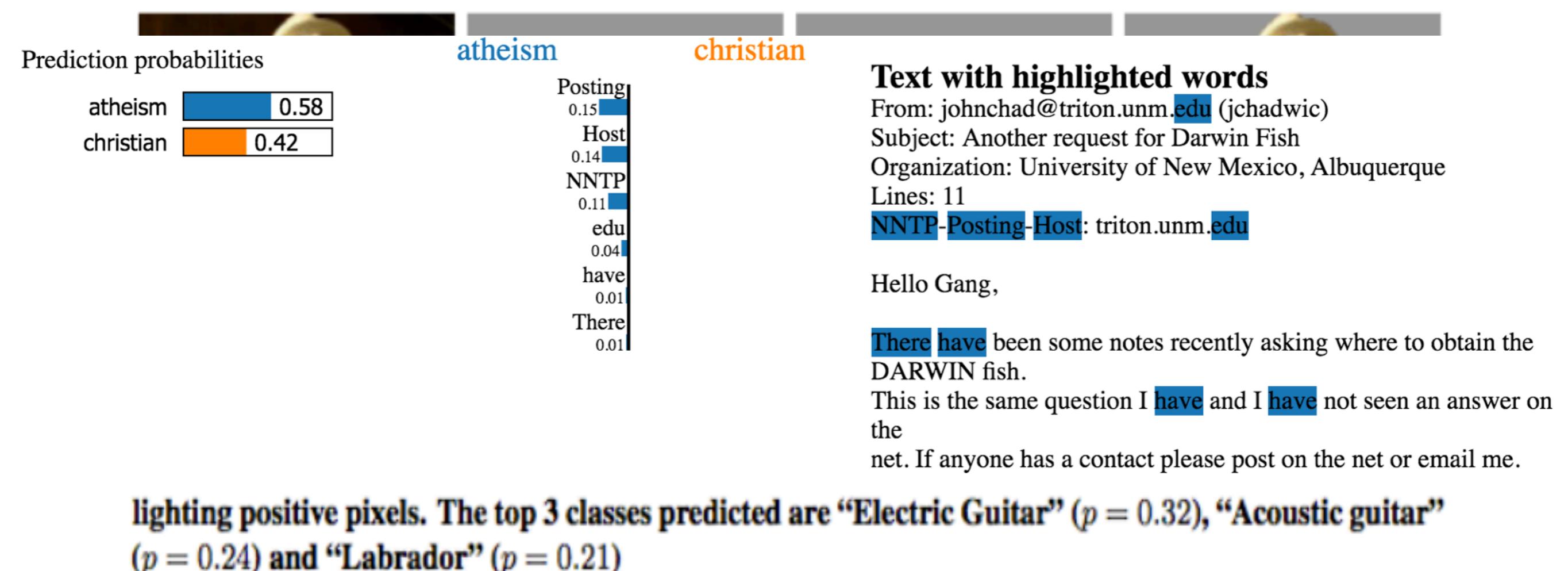
Some $x, f(x), E$ triples



Input x
Predict $f(x)$



Explanation Technique: LIME



Explanation Technique: Influence Functions

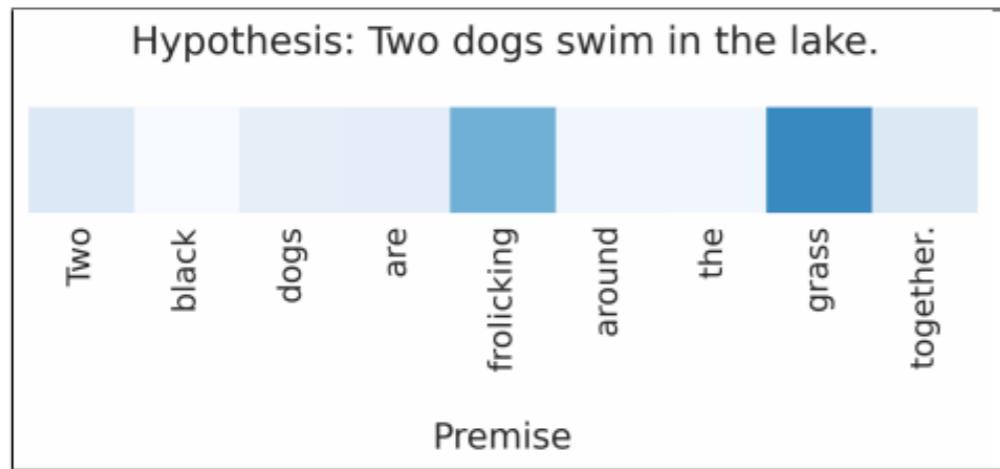
- What would happen if a given training point didn't exist?
- Retraining the network is prohibitively slow, hence approximate the effect using influence functions.



