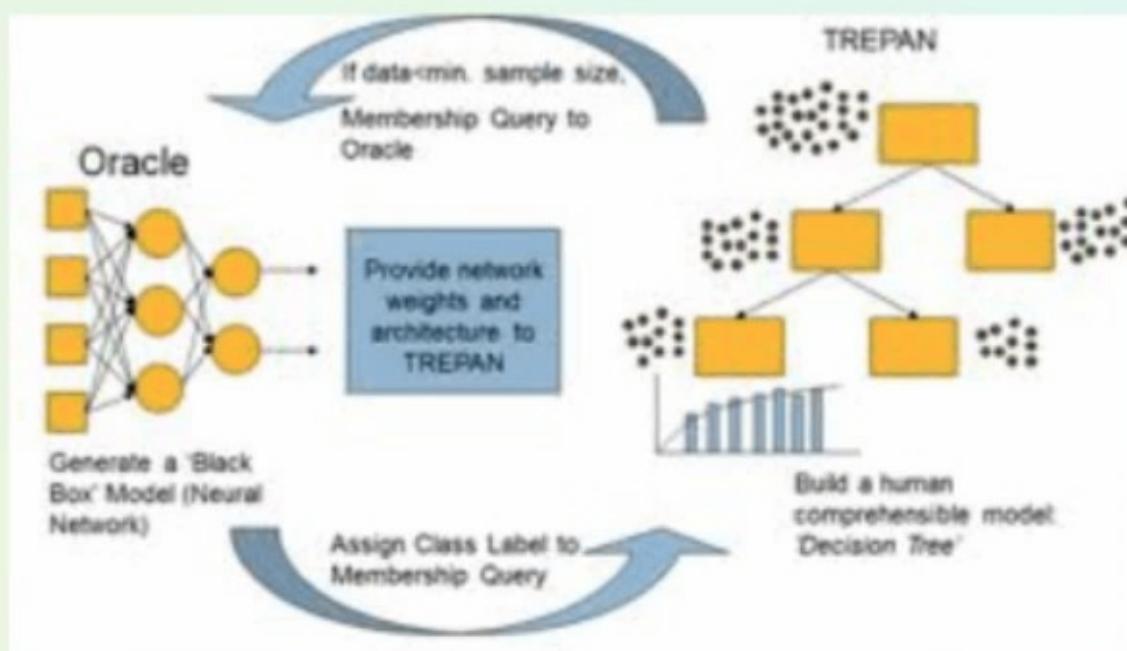
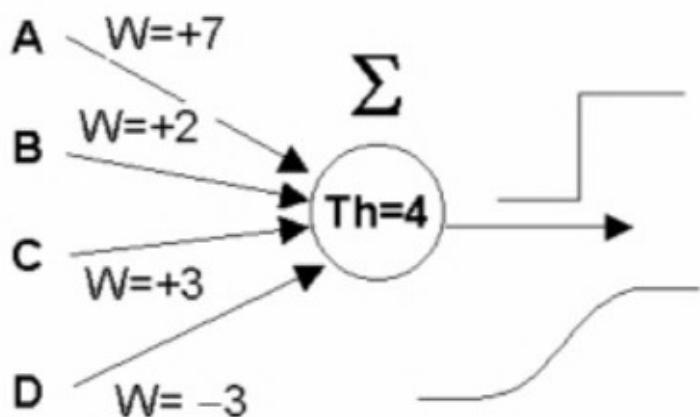


Tworzenie interpretowalnych reprezentacji zastępczych dla złożonych modeli (ANN)



IF{ $(A \text{ and } B) \text{ or } (A \text{ and } C) \text{ or }$
 $(A \text{ and } \text{not } D) \text{ and }$
 $(A \text{ and } B \text{ and } C) \text{ and }$
 $(B \text{ and } C \text{ and } \text{not } D)\}$
THEN True
ELSE False

