

Post Hoc, Local and Model-Agnostic Explanations

AAAI 2023 Tutorial on Trustworthy and Responsible AI

Yilun Zhou MIT CSAIL & Amazon



Outline

- Why model explanations?
- How to compute model explanations? -- Definitions
- How to evaluate model explanations? -- Evaluations
- A definition-evaluation duality





Task for DNN	Caption image	Recognise pneumonia
Problem	Describes green hillside as grazing sheep	Fails on scans from new hospitals
Shortcut	Uses background to recognise primary object	Looks at hospital token, not lung
		(Geirhos et al. 2020)





Why is my model failing? Because ... (debugging and diagnosis)

	(Geirhos et al. 2020)





Computational prediction



How does my model predict this structure? Because ... (scientific discovery)

T1037 / 6vr4 90.7 GDT (RNA polymerase domain) **1049 / 6y4f** 93.3 GDT adhesin tip)

Experimental result

Computational prediction



Global Explains the model "more generally" Local Requires an input instance





Global

Post Hoc, Local and Model-Agnostic Explanations



Post Hoc

Global Explains the model "more generally"

Local Requires an input instance



Intrinsic Generated during model prediction

Generated by an external explainer after model prediction

Yilun Zhou

Post Hoc, Local and Model-Agnostic Explanations









CNN filter visualization (Olah et al., 2017)

Post Hoc Generated by an external explainer after model prediction

Global Explains the model "more generally"





Post Hoc, Local and Model-Agnostic Explanations





CNN filter visualization (Olah et al., 2017)

Post Hoc Generated by an external explainer after model prediction

This presentation

Global Explains the model "more generally"





Post Hoc, Local and Model-Agnostic Explanations



Local Post Hoc Explanations

- What is a local *post hoc* explanation?
 - Some description of the model's local decision making logic
- If this image patch is grayed out, the model prediction will change to this.



Occlusion Saliency (Zeiler and Fergus, 2014)



Local Post Hoc Explanations

- What is a local *post hoc* explanation?
 - Some description of the model's local decision making logic
- If this input feature is changed to this value, the model prediction will be different.





Local Post Hoc Explanations

- What is a local *post hoc* explanation?
 - Some description of the model's local decision making logic
- If this training instance is not present in the dataset, the model will make a different prediction.



Post Hoc, Local and Model-Agnostic Explanations





Occlusion Saliency (Zeiler and Fergus, 2014)



• Vanilla gradient: $s = \nabla_x f(x)$



- Vanilla gradient: $s = \nabla_x f(x)$
- SmoothGrad: $s = \mathbb{E}_{\epsilon}[\nabla_{x} f(x + \epsilon)]; \epsilon_{i}$ independent for every x_{i}





- Vanilla gradient: $s = \nabla_x f(x)$
- SmoothGrad: $s = \mathbb{E}_{\epsilon}[\nabla_{x} f(x + \epsilon)]; \epsilon_{i}$ independent for every x_{i}





- Vanilla gradient: $s = \nabla_x f(x)$
- SmoothGrad: $s = \mathbb{E}_{\epsilon}[\nabla_{x} f(x + \epsilon)]; \epsilon_{i}$ independent for every x_{i}
- Integrated gradient: $s = \int_0^1 \nabla_x f(x_0 + u(x x_0)) du$





- Vanilla gradient: $s = \nabla_x f(x)$
- SmoothGrad: $s = \mathbb{E}_{\epsilon}[\nabla_{x} f(x + \epsilon)]; \epsilon_{i}$ independent for every x_{i}
- Integrated gradient: $s = \int_0^1 \nabla_x f(x_0 + u(x x_0)) du$
- Occlusion: $s_i = f(x) f(x_{-i})$



- Vanilla gradient: $s = \nabla_x f(x)$
- SmoothGrad: $s = \mathbb{E}_{\epsilon}[\nabla_{x} f(x + \epsilon)]; \epsilon_{i}$ independent for every x_{i}
- Integrated gradient: $s = \int_0^1 \nabla_x f(x_0 + u(x x_0)) du$
- Occlusion: $s_i = f(x) f(x_{-i})$
- LIME: $f \sim s^T m + b$

Feature mask m	Masked sentence	$f(\cdot)$
0, 0, 0, 0	<i>um</i>	0.49
0, 0, 0, 1	"beautiful"	0.9
0, 0, 1, 0	"and"	0.52
0, 0, 1, 1	"and beautiful"	0.91
1, 1, 1, 0	"It's good and"	0.89
1, 1, 1, 1	"It's good and beautiful"	0.95



