

---

# Advanced Deep Learning: Explainability

Anders Søgaard

---

coASfal

# Course outline

**Goal 1:** Quick tour of recent developments in deep learning

**Goal 2:** Inspiration for thesis/research projects

Week	Lecturer	Subject	Literature	Assignment
1	Stefan	Introduction to Neural Networks.	d2l 2.1-2.5, 2.7, 11.5.1, <b>slides</b>	
2	Stefan	CNNs; FCNs; U-Nets. Data augmentation; invariance; regularization e.g. dropout	d2l 6, 7, 13.9-13.11, <b>slides</b>	Assignment 1 (May 10)
3	Anders/Phillip	May 9 (A): RNNs May 11 (P): Transformers	d2l 8 Transformers: d2l 10.5-10.7 + <a href="#">Vaswani et al. (2017)</a> <sup>e</sup>	Assignment 2 (May 20)
4	Phillip/Anders	May 16 (P): Representation and Adversarial Learning May 18 (A): A Learning Framework + Self-supervised Learning + Contrastive Learning	Autoencoders: <a href="#">blog post</a> <sup>e</sup> GANs: <a href="#">Goodfellow (2016)</a> <sup>e</sup> Self-supervised learning: <a href="#">blog post</a> <sup>e</sup> Contrastive learning: <a href="#">Dor et al. (2018)</a> <sup>e</sup> Adversarial examples: <a href="#">Goodfellow et al. (2015)</a> <sup>e</sup>	
5	Anders	May 23: General Properties, e.g., Scaling Laws, Lottery Tickets, Bottleneck Phenomena May 25: Applications of Representation, Adversarial and Contrastive Learning	GANs: <a href="#">Lample et al. (2018)</a> <sup>e</sup> Autoencoders: <a href="#">Chandar et al. (2011)</a> <sup>e</sup> Contrastive learning: <a href="#">Yu et al. (2018)</a> <sup>e</sup> DynaBench: <a href="#">Talk by Douwe Kiela</a> <sup>e</sup> (Facebook, now HuggingFace) Scaling laws: <a href="#">Kaplan et al. (2020)</a> <sup>e</sup>	Assignment 3 [ <i>MC on Representation Learning</i> /1p Report on Lottery Ticket extraction] (June 3)
6	Anders	May 30: Interpretability, Transparency, and Trustworthiness & Deep Learning for Scientific Discovery June 1: Interpretability (Feature Attribution), including Guest Lecture by Stephanie Brandl	DL for Scientific Discovery: <a href="#">Sullivan (2022)</a> <sup>e</sup> Interpretability/Background: <a href="#">Segaard (2022)</a> <sup>e</sup>	
7	Anders	June 6: <i>Off (no teaching)</i> June 8: Interpretability (Training Data Influence)	Literature: <a href="#">Feng and Boyd-Graber (2018)</a> <sup>e</sup> ; <a href="#">Jiang and Senge (2021)</a> <sup>e</sup>	
8	Anders	June 13-15: Best Practices	Literature: <a href="#">Dodge et al. (2019)</a> <sup>e</sup> and <a href="#">Raji et al. (2021)</a> <sup>e</sup>	Assignment 4 [ <i>MC on Interpretability</i> ; 1p Report on Best Practices] (June 21)

