

# Explaining nonlinear classification decisions with deep Taylor decomposition

Pattern recognition  
650 citation

Presenter: Jinduk park



# Contents

## **Part1. Introduction to XAI**

- What is explainable AI (XAI) ?
- Taxonomy of XAI

## **Part2. Deep Taylor decomposition**

- Definition of relevance score
- Motivation to use decomposition-based method
- Taylor decomposition
- DeepTaylor decomposition

## **Wrap up**

- Consideration of limitations
- Summary and conclusion

## **Part1. Introduction to XAI**

## What is explainable AI (XAI)?

Deep neural network (DNN) is **successfully** applied to many research area **in terms of its performance**. (e.g. natural language processing, image classification, human action recognition..)

# DEEP LEARNING EVERYWHERE

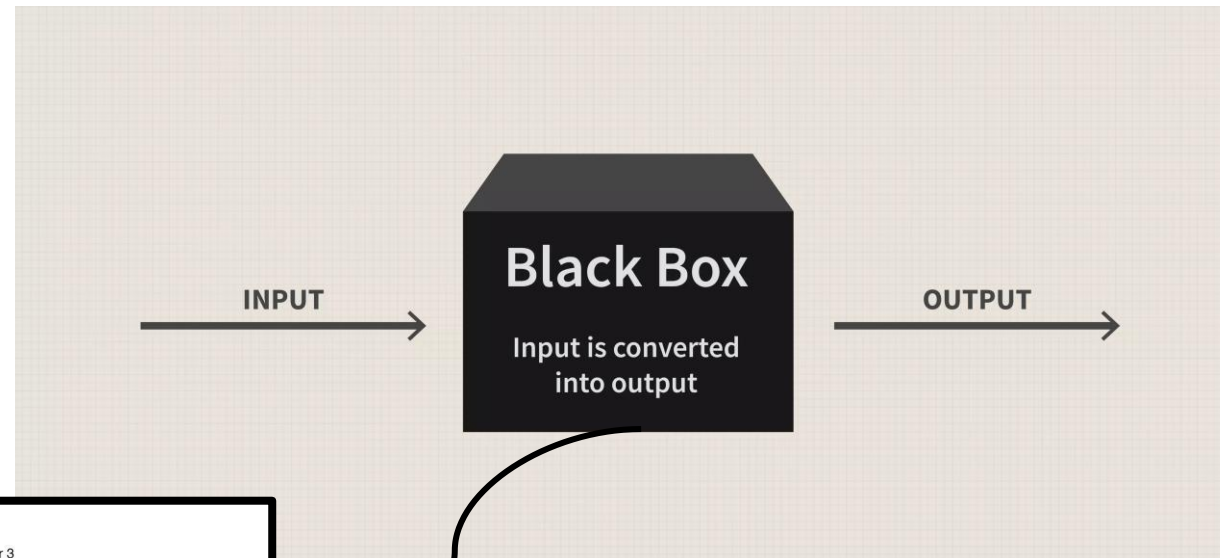


| INTERNET & CLOUD  | MEDICINE & BIOLOGY  | MEDIA & ENTERTAINMENT                                     | SECURITY & DEFENSE  | AUTONOMOUS MACHINES   |
|---|---|---|---|---|
| Image Classification<br>Speech Recognition<br>Language Translation<br>Language Processing<br>Sentiment Analysis<br>Recommendation | Cancer Cell Detection<br>Diabetic Grading<br>Drug Discovery | Video Captioning<br>Video Search<br>Real Time Translation | Face Detection<br>Video Surveillance<br>Satellite Imagery | Pedestrian Detection<br>Lane Tracking<br>Recognize Traffic Sign |

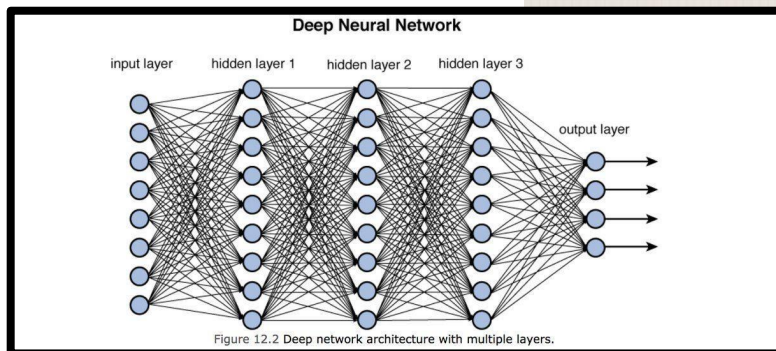
source: [developer.nvidia.com/deep-learning-course](https://developer.nvidia.com/deep-learning-course)

## What is explainable AI (XAI)?

However, because of its complex relations of **nonlinearity**, DNN is regarded as a **black box model**, which means we don't know its internal working.



