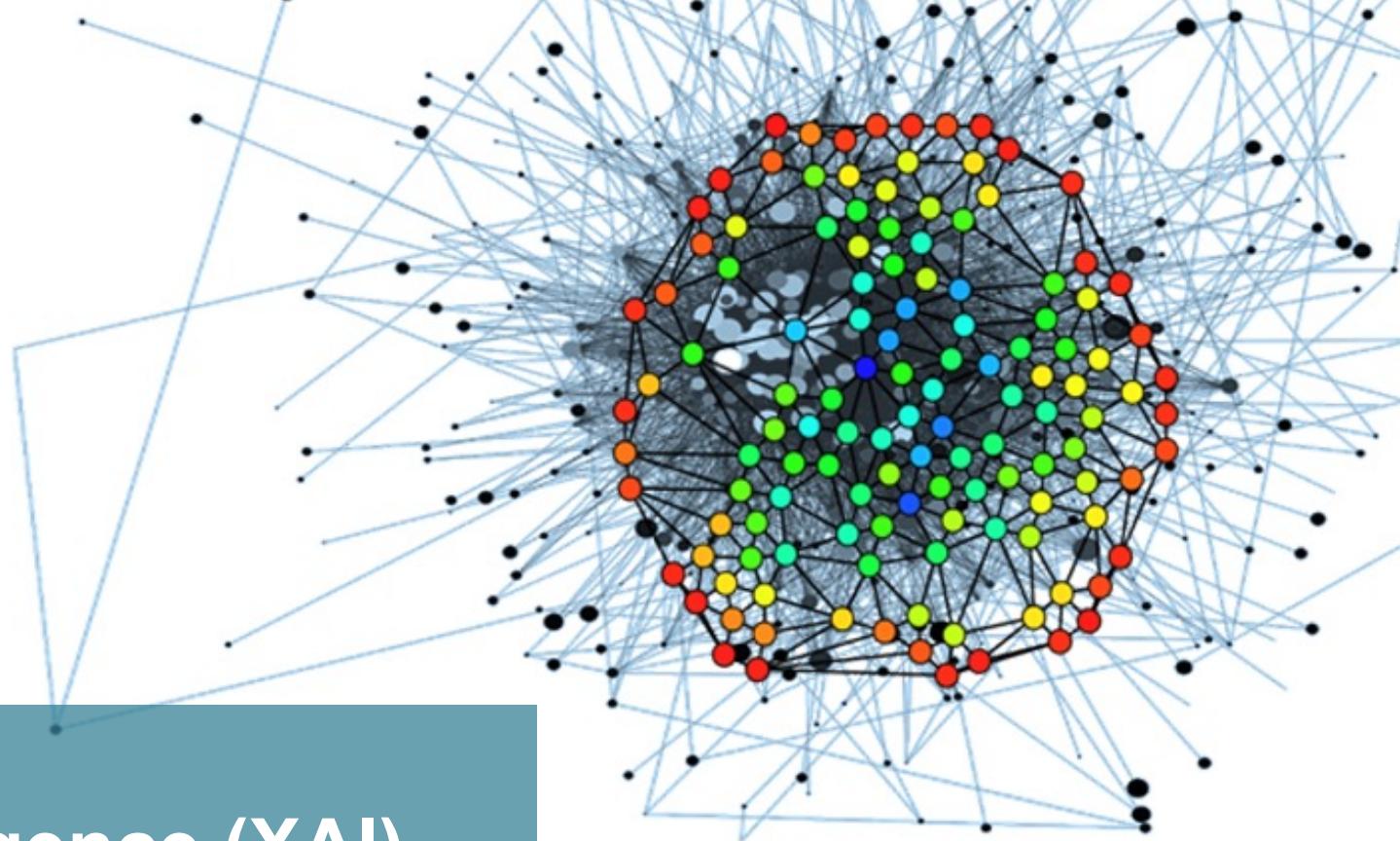


Explainable Artificial Intelligence (XAI)

An introduction



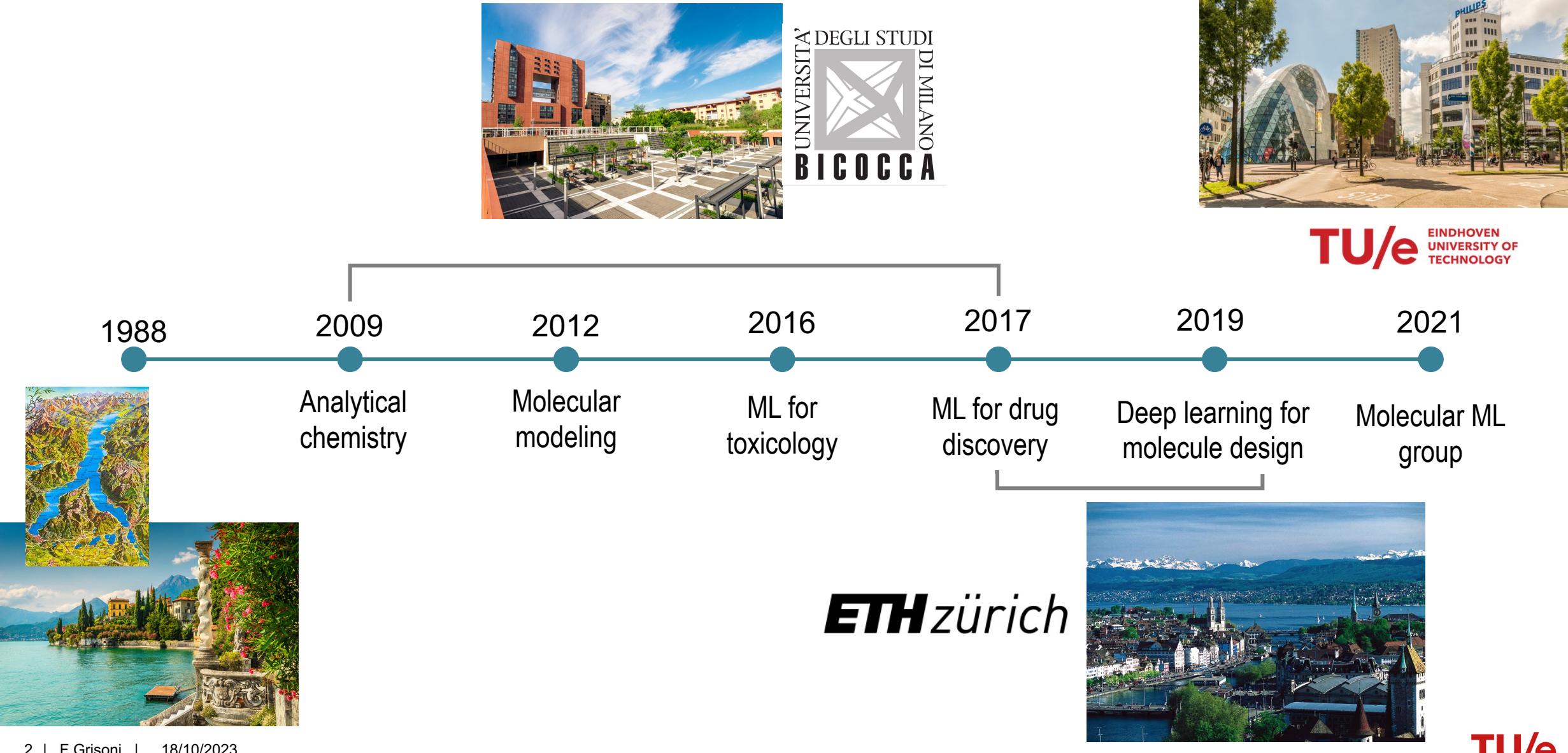
Francesca Grisoni, Assistant Professor

Institute for Complex Molecular Systems & Eindhoven AI Systems Institute

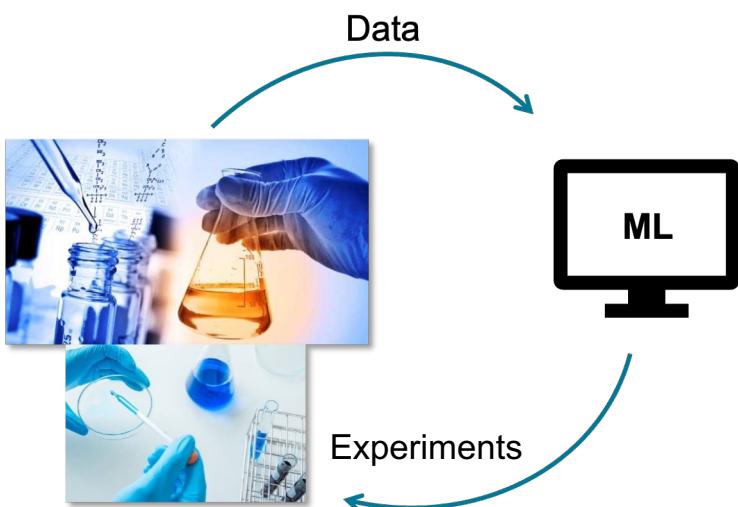
Department of Biomedical Engineering

f.grisoni@tue.nl

About me



Molecular machine learning for drug discovery



X @fra_grisoni
@molecularML

✉ f.grisoni@tue.nl

Reach out for thesis
opportunities!

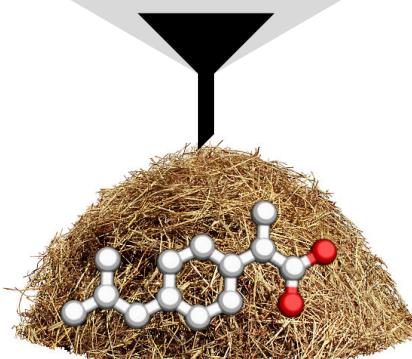
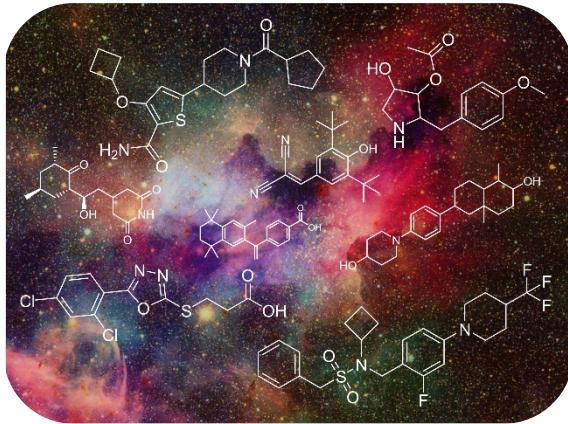


Emanuele Criscuolo
Derek van Tilborg
Rıza Özçelik
Yves Nana Teukam
Helena Brinkmann
Luke Rossen
Cristina Izquierdo-Lozano
Sarah de Ruiter
Meilina Reksoprodjo
Laura van Weesep
Inge Groffen
Sanne van de Vorst

AI for drug discovery

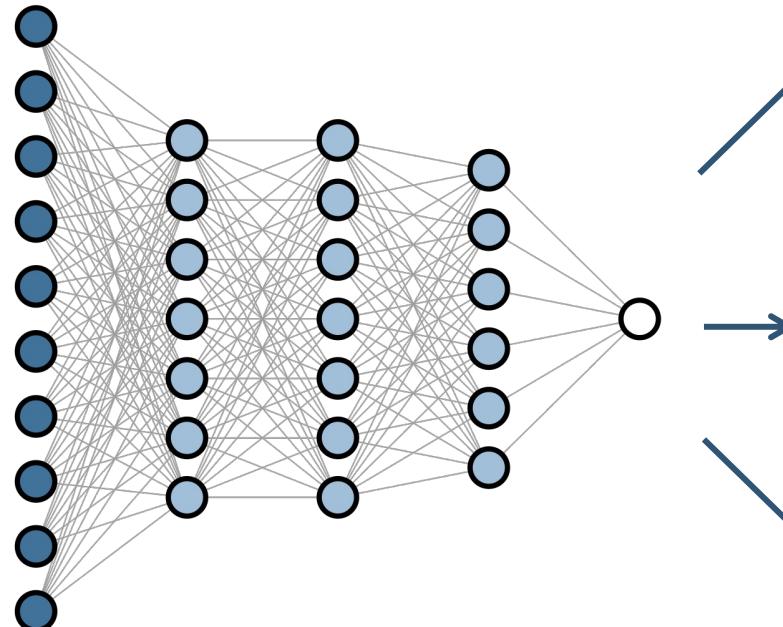
Chemical universe

$10^{23} - 10^{100}$ molecules

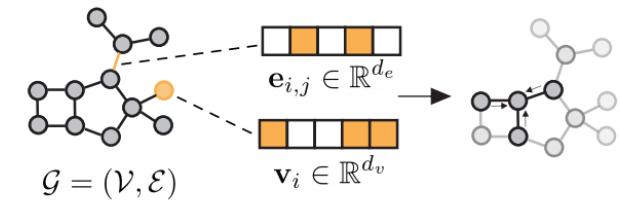


Drug discovery
Finding a *needle in a haystack*

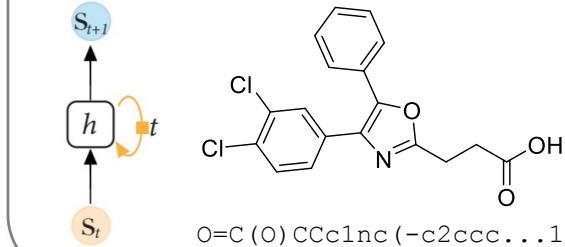
Deep learning



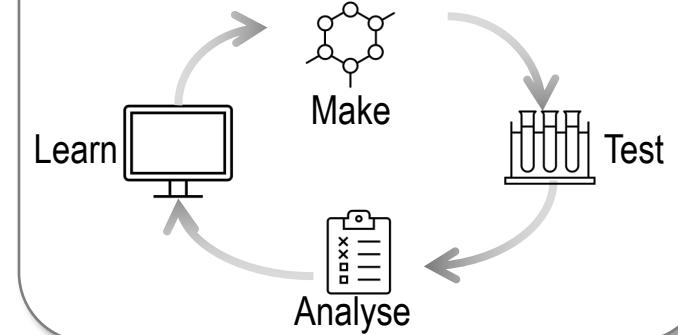
Molecular property prediction



De novo molecule design



Active learning

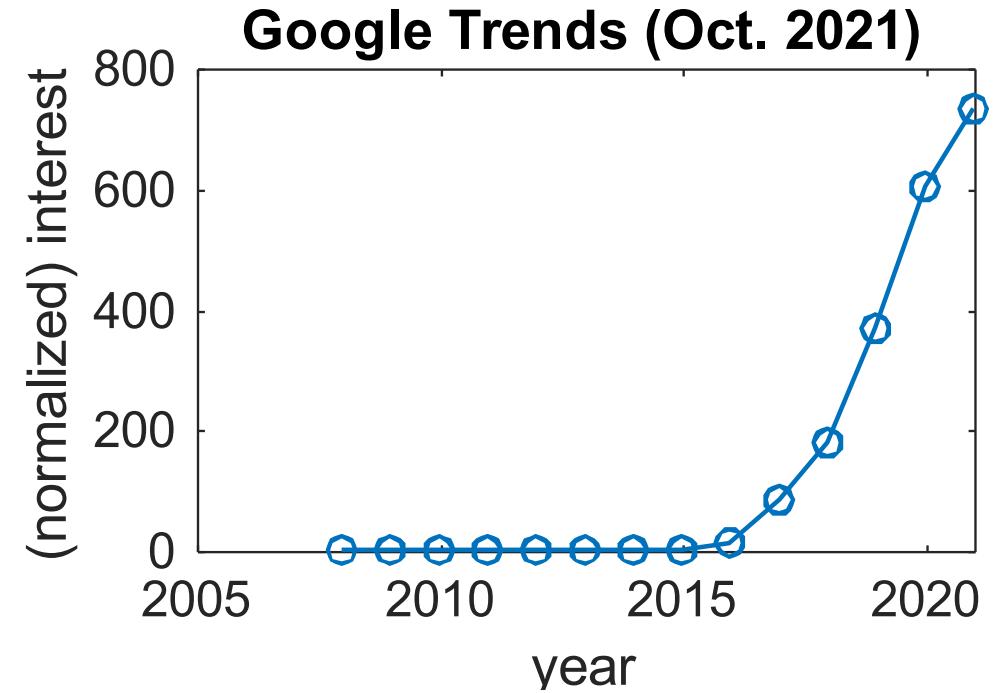


Ice breaker survey

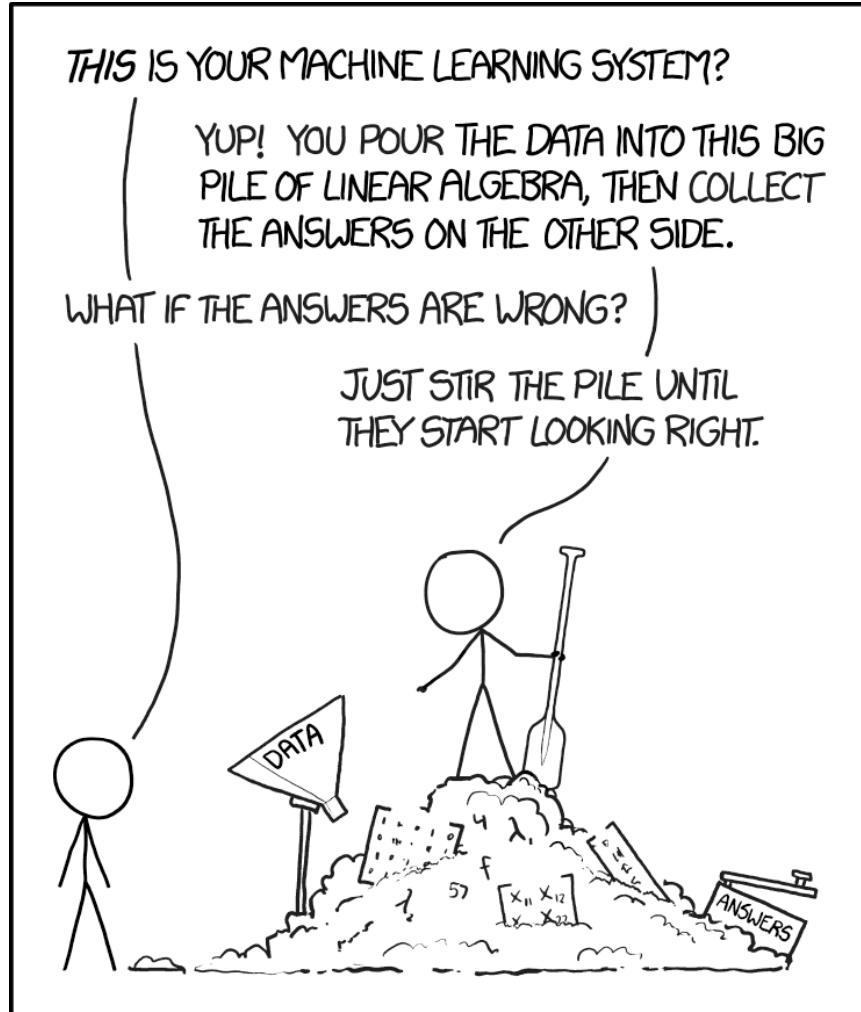


- Have you ever heard of explainable AI?
- What do you think it is?

Explainable AI (XAI)



Explainable AI (XAI)



Extraction of **relevant knowledge** from a machine learning model concerning **relationships contained in data or learned by the model**.