# Course outline

**Goal 1:** Quick tour of recent developments in deep learning

**Goal 2:** Inspiration for thesis/research projects

Architectures

Framework

Fairness /

Explainable AI

Methodology

Week	Lecturer	Subject	Literature	Assignment
1	Stefan	Introduction to Neural Networks.	d2l 2.1-2.5, 2.7, 11.5.1, <a href="#">slides</a>	
2	Stefan	CNNs; FCNs; U-Nets. Data augmentation; invariance; regularization e.g. dropout	d2l 6, 7, 13.9-13.11, <a href="#">slides</a>	Assignment 1 (May 10)
3	Anders/Phillip	May 9 (A): RNNs May 11 (P): Transformers	d2l 8 Transformers: d2l 10.5-10.7 + <a href="#">Vaswani et al. (2017)</a> <sup>e</sup>	Assignment 2 (May 20)
		May 16 (P): Representation and Adversarial Learning	Autoencoders: <a href="#">blog post</a> <sup>e</sup> GANs: <a href="#">Goodfellow (2016)</a> <sup>e</sup>	
4	Phillip/Anders	May 18 (A): A Learning Framework + Self-supervised Learning + Contrastive Learning	Self-supervised learning: <a href="#">blog post</a> <sup>e</sup> Contrastive learning: <a href="#">Dor et al. (2018)</a> <sup>e</sup> Adversarial examples: <a href="#">Goodfellow et al. (2015)</a> <sup>e</sup>	
5	Anders	May 23: General Properties, e.g., Scaling Laws, Lottery Tickets, Bottleneck Phenomena May 25: Applications of Representation, Adversarial and Contrastive Learning	GANs: <a href="#">Lample et al. (2018)</a> <sup>e</sup> Autoencoders: <a href="#">Chandar et al. (2011)</a> <sup>e</sup> Contrastive learning: <a href="#">Yu et al. (2018)</a> <sup>e</sup> DynaBench: <a href="#">Talk by Douwe Kiela</a> <sup>e</sup> (Facebook, now HuggingFace) Scaling laws: <a href="#">Kaplan et al. (2020)</a> <sup>e</sup>	Assignment 3 [MC on Representation Learning/1p Report on Lottery Ticket extraction] (June 3)
6	Anders	May 30: Interpretability, Transparency, and Trustworthiness & Deep Learning for Scientific Discovery June 1: Interpretability (Feature Attribution), including Guest Lecture by Stephanie Brandl	DL for Scientific Discovery: <a href="#">Sullivan (2022)</a> <sup>e</sup> Interpretability/Background: <a href="#">Segaard (2022)</a> <sup>e</sup>	
7	Anders	June 6: Off (no teaching) June 8: Interpretability (Training Data Influence)	Literature: <a href="#">Feng and Boyd-Graber (2018)</a> <sup>e</sup> ; <a href="#">Jiang and Senge (2021)</a> <sup>e</sup>	
8	Anders	June 13-15: Best Practices	Literature: <a href="#">Dodge et al. (2019)</a> <sup>e</sup> and <a href="#">Raji et al. (2021)</a> <sup>e</sup>	Assignment 4 [MC on Interpretability; 1p Report on Best Practices] (June 21)

---

---

# Today

- a) Taxonomies of XAI methods
  - b) Interpretable approximations
  - c) Feature attribution methods
  - d) Training data influence
  - e) Probing and error analyses
-



System Feedback

Individual Feedback

Is System Bad?

Is System Biased?

Has System Learnt Scientific Facts/Hypotheses?

Is System Buggy?



Is Prediction Erroneous?

Is Prediction Biased?

Is Prediction Providing Normative Rationale?

Is the User Provided Useful Feedback?

