Interpreting neural network models

Interpretability is the degree to which a human can understand the cause of a decision.

Why interpret?

Local interpretation: Explain a model's prediction on a single input example.

Global interpretation: Explain how a model makes predictions across all examples.

Example model: Predicting contents of an image



Example model: Predicting gene expression from DNA sequence



Outline

- Model-based interpretation
- Saliency map
- Understanding interactions

Outline

Model-based interpretation

- Saliency map
- Understanding interactions

Understanding simple neural networks



Inspecting individual nodes



https://medium.com/dataseries/visualizing-the-feature-maps-and-filters-by-convolutional-neural-networks-e1462340518e

Inspecting individual nodes



Inspecting individual nodes





What impact does each node have on the final prediction?



Visualizing a hidden node

Idea: Use gradient descent to find the input that maximizes the activation of a given node.

Visualizing a hidden node

Idea: Use gradient descent to find the input that maximizes the activation of a given node.



Maximize

Visualizing a hidden node

Idea: Use gradient descent to find the input that maximizes the activation of a given node.



Maximize

Minimize

Questions about model-based interpretation

Outline

- Model-based interpretation
- Saliency map
 - Understanding interactions

Saliency map

Saliency map: An estimate of the importance of each input feature.



(a) Sheep - 26%, Cow - 17% (b) Importance map of '*sheep*' (c) Importance map of '*cow*'



(d) Bird - 100%, Person - 39% (e) Importance map of 'bird' (f) Importance map of 'person'

Interpretation through perturbation



A model's gradient gives the saliency of each input (backpropagation)



A model's gradient gives the saliency of each input (backpropagation)



Problem: Saturating nonlinearities.

There are many solutions: Integrated Gradients, DeepLIFT, DeepSHAP, ...

Integrated gradients



Integrated gradients

IG(input, base) ::= (input - base) * $\int_{0^{-1}} \nabla F(\alpha * input + (1 - \alpha) * base) d\alpha$

Original image



Integrated Gradients



DeepLIFT

 $x_i \times \frac{\partial Y}{\partial x_i}$

DeepLIFT

$$x_i \times \frac{\partial Y}{\partial x_i} \to (x_i - x_i^{baseline}) \times \frac{Y - Y^{baseline}}{x_i - x^{baseline}}$$

DeepLIFT

$$x_i \times \frac{\partial Y}{\partial x_i} \to (x_i - x_i^{baseline}) \times \frac{Y - Y^{baseline}}{x_i - x^{baseline}}$$

y = ReLU(x) = max(0, x)



Saliency: Local to global



Saliency: Local to global



Saliency: Local to global



DOI: 10.1038/s41588-021-00782-6

Outline

- Model-based interpretation
- Saliency map
- Understanding interactions

Multiple perturbation for identifying interactions



Multiple perturbation for identifying interactions



Limitations of interpretability

Attention

Motivating problem: translation

FRENCH ← ENGLISH	FRENCH ← ENGLISH	FRENCH
Un sourire coûte moins cher que ×	Un sourire coûte moins cher que ×	Un sourire coûte moins cher que ×
l'électricité, mais donne autant	l'électricité, mais donne autant	l'électricité, mais donne autant
de lumière	de lumière	de lumière
A smile costs less expensive than ☆	A smile costs cheaper than	A smile costs less than electricity, ☆
electricity, but gives as many light	electricity, but gives as much light	but gives as much light
✓ offline	🕑 offline 🔹 🔹 🔹 🔹	•) โ :
Idea #1: Recurrent neural network



Idea #1: Recurrent neural network



Idea #1: Recurrent neural network



Idea #2: Give the network access to more hidden states



parameters

Encode $h_i = f(x_i, h_{i-1})$ Context $c = \bigwedge h_{T-M:T}$ Decode $s_j = f(s_{j-1}, y_{j-1}, c)$ $(H+Y+MH)\times H\times L$ Generate $y_j = g(y_{j-1}, s_j, c)$

Aside: What makes one architecture better than another?

- Complexity vs. generalizability tradeoff.
- Inductive bias.
- Computational efficiency.

