CS388: Natural Language Processing

Lecture 21: Interpretation





Kenton Lee talk later today in Eunsol Choi's class

P2 back soon

Administrivia



Recall: Dialogue Systems (Gunrock)



Chen ... Zhou Yu (2018)









- Fluent, on-topic generations, but system still gets into weird situations



Recall: Dialogue Systems (Blender)

Pre-train this model on >300GB of social media data (Meena/Blender)

I'm from the midwest, but I live in the south now. How about you?

Born and raised in the south

What part of the south? I was born in georgia, but raised in tennessee





Source Document (abbreviated)

The researchers examined three types of coral in reefs off the Fisheries off the coast of Fiji are protectcoast of Fiji ... The researchers found when fish were plentiful, ing coral reefs from the effects of global warming, according to a study in the jourthey would eat algae and seaweed off the corals, which appeared to leave them more resistant to the bacterium Vibrio coralliilytinal Science. cus, a bacterium associated with bleaching. The researchers suggested the algae, like warming temperatures, might render the corals' chemical defenses less effective, and the fish were protecting the coral by removing the algae.

Sacoolas, who has immunity as a diplomat's wife, w in a traffic collision ... Prime Minister Johnson was about the case while speaking to the press at a hosp ford. He said, "I hope that Anne Sacoolas will co if we can't resolve it then of course I will be raising personally with the White House."

These look great! But they're not always factual

Recall: Summarization Systems

BART	Summary
------	---------

vas involved	Boris Johnson has said he will raise the is-
s questioned	sue of US diplomat Anne Sacoolas' diplo-
oital in Wat-	matic immunity with the White House.
ome back	
ng it myself	





- care?
- Local explanations: erasure techniques
- Gradient-based methods
- Text-based explanations
- Evaluating explanations

Today

Interpreting neural networks: what does this mean and why should we

Interpreting Neural Networks



- Sentiment w/LSTMs



- Looking at individual neurons usually doesn't tell us much
- which ones actually mean something?

Interpreting Neural Networks

Neural models have complex behavior. How can we understand them?

Sentiment w/BERT: there are hundreds of attention computations...



- Sentiment w/DANs:



- Tells us how these words combine
- How do we know why a neural network model made the prediction it made?

Interpreting Neural Networks

Neural models have complex behavior. How can we understand them?

Ground Truth DAN

negative negative positive positive negative negative positive negative

Left side: predictions the model makes on individual words

lyyer et al. (2015)





- Trust: if we see that models are behaving in human-like ways and making
- tell us that x causes y? Not necessarily, but it might be helpful to know
- **Fairness:** ensure that predictions are non-discriminatory

Why explanations?

human-like mistakes, we might be more likely to trust them and deploy them

• Causality: if our classifier predicts class y because of input feature x, does that

Informativeness: more information may be useful (e.g., predicting a disease diagnosis isn't that useful without knowing more about the patient's situation)

Lipton (2016)





- they do (e.g., a decision tree with <10 nodes)
- Explanations of more complex models
 - Local explanations: highlight what led to this classification decision. predicted a different class) — focus of this lecture
 - **Text explanations:** describe the model's behavior in language
 - understand more about how our model works

Why explanations?

Some models are naturally transparent: we can understand why they do what

(Counterfactual: if these features were different, the model would've

Model probing: auxiliary tasks, challenge sets, adversarial examples to

Lipton (2016); Belinkov and Glass (2018)



Local Explanations

(which parts of the input were responsible for the model's prediction on this particular data point?)



Similar to a DAN model, but (1) extra BiLSTM layer; (2) attention layer instead of just a sum Jain and Wallace (2019)

good









the movie was not

- Attention places most mass on good did the model ignore not? What if we removed *not* from the input? Jain and Wallace (2019)

Attention Analysis

Negative

good





An explanation could help us answer counterfactual questions: if the input were x' instead of x, what would the output be?

that movie was not great, in fact it was terrible !

that movie was not , in fact it was terrible !

that movie was _____ great, in fact it was _____ !

Attention can't necessarily help us answer this!

Local Explanations

Model

+



that movie was not great, in fact it was terrible ! movie was not great, in fact it was terrible ! that _____ was not great, in fact it was terrible ! that movie _____ not great, in fact it was terrible ! that movie was _____ great, in fact it was terrible ! that movie was not _____, in fact it was terrible !

Delete each word one by and one and see how prediction prob changes

- prob = 0.97- prob = 0.97- prob = 0.98- prob = 0.97- prob = 0.8- prob = 0.99







the output

that movie was not great, in fact it was terrible !