Source: <https://www.investopedia.com/terms/b/blackbox.asp>

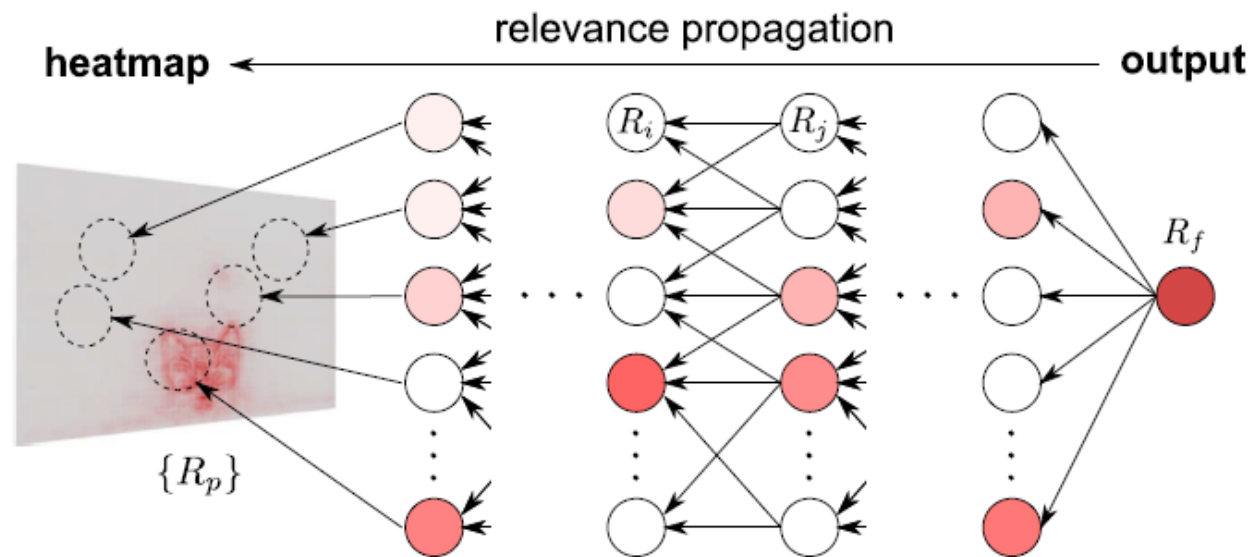


Source: <https://towardsdatascience.com/training-deep-neural-networks-9fdb1964b964>

This property makes **unfaithfulness** of the model prediction. (which is crucial for medical, bio, law domain)

## What is explainable AI (XAI)?

Explainable AI (XAI) explains decisions of a machine learning model, usually in terms of input variables.

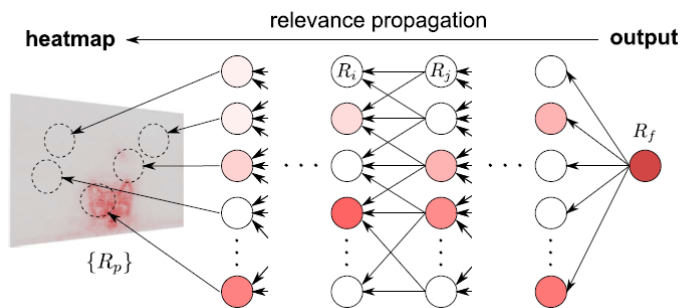


Such a explanation gives users and engineers credibility by providing [transparency](#) to the model.

# Taxonomy of XAI

Whether XAI model **concerns inside of the model or not**

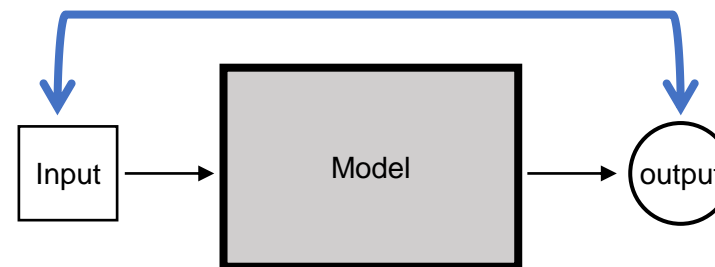
## Model-transparent method



highlight which particular input features triggered key activations within a **model's weights**.

SA [1], GradCAM [2], LRP [3], **DTD**.

## Model-agnostic method



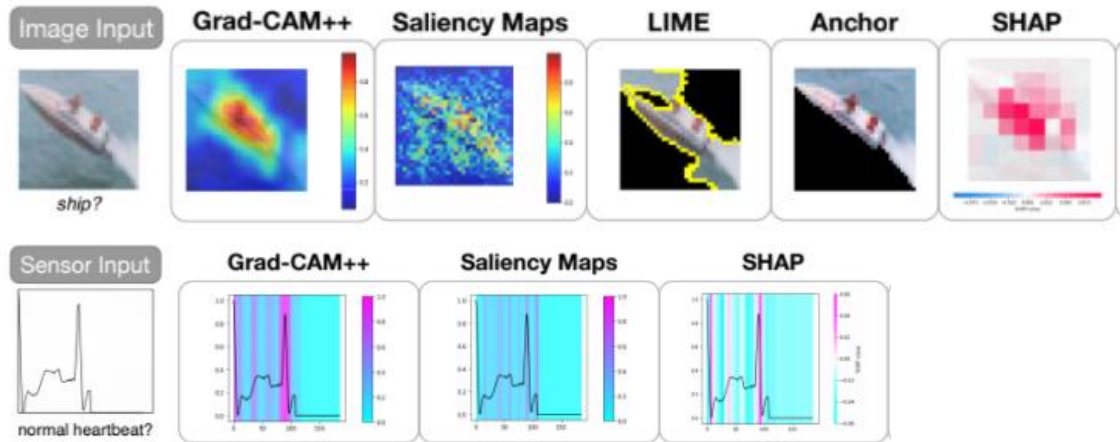
treat the model as totally a black-box and attempt to approximate the **relationship** between the **input and the output prediction**.

LIME [4], SHAP [5], Anchor [6].

# Taxonomy of XAI

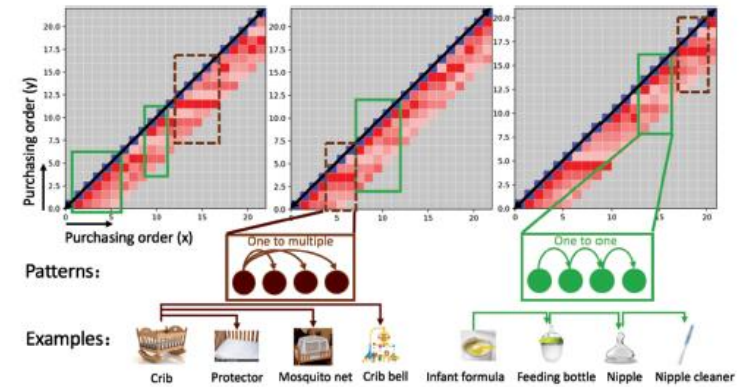
How to provide model explanation to user varies depending on specific domain and XAI method.

[Heatmap visualization]



Source: Jeyakumar, Jeya Vikranth, et al. "How can i explain this to you? an empirical study of deep neural network explanation methods." *Advances in Neural Information Processing Systems* (NIPS) (2020).

[Behavior sequence]



Source: Zhang, Yongfeng, and Xu Chen. "Explainable recommendation: A survey and new perspectives." *arXiv preprint arXiv:1804.11192* (2018).

[Word cloud / sentence generation]



Source: Zhang, Yongfeng, and Xu Chen. "Explainable recommendation: A survey and new perspectives." *arXiv preprint arXiv:1804.11192* (2018).