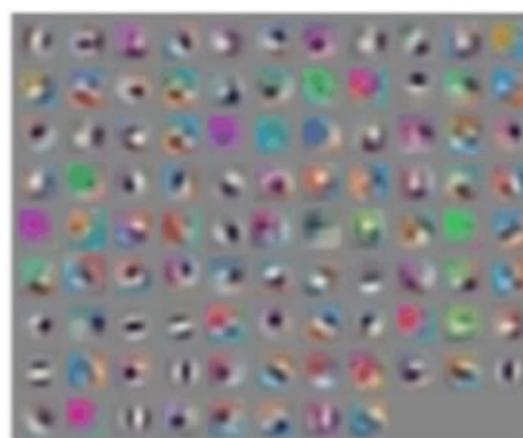


Deconvolution of CNN and other visualizations of internal layers

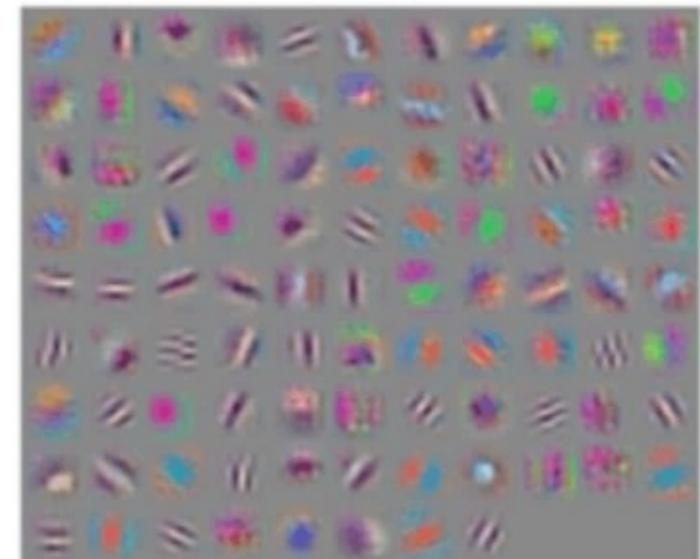
conv1



conv2



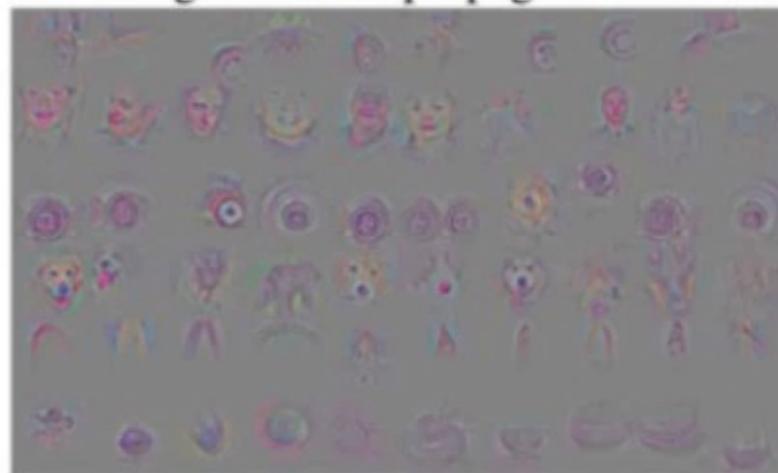
conv3



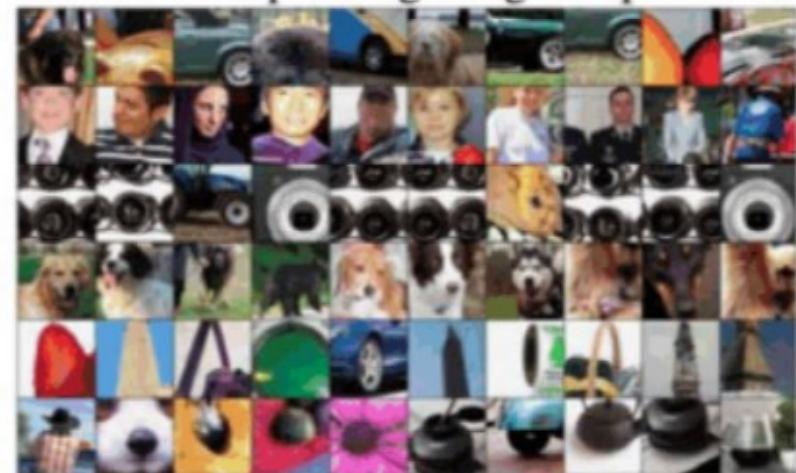
deconv



guided backpropagation



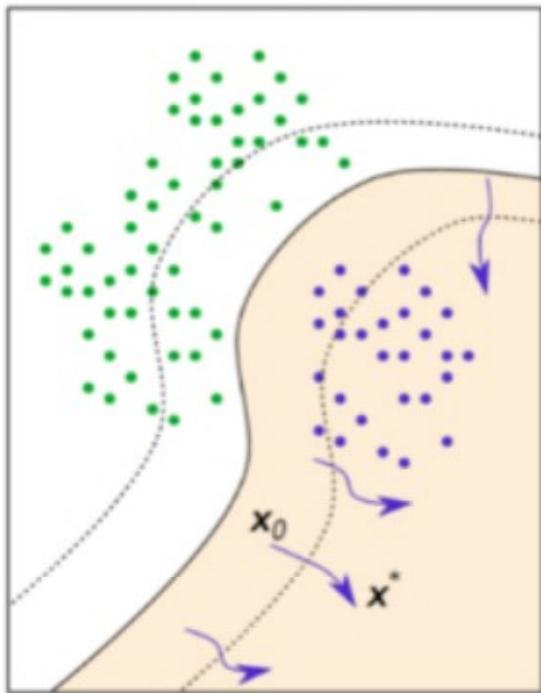
corresponding image crops