Idea #3: Do some computation for each input, then average



parameters

Encode $h_i = f(x_i, h_{i-1})$

Decode
$$s_j = f(s_{j-1}, y_{j-1}, c)$$

$$\sum_{i=1}^T f(s_{j-1}, y_{j-1}, h_i) \quad (H+Y+H) \times H \times L$$



Function	Traditional Encoder- Decoder
Encode	$h_i = f(x_i, h_{i-1})$

Context $c = h_T$

Decode $s_j = f(s_{j-1}, y_{j-1}, c)$ Generate $y_j = g(y_{j-1}, s_j, c)$

Chaudhari, S., Mithal, V., Polatkan, G., & Ramanath, R. (2021). An attentive survey of attention models. ACM *Transactions on Intelligent Systems and Technology (TIST)*, *12*(5), 1-32.



Function	Traditional Encoder-	Encoder-Decoder
	Decoder	with Attention
Encode	$h_i = f(x_i, h_{i-1})$	$h_i = f(x_i, h_{i-1})$
Context	$c = h_T$	

Decode	$s_j = f(s_{j-1}, y_{j-1}, c)$	$s_j = f(s_{j-1}, y_{j-1}, c_j)$
Generate	$y_j = g(y_{j-1}, s_j, c)$	$y_j = g(y_{j-1}, s_j, c_j)$

Chaudhari, S., Mithal, V., Polatkan, G., & Ramanath, R. (2021). An attentive survey of attention models. *ACM Transactions on Intelligent Systems and Technology (TIST)*, *12*(5), 1-32.



Function	Traditional Encoder-	Encoder-Decoder
	Decoder	with Attention
Encode	$h_i = f(x_i, h_{i-1})$	$h_i = f(x_i, h_{i-1})$
Context	$c = h_T$	

$$\alpha_{ij} = p(e_{ij})$$
$$e_{ij} = a(s_{j-1}, h_i)$$

Decode $s_j = f(s_{j-1}, y_{j-1}, c)$ Generate $y_j = g(y_{j-1}, s_j, c)$ y

$$y_j = g(y_{j-1}, s_j, c_j)$$



Function	Traditional Encoder-	Encoder-Decoder
	Decoder	with Attention
Encode	$h_i = f(x_i, h_{i-1})$	$h_i = f(x_i, h_{i-1})$
Context	$c = h_T$	$c_j = \sum_{i=1}^T lpha_{ij} h_i$
		$\alpha_{ij} = p(e_{ij})$
		$e_{ij} = a(s_{j-1}, h_i)$
Decode	$s_j = f(s_{j-1}, y_{j-1}, c)$	$s_j = f(s_{j-1}, y_{j-1}, c_j)$
Generate	$y_j = g(y_{j-1}, s_j, c)$	$y_j = g(y_{j-1}, s_j, c_j)$

Chaudhari, S., Mithal, V., Polatkan, G., & Ramanath, R. (2021). An attentive survey of attention models. *ACM Transactions on Intelligent Systems and Technology (TIST)*, *12*(5), 1-32.



Function	Encoder-Decoder with Attention	# parameters
Encode	$h_i = f(x_i, h_{i-1})$	(X+H)×H×L
Context	$c_j = \sum_{i=1} \alpha_{ij} h_i$	0
	$\alpha_{ij} = p(e_{ij})$	0
	$e_{ij} = a(s_{j-1}, h_i)$	(n+n)×n×L
Decode Generate	$s_j = f(s_{j-1}, y_{j-1}, c_j)$ $y_j = g(y_{j-1}, s_j, c_j)$	(H+Y+H)×H×L

Chaudhari, S., Mithal, V., Polatkan, G., & Ramanath, R. (2021). An attentive survey of attention models. *ACM Transactions on Intelligent Systems and Technology (TIST)*, *12*(5), 1-32.

Outline

- Attention variants
- Self-attention and transformers
- Applications





Function	Equation
similarity	$a(k_i,q) = sim(k_i,q)$
dot product	$a(k_i,q) = q^T k_i$
scaled dot product	$a(k_i, q) = \frac{q^T k_i}{\sqrt{d_k}}$



Chaudhari, S., Mithal, V., Polatkan, G., & Ramanath, R. (2021). An attentive survey of attention models. ACM *Transactions on Intelligent Systems and Technology (TIST)*, *12*(5), 1-32.