- Vanilla gradient: $s = \nabla_x f(x)$
- SmoothGrad: $s = \mathbb{E}_{\epsilon}[\nabla_{x} f(x + \epsilon)]; \epsilon_{i}$ independent for every x_{i}
- Integrated gradient: $s = \int_0^1 \nabla_x f(x_0 + u(x x_0)) du$
- Occlusion: $s_i = f(x) f(x_{-i})$
- LIME: $f \sim s^T m + b$
- SHAP: $s_i = \sum_{S \subseteq F \setminus \{i\}} \frac{(|S|!(|F|-|S|-1)!)}{|F|!} [f_{S \cup \{i\}}(x_{S \cup \{i\}}) f_S(x_S)]$



- Vanilla gradient: $s = \nabla_x f(x)$
- SmoothGrad: $s = \mathbb{E}_{\epsilon}[\nabla_{x} f(x + \epsilon)]; \epsilon_{i}$ independent for every x_{i}
- Integrated gradient: $s = \int_0^1 \nabla_x f(x_0 + u(x x_0)) du$
- Occlusion: $s_i = f(x) f(x_{-i})$
- LIME: $f \sim s^T m + b$
- SHAP: $s_i = \sum_{S \subseteq F \setminus \{i\}} \frac{(|S|!(|F|-|S|-1)!)}{|F|!} [f_{S \cup \{i\}}(x_{S \cup \{i\}}) f_S(x_S)]$ $e = D(x) \in \mathbb{R}^L$





 $e=D(x)\in\mathbb{R}^L$



Problem: don't know how models work



Problem: don't know how models work Solution: develop explanation methods















Feature importance ⇔ model prediction change with feature removal



Feature importance ⇔ model prediction change with feature removal





Feature importance ⇔ model prediction change with feature removal



Gradient





Feature importance ⇔ model prediction change with feature removal



Gradient



Yilun Zhou



Feature importance ⇔ model prediction change with feature removal



Gradient




























Feature importance ⇔ model prediction change with feature removal





Post Hoc, Local and Model-Agnostic Explanations







Feature importance ⇔ model prediction change with feature removal



Post Hoc, Local and Model-Agnostic Explanations



Common Evaluation Metrics

- Comprehensiveness and sufficiency
 - Also known as deletion and insertion metrics
 - Comprehensiveness also known as area over perturbation curve (AoPC)



Common Evaluation Metrics

- Comprehensiveness and sufficiency
 - Also known as deletion and insertion metrics
 - Comprehensiveness also known as area over perturbation curve (AoPC)
- Decision flip rate under most important feature removal
- Number of removals required for decision flip
- Prediction change rank correlation
- Etc.



Common Evaluation Metrics

- Comprehensiveness and sufficiency
 - Also known as deletion and insertion metrics
 - Comprehensiveness also known as area over perturbation curve (AoPC)
- Decision flip rate under most important feature removal
- Number of removals required for decision flip
- Prediction change rank correlation
- Etc. Definition: $e = D(x) \in \mathbb{R}^D$ Evaluation: $q = E(x, e) \in \mathbb{R}$

MIT CSAIL & Amazon



Evaluating Explanations





Evaluating Explanations





Are Explanations (Necessarily) Correct?

If we know that a specific feature is crucial to the model prediction, can feature attribution explanations identify its importance?





X: original input imageY: original output label





Post Hoc, Local and Model-Agnostic Explanations





X: original input image Y: original output label \hat{Y} : modified output label \hat{X} : modified input image

Post Hoc, Local and Model-Agnostic Explanations



Slides and Resources



Watermark

on bottom

Watermark

on top





Aggregate



Yilun Zhou

Post Hoc, Local and Model-Agnostic Explanations

MIT CSAIL & Amazon





















% Attribution = $\frac{\sum A}{\sum A + \sum A}$



Evaluating Image Saliency Maps

SHAP Watermark





Evaluating Image Saliency Maps

SHAP Watermark





– – – – Random saliency map

- - - Optimal saliency map

Post Hoc, Local and Model-Agnostic Explanations



Post Hoc, Local and Model-Agnostic Explanations









"A saliency map illustrates which area of the input image is considered to be responsible for the cancer [...] Figure 4a [...] shows that the image classifier was able to correctly locate the cancerous region on which its decision was based." (Shen et al., Scientific Reports, 2019)







(Shen et al., Scientific Reports, 2019)









Unknown spurious correlation Can we trust these explainers?