Jost Tobias Springenberg, Alexey Dosovitskiy, Thomas Brox & Martin Riedmiller 2014. Striving for simplicity:
The all convolutional net. arXiv:1412.6806.

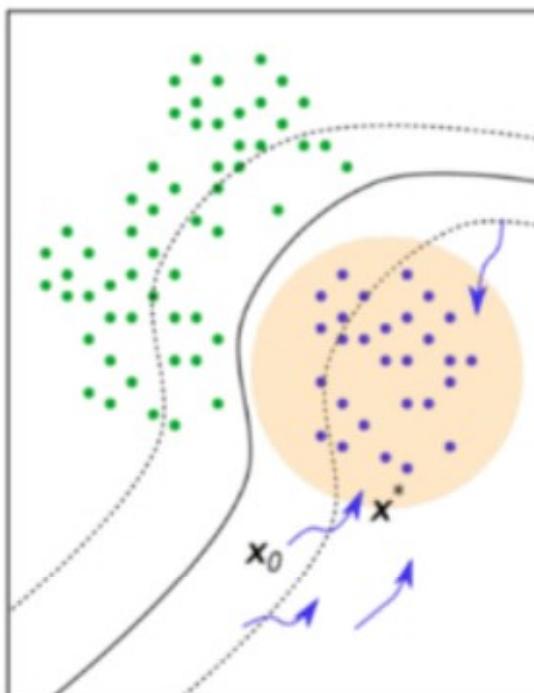
Skupienie uwagi na przykładach

model analysis

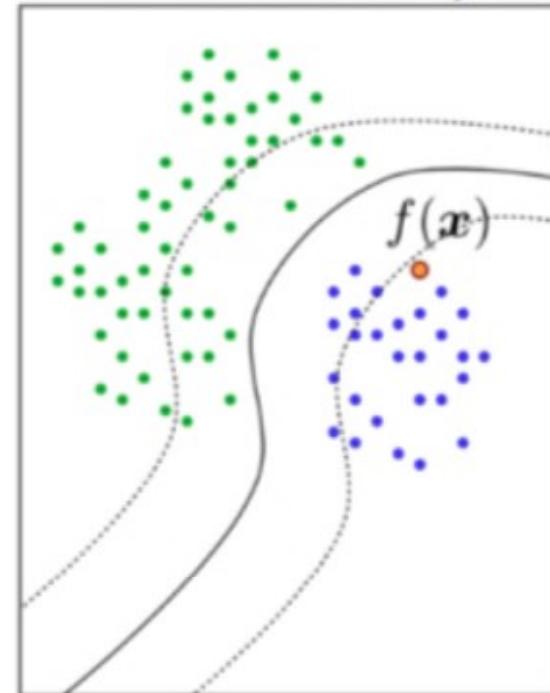


Find the input pattern that maximizes class probability.

