

Today
Interpreting neural networks: what does this mean and why should we care?
Local explanations: erasure techniques
 Gradient-based methods
 Evaluating explanations



Why explanations?	Why explanations?	
 Trust: if we see that models are behaving in human-like ways and making human-like mistakes, we might be more likely to trust them and deploy them 	Some models are naturally transparent: we can understand why they do what they do (e.g., a decision tree with <10 nodes)	
Causality: if our classifier predicts class y because of input feature x, does that tell us that x causes y? Not necessarily, but it might be helpful to know	 Explanations of more complex models Local explanations: highlight what led to this classification decision. 	
 Informativeness: more information may be useful (e.g., predicting a disease diagnosis isn't that useful without knowing more about the patient's situation) 	 (Counterfactual: if these features were different, the model would ve predicted a different class) — focus of this lecture Text explanations: describe the model's behavior in language 	
Fairness: ensure that predictions are non-discriminatory	Model probing: auxiliary tasks, challenge sets, adversarial examples to	

understand more about how our model works

Lipton (2016)

Lipton (2016); Belinkov and Glass (2018)





 Delete each word one by and one and see how prediction prob chang 		
that movie was not great , in fact	it was terrible ! — prob = 0.97	
movie was not great , in fact	t was terrible ! — prob = 0.97	
that was not great , in fact	t was terrible ! — prob = 0.98	
that movie not great, in fact	it was terrible ! — prob = 0.97	
that movie was great, in fact	it was terrible ! — prob = 0.8	
that movie was not, in fact i	t was terrible ! — prob = 0.99	

Erasure Method

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

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

- not contributed to predicting the negative class (removing it made it less negative), great contributed to predicting the positive class (removing it made it more negative)
- Will this work well?

- Inputs are now unnatural, model may behave in "weird" ways
- Saturation: if there are two features that each contribute to negative predictions, removing each one individually may not do much





Problems with LIME	
Lots of moving parts here: what perturbations to use? what model to train? etc.	
Expensive to call the model all these times	Cradiant based Mathads
Linear assumption about interactions may not be reliable	Gradient-based Methods











No positive results on	"human-Al teaming"	with explanations	Bansal et al. (2020)
------------------------	--------------------	-------------------	----------------------



Ongoing Conversation Lots of ongoing research: How do we interpret explanations? How do *users* interpret our explanations? How should *automated systems* make use of explanations? Still a growing area

Packages

- AllenNLP Interpret: https://allennlp.org/interpret
- Captum (Facebook): https://captum.ai/

۲

- LIT (Google): https://ai.googleblog.com/2020/11/the-language-interpretability-tool-lit.html
- Various pros and cons to the different frameworks

٢	Takeaways
Many other w	rays to do explanation:
Probing task tags?	s: do vectors capture information about part-of-speech
Diagnostic t	est sets ("unit tests" for models)
Building mc	dels that are explicitly interpretable (decision trees)
	Wallace, Gardner, Singh Interpretability Tutorial at EMNLP 2020