- **Transparency:** *how did the system reached an answer?*
- **Justification:** *is the answer acceptable?*
- **Informativeness:** *what can I learn from it?*
- **Uncertainty estimation:** *how reliable is a prediction?*

Explainability-performance trade-off

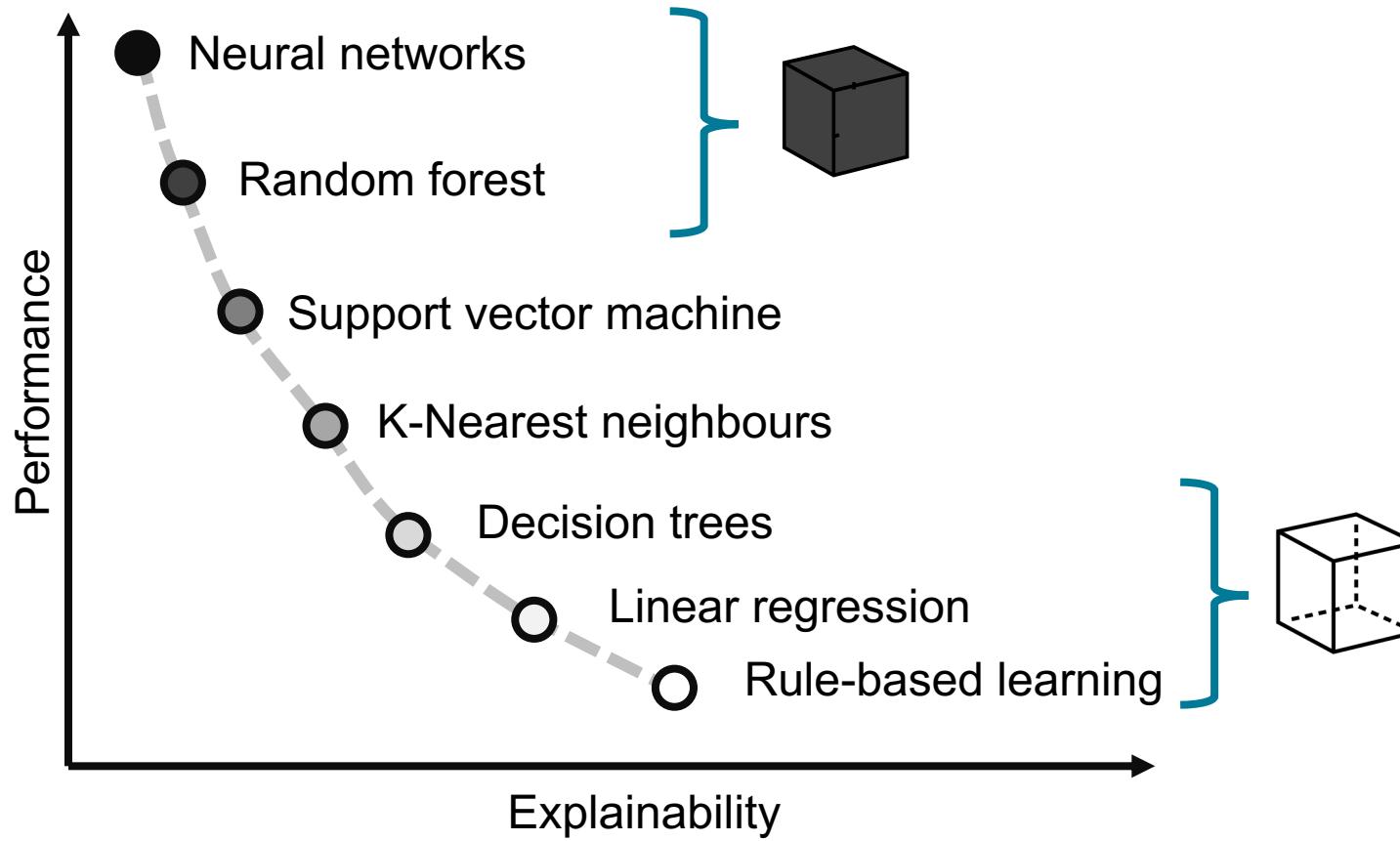


Figure inspired by: Morocho-Cayamcela et al. (2019) *IEEE Access* 7, 137184.

Explainability-performance trade-off

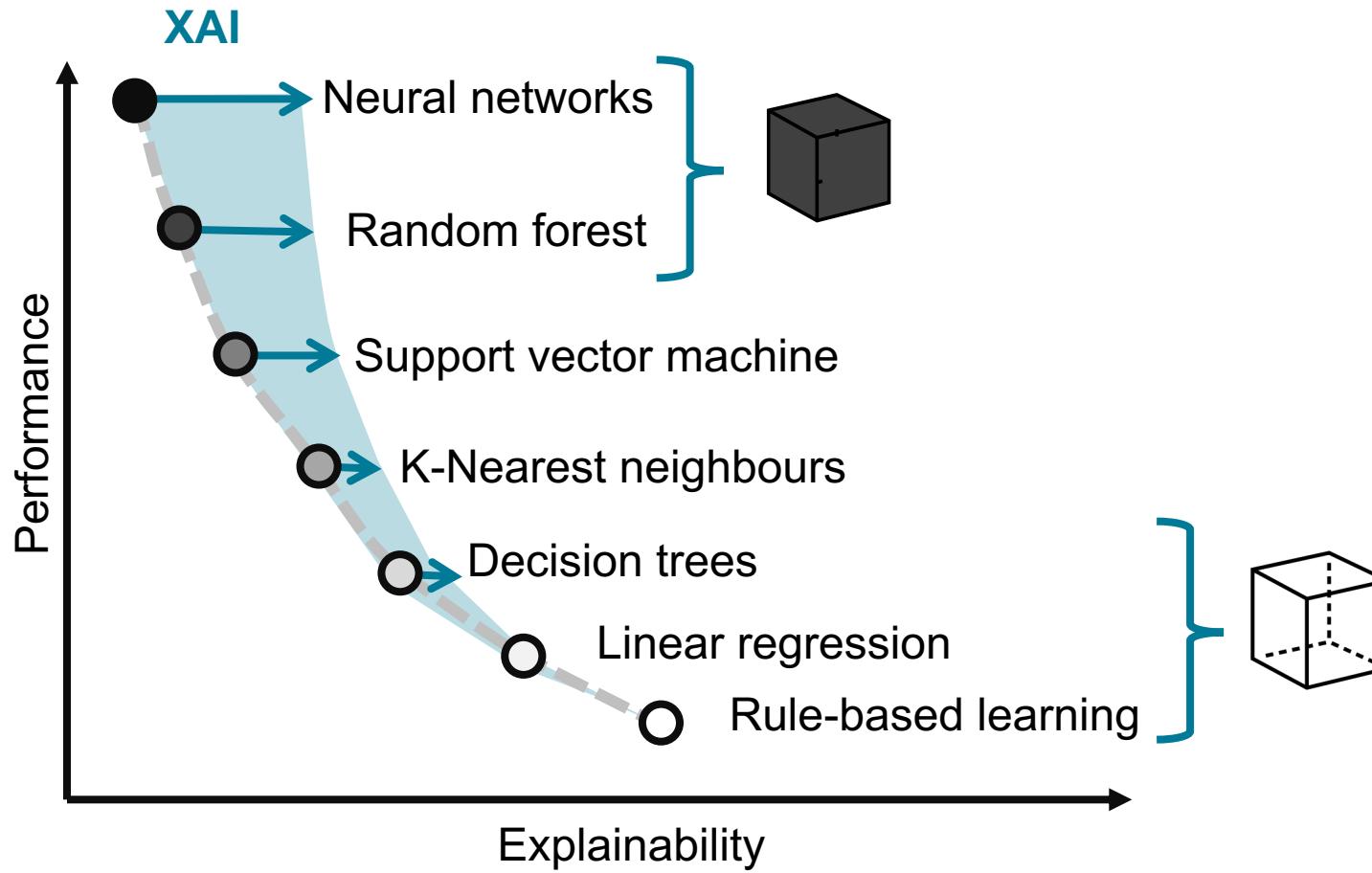


Figure inspired by: Morocho-Cayamcela et al. (2019) *IEEE Access* 7, 137184.

Clever Hans (“der kluge Hans”)



- Berlin, 1900.
- Horse claimed to perform arithmetic (hoof tapping).
- The horse was responding directly to involuntary cues in the body language of the human trainer.

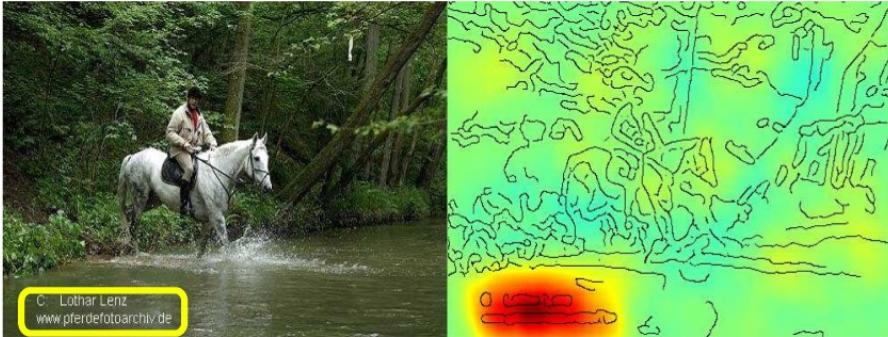


Clever Hans effect
Correct outcome for the wrong reasons

Sebeok and Rosenthal (1981). *Annals of the New York Academy of Sciences*.

“Clever Hans” in AI (shortcut learning)

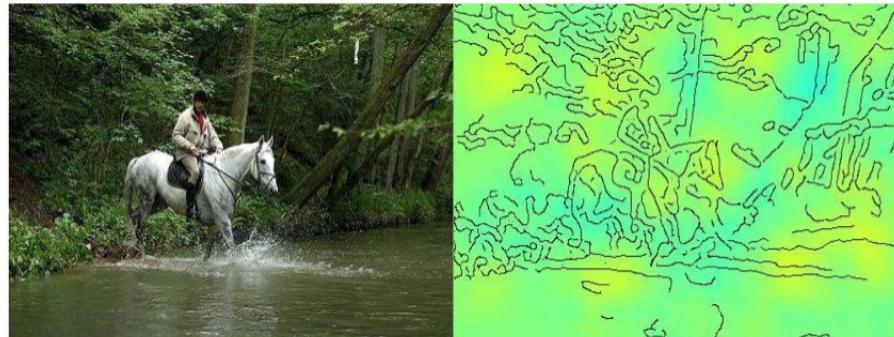
Horse-picture from Pascal VOC data set



Source tag present



Classified as horse

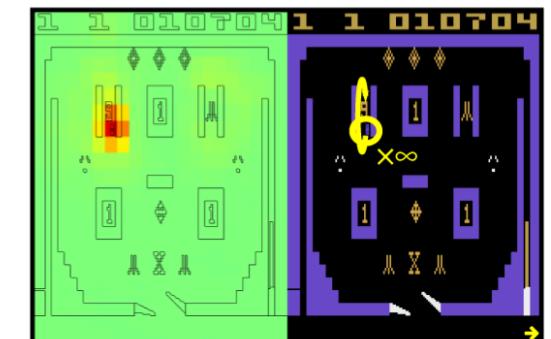
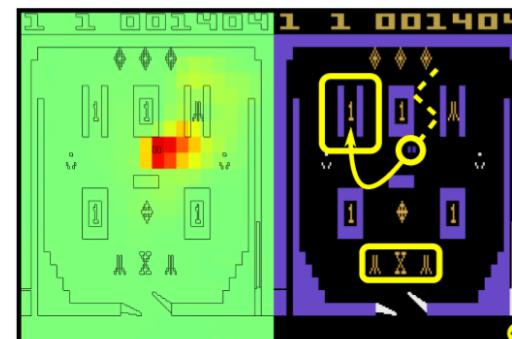
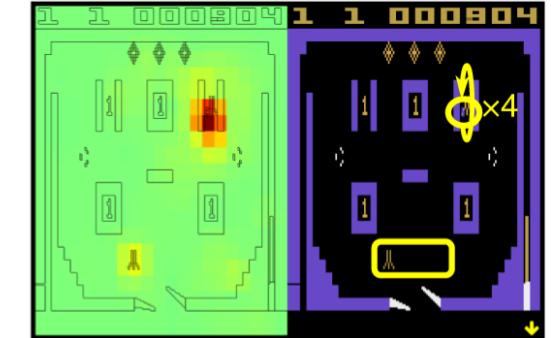
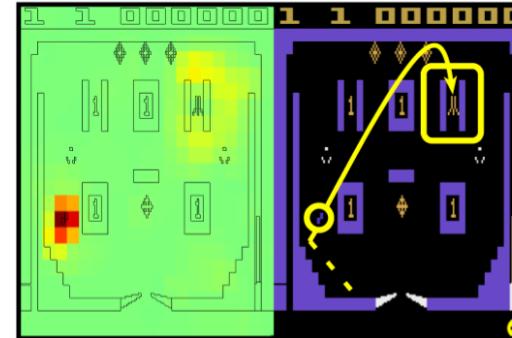


No source tag present



Not classified as horse

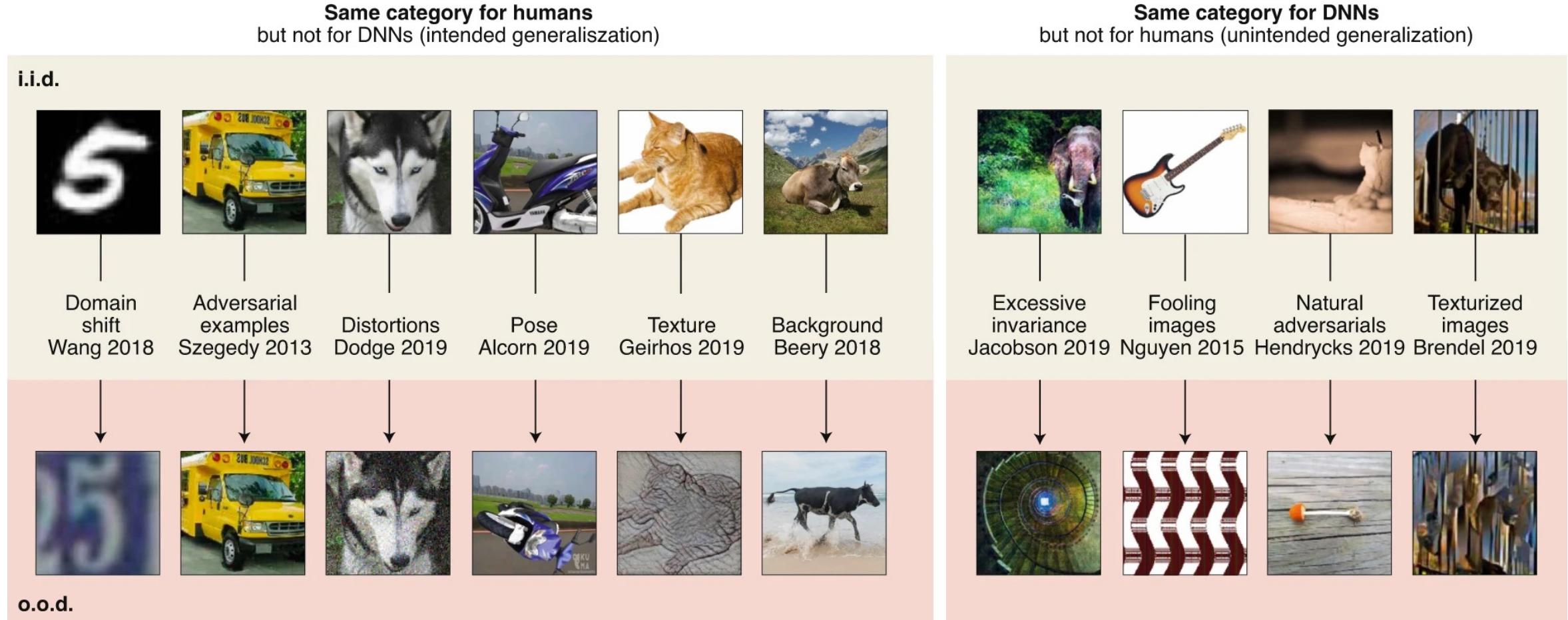
Pinball - relevance during game play



Identify valid vs invalid problem-solving behaviours

Lapuschkin et al. (2019) *Nature Communications* **10**, 1096.

Surface learning and unintended generalization



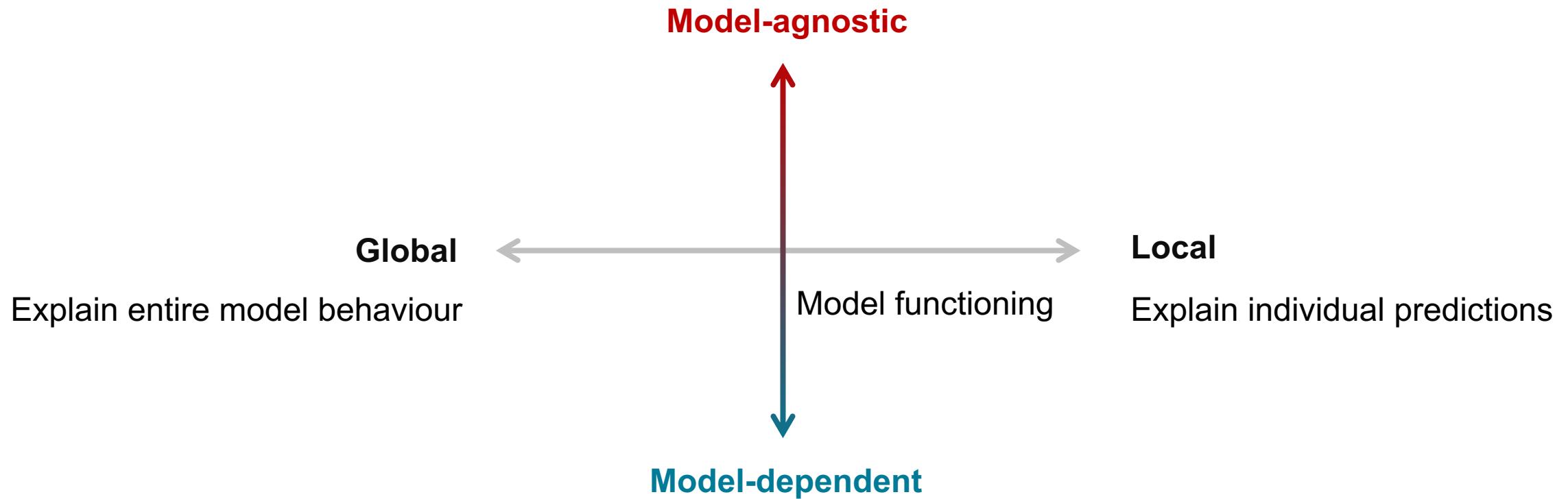
Performance metrics might just not be enough!

Geirhos et al. (2020). *Nature Machine Intelligence* 2, 665.

XAI nomenclature

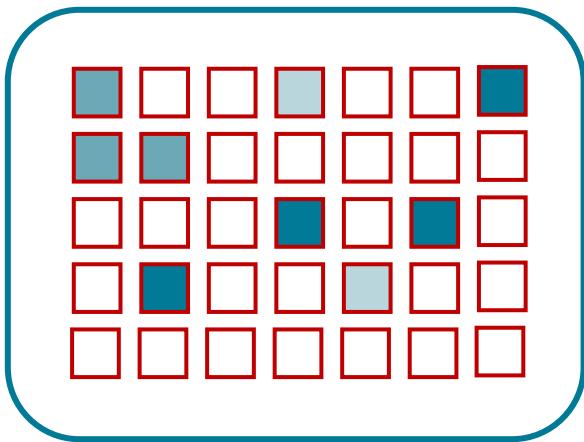


XAI nomenclature



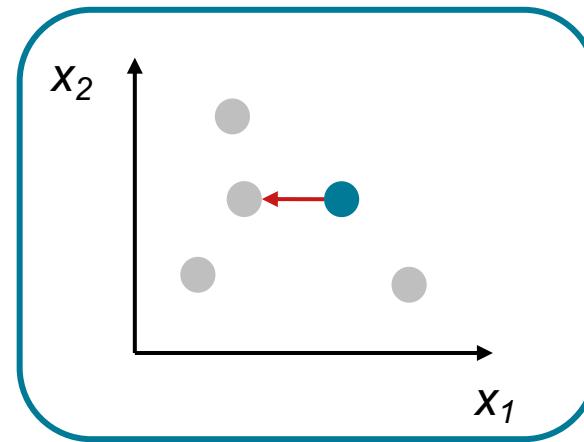
Flavours of XAI

Feature attribution



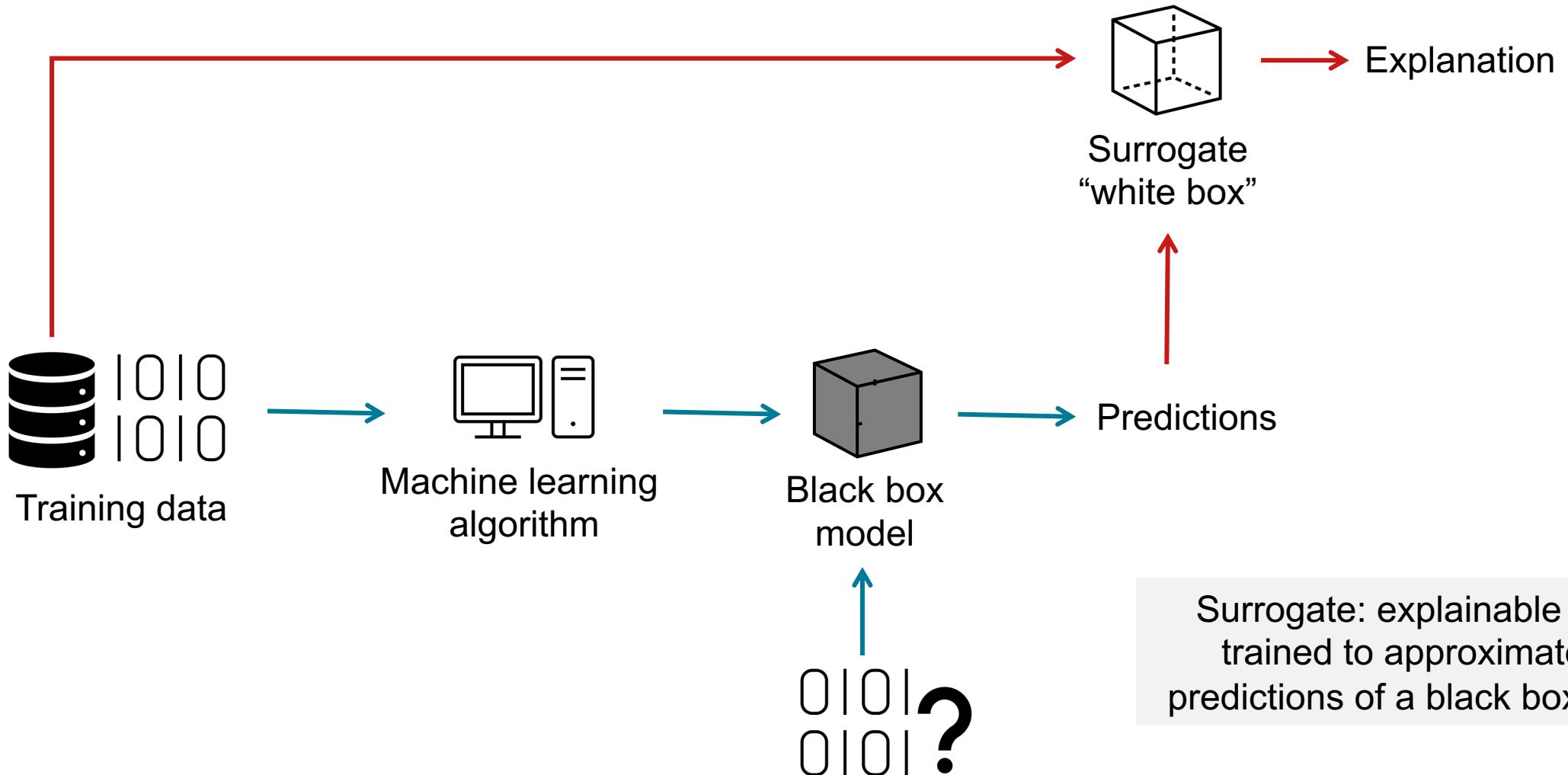
Relevance of the features in
the model behaviour.