## Part2. Deep Taylor decomposition (DTD)

One of the *Model-transparent* XAI method that provides *heatmap* when explain model prediction.  
(specifically assume image classification task)

## Definition of relevance score

- Basic concept & notations

- Model  $f(x)$

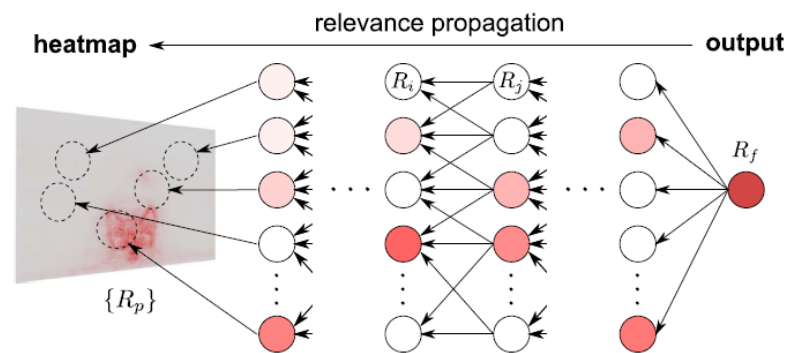
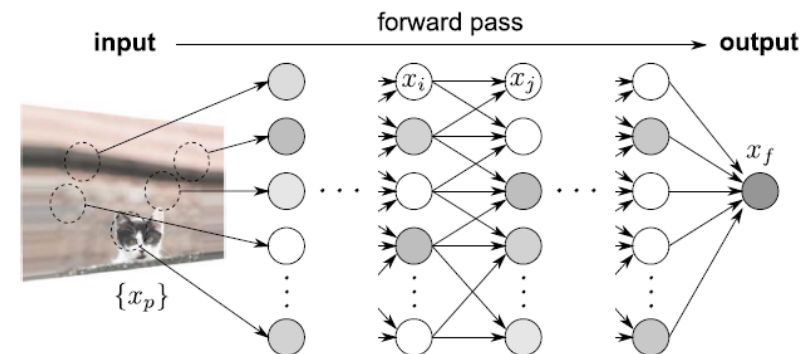
$f: \mathbb{R}^d \rightarrow \mathbb{R}^+$  : Quantifies the presence of certain object.

- Relevance score  $R_p(x)$

To what extent **the pixel  $p$**  contributes to explaining the classification decision in  $f(x)$ .

- Heat mapping  $\mathbf{R}(x)$

Contains set of relevance scores designated to each pixels.



## Definition of relevance score

- Required properties for the relevance score

Definition 1. **conservative**

$$\forall \mathbf{x}: f(\mathbf{x}) = \sum_p R_p(\mathbf{x}).$$

\* These are not sufficient conditions nor strict definition of relevance score though.

Definition 2. **positive**

$$\forall \mathbf{x}, p: R_p(\mathbf{x}) \geq 0$$

Definition 3. **consistent**

- Conservative + positive = consistent
- If the relevance score is consistent, than it is naturally forced to follow:

$$(f(\mathbf{x}) = 0) \Rightarrow (\mathbf{R}(\mathbf{x}) = 0)$$

# Motivation to use decomposition-based method

- Natural decomposition vs SA [1]

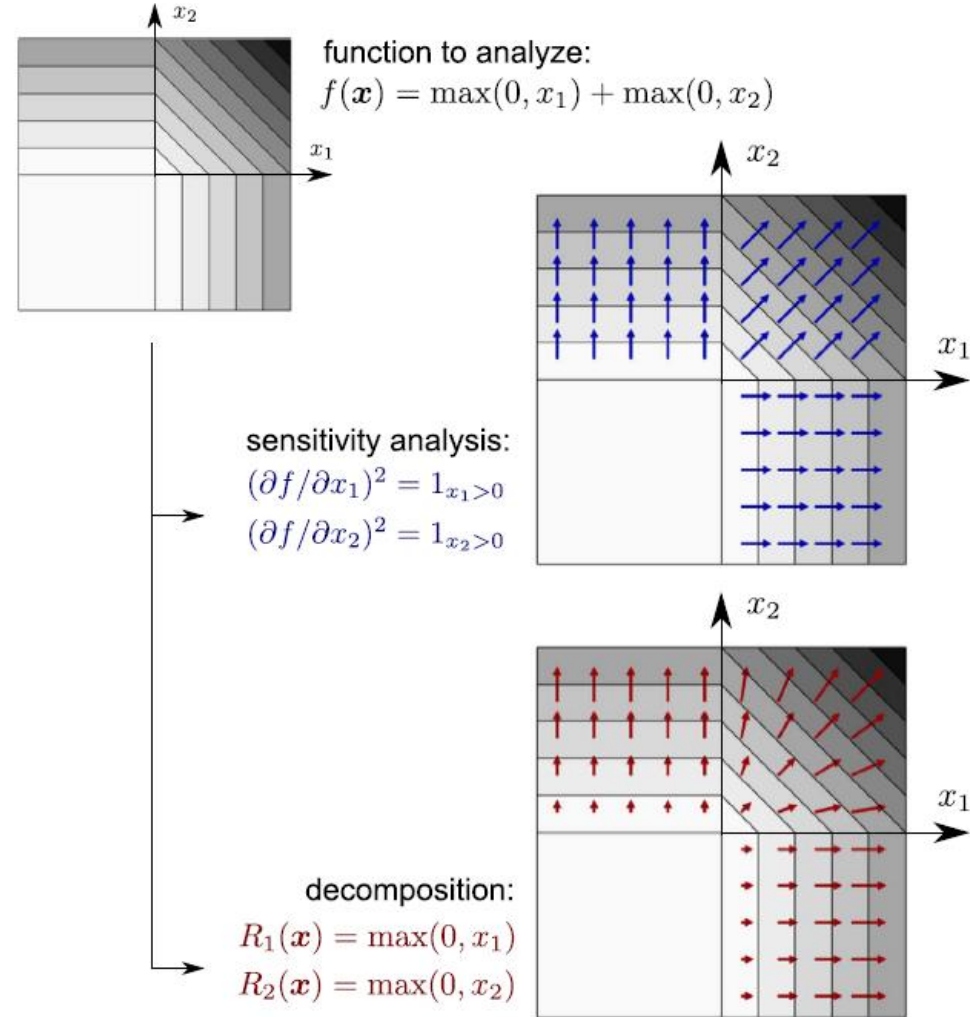
Natural decomposition method that follows definition 1~3.

$$f(\mathbf{x}) = \sum_p \sigma_p(x_p)$$

$$R_p(\mathbf{x}) = \sigma_p(x_p)$$

Even naïve decomposition method **has more expressive** power than SA method.