"A saliency rea of the input image is considered to Known spurious correlation he cancer [...] Figure 4a [...] shows that the explainers don't work shows that the explainers don't work able to correctly locate the cancerous region on which its decision was based." (Shen et al., Scientific Reports, 2019)



Explanation Understandability

Do the explanations correctly explain the model prediction logic?



Explanation Understandability

Do the explanations correctly explain the model prediction logic?

Do people correctly understand the model prediction logic from the explanations?



Correct but Not Understandable Explanations



Computation trace: a fully correct explanation for the prediction



Correct but Not Understandable Explanations

Resources



Computation trace: a fully correct but totally not understandable explanation for the prediction



Yilun Zhou

Post Hoc, Local and Model-Agnostic Explanations













A positive sentiment word contributes very positively to the prediction







A positive sentiment word contributes very positively to the prediction



Opposing instance







A positive sentiment word contributes very positively to the prediction...

unless it is negated

Supporting instance

Opposing instance







A positive sentiment word contributes very positively to the prediction...

unless it is negated

Daring, mesmerizing and exceed-		
ingly hard to lorget.	Deminar	0 1 2
Label = Positive mesmerizing:	0.13	
Prediction = Positive		

Supporting instance

Opposing instance


Understanding a Sentiment Classifier





A positive sentiment word contributes very positively to the prediction...

unless it is negated, or near another positive word

Daring, mesmerizing		
ingly hard to longet.	Doming	0 1 2
Label = Positive	mesmerizing:	0.13
Prediction = Positive		

Supporting instance

Opposing instance



Understanding a Sentiment Classifier





A positive sentiment word contributes very positively to the prediction...

unless it is negated, or near another positive word





Understanding a Sentiment Classifier





A positive sentiment word sometimes contributes very positively to the prediction...

unless it is negated, or near another positive word





76

Understanding a Sentiment Classifier





A positive sentiment word sometimes contributes very positively to the prediction...





Explanation Summary (ExSum)





Explanation Summary (ExSum)





Explanation Summary (ExSum)





GUI for Developing ExSum Rules

ExSum Inspection

Rule Union: ((((R1 > R4) > R3) > R5) > R6) > R7 CF Without Rule 7: (((R1 > R4) > R3) > R5) > R6

Rule Selection

		_
Rule 1: negation	Metrics	C
	Coverage	0.
Rule 2: highly positive adjectives have positive saliency	Validity	0.
Rule 3: highly negative adjectives have	Sharpness	0.
Rule 4: highly positive words have positive saliency	Val	
Rule 5: highly negative words have negative saliency		V
Rule 6: Person names have small saliency Rule 7: Stop words have small saliency	Shp	
Reset	Applicat Par	oility amet
Save	(None)	

cs	CF	→ Full	Selected	
age	0.130	→ 0.605	0.475	
у	0.911	→ 0.915	0.916	
ness	0.539	→ 0.269	0.195	
Val		Cov	Val Shp	

Metric Values

Parameter Values

ty Function Behavior Function Parameters saliency lower range -0.1 saliency upper range <u>AutoTu</u> 0.13

Rule Ur Selecte	iion ↔ ed Rule	Sentence ↔ F	EU	All ↔ Invali	d	New Examples
y=0:0.00 the time re	f you collecte equired to be	ed all the moment oil a four - minute	s of cohere egg .	nt dialogue , <u>t</u>	<u>they</u> still we	ould n't add up
y=1:1.00	Ranks among	g Willams <u>'</u> best so	reen work			
y=1 : 1.00 I , it 's signif	ike the film icant without	's almost anthrop being overstated .	ologically d	letailed realiz	ation of ea	rly - '80s subur
y=0 : 0.00 I and an une	t is a comed asy alliance	y that 's not very f , at that) .	unny and a	n action mov	ie <u>that</u> is n	ot very thrilling
y=1:1.00	A triumph , r	elentless and beau	itiful in its c	lownbeat dar	kness <u>.</u>	
y=0:0.001	hey felt like	the same movie to	me.			
y=1 : 1.00 T Sayles dial	'he story fee ogue <u>.</u>	els more like a serio	ous read , fi	lled with hea	vy doses of	always enticin <u>o</u>
y=1 : 1.00 I picture ho	Sehind <u>the</u> si sts a parka-w	now games and lo vrapped dose of he	vable Siber art .	ian huskies (plus one sh	eep dog) , the
y=1 : 1.00 I others to s ferocious o	t is a film th tand up and lebate for ye	at will have people applaud , and wil ears to come .	e walking ○ I , undoubt	ut halfway th edly , leave <u>b</u>	rough , wil o <u>th</u> camps	l encourage engaged in a
	N. N.:	enjoy the same fre	e ride from	critics afford	led to Clint	Eastwood in th