Instance-based



Producing alternative inputs to
achieve a similar or different result.

Surrogate models (Feature attribution)



Local interpretable model-agnostic explanations (LIME)¹



“Why Should I Trust You?” Explaining the Predictions of Any Classifier

Marco Tulio Ribeiro
University of Washington
Seattle, WA 98105, USA
marcotcr@cs.uw.edu

Sameer Singh
University of Washington
Seattle, WA 98105, USA
sameer@cs.uw.edu

Carlos Guestrin
University of Washington
Seattle, WA 98105, USA
guestrin@cs.uw.edu

ABSTRACT

Despite widespread adoption, machine learning models remain mostly black boxes. Understanding the reasons behind predictions is, however, quite important in assessing *trust*, which is fundamental if one plans to take action based on a prediction, or when choosing whether to deploy a new model. Such understanding also provides insights into the model, which can be used to transform an untrustworthy model or prediction into a trustworthy one.

In this work, we propose LIME, a novel explanation technique that explains the predictions of *any* classifier in an interpretable and faithful manner, by learning an interpretable model locally around the prediction. We also propose a method to explain models by presenting representative individual predictions and their explanations in a non-redundant way, framing the task as a submodular optimization problem. We demonstrate the flexibility of these methods by explaining different models for text (e.g., random forests) and image classification (e.g., neural networks). We show the utility of explanations via novel experiments, both simulated and with human subjects, on various scenarios that require trust: deciding if one should trust a prediction, choosing between models, improving an untrustworthy classifier, and identifying why a classifier should not be trusted.

how much the human understands a model's behaviour, as opposed to seeing it as a black box.

Determining trust in individual predictions is an important problem when the model is used for decision making. When using machine learning for medical diagnosis [6] or terrorism detection, for example, predictions cannot be acted upon on blind faith, as the consequences may be catastrophic.

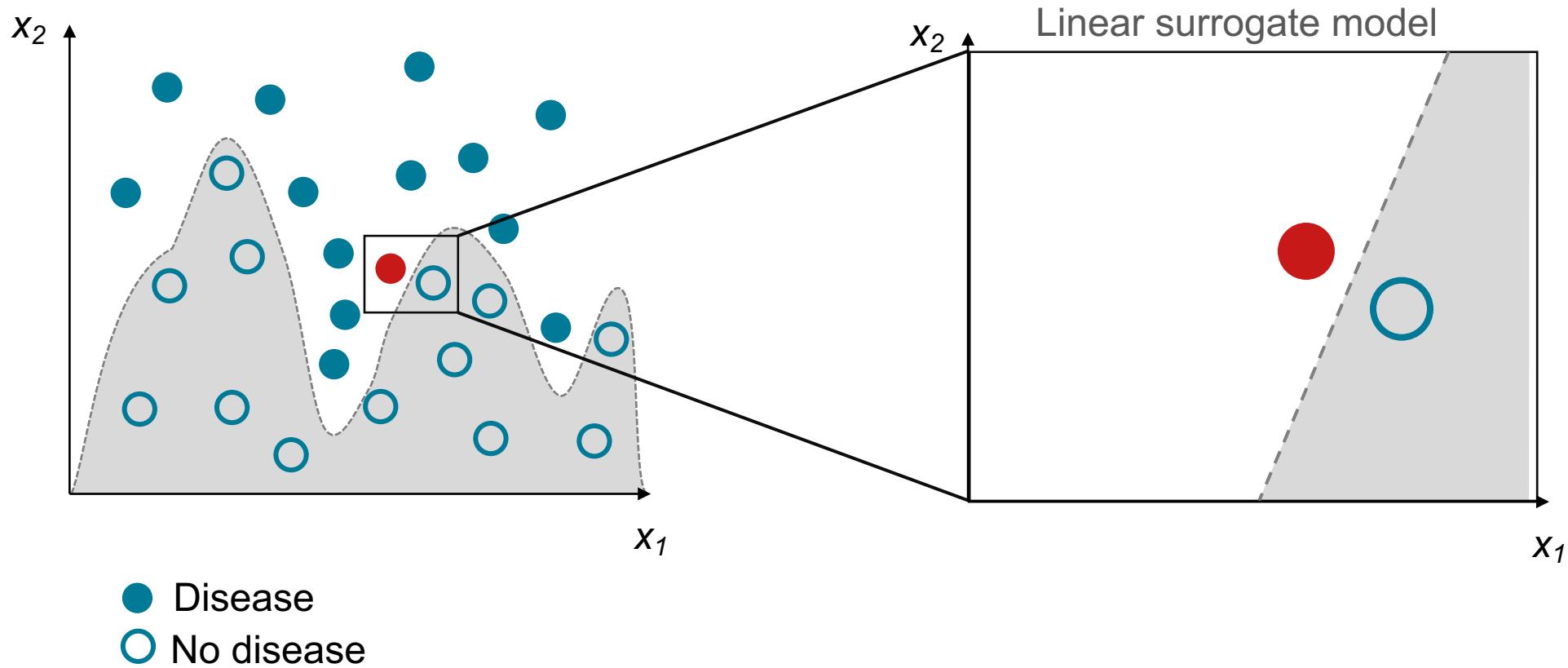
Apart from trusting individual predictions, there is also a need to evaluate the model as a whole before deploying it “in the wild”. To make this decision, users need to be confident that the model will perform well on real-world data, according to the metrics of interest. Currently, models are evaluated using accuracy metrics on an available validation dataset. However, real-world data is often significantly different, and further, the evaluation metric may not be indicative of the product’s goal. Inspecting individual predictions and their explanations is a worthwhile solution, in addition to such metrics. In this case, it is important to aid users by suggesting which instances to inspect, especially for large datasets.

In this paper, we propose providing explanations for individual predictions as a solution to the “trusting a prediction” problem, and selecting multiple such predictions (and explanations) as a solution to the “trusting the model” problem. Our main contributions are summarized as follows.

- Works on any “black-box” model
- Does not address the model internal functioning
- “Locally-faithful” explanations

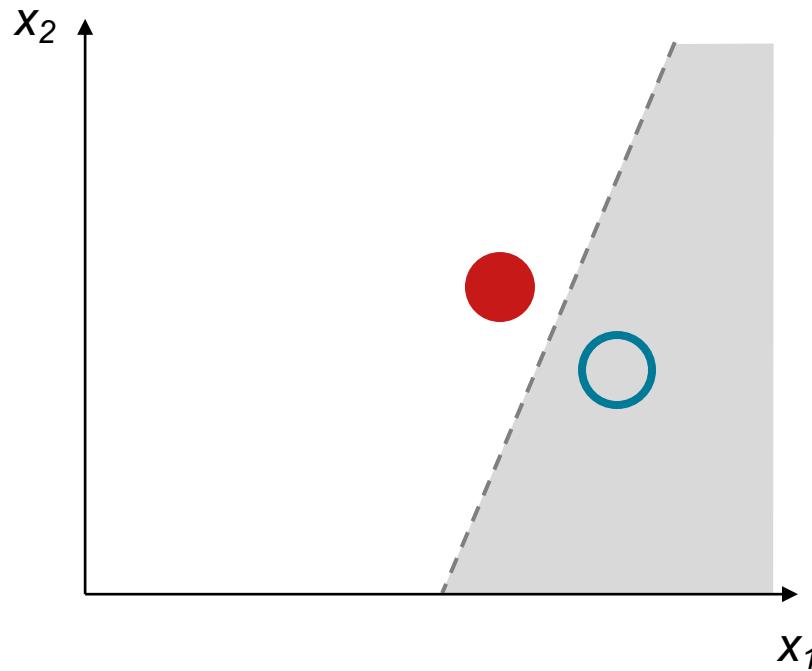
¹Ribeiro et al. (2016). Proc. ACM SIGKDD Int, 1135.

Local interpretable model-agnostic explanations (LIME)¹



¹Ribeiro et al. (2016). Proc. ACM SIGKDD Int, 1135.

Local interpretable model-agnostic explanations (LIME)¹

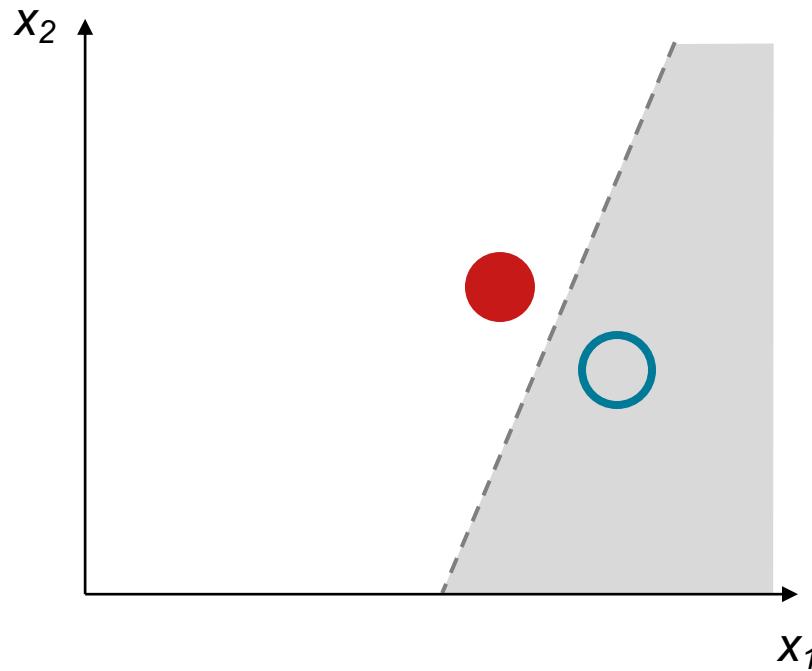


$$\xi(x) \triangleq \arg \min_{g \in G} L(f, g, \pi_x) \triangleq \Omega(g)$$

Approximation Complexity

¹Ribeiro et al. (2016). Proc. ACM SIGKDD Int, 1135.

Local interpretable model-agnostic explanations (LIME)¹



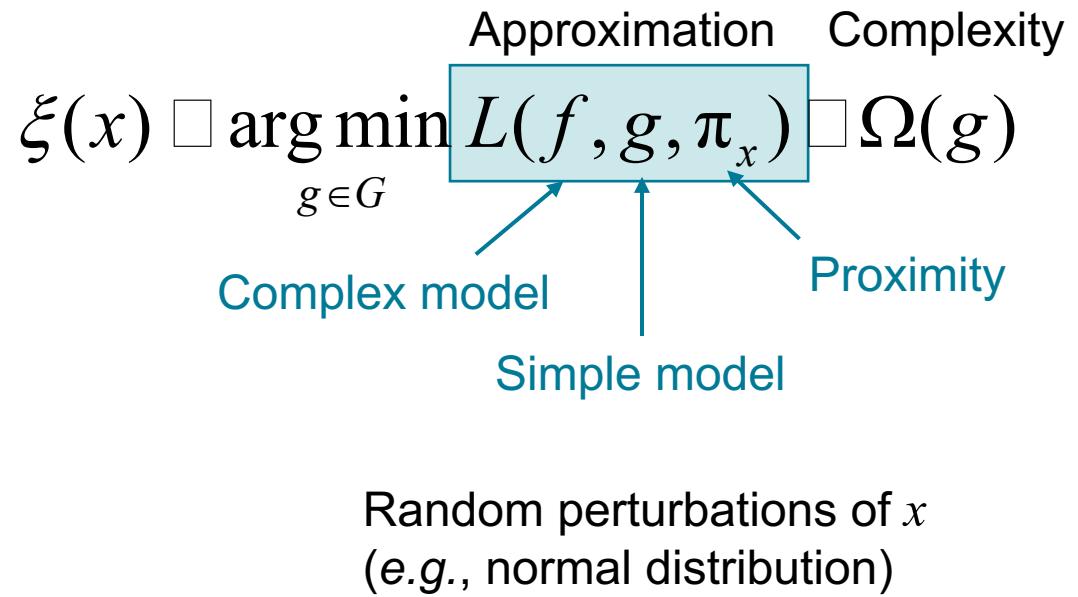
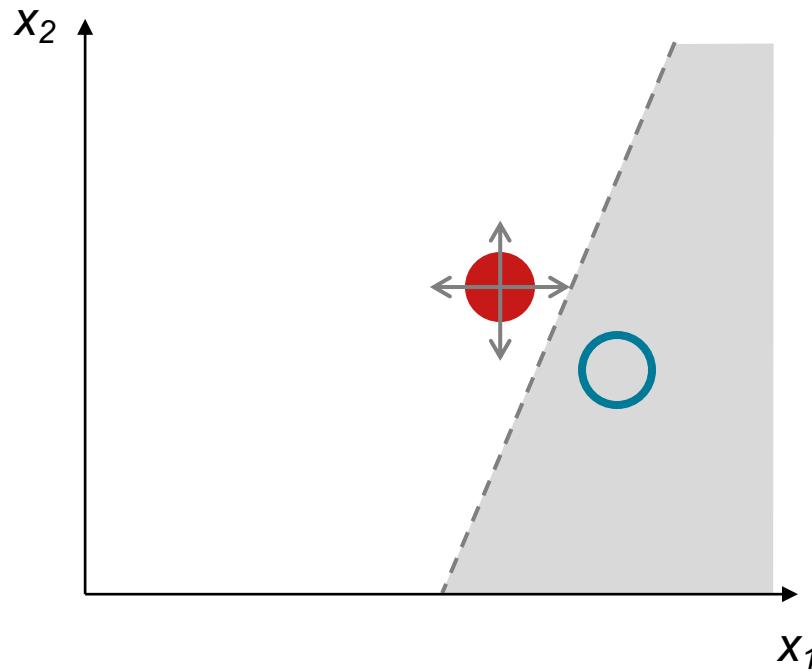
$$\xi(x) \triangleq \arg \min_{g \in G} L(f, g, \pi_x) \triangleq \Omega(g)$$

Approximation Complexity
Complex model Simple model
Proximity

A diagram illustrating the LIME optimization process. The equation $\xi(x) \triangleq \arg \min_{g \in G} L(f, g, \pi_x) \triangleq \Omega(g)$ is shown. The term $L(f, g, \pi_x)$ is highlighted with a blue box and labeled "Approximation". The term $\Omega(g)$ is also highlighted with a blue box and labeled "Complexity". Arrows point from the labels "Complex model" and "Simple model" to the terms $L(f, g, \pi_x)$ and $\Omega(g)$ respectively. An arrow points from the label "Proximity" to the variable π_x .

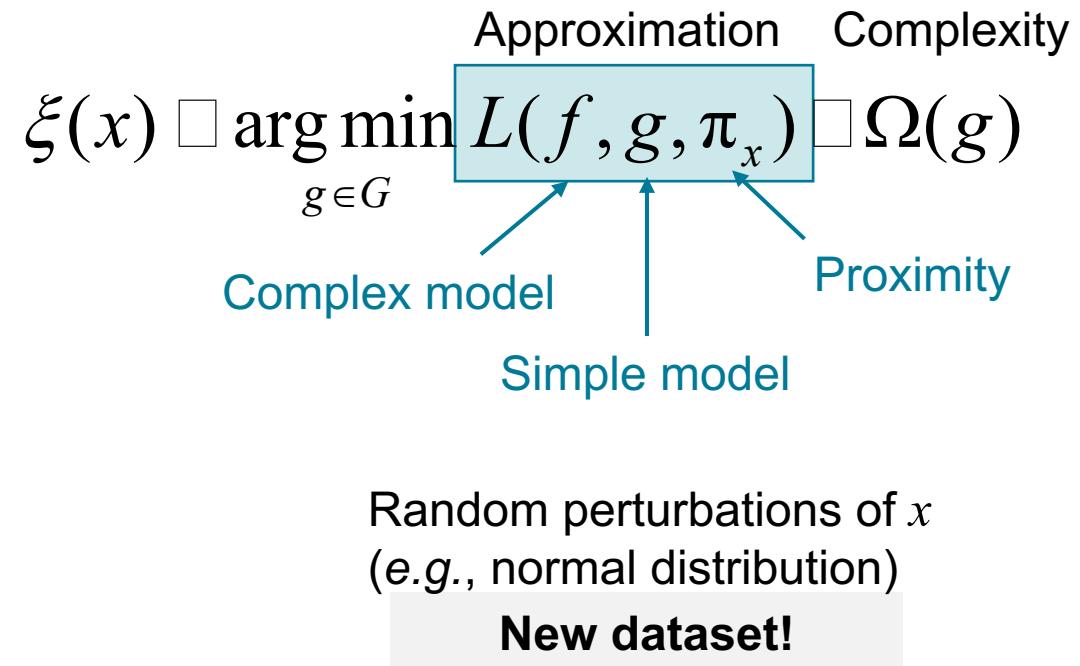
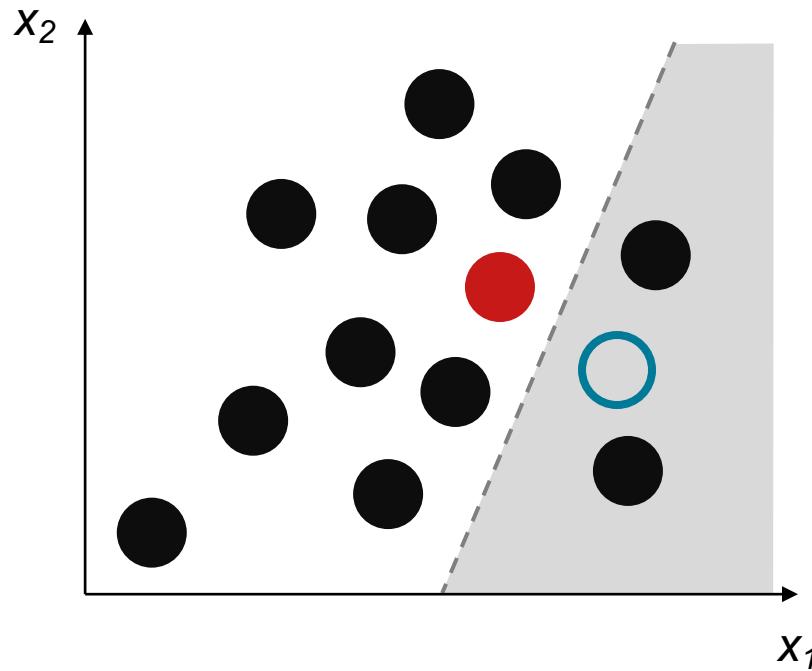
¹Ribeiro et al. (2016). Proc. ACM SIGKDD Int, 1135.