However, It is still not enough.



Taylor decomposition

A decomposition method based on the **Taylor expansion** of the function at some **well-chosen root point  $\tilde{x}$**

root point  $\tilde{x}$  where  $f(\tilde{x}) = 0$

The Taylor expansion gives:

$$f(x) = f(\tilde{x}) + \left( \frac{\partial f}{\partial x} \Big|_{x=\tilde{x}} \right)^T \cdot (x - \tilde{x}) + \varepsilon$$

$$= 0 + \sum_p \underbrace{\frac{\partial f}{\partial x_p} \Big|_{x=\tilde{x}} \cdot (x_p - \tilde{x}_p)}_{R_p(x)} + \varepsilon$$


★★★★

Higher order terms are complex and hard to redistribute

We takes **1-st order term** as a relevance score  $R_p(x)$

Taylor decomposition

“ What is the philosophy behind the formulation? “

$$f(\mathbf{x}) = \boxed{f(\tilde{\mathbf{x}})} + \overbrace{\left( \frac{\partial f}{\partial \mathbf{x}} \Big|_{\mathbf{x}=\tilde{\mathbf{x}}} \right)^{\top} \cdot (\mathbf{x} - \tilde{\mathbf{x}})}^{\mathbf{R}(\mathbf{x})} + \varepsilon$$


Transposing the term gives:

$$f(\mathbf{x}) - \overset{0}{f(\tilde{\mathbf{x}})} = \mathbf{R}(\mathbf{x})$$

“ From the **absent of the object** ( $f(\tilde{\mathbf{x}}) = 0$ ),  
**how much  $x$  contributes** to the model classification  $f(x)$ . ”

## Taylor decomposition

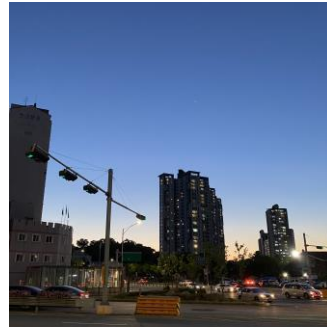
“ Then, **how to find such a root point  $\tilde{x}$**  that satisfying  $(f(\tilde{x}) = 0)$ ? ”

+ The root point should be **admissible**:  
nearest in the Euclidean sense to the actual data point  $x$

$f(.)$  : white car classifier

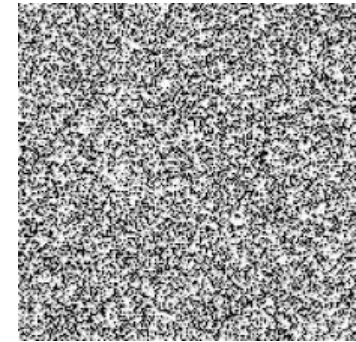


$x$



$\tilde{x}$  (good !)

It will not exactly look like this though.. 😊



$\tilde{x}$  (bad !)

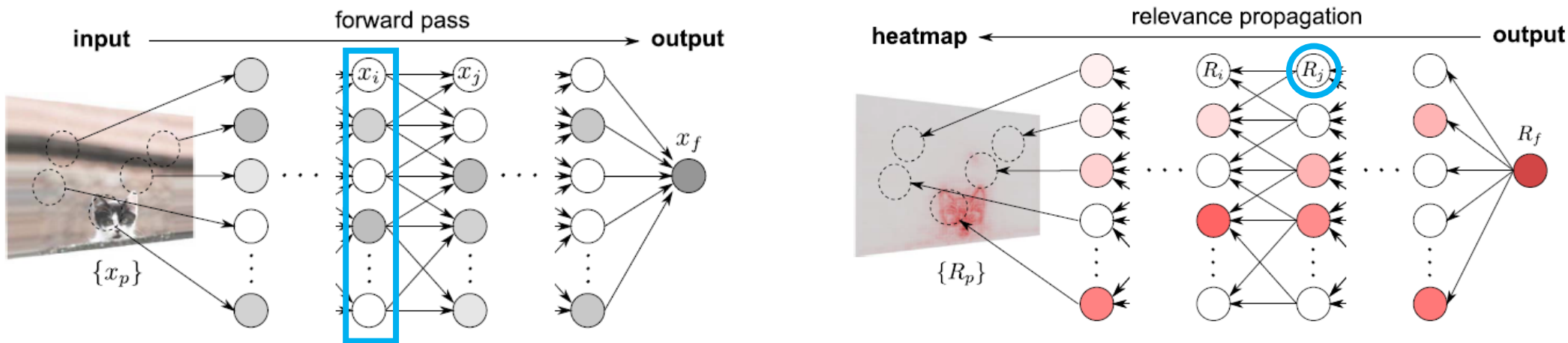
One possible way is to solve **optimization** problem with the minimization objective:

$$\min_{\xi} \|\xi - x\|^2 \quad \text{subject to } f(\xi) = 0 \quad \text{and } \xi \in \mathcal{X},$$

But It is **time-consuming** and thus undesirable.

# Deep Taylor decomposition (DTD)

Deep Taylor decomposition, specifically designed Taylor-decomposition-based redistribution method to apply to deep neural network (DNN).



because of its top-down dependency property of DNN, we can always decompose  $R_j$  in terms of previous layer's units  $\{x_i\}$ .