→
Most influential train images



Explanation Technique: Attention



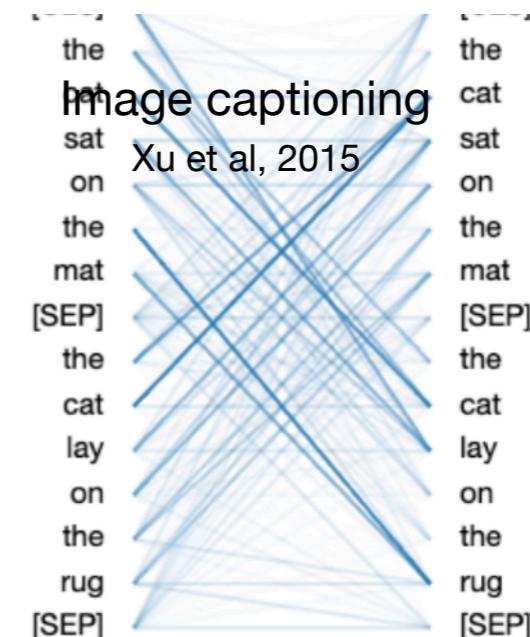
Entailment
Rocktäschel et al, 2015

why does zebras have stripes ?
what is the purpose or those stripes ?
who do they serve the zebras in the wild life ?
this provides camouflage - predator vision is such that it is usually difficult for them to see complex patterns

Document classification
Yang et al, 2016



A stop sign is on a road with a mountain in the background.



BERTViz
Vig et al, 2019

Explanation Technique: Attention

Attention is not Explanation

Sarthak Jain

Northeastern University

jain.sar@husky.neu.edu

Byron C. Wallace

Northeastern University

b.wallace@northeastern.edu

1. Attention is only mildly correlated with other importance score techniques
2. Counterfactual attention weights should yield different predictions, but they do not

Attention is not not Explanation

Sarah Wiegreffe*

School of Interactive Computing
Georgia Institute of Technology
saw@gatech.edu

Yuval Pinter*

School of Interactive Computing
Georgia Institute of Technology
uvp@gatech.edu

"Attention *might* be an explanation."

- Attention scores can provide a (plausible) explanation not the explanation.
- Attention is not explanation if you don't need it
- Agree that attention is indeed manipulable,

"this should provide pause to researchers who are looking to attention distributions for one true, faithful interpretation of the link their model has established between inputs and outputs."

Learning to Deceive with Attention-Based Explanations

Danish Pruthi[†], Mansi Gupta[‡], Bhuwan Dhingra[†], Graham Neubig[†], Zachary C. Lipton[†]

[†]Carnegie Mellon University, Pittsburgh, USA

[‡]Twitter, New York, USA

ddanish@cs.cmu.edu, mansig@twitter.com,
{bdhingra, gneubig, zlipton}@cs.cmu.edu

| Attention | Biography | Label |
|-----------|--|-----------|
| Original | Ms. X practices medicine in Memphis, TN and is affiliated ... Ms. X speaks English and Spanish. | Physician |
| Ours | Ms. X practices medicine in Memphis , TN and is affiliated ... Ms. X speaks English and Spanish. | Physician |

- Manipulated models perform better than no-attention models
- Elucidate some workarounds (what happens behind the scenes)