Local interpretable model-agnostic explanations (LIME)¹



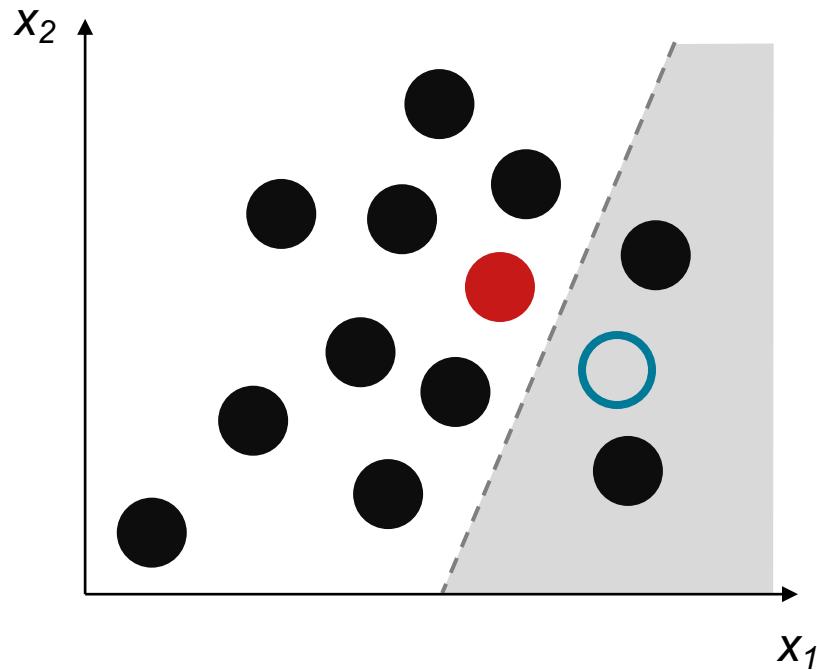
¹Ribeiro et al. (2016). Proc. ACM SIGKDD Int, 1135.

Local interpretable model-agnostic explanations (LIME)¹



¹Ribeiro et al. (2016). Proc. ACM SIGKDD Int, 1135.

Local interpretable model-agnostic explanations (LIME)¹



$$\xi(x) \triangleq \arg \min_{g \in G} L(f, g, \pi_x)$$

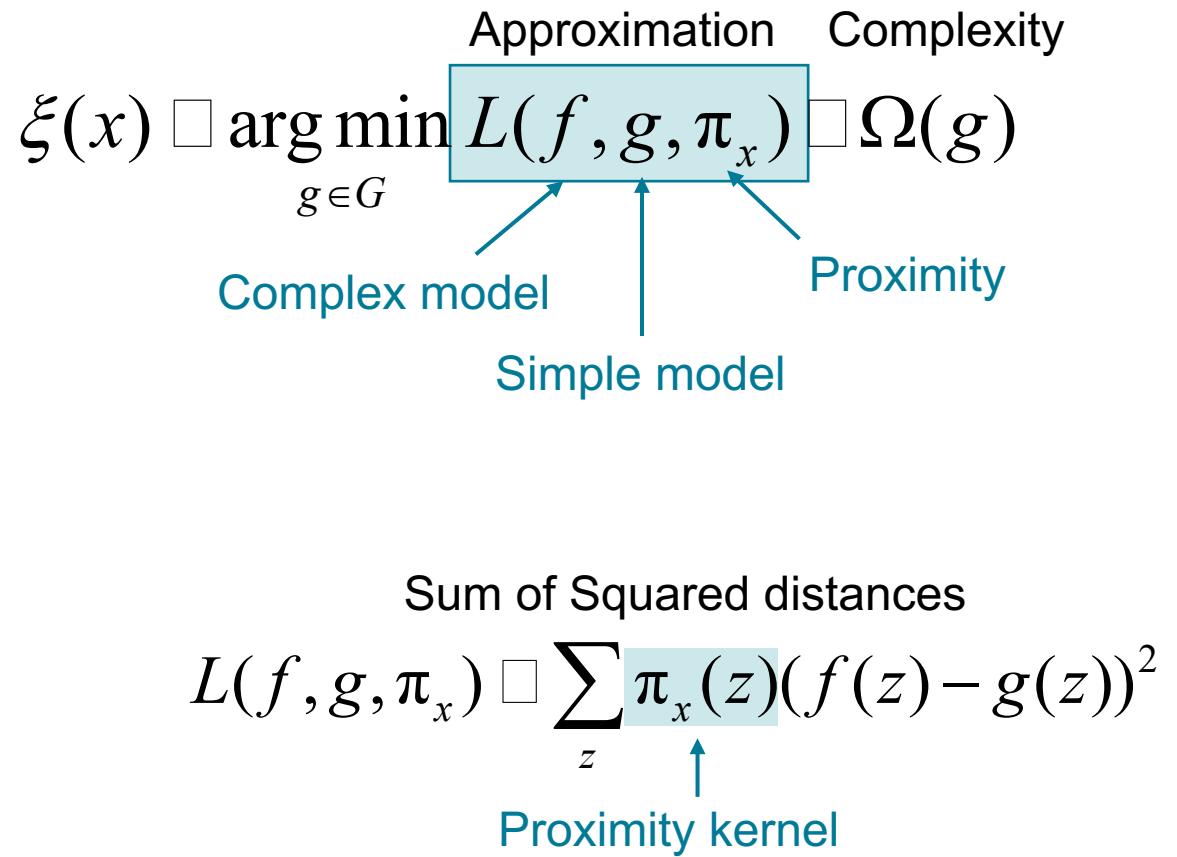
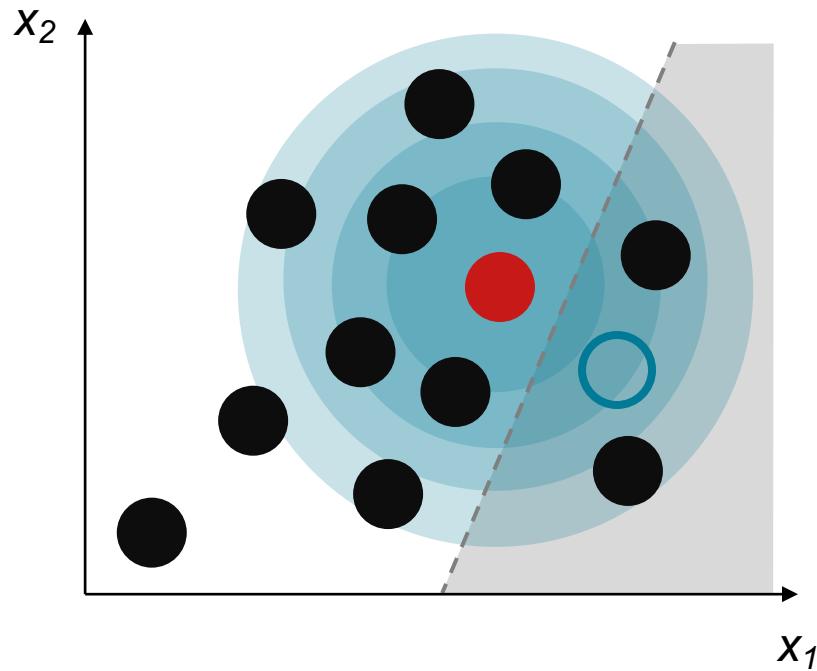
Approximation Complexity
Complex model Simple model Proximity

Sum of Squared distances

$$L(f, g, \pi_x) \triangleq \sum_z \pi_x(z)(f(z) - g(z))^2$$

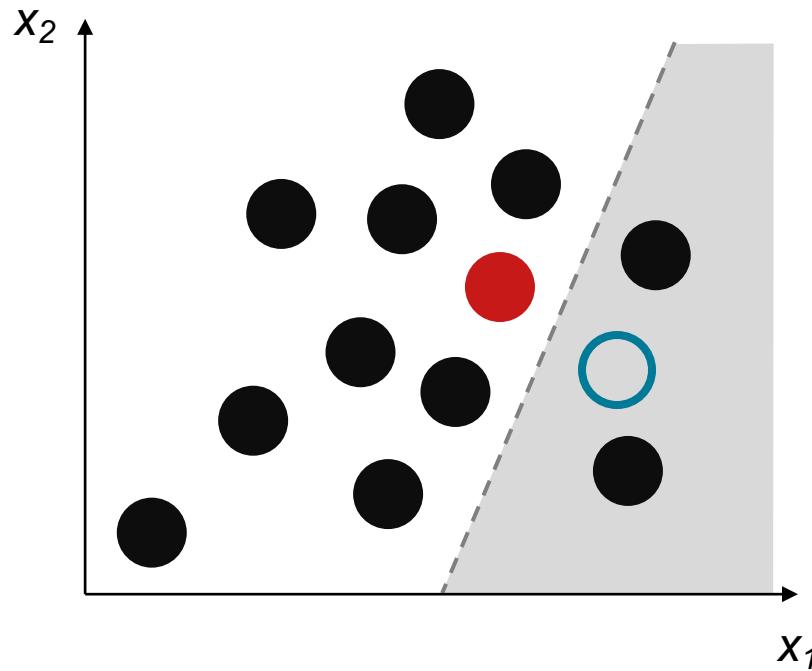
¹Ribeiro et al. (2016). Proc. ACM SIGKDD Int, 1135.

Local interpretable model-agnostic explanations (LIME)¹



¹Ribeiro et al. (2016). Proc. ACM SIGKDD Int, 1135.

Local interpretable model-agnostic explanations (LIME)¹

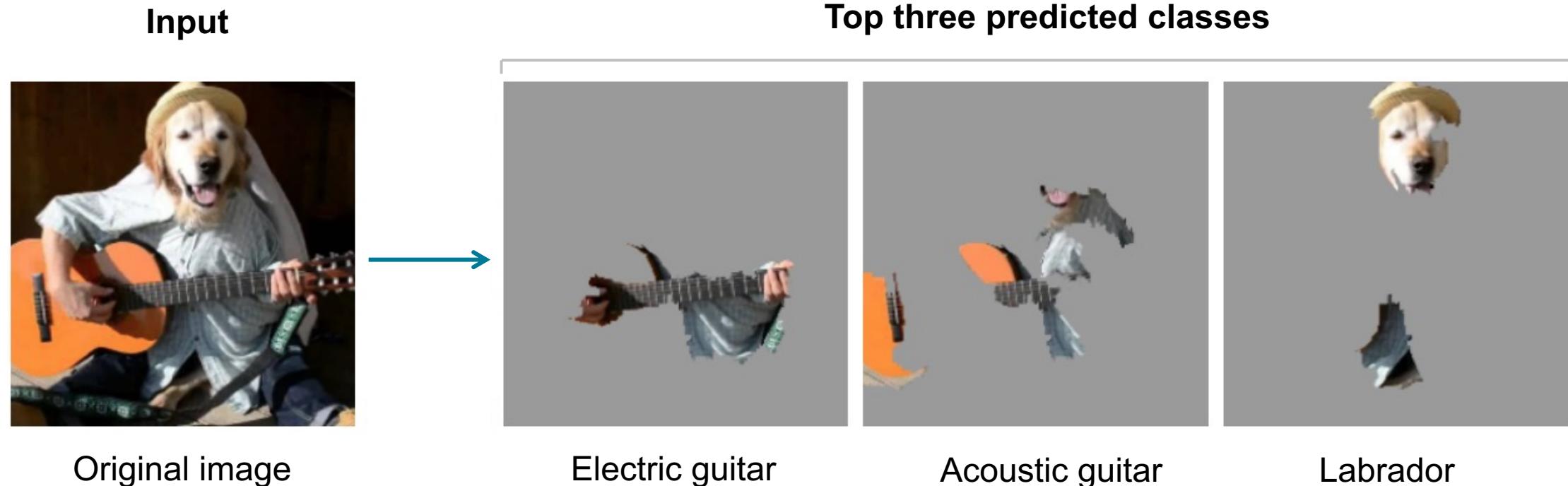


$$\xi(x) \triangleq \arg \min_{g \in G} L(f, g, \pi_x) \triangleq \Omega(g)$$

Regularization
Sparse linear models
(e.g., LASSO)

¹Ribeiro et al. (2016). Proc. ACM SIGKDD Int, 1135.

Local interpretable model-agnostic explanations (LIME)¹



¹Ribeiro et al. (2016). Proc. ACM SIGKDD Int, 1135.

Surrogate models

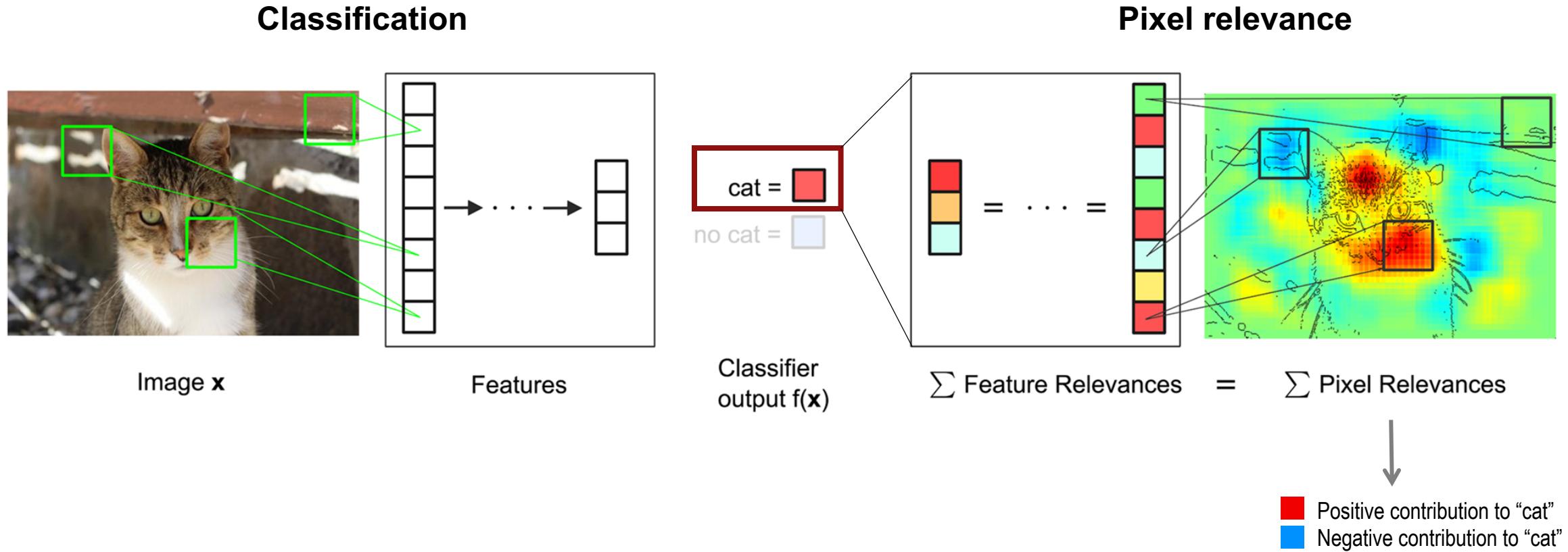
Pros

- Easy to understand
- Can use *some* known information
- Model agnostic

Cons

- The surrogate model is not the same as our original model
- Approximated explanations
- Conclusions about the black box model, not the data

Pixel attribution methods (Feature attribution)



Bach et al. (2015). PLoS one **10**, e0130140.

Layer-wise relevance propagation (LRP)



RESEARCH ARTICLE

On Pixel-Wise Explanations for Non-Linear Classifier Decisions by Layer-Wise Relevance Propagation

Sebastian Bach^{1,2*}, Alexander Binder^{2,5*}, Grégoire Montavon², Frederick Klauschen³, Klaus-Robert Müller^{2,4*}, Wojciech Samek^{1,2}

1 Machine Learning Group, Fraunhofer Heinrich Hertz Institute, Berlin, Germany, **2** Machine Learning Group, Technische Universität Berlin, Berlin, Germany, **3** Charité University Hospital, Berlin, Germany, **4** Department of Brain and Cognitive Engineering, Korea University, Seoul, Korea, **5** ISTD Pillar, Singapore University of Technology and Design (SUTD), Singapore

* These authors contributed equally to this work.

* sebastian.bach@hhi.fraunhofer.de (SB), klaus-robert.mueller@tu-berlin.de (KM), wojciech.samek@hhi.fraunhofer.de (WS)



OPEN ACCESS

Citation: Bach S, Binder A, Montavon G, Klauschen F, Müller K-R, Samek W (2015) On Pixel-Wise Explanations for Non-Linear Classifier Decisions by Layer-Wise Relevance Propagation. PLoS ONE 10(7): e0130140. doi:10.1371/journal.pone.0130140

Editor: Oscar Deniz Suarez, Universidad de Castilla-La Mancha, SPAIN

Received: May 19, 2014

Accepted: May 15, 2015

Published: July 10, 2015

Copyright: © 2015 Bach et al. This is an open access article distributed under the terms of the [Creative Commons Attribution License](#), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

Data Availability Statement: All relevant data are

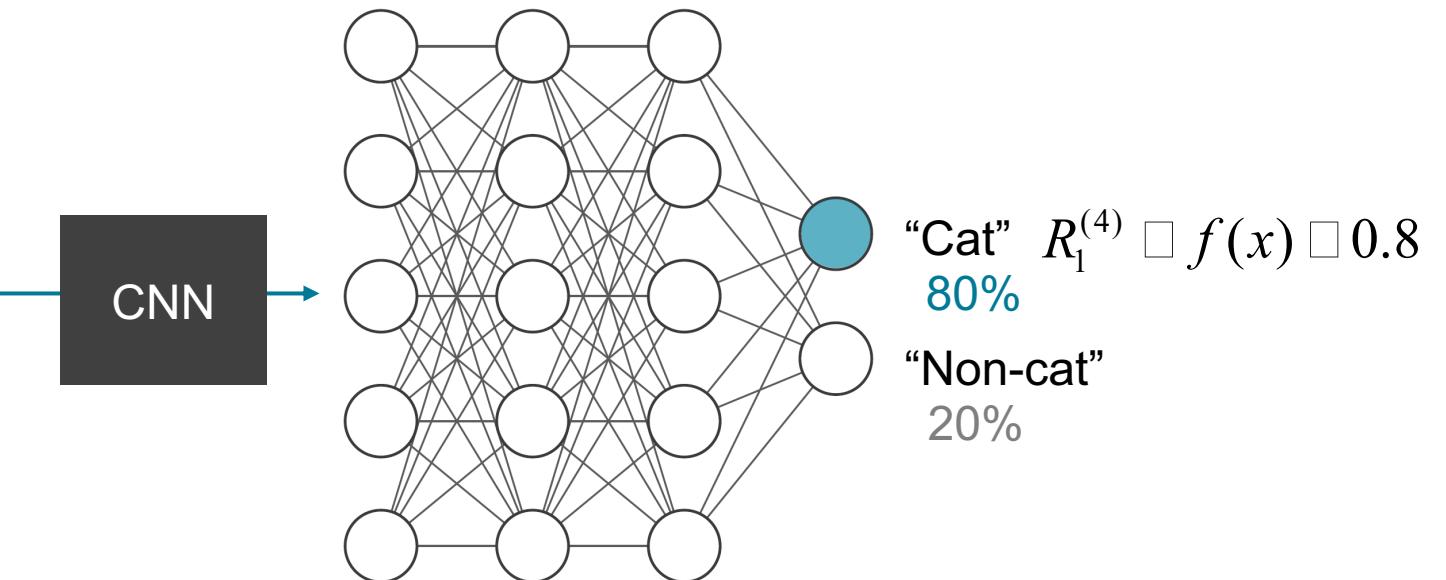
Abstract

Understanding and interpreting classification decisions of automated image classification systems is of high value in many applications, as it allows to verify the reasoning of the system and provides additional information to the human expert. Although machine learning methods are solving very successfully a plethora of tasks, they have in most cases the disadvantage of acting as a black box, not providing any information about what made them arrive at a particular decision. This work proposes a general solution to the problem of understanding classification decisions by pixel-wise decomposition of nonlinear classifiers. We introduce a methodology that allows to visualize the contributions of single pixels to predictions for kernel-based classifiers over Bag of Words features and for multilayered neural networks. These pixel contributions can be visualized as heatmaps and are provided to a human expert who can intuitively not only verify the validity of the classification decision, but also focus further analysis on regions of potential interest. We evaluate our method for classifiers trained on PASCAL VOC 2009 images, synthetic image data containing geometric shapes, the MNIST handwritten digits data set and for the pre-trained ImageNet model available as part of the Caffe open source package.