Post Hoc, Local and Model-Agnostic Explanations

MIT CSAIL & Amazon



GUI for Developing ExSum Rules

ExSum Inspection

Rule Selection	Rule Selection		alues		Example Visualization		
	Metrics	$CF \longrightarrow Full$	Selected	Rule Union ↔	Sentence \leftrightarrow FEU	All ↔ Invalid	
ule 2: highly positive saliency	oip install	exsum					ould n't add up to
Jle 3: highly negative \$ g	it clone l	nttps://g	jithub.co	m/YilunZho	u/exsum-d	emos	
le 4: highly positive	d exsum	-demos					rly - '80s suburb
Jle 5: highly negative	xsum sst	_rule_u	nion.py				ot very thrilling (
ule 6: Person names h ule 7: Stop words hav	en up a l	orowser	to localh	.ost:5000 to i	interact wi	th the GU	always enticing
lle 6: Person names h lle 7: Stop words hav Rese Mo	en up a l ore inforn	orowser nation a	to localh t https://y	ost:5000 to i yilunzhou.gi	interact witht. .thub.io/ex	th the GU	always enticing
ule 6: Person names h ule 7: Stop words hav Rese Save	oen up a l ore inform	orowser nation at	to localh t https://y	ost:5000 to i yilunzhou.gi	interact with thub.io/ex	th the GU	always enticing eep dog) , the
lle 6: Person names h lle 7: Stop words hav Rese Save	oen up a l ore inforn _(None)	orowser nation at	to localh thttps://y saliency lower range AutoTu -0.1	ost:5000 to i yilunzhou.gi picture hosts a parka- y=1 : 1.00 It is a film t others to stand up an ferocious debate for	interact with thub.io/ex wrapped dose of heart. hat will have people walki id applaud, and will, under years to come.	th the GU sum/ ng out halfway through pubtedly , leave <u>both</u> ca	always enticing eep dog) , the n , will encourage amps engaged in a
ule 6: Person names r ule 7: Stop words hav Rese Save	oen up a l ore inform	orowser nation a	to localh thttps://y saliency lower range AutoTu -0.1	picture hosts a parka- y=1 : 1.00 It is a film t others to stand up an ferocious debate for y=0 : 0.05 De Niro ma lazy Bloodwork .	interact with thub.io/ex wrapped dose of heart. hat will have people walki d applaud, and will, under y enjoy the same free ride	th the GU sum/ ng out halfway through pubtedly , leave <u>both</u> ca	always enticing eep dog) , the en , will encourage amps engaged in a Clint Eastwood in the

Yilun Zhou

Post Hoc, Local and Model-Agnostic Explanations



The Many Faces of Understandability





The Many Faces of Understandability





The Many Faces of Understandability



Definition vs. Evaluation



Definition

Gradient $D_g(x) = \nabla_x f(x)$ Evaluation



Definition vs. Evaluation



Definition

Gradient $D_g(x) = \nabla_x f(x)$ Evaluation



Gradient $E_g(x, e) = -||\nabla_x f(x) - e||$

Definition vs. Evaluation



Definition

Gradient $D_g(x) = \nabla_x f(x)$ Evaluation



Comprehensiveness $D_{\kappa}(x) = \underset{e}{\operatorname{argmax}} \frac{1}{L+1} \sum_{l=0}^{L} f(x) - f\left(\bar{x}_{e}^{(l)}\right)$ Gradient $E_g(x, e) = -||\nabla_x f(x) - e||$











Beam Search Results

$$\kappa(x, e) = \frac{1}{L+1} \sum_{l=0}^{L} f(x) - f\left(\bar{x}_{e}^{(l)}\right)$$

$$e^* = \operatorname*{argmax}_{e} \kappa(x, e)$$

Algorithm 1: Beam search for finding e^* .