Explanation Techniques: gradient based importance scores

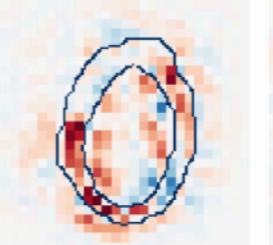
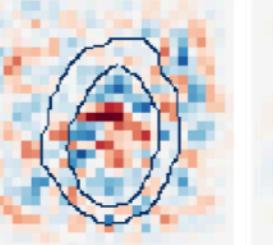
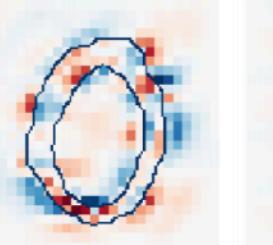
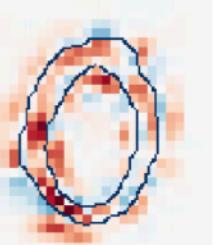
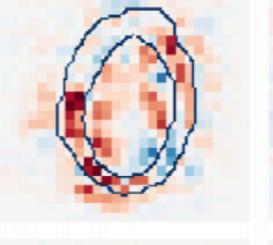
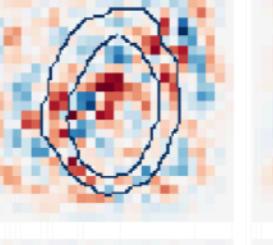
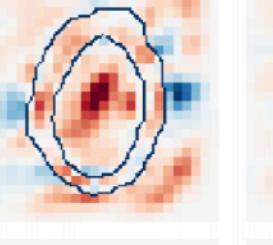
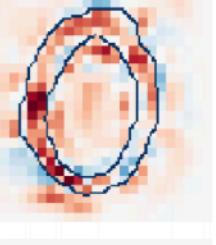
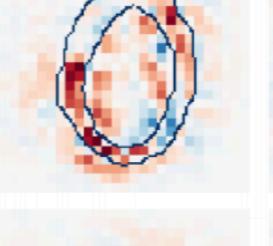
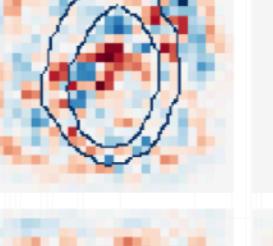
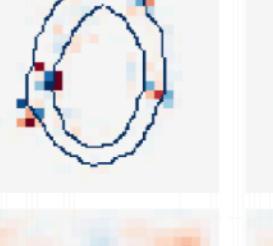
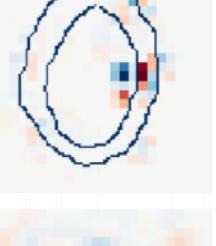
| Method | Attribution $R_i^c(x)$ | Example of attributions on MNIST | | | |
|----------------------------------|--|---|---|---|---|
| Gradient * Input | $x_i \cdot \frac{\partial S_c(x)}{\partial x_i}$ | ReLU | Tanh | Sigmoid | Softplus |
| Integrated Gradient | $(x_i - \bar{x}_i) \cdot \int_{\alpha=0}^1 \frac{\partial S_c(\tilde{x})}{\partial (\tilde{x}_i)} \Big _{\tilde{x}=\bar{x}+\alpha(x-\bar{x})} d\alpha$ |  |  |  |  |
| <u>ϵ-LRP</u> | $x_i \cdot \frac{\partial^g S_c(x)}{\partial x_i}, \quad g = \frac{f(z)}{z}$ |  |  |  |  |
| <u>DeepLIFT</u> | $(x_i - \bar{x}_i) \cdot \frac{\partial^g S_c(x)}{\partial x_i}, \quad g = \frac{f(z) - f(\bar{z})}{z - \bar{z}}$ |  |  |  |  |

Figure from Ancona et al, ICLR 2018

Explanation Technique: Extractive Rationale Generation

Key idea: find minimal span(s) of text that can (by themselves) explain the prediction

- Generator (x) outputs a probability distribution of each word being the rational
- Encoder (x) predicts the output using the snippet of text x
- Regularization to support contiguous and minimal spans

Review

the beer was n't what i expected, and i'm not sure it's "true to style", but i thought it was delicious. **a very pleasant ruby red-amber color** with a relatively brilliant finish, but a limited amount of carbonation, from the look of it. aroma is what i think an amber ale should be - a nice blend of caramel and happiness bound together.

Ratings

Look: 5 stars

Smell: 4 stars

Figure 1: An example of a review with ranking in two categories. The rationale for Look prediction is shown in bold.

Future Directions

- Need automatic methods to evaluate interpretations
- Complete the feedback loop: update the model based on explanations

Thank You!

Questions?