By Taylor expansion,

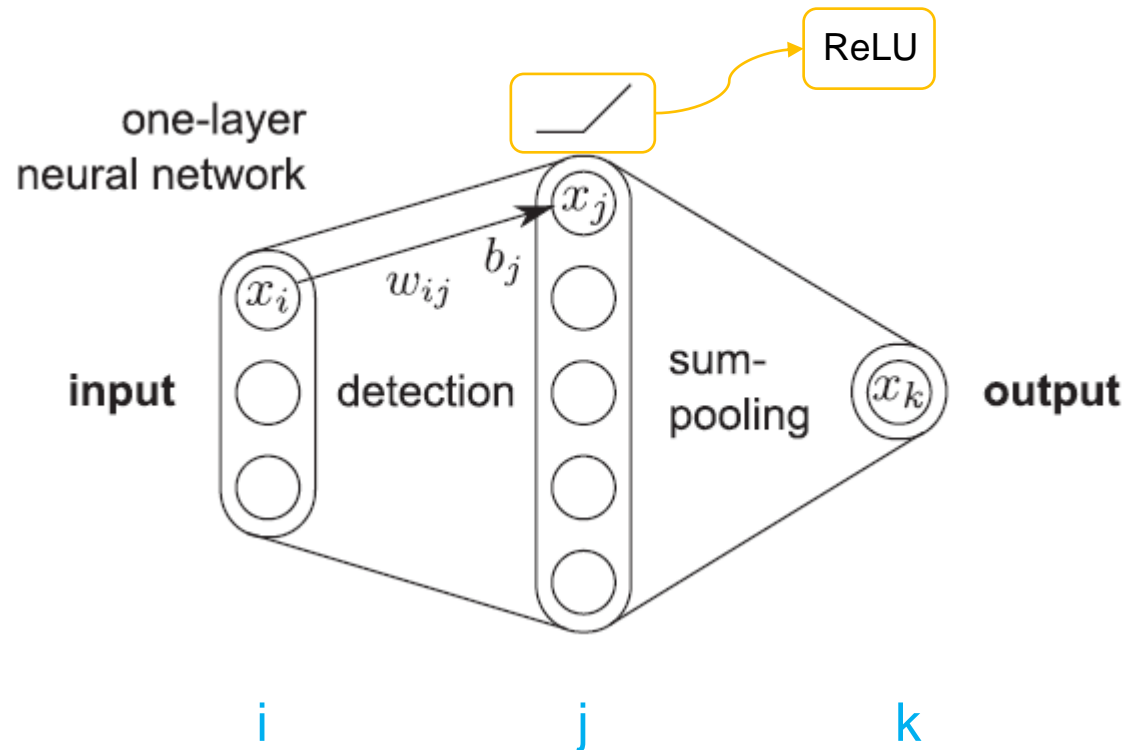
$$R_j = \left( \frac{\partial R_j}{\partial \{x_i\}} \Big|_{\{\tilde{x}_i\}^{(j)}} \right)^T \cdot (\{x_i\} - \{\tilde{x}_i\}^{(j)}) + \varepsilon_j = \sum_i \underbrace{\frac{\partial R_j}{\partial x_i} \Big|_{\{\tilde{x}_i\}^{(j)}}}_{R_{ij}} \cdot (x_i - \tilde{x}_i^{(j)}) + \varepsilon_j$$

Taylor residual



## DTD to one-layer networks

Application DTD on **one-layer network** with **ReLU** non-linear activation.



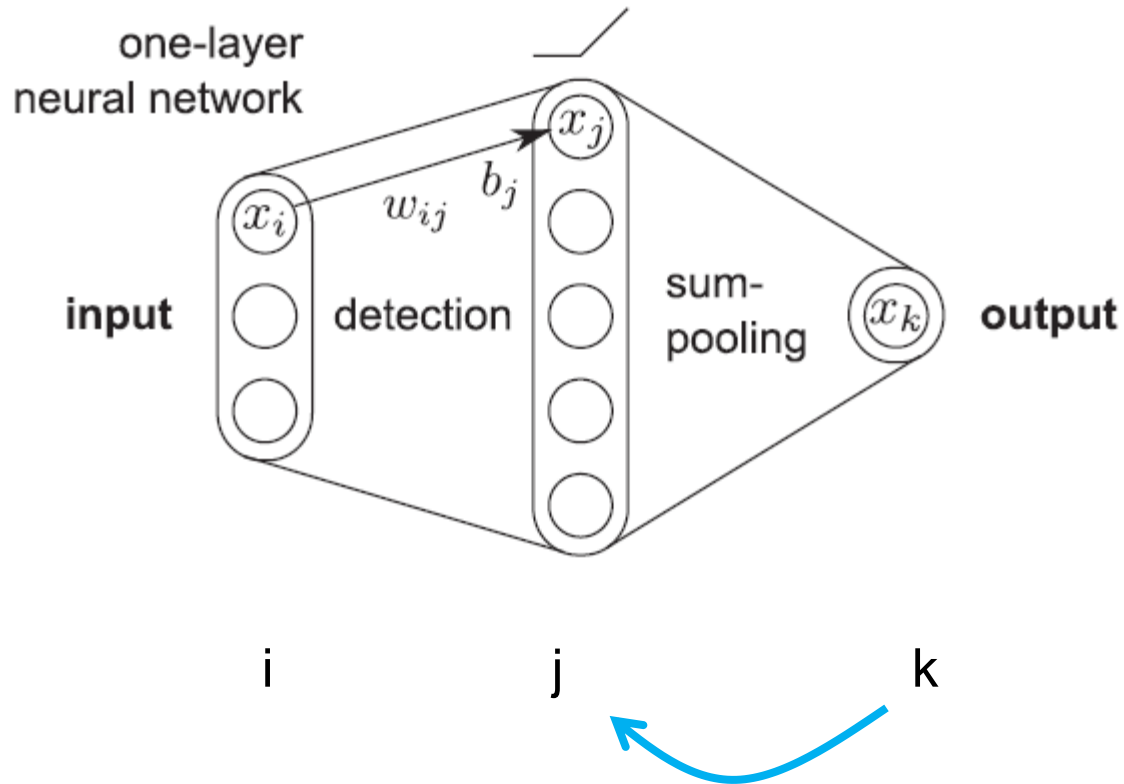
$$x_j = \max\left(0, \sum_i x_i w_{ij} + b_j\right) \quad \text{and} \quad x_k = \sum_j x_j$$

\* **Logic flow** for deriving DTD method.

1. Find **root point**  $\tilde{x}$ .
2. Based on the found root point, conduct **Taylor decomposition w.r.t previous layer**.
3. Derive relevance **redistribution rule**.