## Motivation

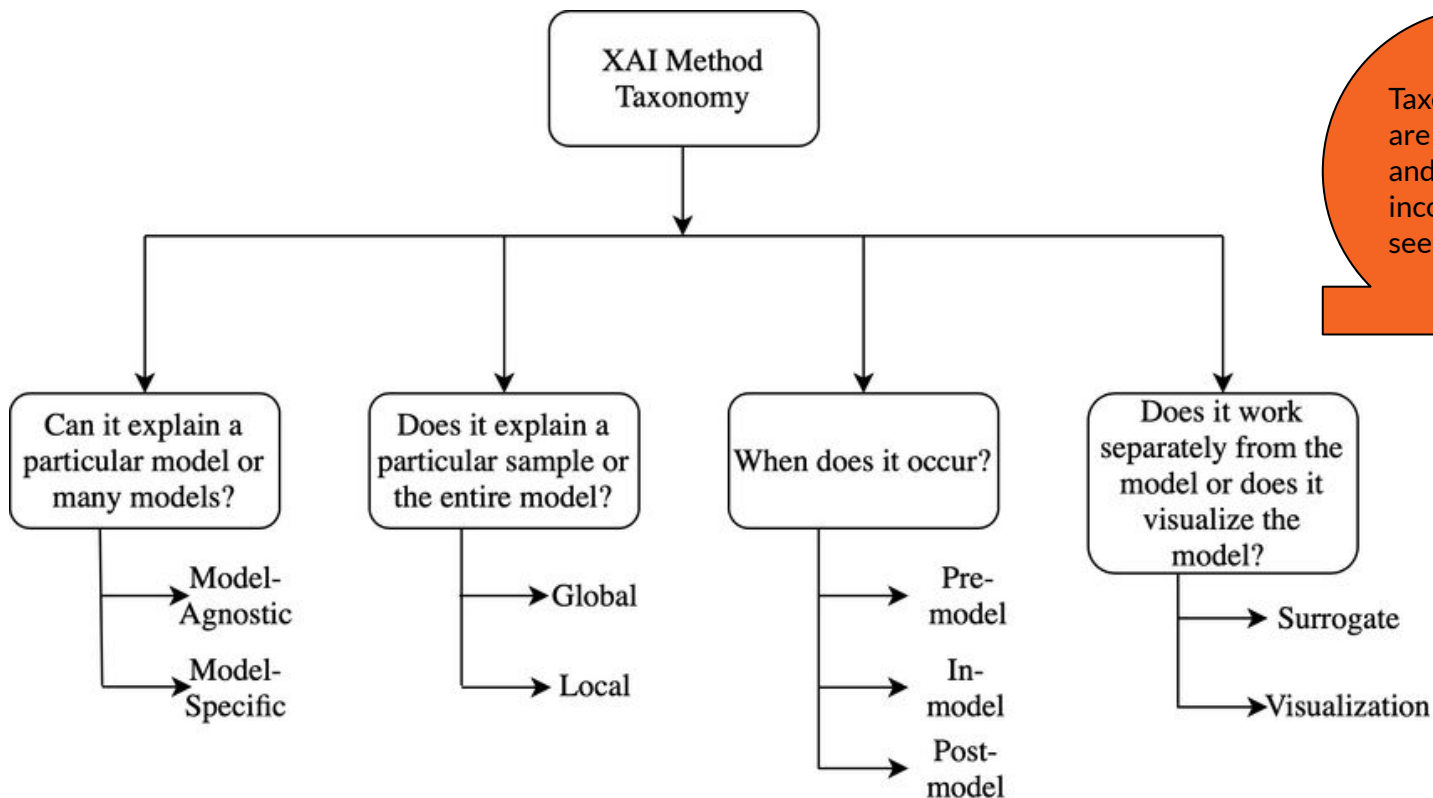
---

---

# Taxonomies

---

# Standard taxonomy: loc-glob + intr-posthoc



Taxonomies are redundant and inconsistent; see book.