decision analysis



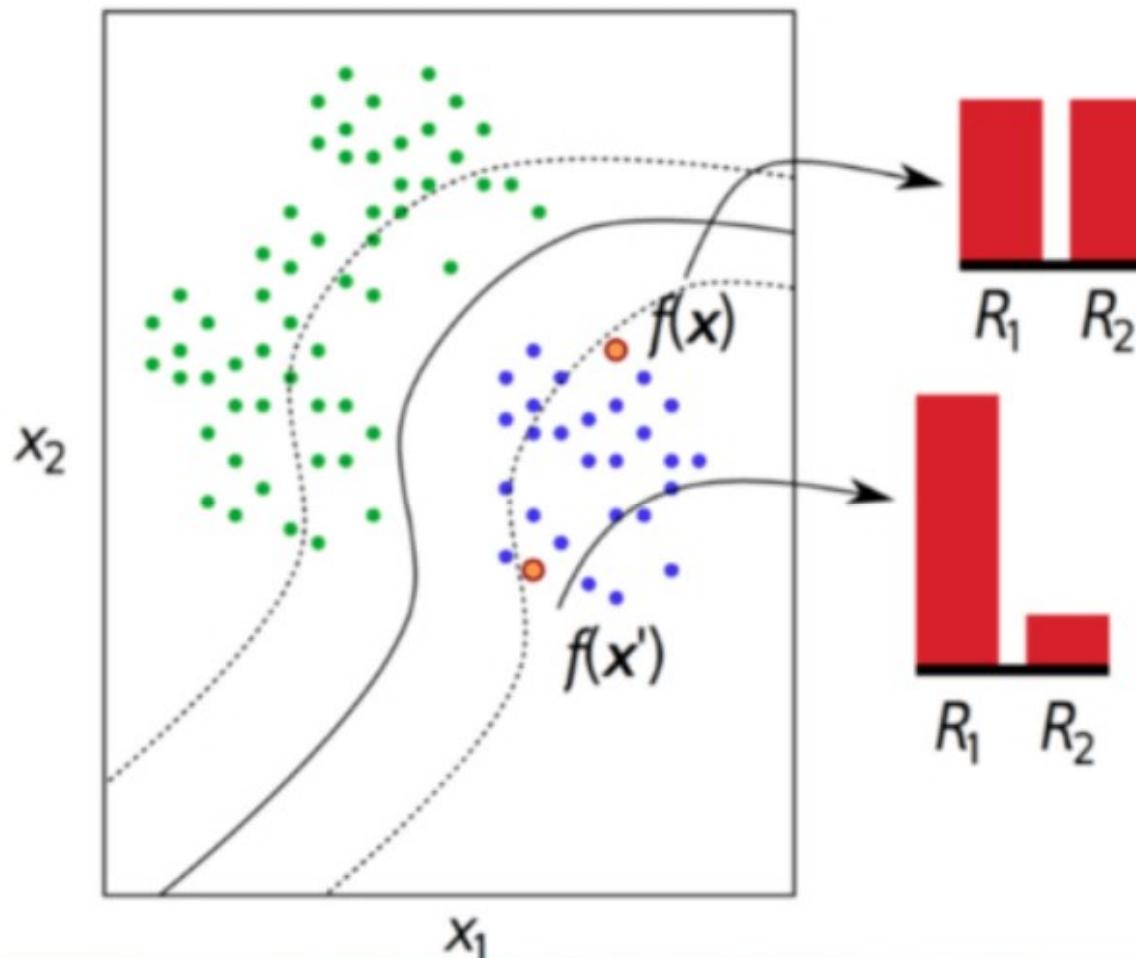
Find the most likely input pattern for a given class.



Explain individual prediction.

Explaining Individual Decisions

Goal: Determine the relevance of each input variable for a given decision $f(x_1, x_2, \dots, x_d)$, by assigning to these variables *relevance scores* R_1, R_2, \dots, R_d .

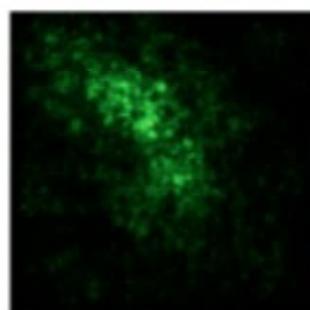


Poszukiwanie krytycznych elementów w danych

Techniques of Interpretation

Sensitivity Analysis

(Simonyan et al. 2014)