 Input: beam size B, metric m, sentence x of length L;
 Let e⁽⁰⁾ be an empty length-L explanation;
 beams ← {e⁽⁰⁾};

4 for
$$l = 1, ..., L$$
 do

5 beams
$$\leftarrow \bigcup_{e \in \text{beams}} \text{ext}(e, L - l + 1);$$

6 beams
$$\leftarrow$$
 choose_best(beams, B);

7 end

s
$$e^* \leftarrow \text{choose_best(beams, 1)};$$

9
$$e^* \leftarrow \operatorname{shift}(e^*);$$

10 return
$$e^*$$
;



Beam Search Results

$$\kappa(x,e) = \frac{1}{L+1} \sum_{l=0}^{L} f(x) - f\left(\bar{x}_{e}^{(l)}\right)$$

 $e^* = \operatorname*{argmax}_{e} \kappa(x, e)$

c(0)

- Algorithm 1: Beam search for finding e^* .
- Input: beam size B, metric m, sentence x of length L;
 Let e⁽⁰⁾ be an empty length-L explanation;

3 beams
$$\leftarrow \{e^{(0)}\};$$

4 for $l = 1, ..., L$ do

s beams
$$\leftarrow \bigcup_{e \in \text{beams}} \text{ext}(e, L - l + 1);$$

6 beams
$$\leftarrow$$
 choose_best(beams, B);

7 end

$$s e^* \leftarrow \text{choose_best(beams, 1);}$$

9
$$e^* \leftarrow \operatorname{shift}(e^*);$$

10 return
$$e^*$$
;

A worthy tribute to a great humanitarian and her vibrant ' co-stars . '

So stupid, so ill-conceived, so badly drawn, it created whole new levels of ugly.



Beam Search Results

$$\kappa(x, e) = \frac{1}{L+1} \sum_{l=0}^{L} f(x) - f\left(\bar{x}_{e}^{(l)}\right)$$

 $e^* = \operatorname*{argmax}_{e} \kappa(x, e)$

Algorithm 1: Beam search for finding e^* .

1 **Input**: beam size *B*, metric *m*, sentence *x* of length *L*;

2 Let $e^{(0)}$ be an empty length-*L* explanation; 3 beams $\leftarrow \{e^{(0)}\};$ 4 for $l = 1, \dots, L$ do

s beams
$$\leftarrow \bigcup_{e \in \text{beams}} \text{ext}(e, L - l + 1);$$

6 beams
$$\leftarrow$$
 choose_best(beams, B);

7 end

$$s e^* \leftarrow \text{choose_best(beams, 1);}$$

9
$$e^* \leftarrow \operatorname{shift}(e^*);$$

10 **return**
$$e^*$$
;

A worthy tribute to a great humanitarian and her vibrant ' co-stars . '

So stupid, so ill-conceived, so badly drawn, it created whole new levels of ugly.

Explainer	Comp $\kappa \uparrow$	$\operatorname{Suff} \sigma \downarrow$	$\operatorname{Diff}\Delta\uparrow$
Grad	0.327	0.108	0.218
IntG	0.525	0.044	0.481
LIME	0.682	0.033	0.649
SHAP	0.612	0.034	0.578
Occl	0.509	0.040	0.469
E*	0.740	0.020	0.720
Random	0.218	0.212	0.006



Proxy Metrics for Explanation Quality



Post Hoc, Local and Model-Agnostic Explanations



Evaluation on Other Aspects

• E* is competitive on other metrics

Explainer	DF_{MIT}^{\uparrow}	$\text{DF}_{\text{Frac}}\downarrow$	$Rank_{Del} \uparrow$
Grad	10.5%	54.5%	0.162
IntG	16.9%	39.6%	0.369
LIME	25.5%	28.1%	0.527
SHAP	23.0%	36.1%	0.369
Occl	26.4%	40.6%	1.000
E*	25.0%	25.2%	0.438
Random	3.4%	72.3%	0.004



Evaluation on Other Aspects

- E* is competitive on other metrics
- E* is robust to perturbations

Explainer	$ DF_{MIT} \uparrow$	$\mathrm{DF}_{\mathrm{Frac}}\downarrow$	$Rank_{Del}\uparrow$
Grad	10.5%	54.5%	0.162
IntG	16.9%	39.6%	0.369
LIME	25.5%	28.1%	0.527
SHAP	23.0%	36.1%	0.369
Occl	26.4%	40.6%	1.000
E*	25.0%	25.2%	0.438
Random	3.4%	72.3%	0.004





Time Efficiency

B	1	5	10	20	50	100	LIME
κ	0.717	0.731	0.734	0.736	0.739	0.740	0.682
σ	0.020	0.020	0.020	0.020	0.020	0.020	0.033
Δ	0.697	0.711	0.714	0.716	0.719	0.720	0.649
T	0.38	0.77	1.15	1.72	2.85	4.37	4.75

















102



pip install solvex

https://yilunzhou.github.io/solvability/

Summary NI

NLP - Word

NLP - Sentence

e CV - Grid Superpixel

CV - Custom Superpixel Tabular

Read this tutorial as a Jupyter notebook here.