\* Such root point ( $f(\tilde{x}) = 0$ ) may not exist in some DNN:  
Give constraint:  $b_j \leq 0$  to ensure *existence* of the root point

## DTD to one-layer networks



- First rule for relevance redistribution.

$$R_k = \sum_j x_j \quad (1)$$

$$R_j = \frac{\partial R_k}{\partial x_j} \Big|_{\{\tilde{x}_j\}} \cdot (x_j - \tilde{x}_j) \quad (2)$$

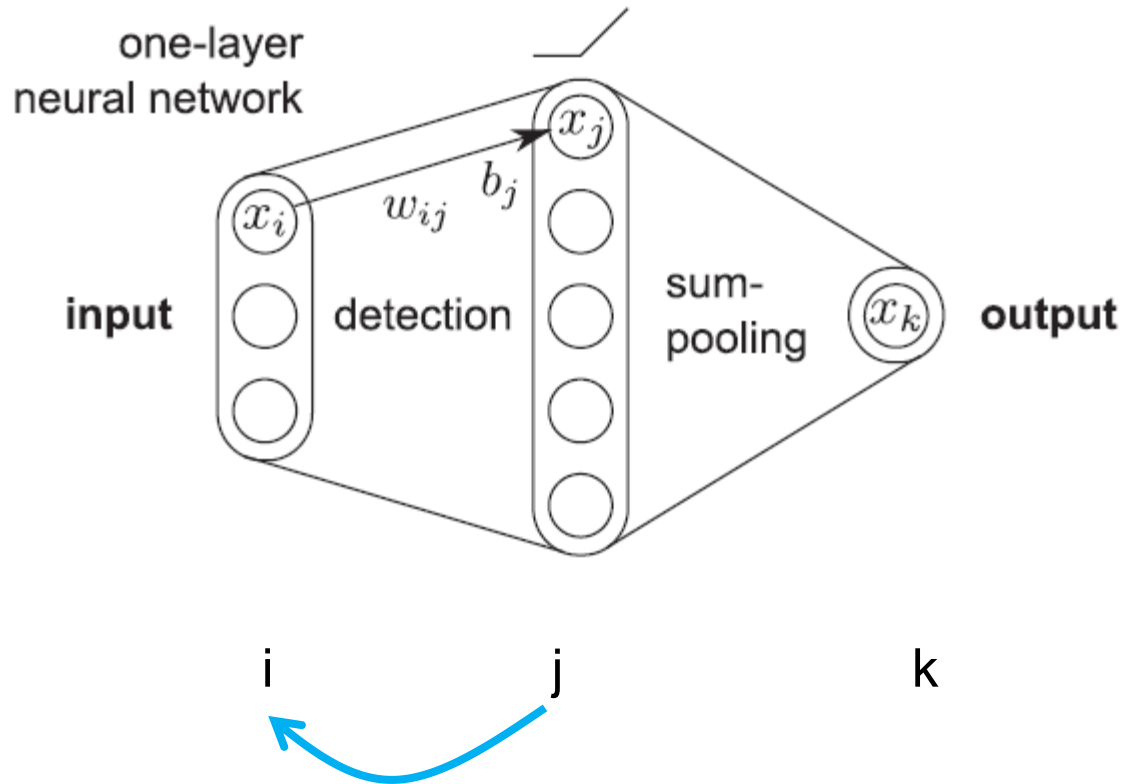
Let's find **root point**  $\{\tilde{x}_j\}$ .

- From (1),  $f(\tilde{x}_j) = \sum_j \tilde{x}_j = 0$
  - Admissible ( $\forall j: \tilde{x}_j \geq 0$ ) : ReLU
- $\tilde{x}_j = 0$

From (2),  $\frac{\partial R_k}{\partial x_j} = 1$ ,  $\tilde{x}_j = 0$  gives  
**first rule for relevance redistribution.**

$$R_j = x_j$$

## DTD to one-layer networks



- **Second rule** for relevance redistribution.

$$\textcircled{R_j} = \max \left( 0, \sum_i x_i w_{ij} + b_j \right), \quad (\text{from first rule})$$

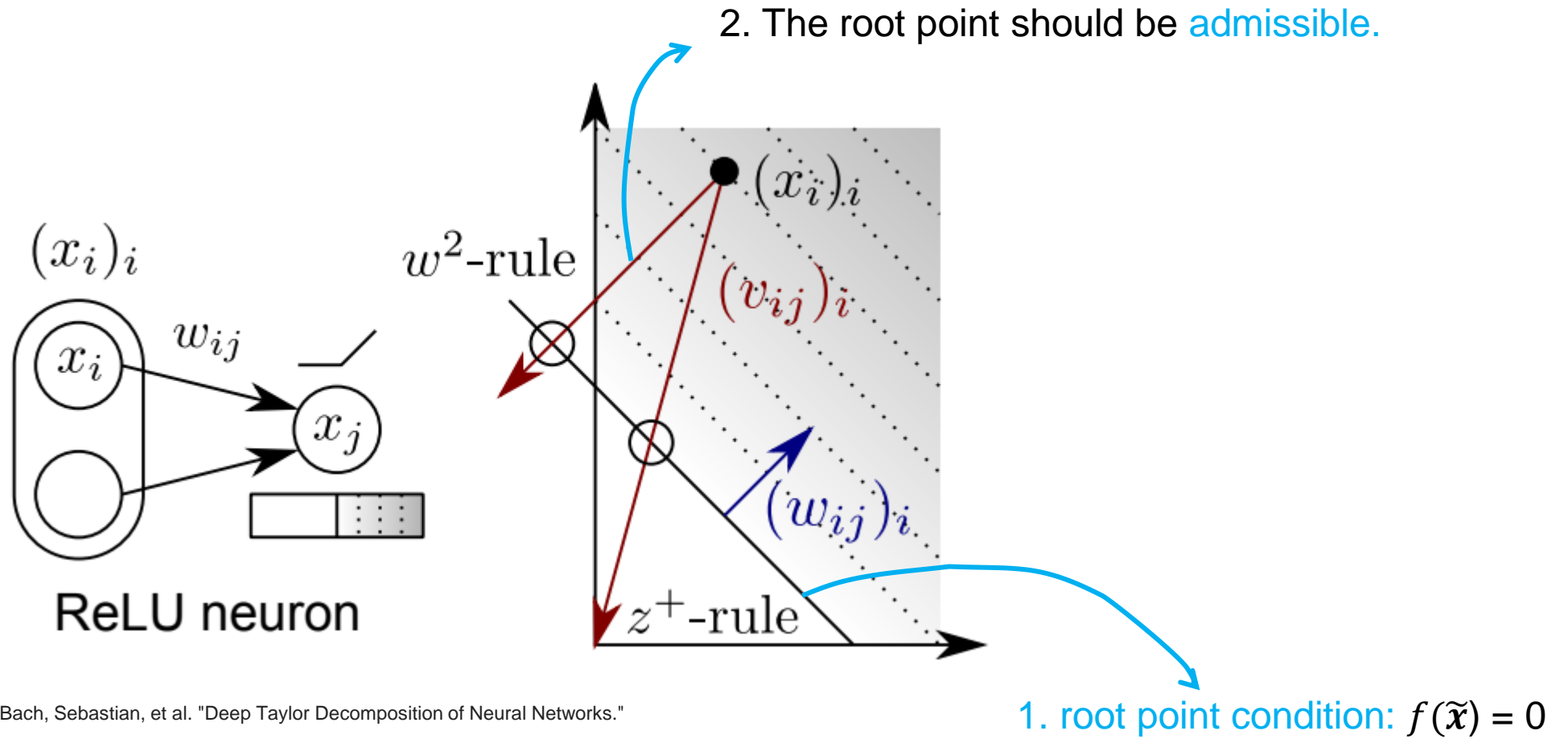
Establish mapping between  $\{x_i\}$  to  $R_j$  by Taylor expansion