---

# Feature attribution

---

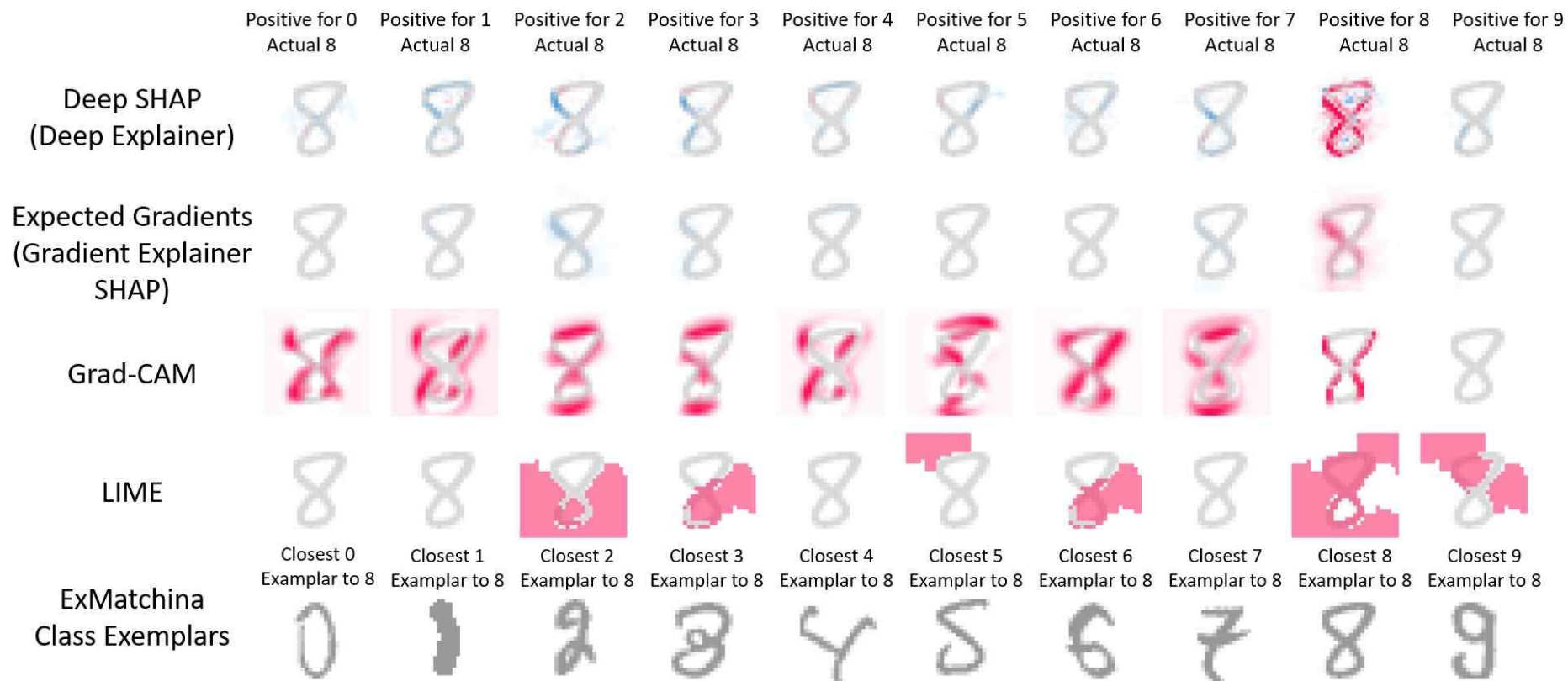


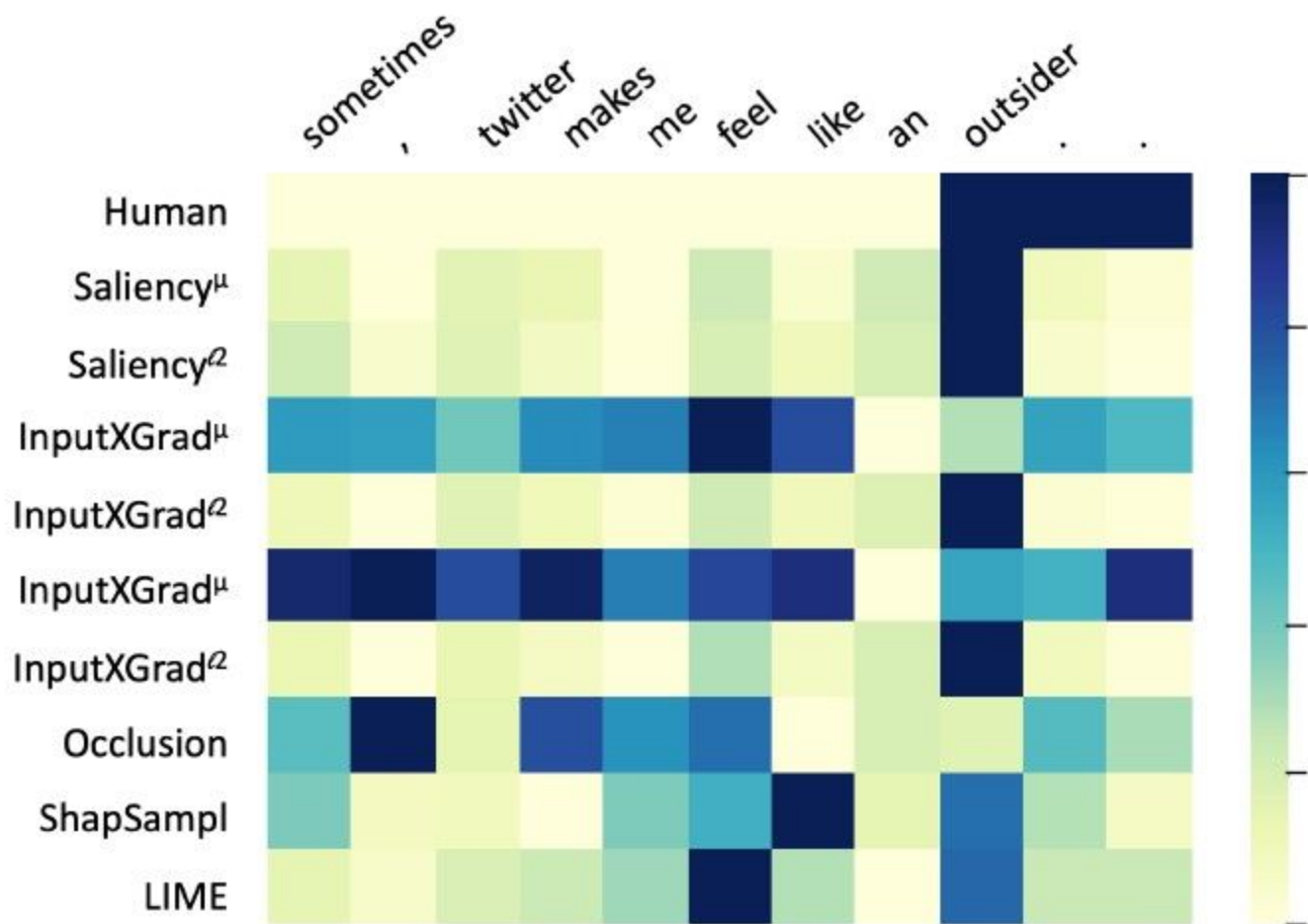
---

# Feature attribution

Using feature attribution methods to 'explain' deep neural networks took off in 2014-15.

Method	Year	Reference
Vanilla gradients	2014	<a href="#">Denil et al. (2014)</a>
Guided back-propagation	2015	<a href="#">Springenberg et al. (2015)</a>
Layer-wise relevance propagation	2015	<a href="#">Bach et al. (2015)</a>
Deep Taylor decomposition	2017	<a href="#">Montavon et al. (2017)</a>
Integrated gradients	2017	<a href="#">Sundararajan et al. (2017)</a>
DeepLift	2017	<a href="#">Shrikumar et al. (2017)</a>





---

# Vanilla gradients

- Compute the gradient of the loss function or the logit of the predicted class with respect to the input embeddings given model parameters.
- We cannot compute the gradient with respect to tokens, only with respect to their embeddings. We therefore must reduce the  $d$ -dimensional gradients to a scalar value, e.g., by computing their  $L_n$  norm.