- Which pixels contributed to a certain prediction
- Model-specific (neural networks)

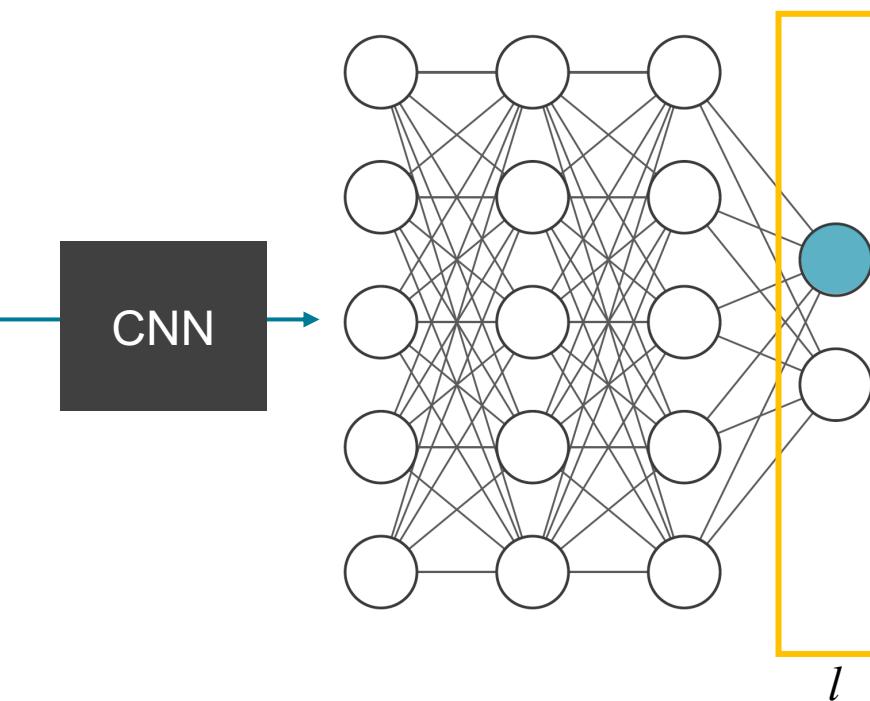
Bach et al. (2015). *PLoS one* **10**, e0130140.

Layer-wise relevance propagation (LRP)



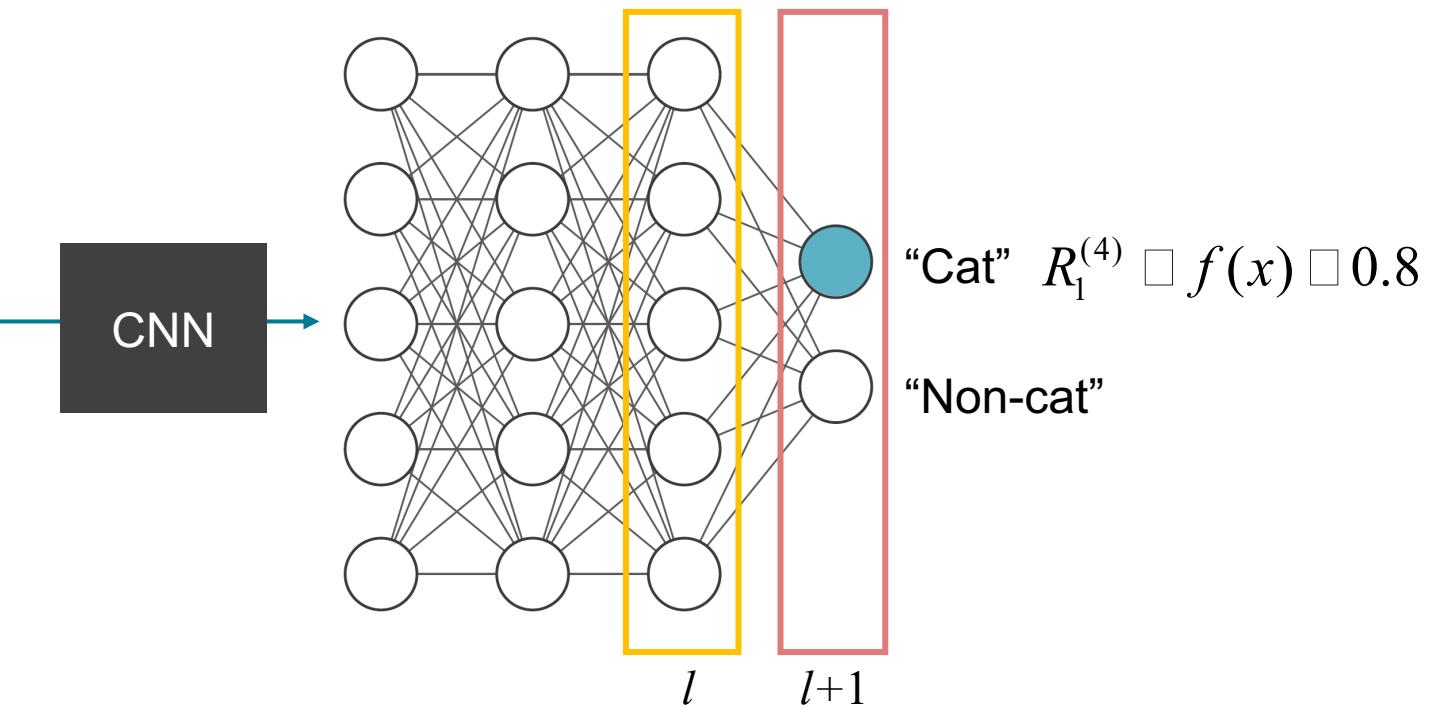
Example from: <https://www.youtube.com/watch?v=PDRewtcqmal>
Bach et al. (2015). PloS one 10, e0130140.

Layer-wise relevance propagation (LRP)



Example from: <https://www.youtube.com/watch?v=PDRewtcqmal>
Bach et al. (2015). PloS one 10, e0130140.

Layer-wise relevance propagation (LRP)

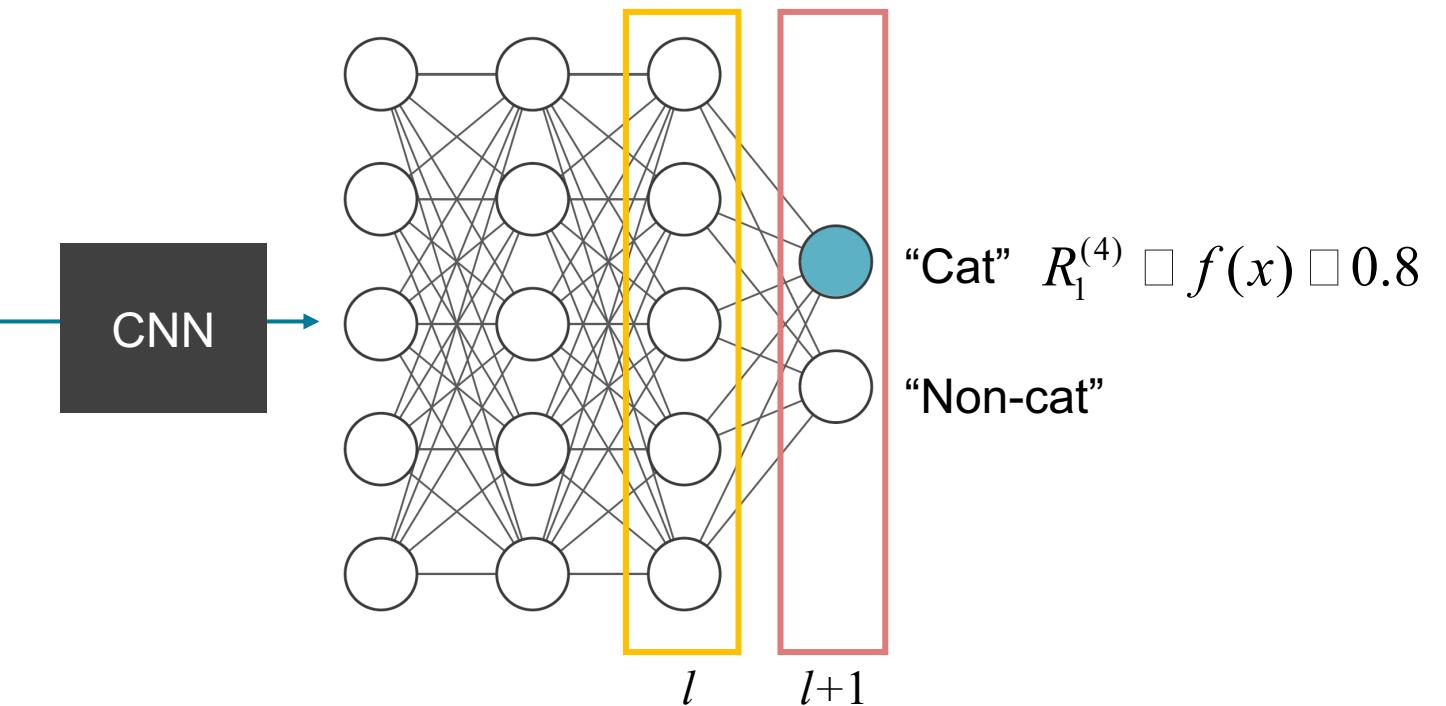


$$R_i^{(l)} \square \sum_j \frac{z_{ij}}{\sum_k z_{kj}} R_j^{(l+1)} \text{ with } z_{ik} \square x_i^{(l)} w_{ij}^{(l,l+1)}$$

Neurons in layer l
Neurons in layer $l+1$

Example from: <https://www.youtube.com/watch?v=PDREwtcqmal>
Bach et al. (2015). PloS one 10, e0130140.

Layer-wise relevance propagation (LRP)

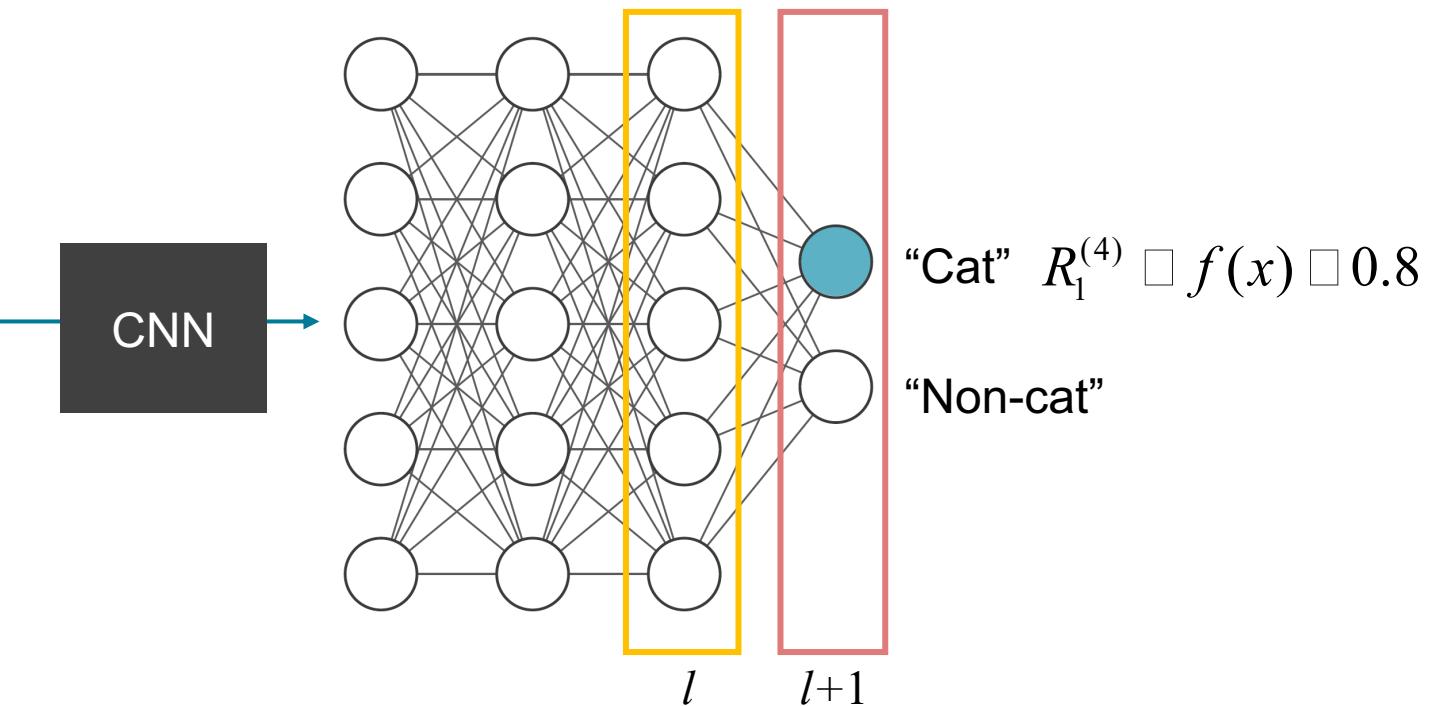


$$R_i^{(l)} \square \sum_j \frac{z_{ij}}{\sum_k z_{kj}} R_j^{(l+1)} \text{ with } z_{ik} \square x_i^{(l)} w_{ij}^{(l,l+1)}$$

Relative activation

Example from: <https://www.youtube.com/watch?v=PDREwtcqmal>
Bach et al. (2015). PloS one 10, e0130140.

Layer-wise relevance propagation (LRP)

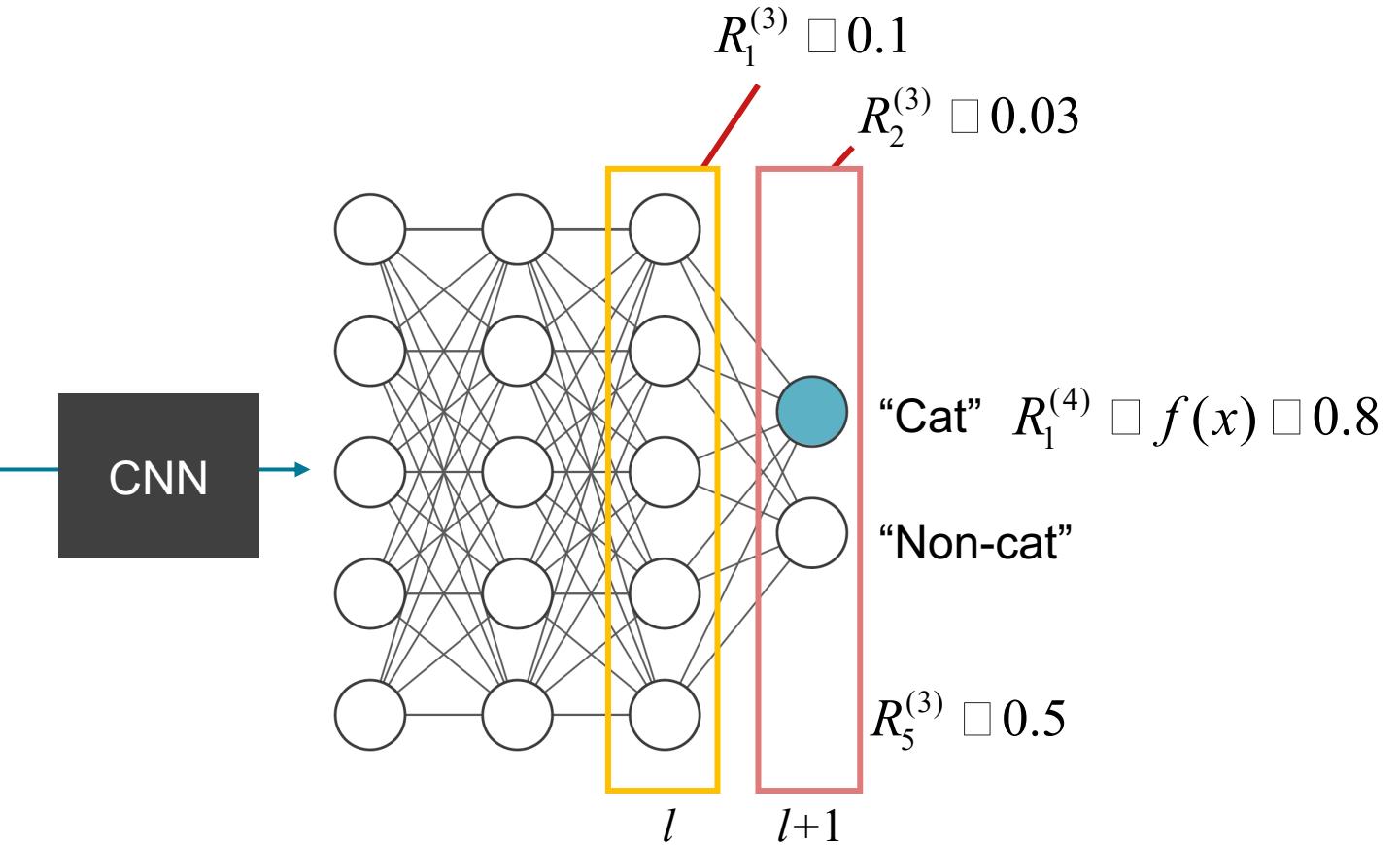


$$R_i^{(l)} \square \sum_j \frac{z_{ij}}{\sum_k z_{kj}} R_j^{(l+1)} \quad \text{with} \quad z_{ik} \square x_i^{(l)} w_{ij}^{(l, l+1)}$$

Propagate the previous relevance backwards

Example from: <https://www.youtube.com/watch?v=PDREwtcqmal>
Bach et al. (2015). PloS one 10, e0130140.

Layer-wise relevance propagation (LRP)



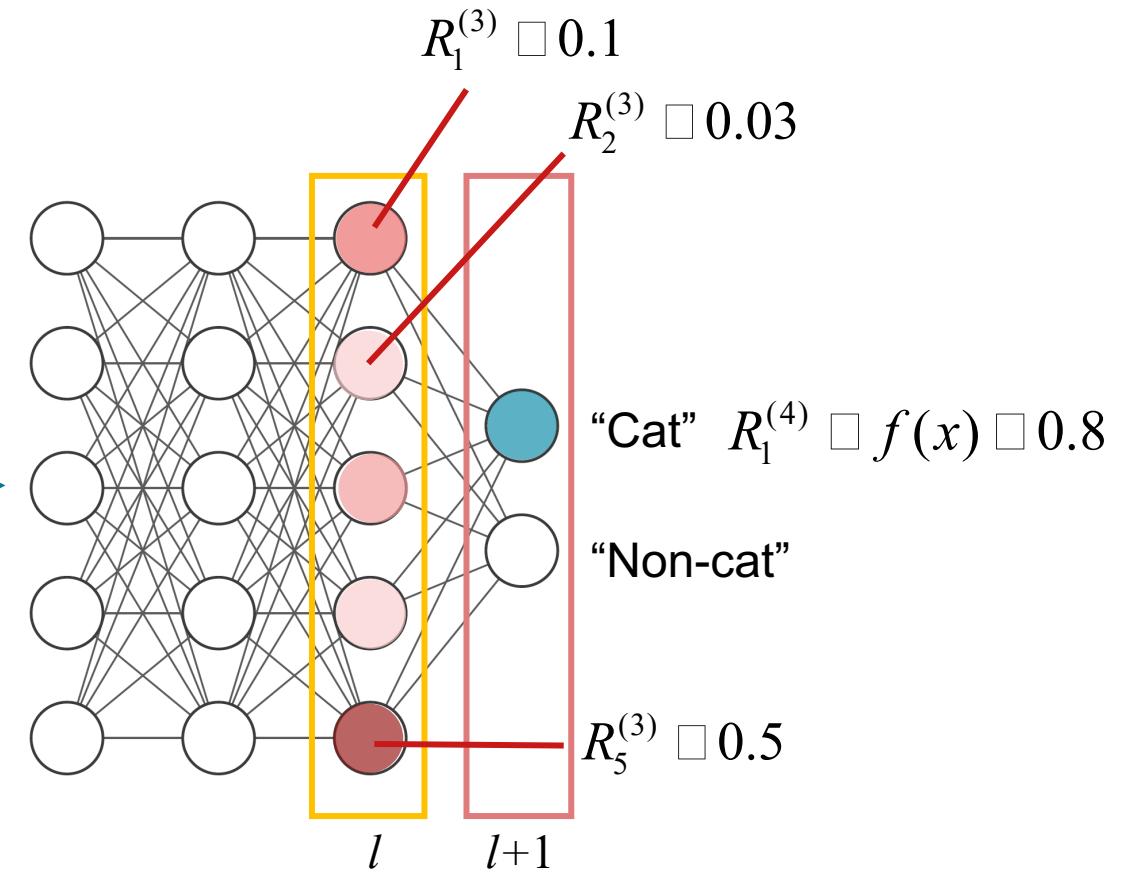
$$R_i^{(l)} \square \sum_j \frac{z_{ij}}{\sum_k z_{kj}} R_j^{(l+1)} \quad \text{with } z_{ik} \square x_i^{(l)} w_{ij}^{(l,l+1)}$$

Example from: <https://www.youtube.com/watch?v=PDREwtcqmal>
Bach et al. (2015). PloS one 10, e0130140.