$$R_i = \sum_j \frac{\textcircled{R_j}}{\partial x_i} \Big|_{\{\tilde{x}_i\}^{(j)}} \cdot (x_i - \tilde{x}_i^{(j)}).$$

Let's find **root point**  $\{\tilde{x}_i\}^j$  again.

Here, each **choice of input domain** will lead to **different rule** for propagating relevance.

## DTD to one-layer networks



Bach, Sebastian, et al. "Deep Taylor Decomposition of Neural Networks."

Illustration of a root point search in the two-dimensional input space of a ReLU neuron.

We solve the **intersection** of the two condition to find admissible root point.

DTD to one-layer networks1.  $w^2$ -rule

Unconstrained input search space

The admissible root point is **the intersection** of the plane equation  $\sum_i \tilde{x}_i^j w_{ij} + b_j = 0$  and the line of maximum descent  $\{\tilde{x}_i\}^{(j)} = \{x_j\} + t\mathbf{w}_j$

root point  
condition

Admissible : be in vicinity



$$\{\tilde{x}_i\}^{(j)} = \left\{ x_j - \frac{w_{ij}}{\sum_i w_{ij}^2} \left( \sum_i x_i w_{ij} + b_j \right) \right\}$$



Inject this into the equation  $R_i = \sum_j \frac{\partial R_j}{\partial x_i} \Big|_{\{\tilde{x}_i\}^{(j)}} (x_i - \tilde{x}_i^{(j)})$ .

$$R_i = \sum_j \frac{w_{ij}^2}{\sum_{i'} w_{i'j}^2} R_j$$

DTD to one-layer networks

## 2. z-rules

We can give some constraints, which leads to different propagation rules

Positive input search space

$$\{\{x_i\} = 0 \leq x_i\}$$



$z^+$ -rules

$$R_i = \sum_j \frac{z_{ij}^+}{\sum_{i'} z_{i'j}^+} R_j$$

$$z_{ij}^+ = x_i w_{ij}^+$$

Bounded input search space

$$\{\{x_i\} = l_i \leq x_i \leq h_i\} \quad \begin{array}{l} l: \text{lower bound} \\ h: \text{upper bound} \end{array}$$



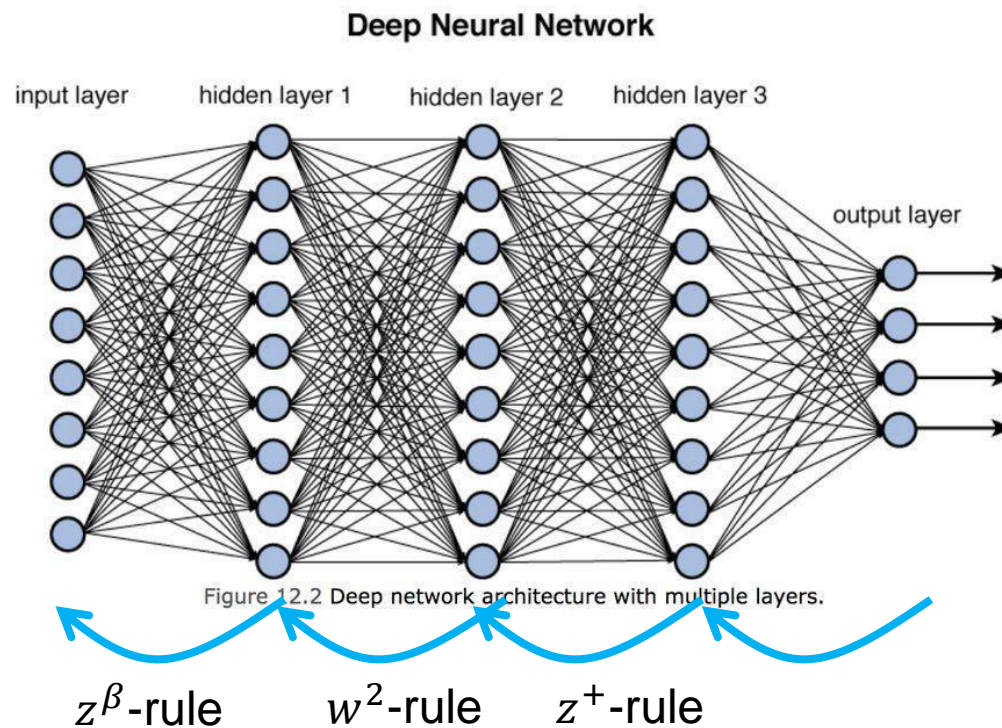
$z^\beta$ -rules

$$R_i = \sum_j \frac{z_{ij} - l_i w_{ij}^+ - h_i w_{ij}^-}{\sum_{i'} z_{i'j} - l_i w_{i'j}^+ - h_i w_{i'j}^-} R_j$$

$$z_{ij} = x_i w_{ij}$$

## DTD to deep networks

DNN is basically constructed with stacking such a simple layer by layer, so we can **serially apply redistribution rule** for different layer considering its activation and constraint.