---

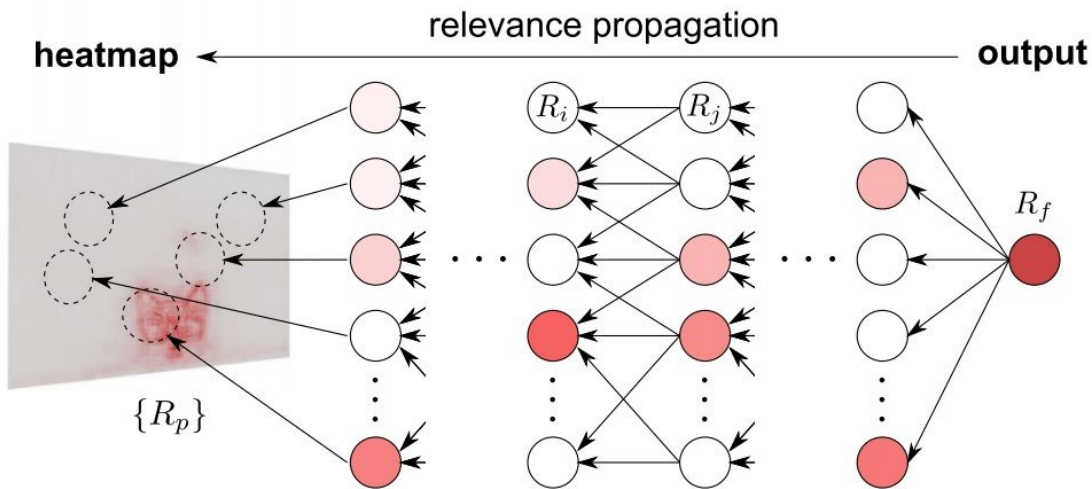
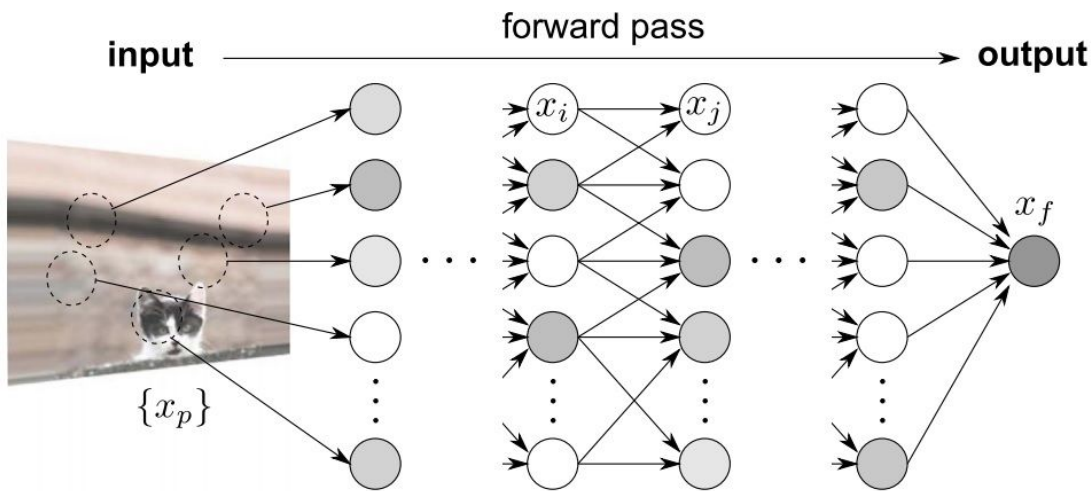
# Vanilla gradients + Guided BP

- Compute the gradient of the loss function or the logit of the predicted class with respect to the input embeddings given model parameters.
  - We cannot compute the gradient with respect to tokens, only with respect to their embeddings. We therefore must reduce the  $d$ -dimensional gradients to a scalar value, e.g., by computing their  $L_n$  norm.
  - While using vanilla gradients relies on actual gradients, guided back-propagation only back-propagates positive error signals, setting negative gradients to zero, reflecting the intuition that positive gradients provide more direct explanations of model decisions.
-