In this demo, we compute word-level explanations for the Huggingface <u>textattack/roberta-base-SST-2</u> model, also the setup presented in the paper.

We first load required packages and the RoBERTa model. Two classes are needed to compute the explanations. **BeamSearchExplainer** implements the beam search algorithm, and ***Masker** implements the feature masking. In this demo, we use **TextWordMasker** since we need to mask out individual words from a text input. The other demos showcase other ***Masker** s.

from solvex import BeamSearchExplainer, TextWordMasker import torch from transformers import AutoTokenizer, AutoModelForSequenceClassification



pip install solvex

Explained label: 1
Function value for label 1: 1.000
Contrary to other reviews, I have zero complaints about the service or the prices. I have been getting tire
service here for the past 5 years now, and compared to my experience with places like Pep Boys, these guys are
experienced and know what they're doing. Also, this is one place that I do not feel like I am being taken
advantage of, just because of my gender. Other auto mechanics have been notorious for capitalizing on my ignorance
of cars, and have sucked my bank account dry. But here, my service and road coverage has all been well explained -
and let up to me to decide. And they just renovated the waiting room. It looks a lot better than it did in
previous years.



pip install solvex

Explained label: 1 Function value for label 1: 1.000 Contrary to other reviews, I have zero complaints about the service or the prices. I have been getting tire service here for the past 5 years now, and compared to my experience with places like Pep Boys, these guys are experienced and know what they're doing. Also, this is one place that I do not feel like I am being taken advantage of, just because of my gender. Other auto mechanics have been notorious for capitalizing on my ignorance of cars, and have sucked my bank account dry. But here, my service and road coverage has all been well explained and let up to me to decide. And they just renovated the waiting room. It looks a lot better than it did in previous years.

Explained label: 232. Function value: 0.159



20

10

0

-10

-20



pip install solvex

Explained label: 1 Function value for label 1.1000 Contrary to other reviews, I have zero complaints about the service or the prices. I have been getting tire service here for the past 5 years now, and compared to my experience with places like Pep Boys, these guys are experienced and know what they're doing. Also, this is one place that I do not feel like I am being taken advantage of, just because of my gender. Other auto mechanics have been notorious for capitalizing on my ignorance of cars, and have sucked my bank account dry. But here, my service and road coverage has all been well explained and let up to me to decide. And they just renovated the waiting room. It looks a lot better than it did in previous years.





Aae = 27.0

-2

Yilun Zhou

Post Hoc, Local and Model-Agnostic Explanations



References

• Definitions

- Simonyan et al. Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps. arXiv 2014
- Zeiler and Fergus. Visualizing and Understanding Convolutional Networks. ECCV 2014
- Ribeiro et al. "Why Should I Trust You?" Explaining the Predictions of Any Classifier. KDD 2016
- Smilkov et al. SmoothGrad: Removing Noise by Adding Noise. arXiv 2017
- Sundararajan et al. Axiomatic Attribution for Deep Networks. ICML 2017
- Lundberg et al. A Unified Approach to Interpreting Model Predictions. NIPS 2017
- Evaluations
 - Samek et al. Evaluating the Visualization of What a Deep Neural Network Has Learned. T-NNLS 2016
 - Petsiuk et al. RISE: Randomized Input Sampling for Explanation of Black-box Models. BMVC 2018
 - Ghorbani et al. Interpretation of Neural Networks is Fragile. AAAI 2019
 - DeYoung et al. ERASER: A Benchmark to Evaluate Rationalized NLP Models. ACL 2020
 - Ross et al. Explaining NLP Models via Minimal Contrastive Editing (MiCE). ACL 2021 (Findings)
 - Zhou et al. Do Feature Attribution Methods Correctly Attribute Features? AAAI 2022
 - Zhou et al. ExSum: From Local Explanations to Model Understanding. NAACL 2022
- Duality
 - Zhou and Shah. The Solvability of Interpretability Evaluation Metrics. EACL 2023 (Findings)