Classifier

Bird

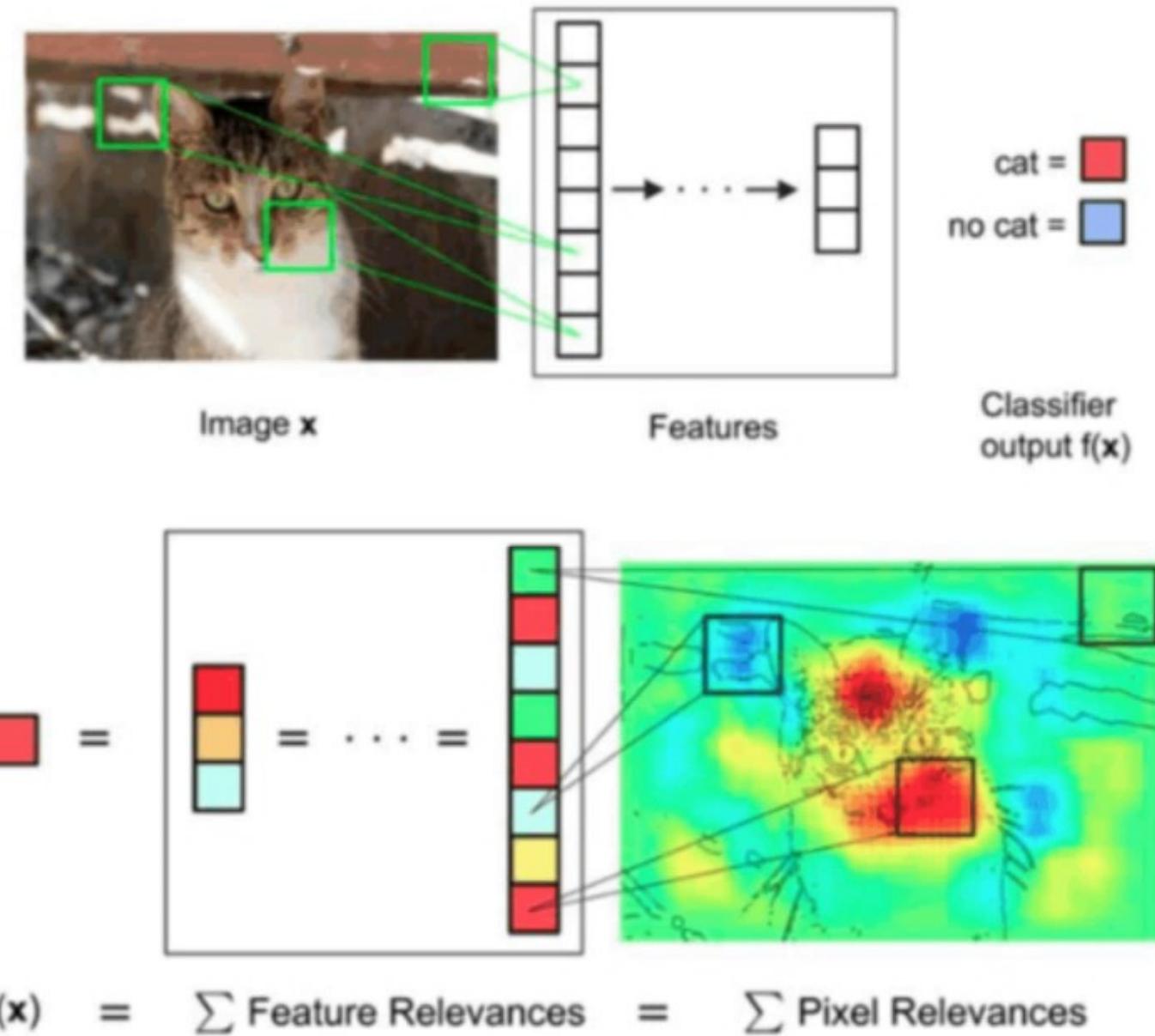
prediction $f(x)$

Explain prediction
*(which pixels lead to decrease
of prediction score when changed)*