# Layerwise Relevance Propagation

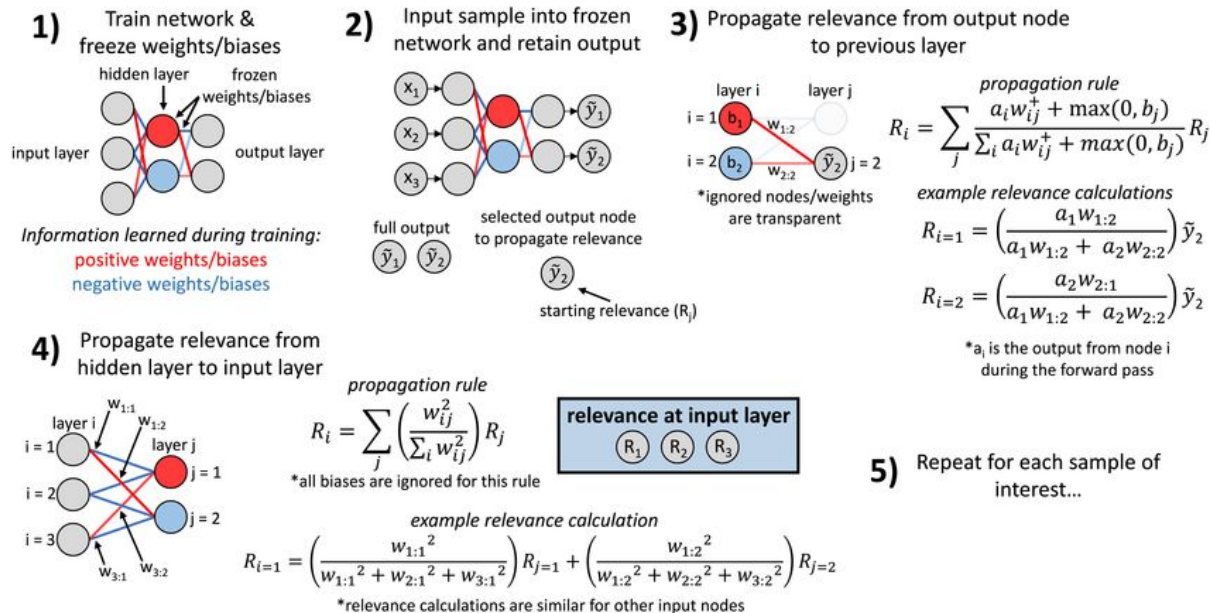
LRP can be seen as an instance of Deep Taylor Expansion, and equivalent to a restricted version of DeepLift.

**Works as follows:** Associate the actual prediction with a relevance score  $R_f$ . Going backwards you compute subsequent relevance scores as the sum of the incoming relevance scores multiplied by (normalized) activation times weight.



## Illustration of Layerwise Relevance Propagation

- LRP back-propagates relevance recursively from the output layer to the input layer.
- Deep Taylor Decomposition (see video on Absalon) is its theoretical motivation.



---

## Other methods

- **Input x Gradient** (like Vanilla Gradients, but multiplied by input; near-equivalent to simple LRP)
  - **DeepLIFT** (normalization by reference point)
  - **Integrated Gradients** (integrating gradients from reference point to data point)
-



---

# Evaluating Feature attribution

---



# Human rationales

- F1 agreement with human rationales.
- The higher the overlap the better.

## Movie Reviews

In this movie, ... Plots to take over the world. The acting is great! The soundtrack is run-of-the-mill, but the action more than makes up for it

(a) Positive (b) Negative

## e-SNLI

H A man in an orange vest leans over a pickup truck  
P A man is touching a truck

(a) Entailment (b) Contradiction (c) Neutral

## Commonsense Explanations (CoS-E)

Where do you find the most amount of leaves?

(a) Compost pile (b) Flowers (c) Forest (d) Field (e) Ground

## Evidence Inference

**Article** Patients for this trial were recruited ... Compared with 0.9% saline, 120 mg of inhaled nebulized furosemide had no effect on breathlessness during exercise.

**Prompt** With respect to *breathlessness*, what is the reported difference between patients receiving *placebo* and those receiving *furosemide*?

(a) Sig. decreased (b) No sig. difference (c) Sig. increased

# Exercise

What's the problem with input reduction and overlap with human rationales?

## Hints:

- 1) Think about i.i.d.
- 2) Think about what we're trying to explain.

---

---

# Training data influence

---

# Influence functions

An old technique for quantifying how the model parameters change as we upweight a training point by an infinitesimal amount.

**Problems:** Expensive and only works for convex models.

**Solution:** Approximations.



Figure 6: Top 5 influential points for the test point: 1479 (CIFAR-10). The model is a ResNet-18 trained with a weight-decay regularization; Only 3 out of the 5 points are semantically similar to the test-point with class "Bird".