Layer-wise relevance propagation (LRP)



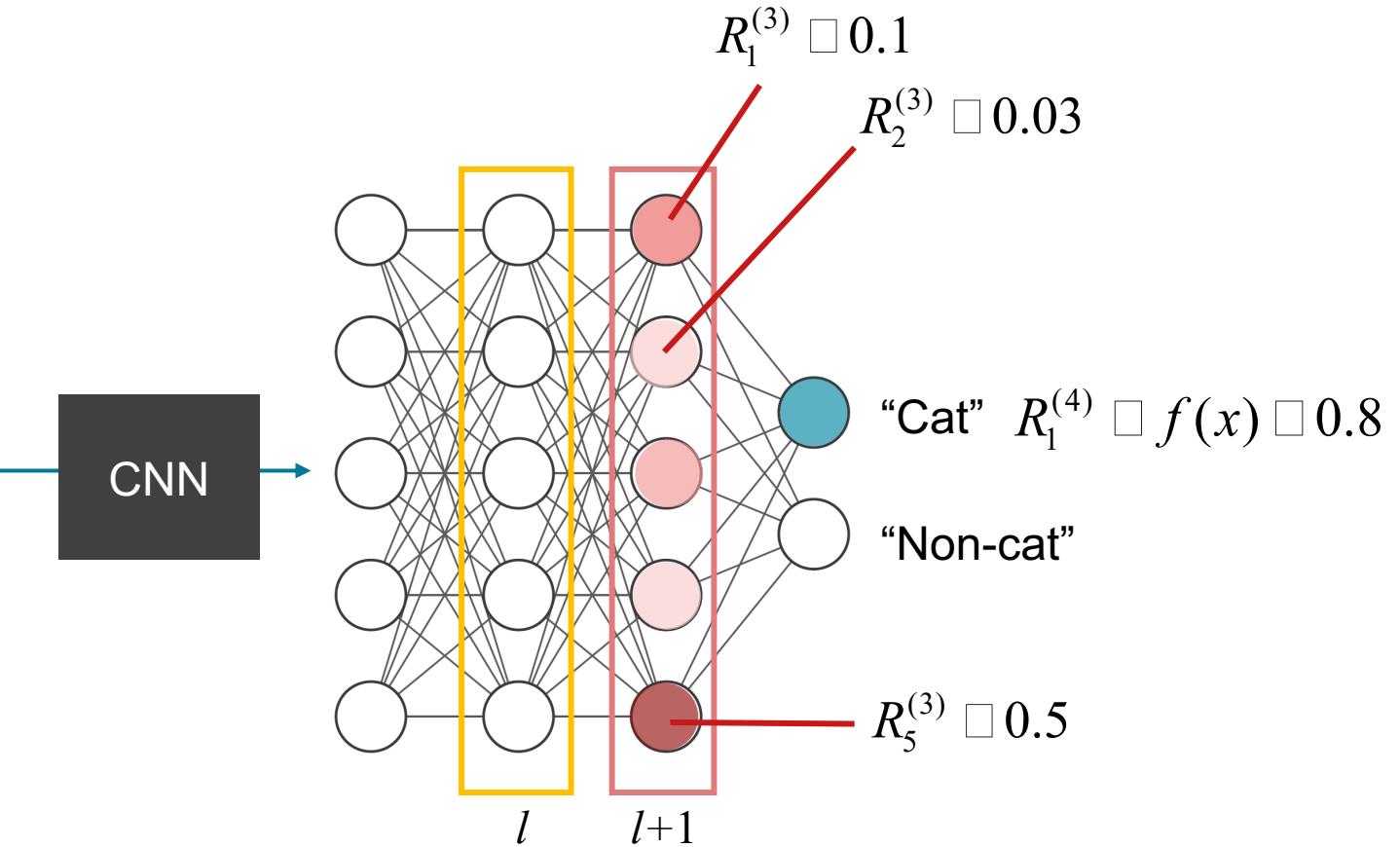
CNN



$$R_i^{(l)} \square \sum_j \frac{z_{ij}}{\sum_k z_{kj}} R_j^{(l+1)} \quad \text{with} \quad z_{ik} \square x_i^{(l)} w_{ij}^{(l,l+1)}$$

Example from: <https://www.youtube.com/watch?v=PDREwtcqmal>
Bach et al. (2015). PloS one 10, e0130140.

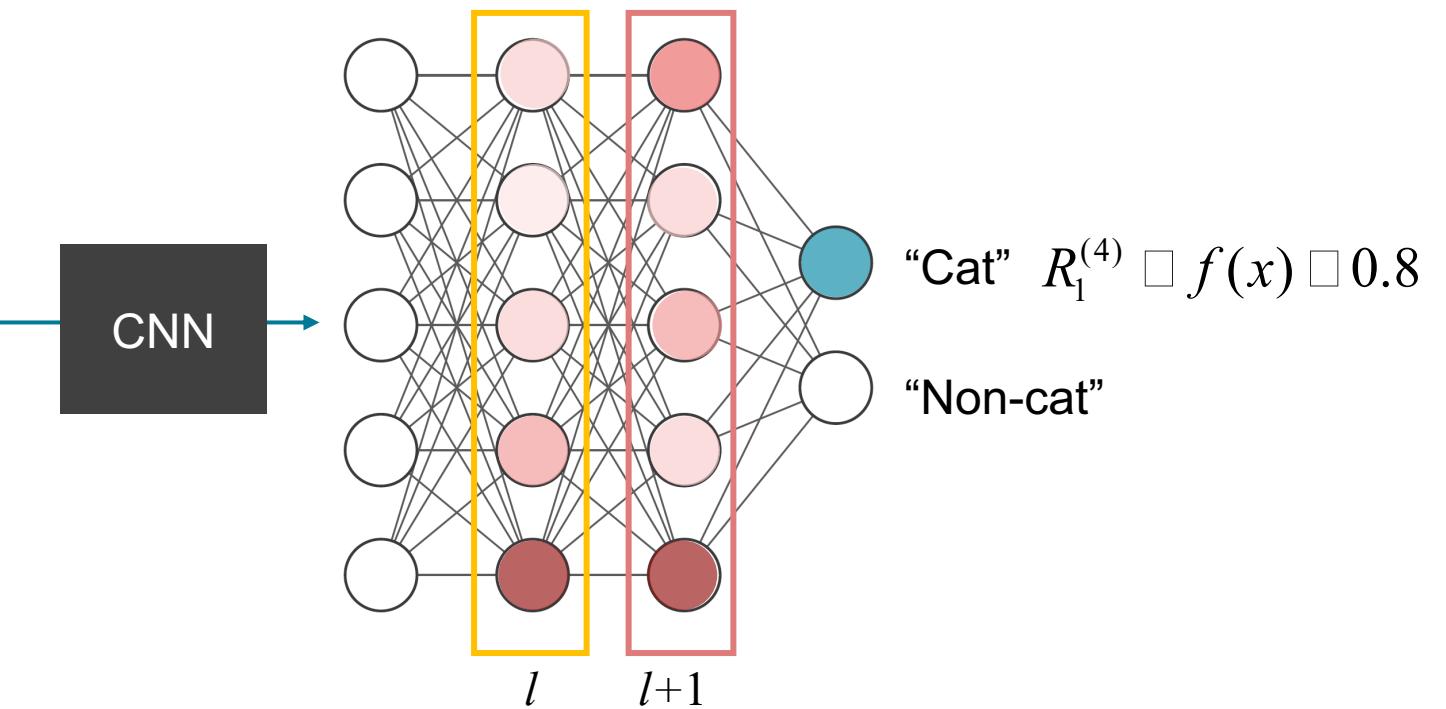
Layer-wise relevance propagation (LRP)



$$R_i^{(l)} \square \sum_j \frac{z_{ij}}{\sum_k z_{kj}} R_j^{(l+1)} \quad \text{with} \quad z_{ik} \square x_i^{(l)} w_{ij}^{(l,l+1)}$$

Example from: <https://www.youtube.com/watch?v=PDREwtcqmal>
Bach et al. (2015). PloS one 10, e0130140.

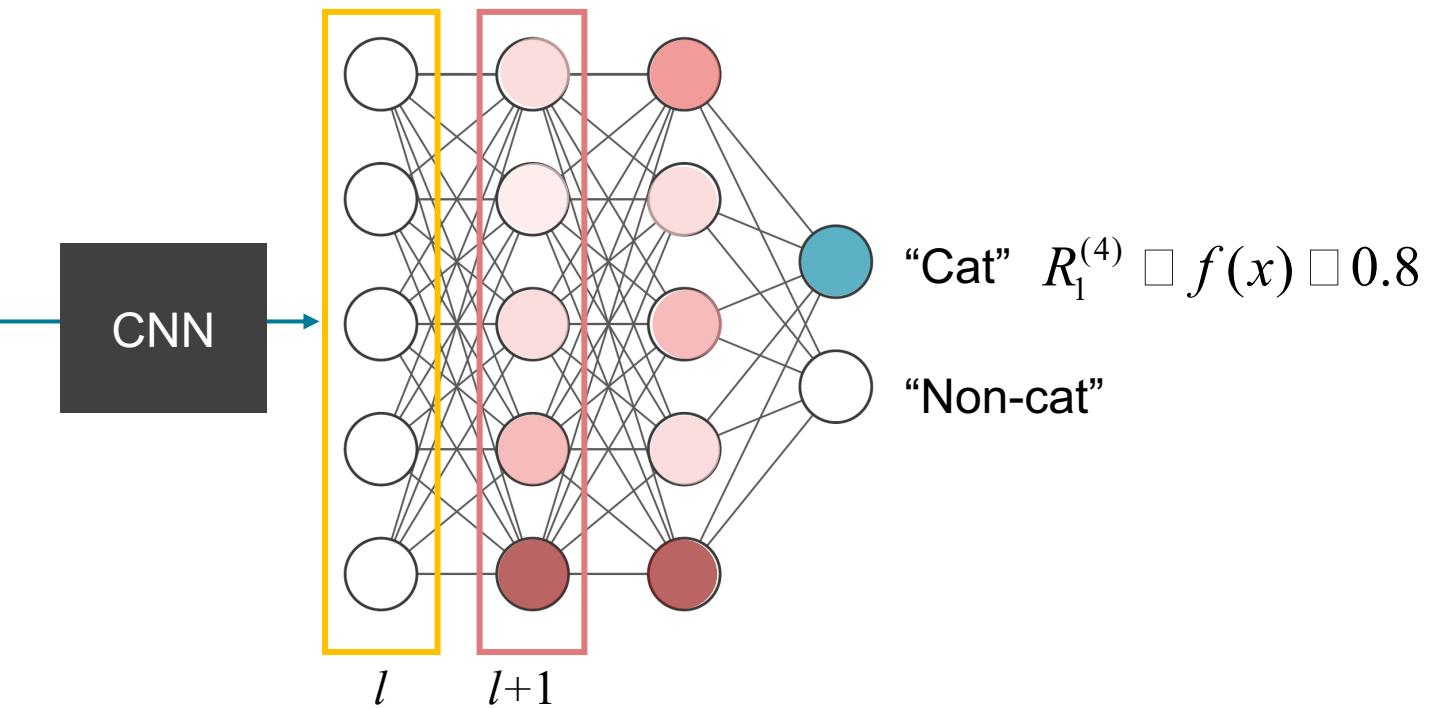
Layer-wise relevance propagation (LRP)



$$R_i^{(l)} \square \sum_j \frac{z_{ij}}{\sum_k z_{kj}} R_j^{(l+1)} \quad \text{with} \quad z_{ik} \square x_i^{(l)} w_{ij}^{(l, l+1)}$$

Example from: <https://www.youtube.com/watch?v=PDREwtcqmal>
Bach et al. (2015). PloS one 10, e0130140.

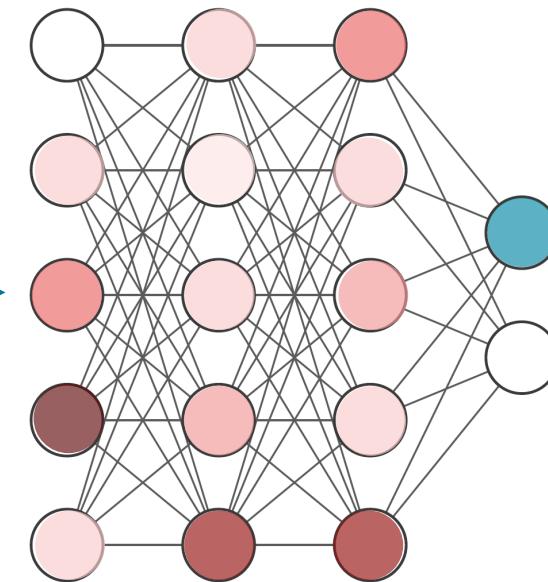
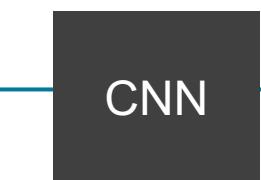
Layer-wise relevance propagation (LRP)



$$R_i^{(l)} \square \sum_j \frac{z_{ij}}{\sum_k z_{kj}} R_j^{(l+1)} \quad \text{with} \quad z_{ik} \square x_i^{(l)} w_{ij}^{(l,l+1)}$$

Example from: <https://www.youtube.com/watch?v=PDREwtcqmal>
Bach et al. (2015). PloS one 10, e0130140.

Layer-wise relevance propagation (LRP)



"Cat" $R_1^{(4)} \square f(x) \square 0.8$
"Non-cat"

$$R_i^{(l)} \square \sum_j \frac{z_{ij}}{\sum_k z_{kj}} R_j^{(l-1)} \quad \text{with} \quad z_{ik} \square x_i^{(l)} w_{ij}^{(l,l-1)}$$

Example from: <https://www.youtube.com/watch?v=PDRewtcqmal>
Bach et al. (2015). *PLoS one* **10**, e0130140.

Layer-wise relevance propagation (LRP)

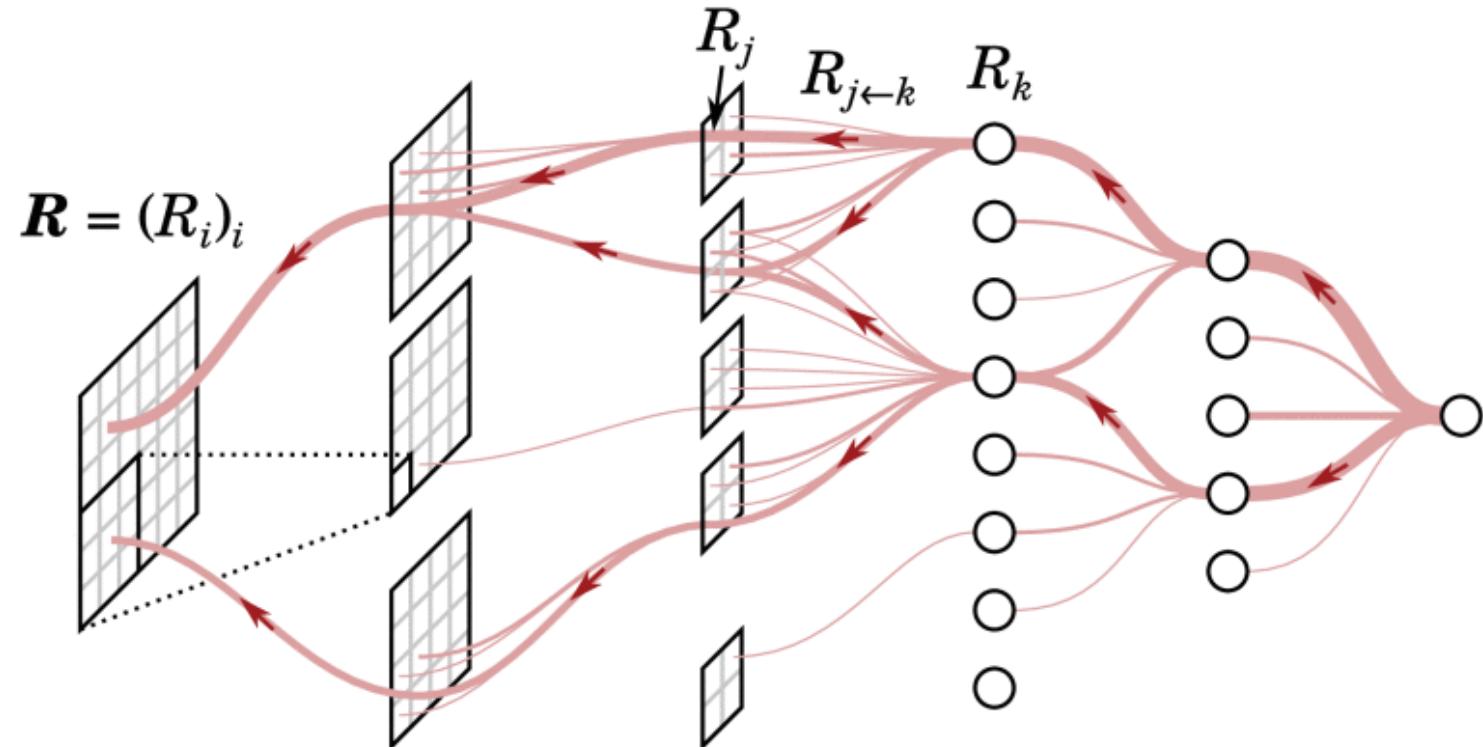
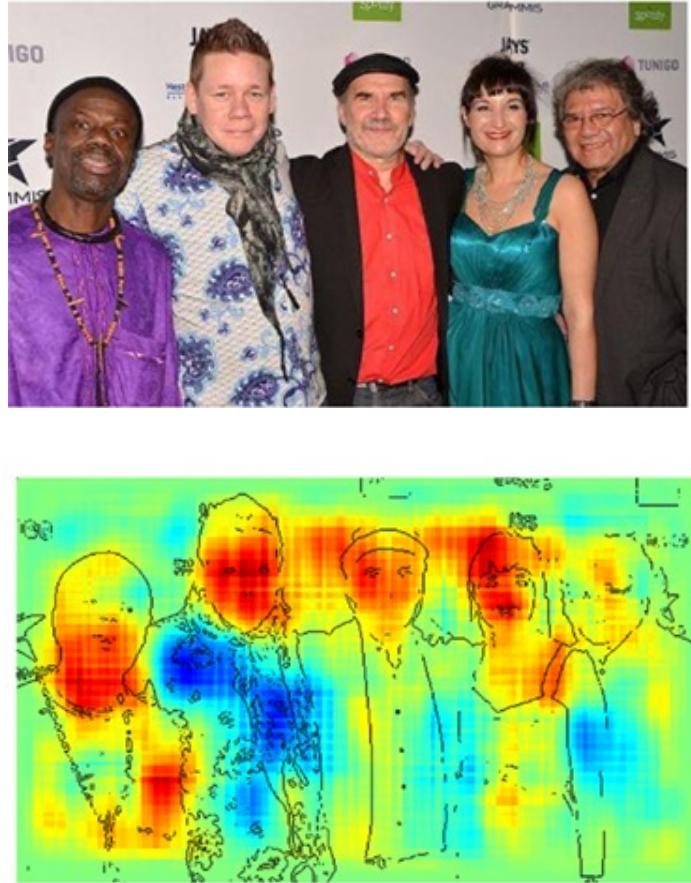
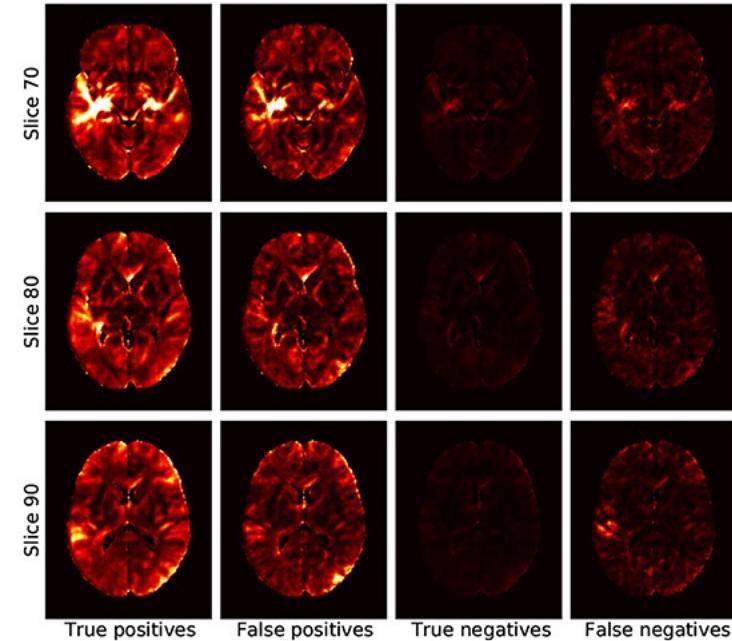


Figure from: Samek et al. (2021). *IEEE Proceedings* **109**, 247.
Bach et al. (2015). *PLoS one* **10**, e0130140.