- made it more negative)
- Will this work well?
 - Inputs are now unnatural, model may behave in "weird" ways

Output: highlights of the input based on how strongly each word affects

In not contributed to predicting the negative class (removing it made it less) negative), great contributed to predicting the positive class (removing it

Saturation: if there are two features that each contribute to negative predictions, removing each one individually may not do much







- Locally-interpretable, model-agnostic explanations (LIME)
- at once
 - words with it)
 - More scalable to complex settings

LIME

Similar to erasure method, but we're going to delete collections of things

Can lead to more realistic input (although people often just delete

Ribeiro et al. (2016)









Components

Break input into components (for text: could use words, phrases, sentences, ...)

https://www.oreilly.com/learning/introduction-to-localinterpretable-model-agnostic-explanations-lime

LIME



Check predictions on > Now we have model subsets of those predictions on

perturbed examples









LIME (cont'd)

- This is what the model is doing on perturbed examples of the input
- Now we train a classifier to predict the model's behavior based on what subset of the input it sees
- The weights of that classifier tell us which parts of the input are important





The movie is mediocre, maybe even bad.

The movie is mediocre, maybe even bad. The movie is mediocre, maybe even bad. The movie is mediocre, maybe even bad.

The movie is mediocre, maybe even bad.

The movie is mediocre, maybe even bad.

The movie is mediocre, maybe even bad.

LIME (cont'd)

This secondary classifier's weights now give us highlights on the input

Negative 99.8%

Negative 98.0%

Negative 98.7%

Positive 63.4%

Positive 74.5%

Negative 97.9%

Wallace, Gardner, Singh Interpretability Tutorial at EMNLP 2020





- to train? etc.
- Expensive to call the model all these times
- Linear assumption about interactions may not be reliable

Problems with LIME

Lots of moving parts here: what perturbations to use? what model



Problems with LIME

Problem: fully removing pieces unnatural

LIME/erasure zeroes out certain features

Alternative approach: look at what this perturbation does locally right around the data point using gradients

Problem: fully removing pieces of the input may cause it to be very

data manifold (points we observe in practice)



Learning a model

Compute derivative of score with respect to weights: how can changing weights improve score of correct class?

score = weights * features (or an NN)

> **Gradient-based** Explanations Compute derivative of score with respect to *features*: how can changing *features* improve score of correct class?



Originally used for images

 S_c = score of class c I_0 = current image $w = \frac{\partial S_c}{\partial J} \Big|$

- change in prediction
- up to get the importance of that word



Higher gradient magnitude = small change in pixels leads to large

For words: "pixels" are coordinates of each word's vector, sum these Simonyan et al. (2013)

















Simonyan et al. (2013)





Integrated Gradients

- would. Gradient-based method says neither is important
- Integrated gradients: compute gradients along a path from the origin to the current data point, aggregate these to learn feature importance
- Intermediate points can reveal new info about features

Suppose you have prediction = A OR B for features A and B. Changing either feature doesn't change the prediction, but changing both







IntegratedGrads_i^{approx}(x) ::=
$$(x_i - x'_i)$$

Scale by total distance

better

steps along the way

Integrated Gradients

 $(x_{i}) \times \sum_{k=1}^{m} \frac{\partial F(x' + \frac{k}{m} \times (x - x')))}{\partial x_{i}} \times \frac{1}{m}$

Compute gradient at the *k*th point along the way w.r.t. the ith feature

Average over the m steps

 x'_i = "baseline" — all PAD or MASK tokens (MASK usually works)

Can be expensive: requires calling forward() and backward() at m

Sundararajan et al. (2017)





Integrated Gradients



Question type classification task:

how many townships have a population above 50? [prediction: NUMERIC] how many athletes are not ranked ? [prediction: NUMERIC] what is the total number of points scored ? [prediction: NUMERIC] which film was before the audacity of democracy ? [prediction: STRING] which year did she work on the most films ? [prediction: DATETIME] what year was the last school established ? [prediction: DATETIME] did charles oakley play more minutes than robert parish? [prediction: YESNO]

what is the difference in population between fora and masilo [prediction: NUMERIC] when did ed sheeran get his first number one of the year ? [prediction: DATETIME]

Sundararajan et al. (2017)





Comparison



(Answer = Stanford University)

Question: Where did the Broncos practice for the Super Bowl? **Passage:** The Panthers used the San Jose State practice facility and stayed at the San Jose Marriott. The Broncos practiced at Stanford University and stayed at the Santa Clara Marriott.