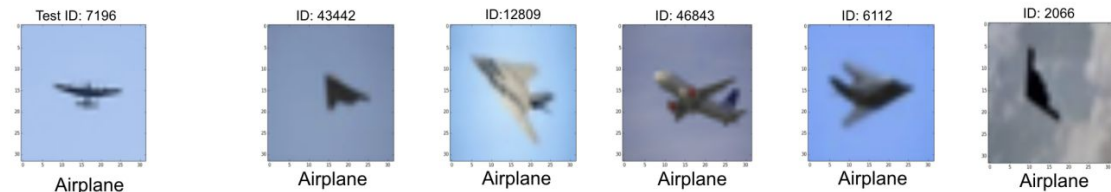


Figure 7: Top 5 influential points for the test point: 7196 (CIFAR-10). The model is a ResNet-18 trained with a weight-decay regularization; All the 5 training points are semantically similar to the test-point from the class "Airplane".

# TracInCP

Store check-points. Make influence of a training data point ( $z$ ) on a test data point ( $z'$ ):

$$\text{TracInCP}(z, z') = \sum_{i=1}^k \eta_i \nabla \ell(w_{t_i}, z) \cdot \nabla \ell(w_{t_i}, z')$$



microphone



microphone



microphone



microphone



acousticguitar



oboe



stage



church



church



church



church



castle



castle



castle



af-chameleon



af-chameleon



af-chameleon



af-chameleon



brocoli



agama



jackfruit



bostonbull



bostonbull



bostonbull



bostonbull



fr-bulldog



fr-bulldog



fr-bulldog



carwheel



carwheel



carwheel



candle



spotlight



loupe



bathtowel

---

# Grad-Cos

**Baseline method:** Simply returns the cosine distance of the gradients of  $z$  and  $z'$ .

**Note:** Common alternative is Grad-Dot.



---

# Evaluating Training Data Influence

---

---

# Leave-one-out influence

- Train a model for all  $n-1$  subsets of your  $n$ -sized training data
  - The influence of  $z$  is the difference in output on  $z'$  between the model trained on all data - and the model trained on all data but  $\{z\}$
  - Often considered **gold** standard
-

---

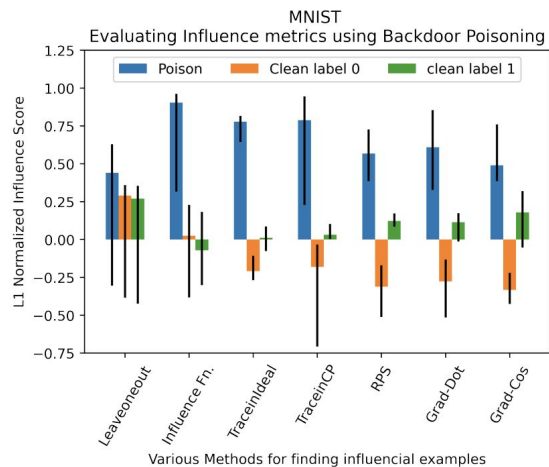
# Heuristics

- The amount of the training data points that have themselves as most influential
- The amount of the test data points in class agreement with their most influential

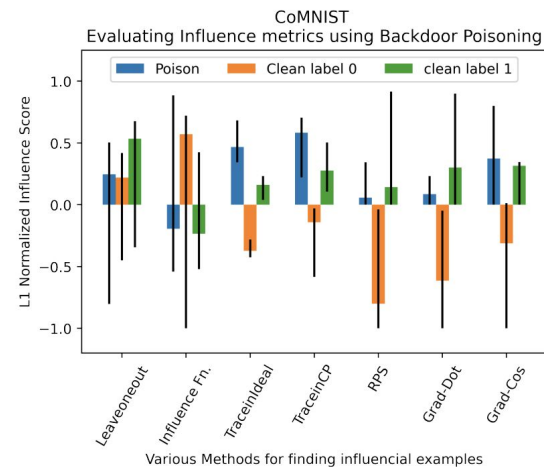
*Exercise:* What's the problem with these heuristics? Think of a counter-example to both.

---

# Backdoor poisoning attacks



(a) MNIST



(b) CoMNIST

# Take-home message

No good automated evaluation  
protocols

**Next Lecture:** The intricacies of  
human evaluation of  
rationales/explanations. As well as  
an equally two-sided story about  
**probing.**

---