Layer-wise relevance propagation (LRP)

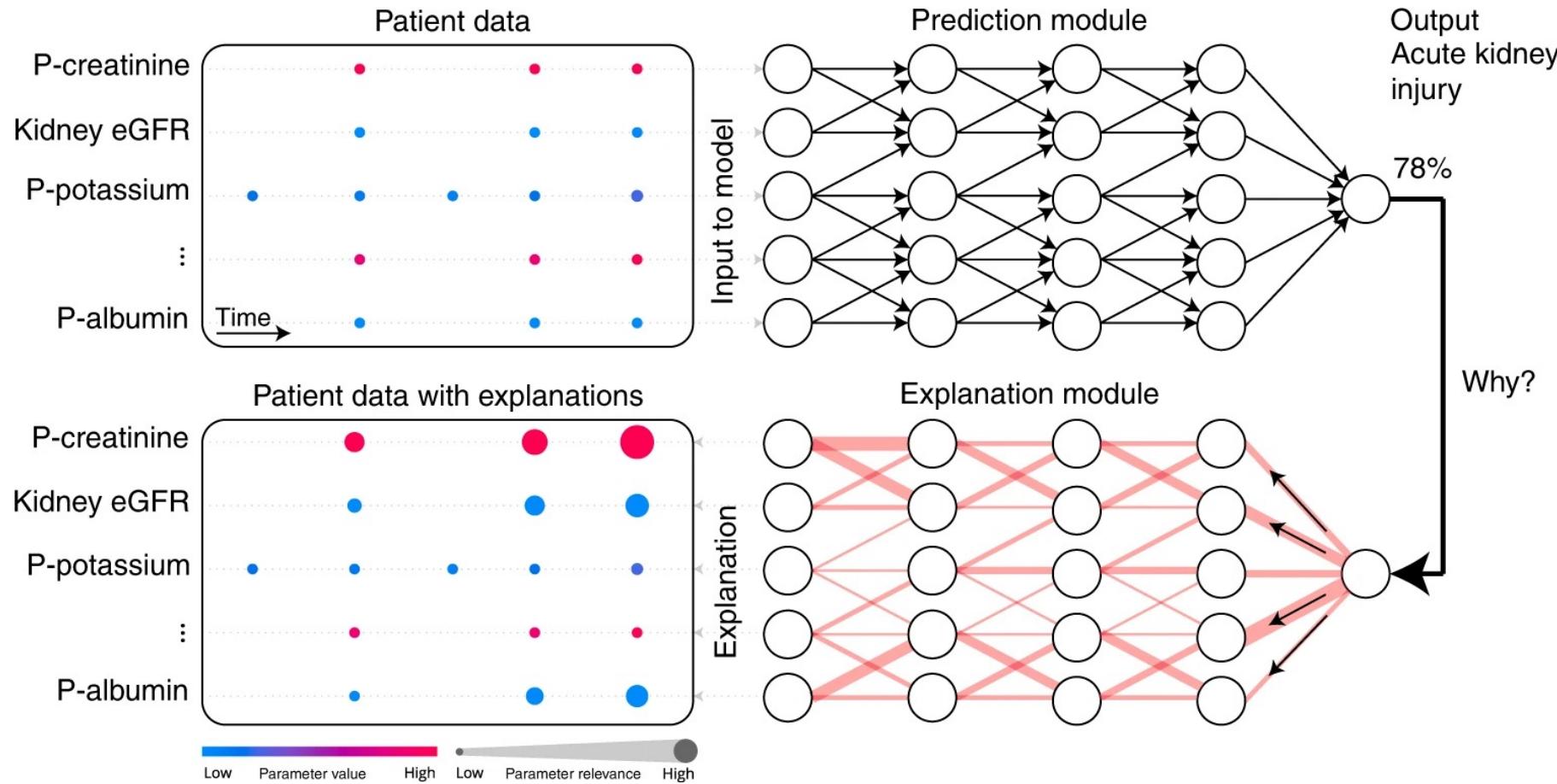


Alzheimer's disease classification



Bach et al. (2015). *PLoS one* **10**, e0130140.
Böhle et al. (2019). *Frontiers Aging Neurosci.* **11**, 194.

LRP on tabular data



Lauritsen et al. (2020). *Nature communications* 11, 11.

Pixel (feature) attribution methods

Pros

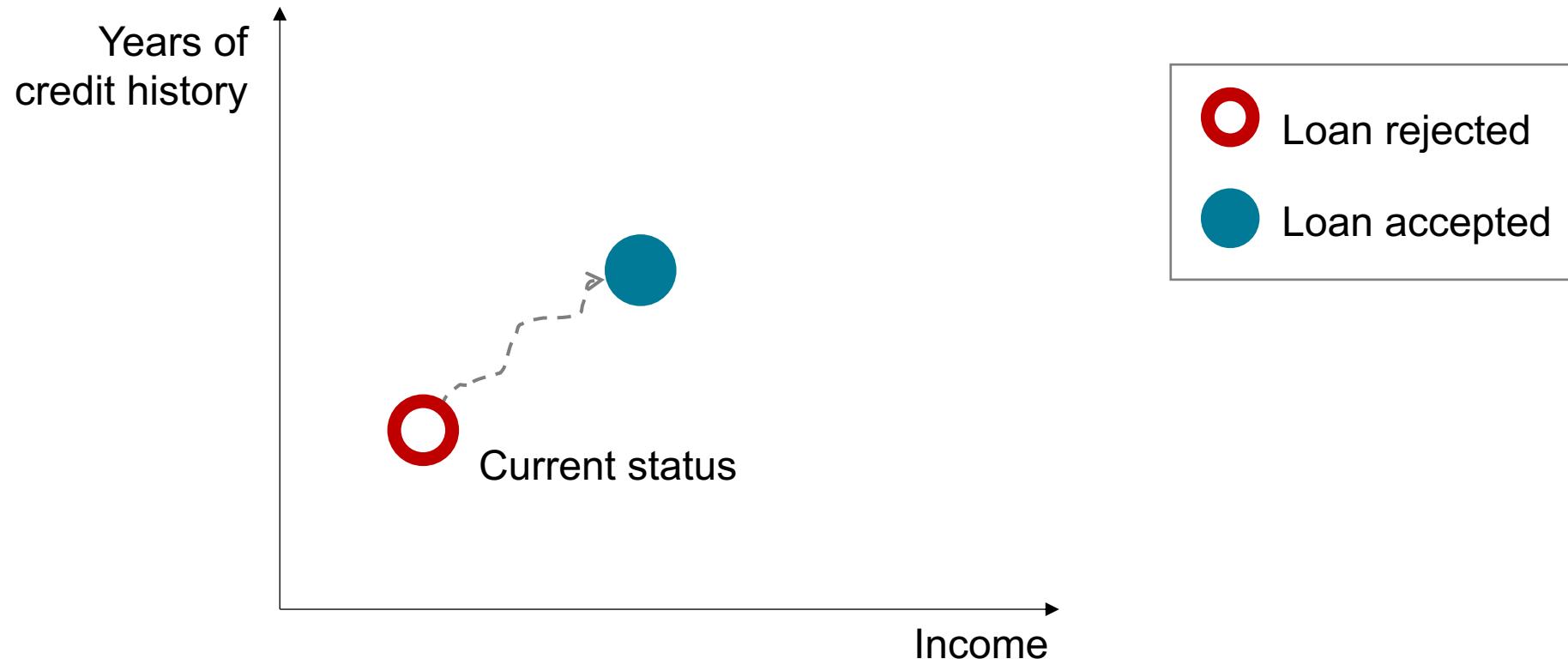
- Visual and quick to inspect
- Important regions in the image
- Allow to detect certain shortcomings

Cons

- Difficult to know whether an explanation is correct
- Small perturbations to an image can lead to very different pixels highlighted
- Often very qualitative

Counterfactual explanations (Instance based)

Counterfactual: smallest change in the input features that changes the predictions to another output.



Adversarial attacks

Adversarial examples: generate false predictions by leveraging the shortages of the algorithm.

- AI safety
- Make machine learning models more secure against manipulations

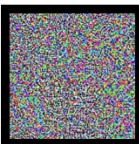


Minimal change in the input to generate a different output



"panda"

Adversarial Noise



+

=

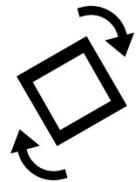


"gibbon"



"vulture"

Adversarial Rotation



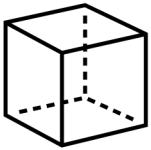
+

=

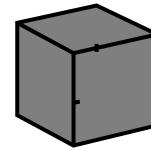
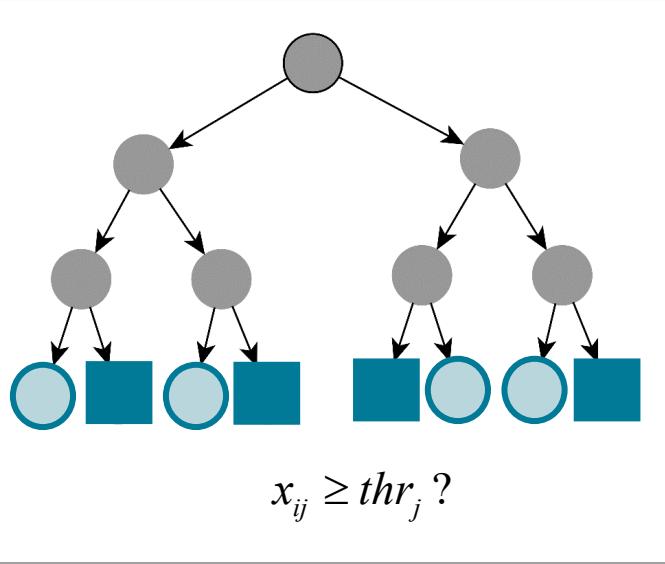


"orangutan"

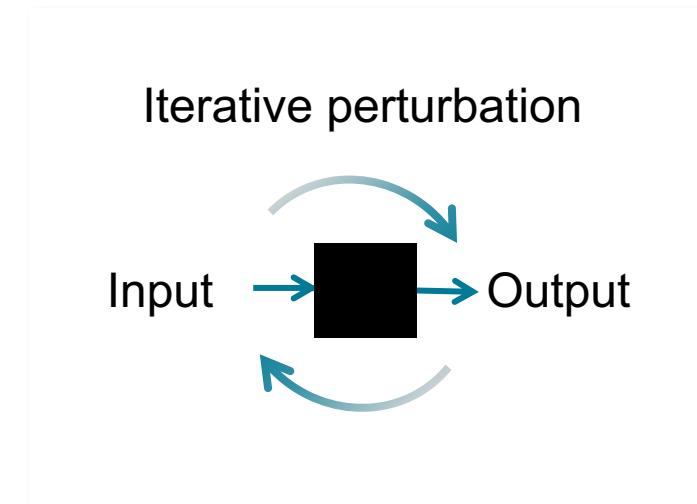
Counterfactuals



White box → model-dependent

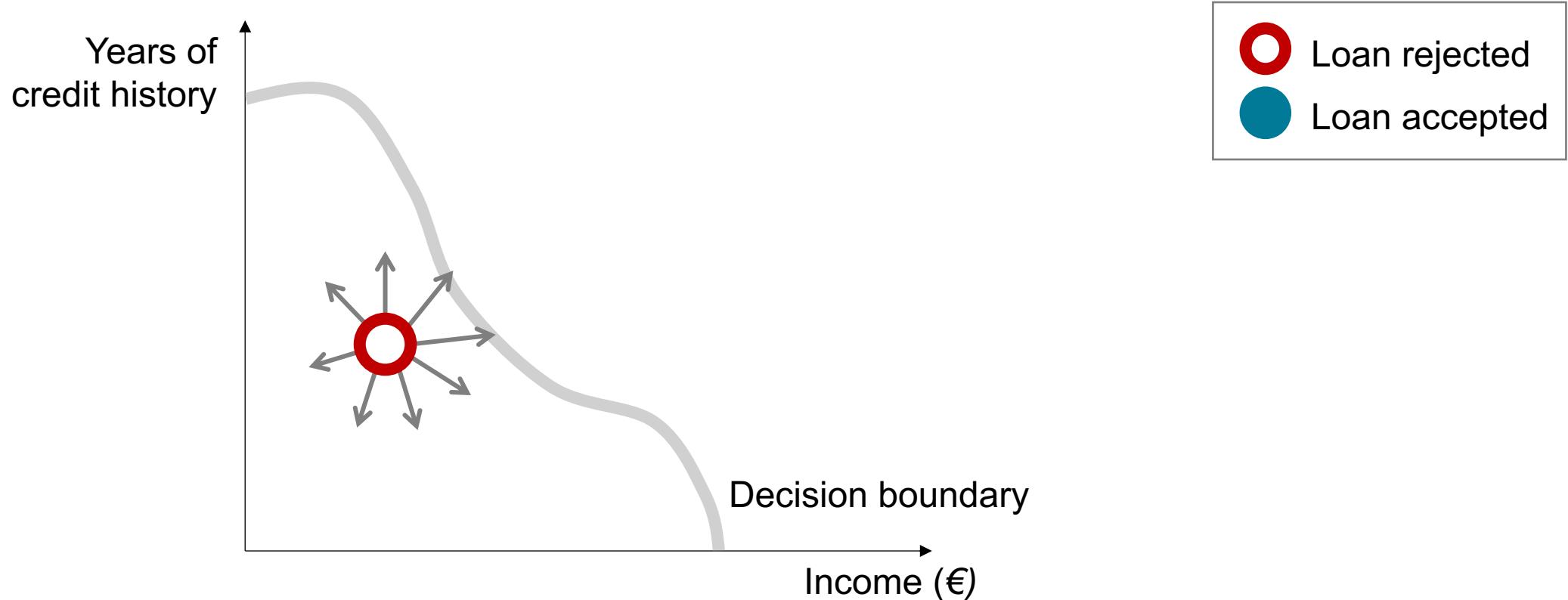


Black box → model-agnostic



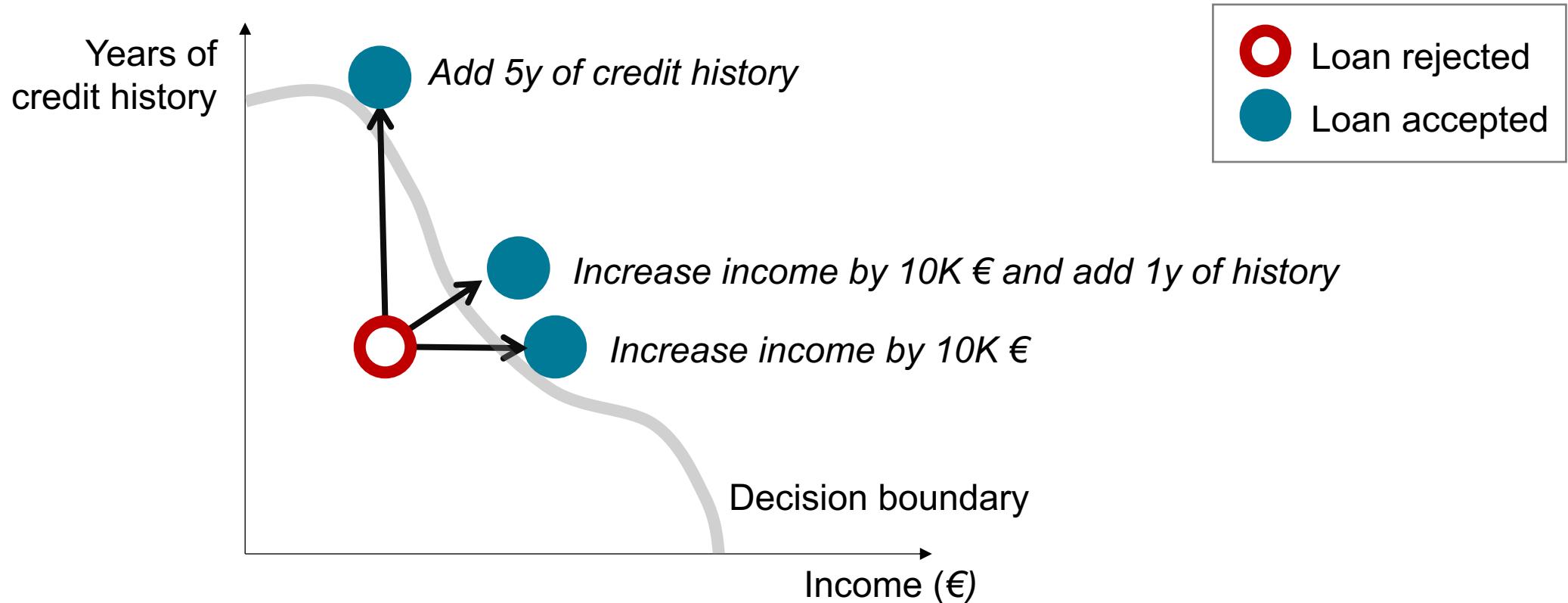
Counterfactual explanations

Counterfactual: smallest change in the input features that changes the predictions to another output.

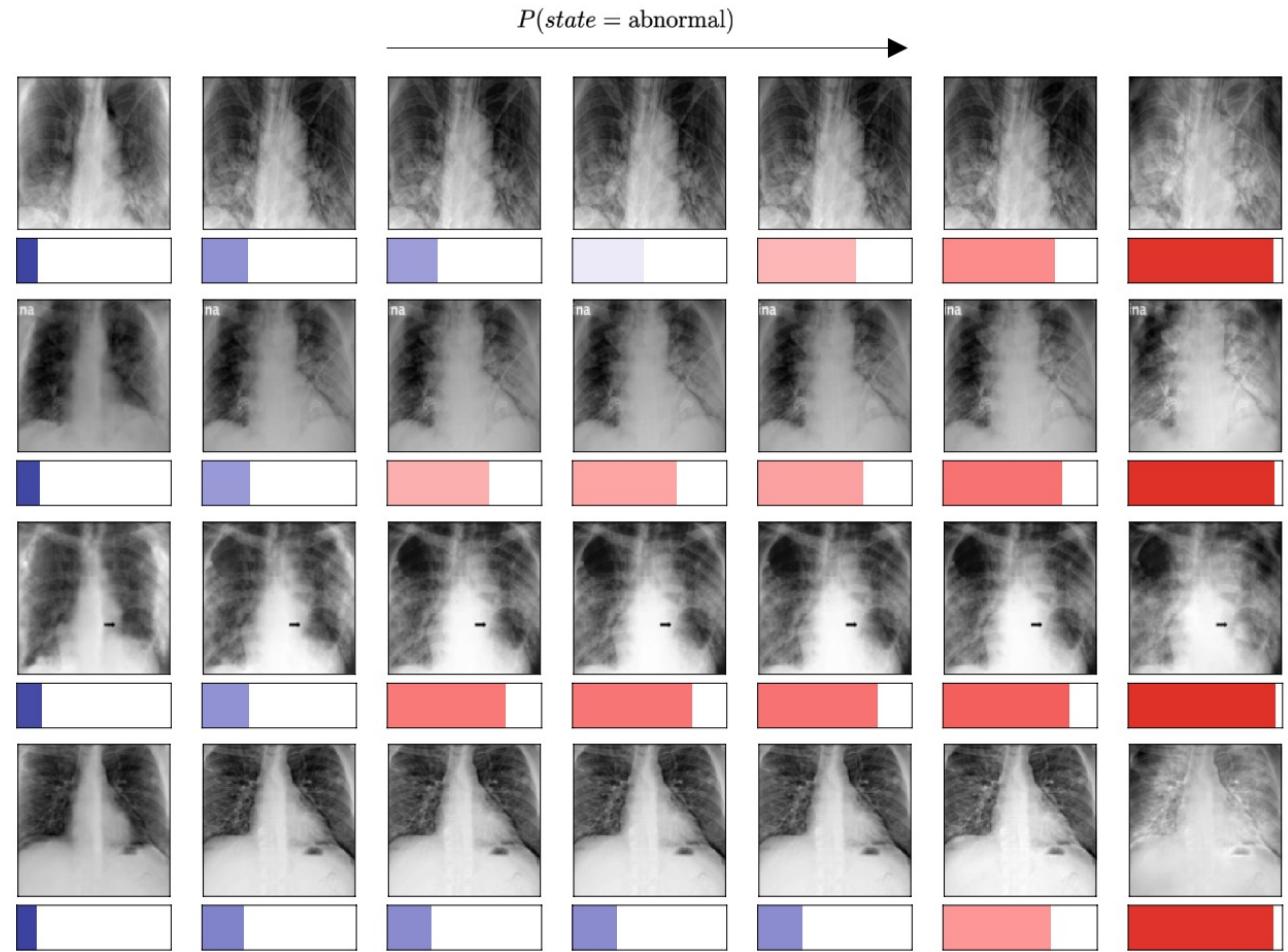
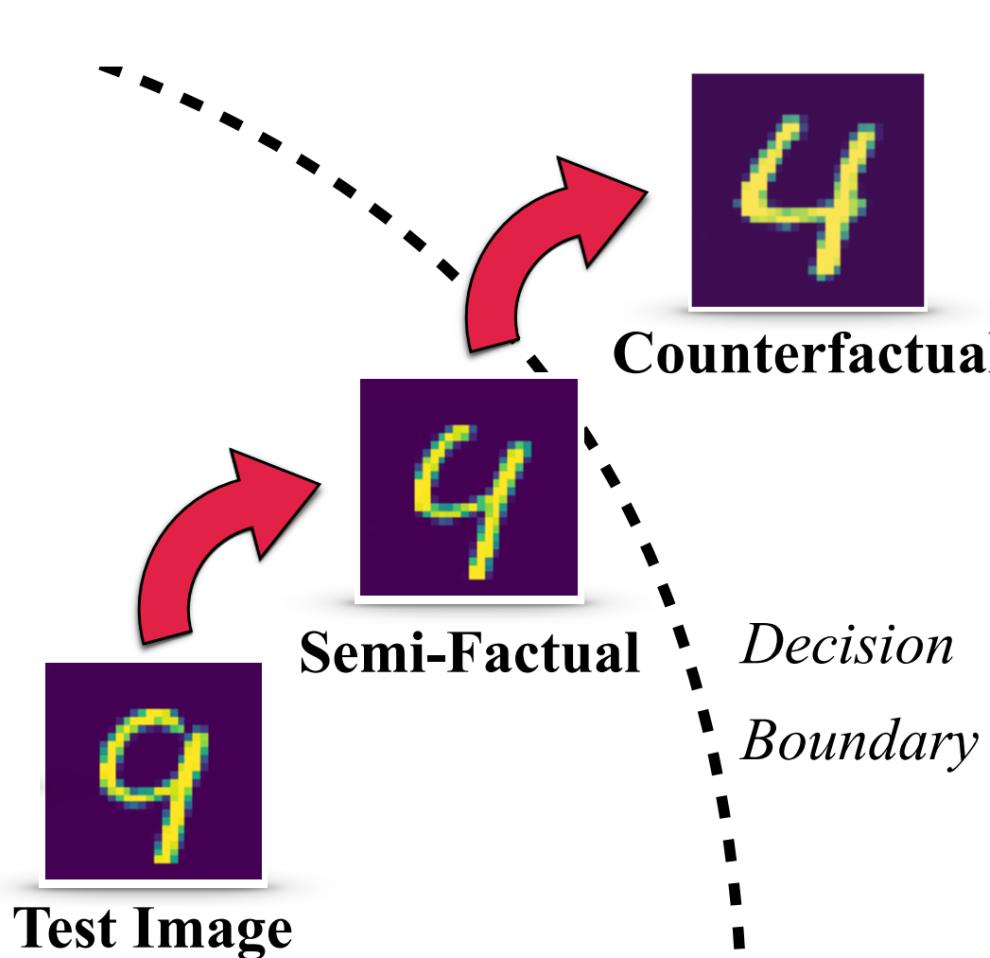


Counterfactual explanations

Counterfactual: smallest change in the input features that changes the predictions to another output.