Function	Equation
similarity	$a(k_i,q) = sim(k_i,q)$
dot product	$a(k_i,q) = q^T k_i$
scaled dot product	$a(k_i, q) = \frac{q^T k_i}{\sqrt{d_k}}$
general	$a(k_i,q) = q^T \tilde{W} k_i$
biased general	$a(k_i, q) = k_i(Wq + b)$







Function	Equation
similarity	$a(k_i,q) = sim(k_i,q)$
dot product	$a(k_i,q) = q^T k_i$
scaled dot product	$a(k_i, q) = \frac{q^T k_i}{\sqrt{d_k}}$
general	$a(k_i,q) = q^T \tilde{W} k_i$
biased general	$a(k_i, q) = k_i(Wq + b)$
activated general	$a(k_i, q) = act(q^T W k_i + b)$
	T^{-} (I () I
deep	$a(k_i, q) = w_{imp}^I E^{(L-1)} + b^L$
	$E^{(l)} = act(W_l E^{(l-1)} + b^l)$
	$E^{(1)} = act(W_1k_i + W_0q) + b^l$

Hard vs. soft attention

Outline

- Attention variants
- Self-attention and transformers
- Applications

Self-attention





How to represent relative position:

• Idea #1: One-hot encoding of distance.





How to represent relative position:

- Idea #1: One-hot encoding of distance.
- Idea #2: Binary encoding.





How to represent relative position:

- Idea #1: One-hot encoding of distance.
- Idea #2: Binary encoding.
- Idea #3: Sinusoids of different frequencies.





Transformer



Vaswani et al, 2017

Outline

- Attention variants
- Self-attention and transformers
- Applications

Self-supervised pre-training of language models

Text Corpus

Nothing is impossible. Even the word impossible says I'm possible Task: Predict from past Nothing Nothing is Nothing is impossible

Predicting gene expression



Predicting gene expression



Avsec, Ž., Agarwal, V., Visentin, D. *et al.* Effective gene expression prediction from sequence by integrating long-range interactions. *Nat Methods* **18**, 1196–1203 (2021). <u>https://doi.org/10.1038/s41592-021-01252-x</u>

Attention can be interpretable



Questions and future directions: Attention

- How can we trade off complexity, generalizability and computation?
- In each application, what attention architectures are effective?
- In what settings are attention weights provide a useful explanation?

Questions and future directions: Interpretation

- Do interpretations report spurious associations present in the training set?
- How can interpretation account for un-identifiability of machine learning?
- What interpretation approaches are relevant for a given goal?
- How can we aggregate gradient-based saliency into a global (or large-scale local) interpretation?

Questions and future directions: Interpretation

- Do interpretations report spurious associations present in the training set?
- How can interpretation account for un-identifiability of machine learning?
- What interpretation approaches are relevant for a given goal?
- How can we aggregate gradient-based saliency into a global (or large-scale local) interpretation?

Speaking

- Content
- Visual aids
 - Slides
 - Schematics
 - Data figures
- · Delivery

Presentations should have a main thesis
Err on the side of too much background

Re-engage the audience using a home slide

Outline

- Introduction
- Methods
- Results
- Conclusion

Machine learning methods for the genotype-phenotype relationship, gene regulation and epigenomics



Text should be in full sentences

Text should be in full sentences

"Method A crash rate too high"

Slide titles: Use a full sentence explaining the main point of the slide.

- Notation, acronyms: Re-introduce every time you use it.
- Animate in each slide element as you present it.
- Make every slide element legible (font size 18+)

Reduce text by translating it to schematics

Reduce text by translating it to schematics

A class of methods known as *semi-automated genome* annotation (SAGA) algorithms are widely used to perform such integrative modeling of diverse genomics data sets. These algorithms take as input a collection of genomics data sets from a particular cell type. They output (1) a set of integer *state labels*, such that each state label putatively corresponds to a type of genomic activity (such as active promoter, active transcription or repressed region), and (2) a partition of the genome and annotation of each genomic segment with one state label. These methods are "semiautomated" because a human performs a functional interpretation of the state labels after the annotation process. In this interpretation step, the human assigns an *interpretation term* to each state label, such as "Promoter" or "Repressed", indicating its putative function.

Segmentation and genome annotation (SAGA) algorithms partition and label the genome on the basis of genomics data sets



ChromHMM: Ernst, J. and Kellis, M. *Nature Biotechnology*, 2010 Segway: Hoffman, M et al. *Nature Methods*, 2012

What biological phenomenon does each unsupervised label correspond to?



Figures should make a point

Walk the audience through each figure









Practice your delivery

Administrivia