Question: Where did the Broncos practice for the Super Bowl? **Passage:** The Panthers used the San Jose State practice facility and stayed at the San Jose Marriott. The Broncos practiced at Stanford University and stayed at the Santa Clara Marriott.

(a) Integrated Gradient (Sundararajan et al., 2017).

Are these good explanations?

(d) Erasure exact search optima.

De Cao et al. (2020)



Text Explanations



Explanations of Bird Classification

Laysan Albatross



and white belly.

yellow beak, and white belly.

Laysan Albatross Description: This is a large bird with a white neck and a black back in the water. and white belly. neck and black back.

Are these features really what the model used?

- **Description:** This is a large flying bird with black wings and a white belly.
- Class Definition: The Laysan Albatross is a large seabird with a hooked yellow beak, black back
- Visual Explanation: This is a *Laysan Albatross* because this bird has a large wingspan, hooked
- Class Definition: The Laysan Albatross is a large seabird with a hooked yellow beak, black back
- Visual Explanation: This is a *Laysan Albatross* because this bird has a hooked yellow beak white

- An explanation should be relevant to both the class and the image

Hendricks et al. (2016)











Explanations of Bird Classification

This is a cardinal because ...



- decision-making
- explanations!

Are these features really what the model used? The decoder looks at the image, but what it reports may not truly reflect the model's

More likely to produce plausible (look good to humans) but unfaithful

Hendricks et al. (2016)



Premise: An adult dressed in black holds a stick. Hypothesis: An adult is walking away, empty-handed. Label: contradiction Explanation: Holds a stick implies using hands so it is not empty-handed.

How do we use this information? If we produce a network to predict it, does that make it an actual explanation of what's happening?

Explanations of NLI

Camburu et al. (2019)

Explanations of NLI

Information from f is fed into the explanation LSTM, but no constraint that this must be used. Different coordinates from f could predict label and explanations

Evaluating Explanations

Faithfulness vs. Plausibility

- Suppose our model is a bag-of-words model with the following:
 - the = -1, movie = -1, good = +3, bad =0
 - the movie was good prediction score=+1
 - the movie was bad prediction score=-2
- Suppose explanation returned by LIME is:
 - the movie was good
 - the movie was bad
- Is this a "correct" explanation?

Plausible explanation: matches what a human would do

the movie was **good** the movie was **bad**

Maybe useful to explain a task to a human, but it's not what the model is really doing!

Faithful explanation: actually reflects the behavior of the model

the movie was good

- and Use Interpretable Models Instead

Faithfulness vs. Plausibility

the movie was bad

We usually prefer faithful explanations; non-faithful explanations are actually deceiving us about what our models are doing!

Rudin: Stop Explaining Black Box Models for High-Stakes Decisions

- Nguyen (2018): delete words from the input and see how quickly the model flips its prediction?
 - Downside: not a "real" use case
- Hase and Bansal (2020): counterfactual simulatability: user should be able to predict what the model would do in another situation
 - Hard to evaluate

Evaluating Explanations

Evaluating Explanations

C I, like others was very excited to read this book. I thought it would show another side to how the Tate family dealt with t he murder of thier daughter Sharon. I didn't have to read mu ch to realize however that the book is was not going to be w hat I expected. It is full of added dialog and assumptions. It makes it hard to tell where the truth ends and the embellish ments begin. It reads more like fan fiction than a true accou nt of this family's tragedy. I did enjoy looking at the early pic tures of Sharon that I had never seen before but they were hardly worth the price of the book.

- Al provides both an explanation for its prediction (blue) and also a possible counterargument (red)

62.7% CONFIDENT 100

Human is trying to label the sentiment. The AI provides its prediction to try to help. Does the human-AI team beat human/AI on their own?

Do these explanations help the human? Slightly, but Al is still better No positive results on "human-Al teaming" with explanations Bansal et al. (2020)

AllenNLP Interpret: https://allennlp.org/interpret

Captum (Facebook): https://captum.ai/

LIT (Google): https://ai.googleblog.com/2020/11/the-language-interpretability-tool-lit.html

Packages

- Many other ways to do explanation:
 - Probing tasks: we looked at these for ELMo, do vectors capture information about part-of-speech tags?
 - Diagnostic test sets ("unit tests" for models)
 - Building models that are explicitly interpretable (decision trees)
- Input attribution methods can be useful for visualization (consider) using these for your final project!)

Wallace, Gardner, Singh Interpretability Tutorial at EMNLP 2020