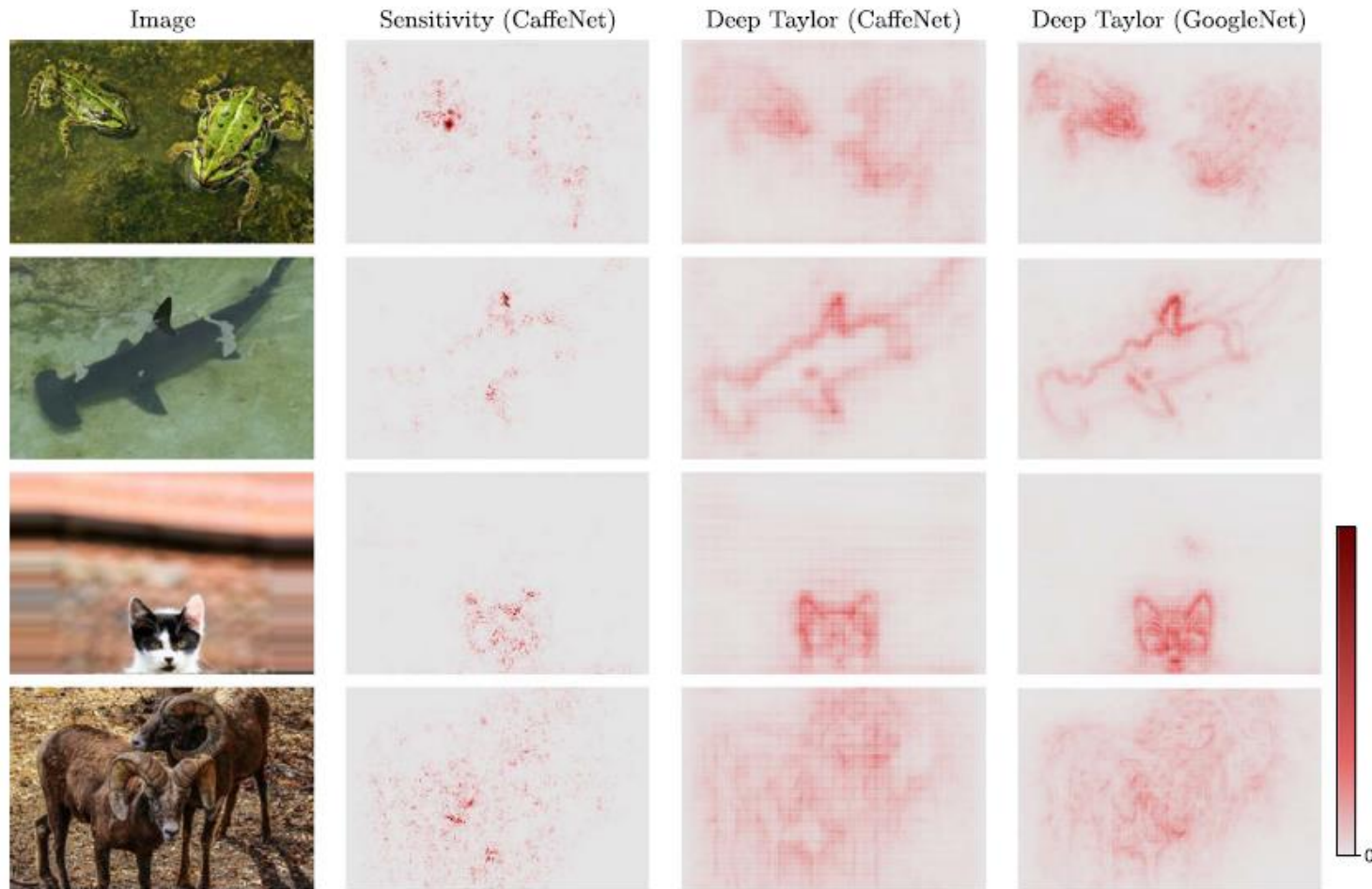


Source: <https://towardsdatascience.com/training-deep-neural-networks-9fdb1964b964>

Stacking redistribution rules makes *training-free* decomposition-based XAI method

## DTD to deep networks

The results (heatmap visualization)



Qualitative evaluation of the method

\* There is no ground-truth for the explanation.



## Wrap up

## Consideration of limitations

1. **Deriving** every relevance redistribution rules for different conditions might be **expensive**.
2. Plus, It seems to be **difficult** to apply DTD to the **more complex network architecture** like LSTM cell.
3. We could infer that **explanation error** occur because of **some assumptions** and **Taylor residuals** in the process of rule derivation.

## Summary and conclusion

- DTD is theoretically well-established model-transparent XAI method that can be applicable for DNN architecture.
- The redistribution rules in DNN are vary depending on the model and data constraints

(Next talk)

- Introduction to some model-agnostic XAI methods
- How to quantitatively evaluate XAI model instead of visual one.

## References

[1] Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. Deep inside convolutional networks: Visualising image classification models and saliency maps. arXiv preprint arXiv:1312.6034, 2013.

[2] Aditya Chattopadhyay, Anirban Sarkar, Prantik Howlader, and Vineeth N Balasubramanian. Grad-cam++: Generalized gradient-based visual explanations for deep convolutional networks. In 2018 IEEE Winter Conference on Applications of Computer Vision (WACV), pages 839–847. IEEE, 2018.

[3] Bach, Sebastian, et al. "On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation." *PloS one* 10.7 (2015): e0130140.

[4] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. " why should i trust you?" explaining the predictions of any classifier. In Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining, pages 1135–1144, 2016.

[5] Scott M Lundberg and Su-In Lee. A unified approach to interpreting model predictions. In Advances in neural information processing systems, pages 4765–4774, 2017.

[6] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. Anchors: High-precision model-agnostic explanations. In AAAI Conference on Artificial Intelligence (AAAI), 2018.

Thanks for your listening.