$$\left\| \frac{\partial}{\partial x_p} f(x) \right\|$$

Specyfika wyjaśnień obrazowych

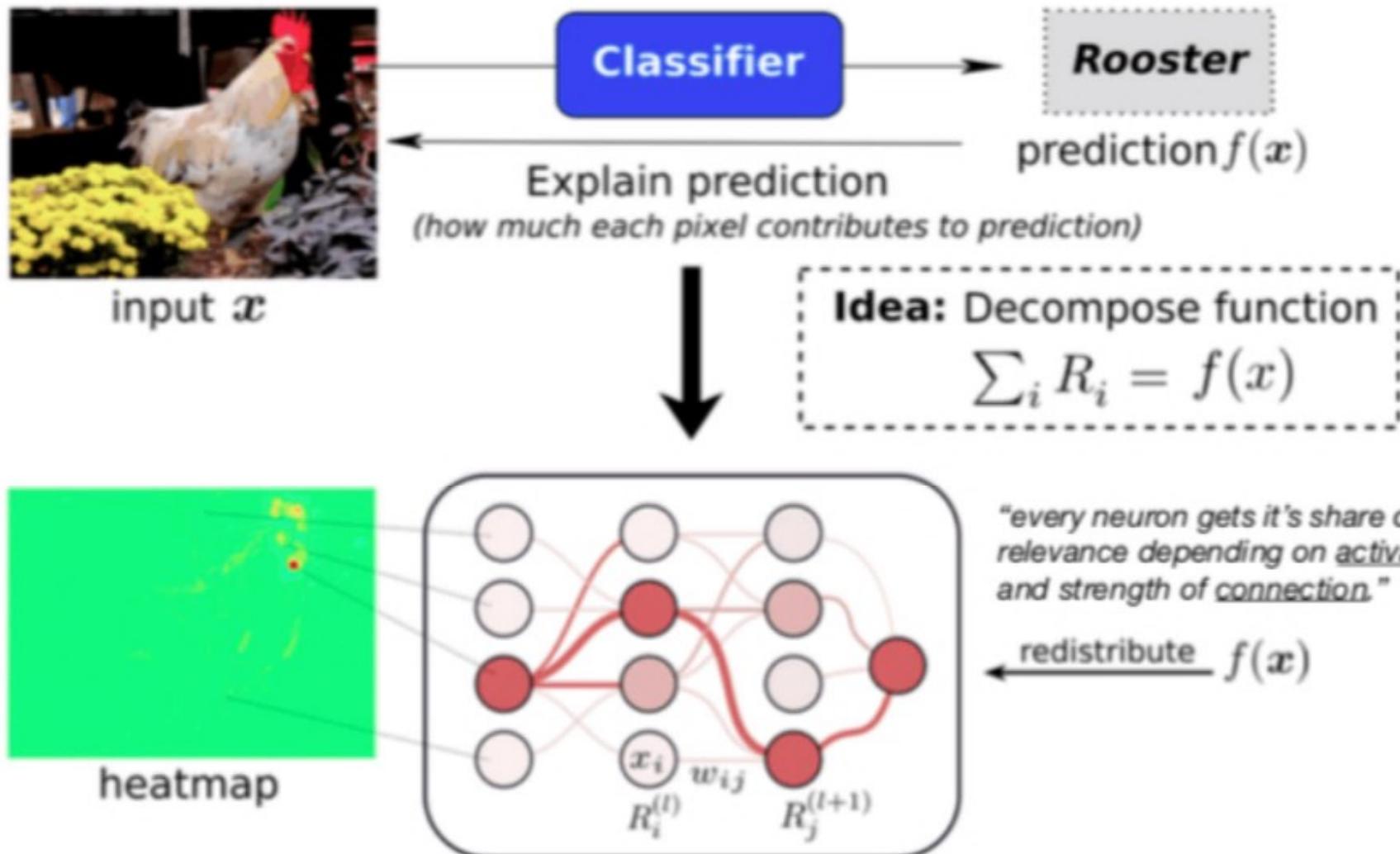
Sebastian Bach, Alexander Binder, Grégoire Montavon, Frederick Klauschen, Klaus-Robert Müller & Wojciech Samek 2015. On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation. *PLoS One*, 10, (7), e0130140, doi:10.1371/journal.pone.0130140.



Techniques of Interpretation

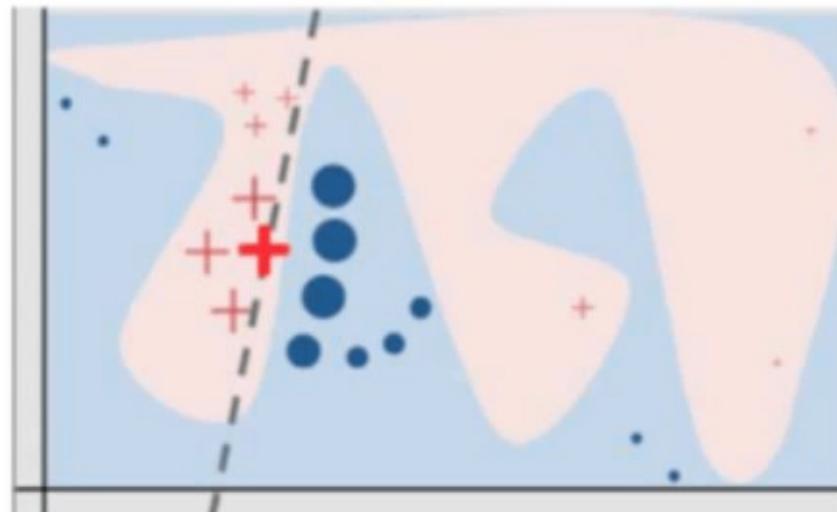
Layer-wise Relevance Propagation (LRP)

(Bach et al. 2015)

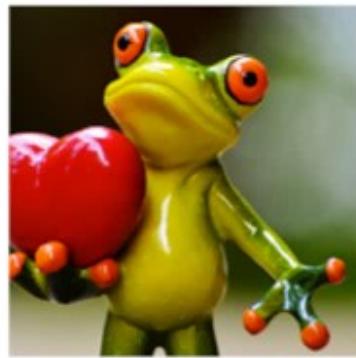


What are explanations in the sense of LIME ?