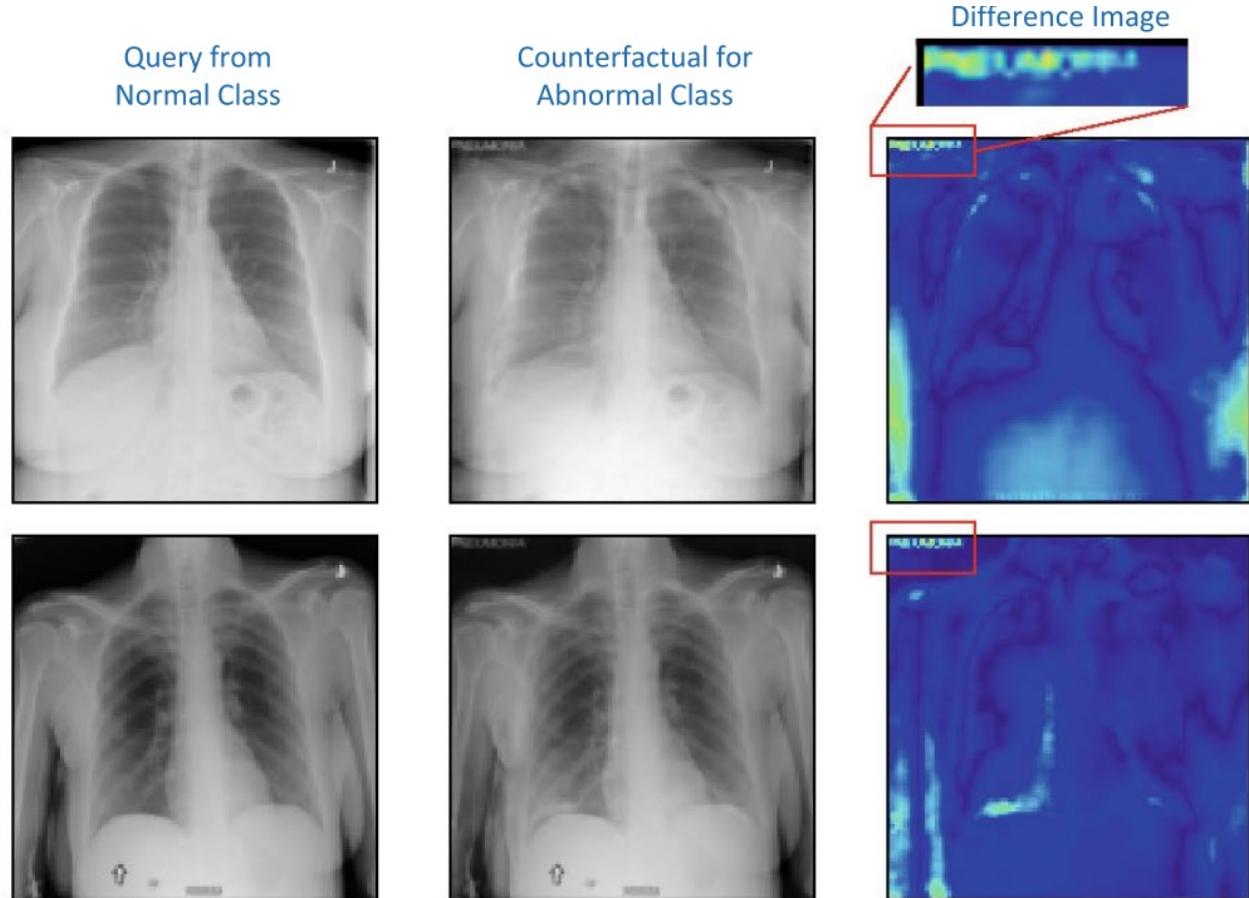
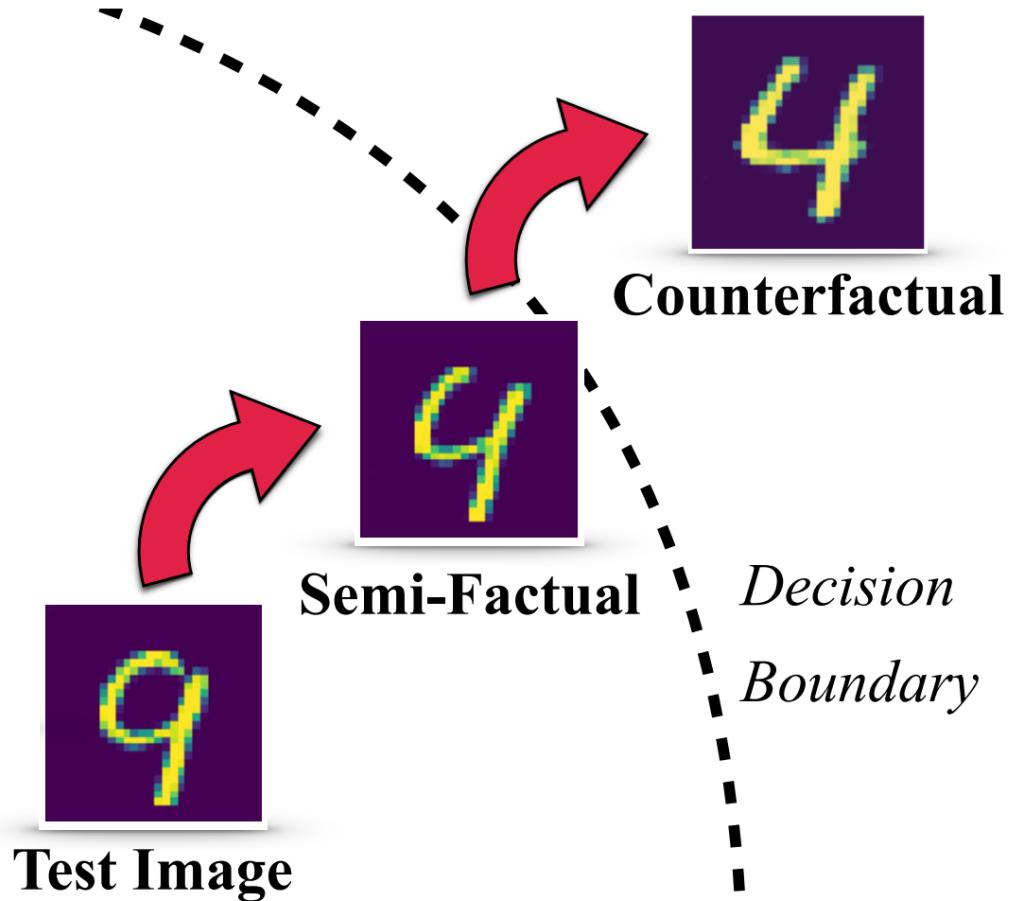
Counterfactual explanations



Kenn and Keane (2021), Proceedings of the AAAI Conference on Artificial Intelligence 35, 11575.

Thiagarajan et al. (2022) Sci Rep 12, 597.

Counterfactual explanations



Kenn and Keane (2021, *Proceedings of the AAAI Conference on Artificial Intelligence* 35, 11575).

Thiagarajan et al. (2022) *Sci Rep* 12, 597.

Feature attribution vs feature visualization

Feature attribution

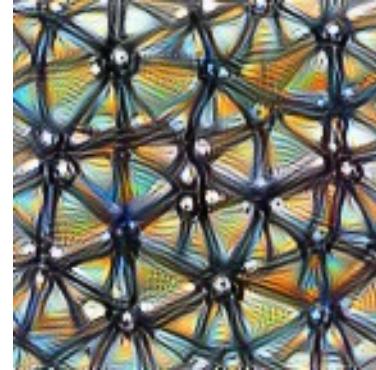


What part of the example is
responsible for the prediction?



Image in

Feature visualization



What is a network looking for?



Image out

Images from: <https://distill.pub/2017/feature-visualization/>

Feature attribution vs feature visualization

Feature Visualization

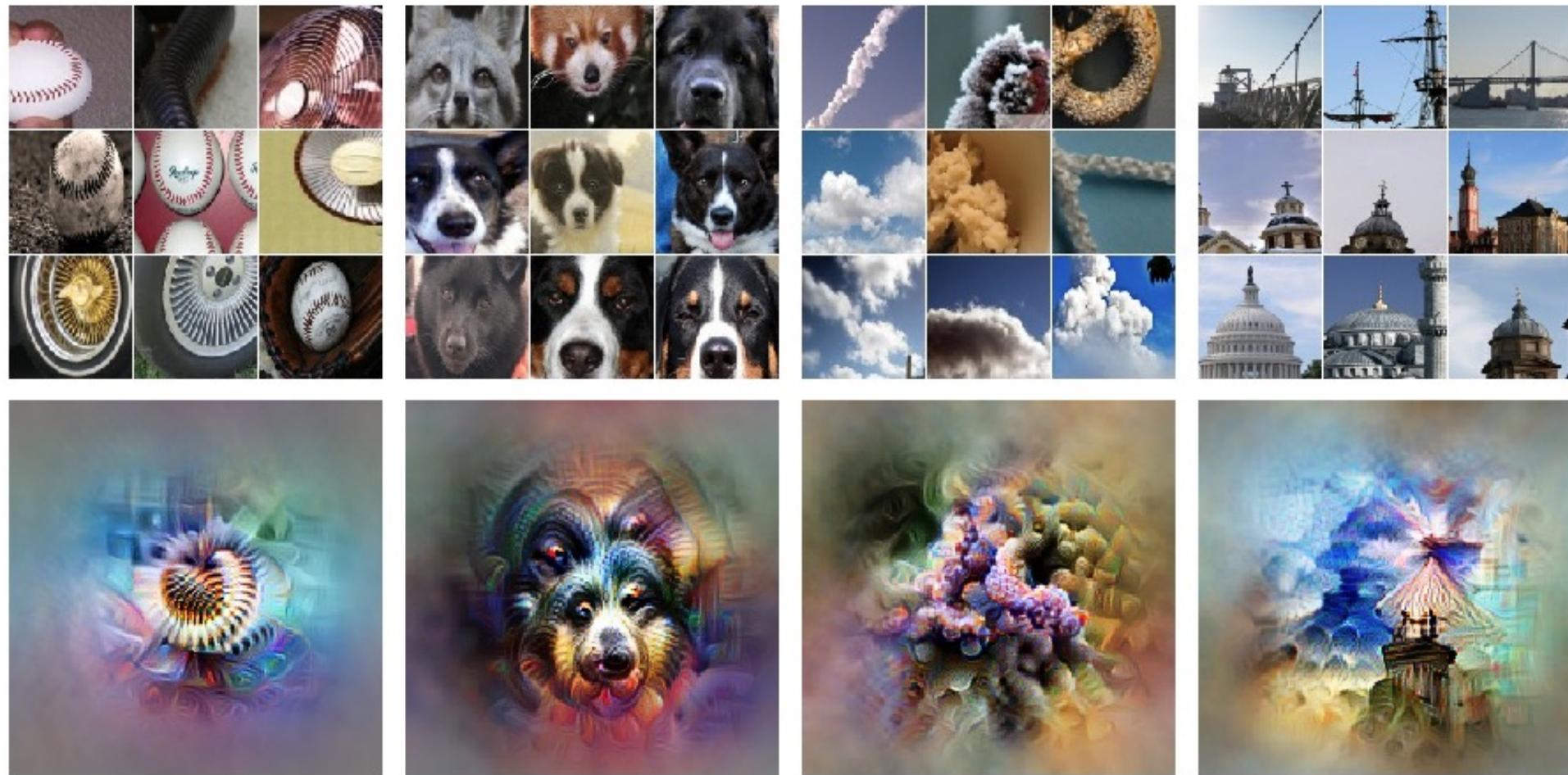
How neural networks build up their understanding of images

Edges (layer conv2d0) **Textures** (layer mixed3a) **Patterns** (layer mixed4a) **Parts** (layers mixed4b & mixed4c) **Objects** (layers mixed4d & mixed4e)

Feature visualization allows us to see how GoogLeNet [1], trained on the ImageNet [2] dataset, builds up its understanding of images over many layers. Visualizations of all channels are available in the [appendix](#).

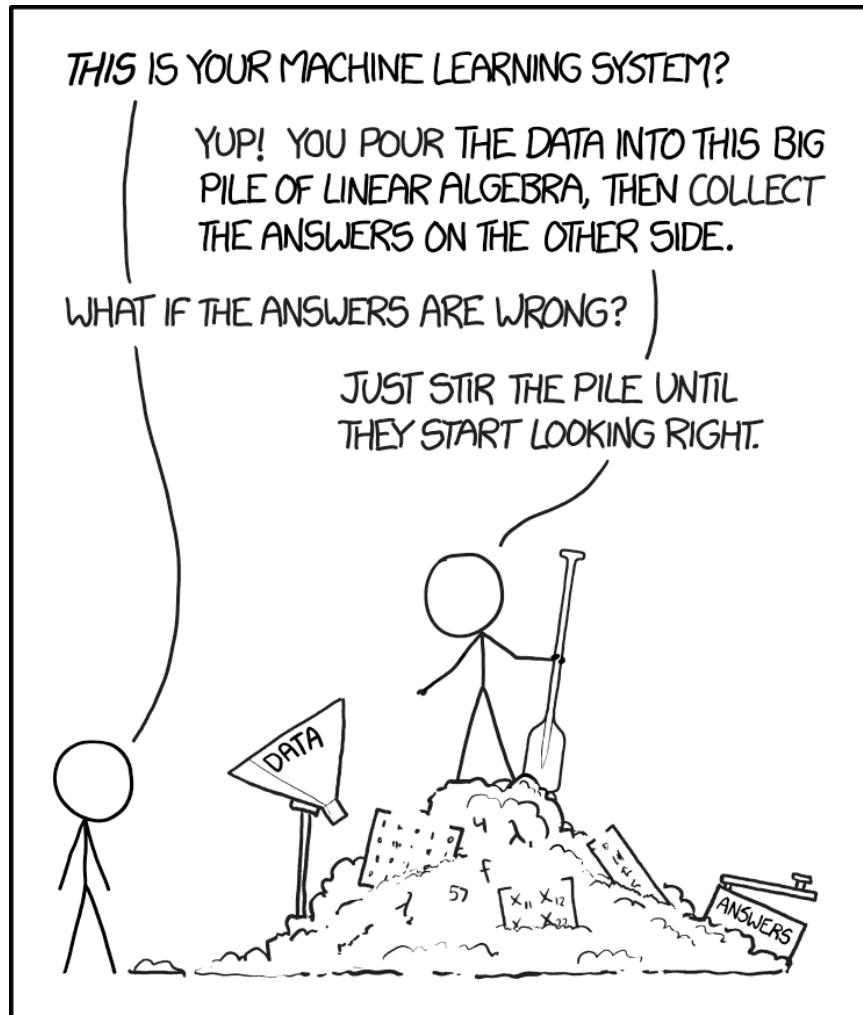
<https://distill.pub/2017/feature-visualization/>

Visualizing neurons



Images from: <https://distill.pub/2017/feature-visualization/>

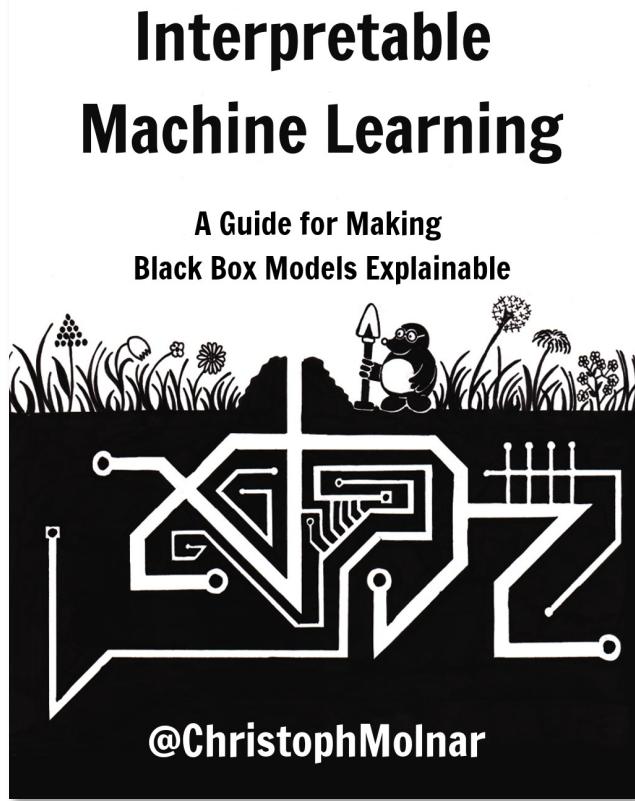
Explainable AI (XAI) – A recap



Extraction of **relevant knowledge** from a machine learning model concerning **relationships contained in data or learned by the model**.

- **Transparency:** *how did the system reached an answer?*
→ Feature attribution (surrogate and layer-wise propagation).
- **Informativeness:** *what can I learn from it?*
→ Instance-based approaches and domain-expertise
- **Justification:** *is the answer acceptable?*
→ Domain-expertise.
- **Uncertainty estimation:** *how reliable is a prediction?*

Want to know more?



<https://christophm.github.io/interpretable-ml-book>

This image is a screenshot of a research article from the journal 'nature machine intelligence'. The article is a 'REVIEW ARTICLE' titled 'Drug discovery with explainable artificial intelligence'. It is authored by José Jiménez-Luna, Francesca Grisoni, and Gisbert Schneider. The abstract discusses the potential of deep learning for drug discovery, highlighting the need for explainability to address ethical concerns and facilitate scientific validation. The page includes a 'Check for updates' button and a DOI link (<https://doi.org/10.1038/s42256-020-00236-4>).

<https://www.nature.com/articles/s42256-020-00236-4>

This image shows a screenshot of a video player interface for a series titled 'EXPLAINABLE AI EXPLAINED!'. The video player has a purple and orange gradient background with a brain icon. The main title 'EXPLAINABLE AI EXPLAINED!' is displayed in large, bold, white letters. Below it, a smaller line of text says 'Introduction'. The video player shows a play button, a progress bar at 0:04 / 52, and other standard controls. Below the video player, there is a list of six video thumbnails, each with a title, a small image, and a timestamp. The titles include 'Explainable AI explained! | #3 LIME', 'Explainable AI explained! | #4 SHAP', 'Explainable AI explained! | #5 Counterfactual explanations and adversarial attacks', 'Explainable AI explained! | #6 Layerwise Relevance Propagation with MRI data', and 'How to explain Graph Neural Networks (with XAI)'. Each thumbnail also includes the word 'DeepFindr'.

Lectures and code!