- Explanation := local linear approximation of the model's behaviour. While the model may be very complex globally, it is easier to approximate it around the vicinity of a particular instance. While treating the model as a black box, we perturb the instance we want to explain and learn a sparse linear model around it -> used as explanation.
- Look at the image: The model's decision function is represented by the blue/pink background = clearly nonlinear. The bright red cross is the instance being explained (let's call it X). We sample instances around X, and weight them according to their proximity to X (weight here is indicated by size). We then learn a linear model (dashed line) that approximates the model well in the vicinity of X, but not necessarily globally!



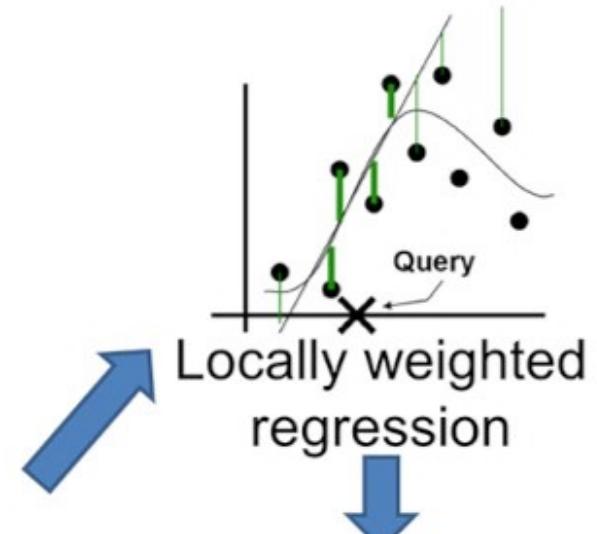
LIME – rozpoznawanie obrazów



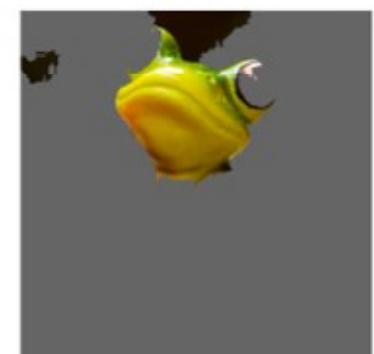
Original Image
 $P(\text{tree frog}) = 0.54$



Perturbed Instances	$P(\text{tree frog})$
A photograph of a tree frog with several red spots removed.	0.85
A photograph of a tree frog with its body colored yellow.	0.00001
A photograph of a tree frog with its eyes colored red.	0.52

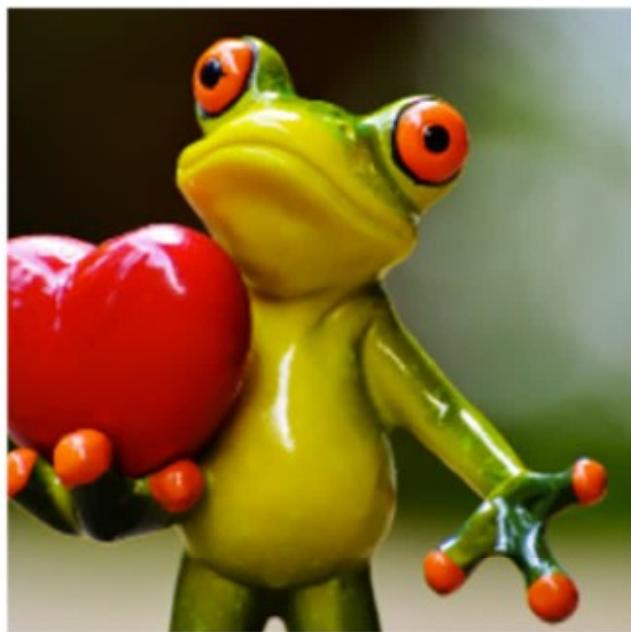


Locally weighted regression

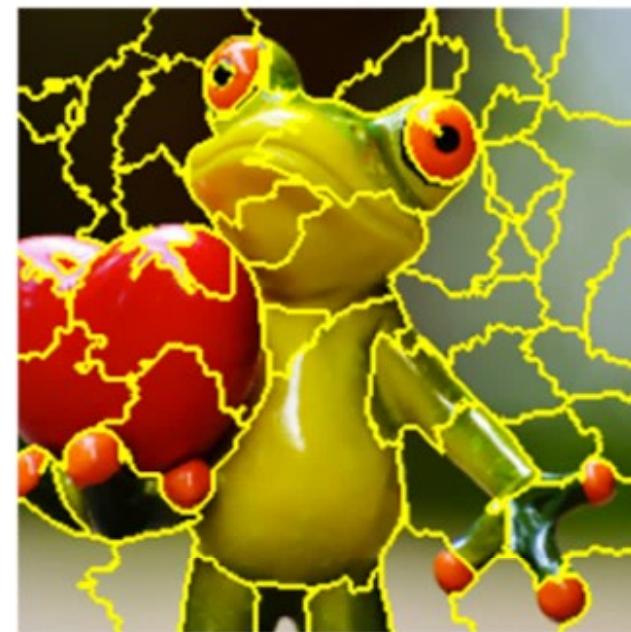


Explanation

LIME – rozpoznawanie obrazów



Original Image



Interpretable
Components

LIME - rozpoznawanie obrazów



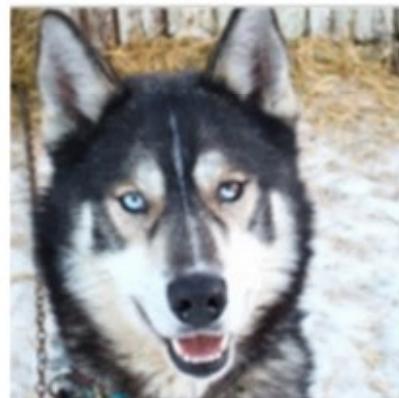
Predicted: **wolf**
True: **wolf**



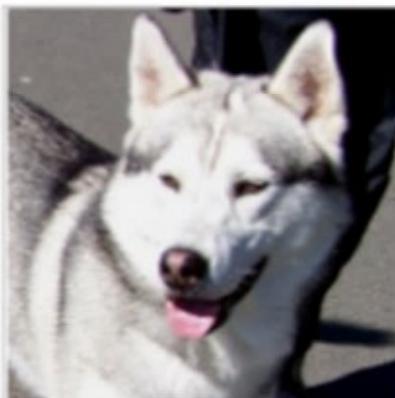
Predicted: **husky**
True: **husky**



Predicted: **wolf**
True: **wolf**



Predicted: **wolf**
True: **husky**



Predicted: **husky**
True: **husky**



Predicted: **wolf**
True: **wolf**

LIME – rozpoznawanie obrazów



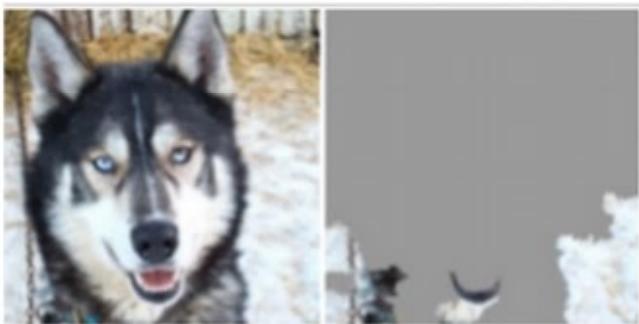
Predicted: **wolf**
True: **wolf**



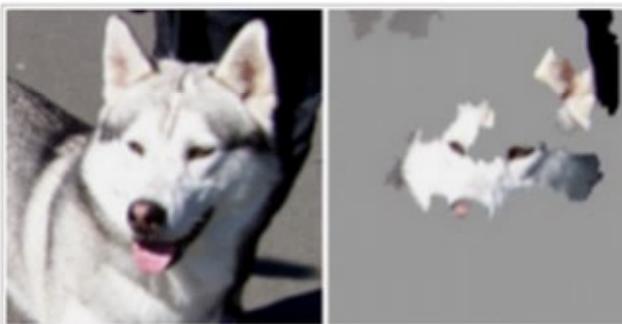
Predicted: **husky**
True: **husky**



Predicted: **wolf**
True: **wolf**



Predicted: **wolf**
True: **husky**



Predicted: **husky**
True: **husky**



Predicted: **wolf**
True: **wolf**

LIME

