

Deep Learning (Explainable AI)

why? \rightarrow open up black box

Activation maximization framework:

$$\max_{x \in X} \log p(w_c | x) - \lambda \|x\|^2$$

\hookrightarrow class prototype

- find x with highest score
- $\lambda \|x\|^2$: regularization

In code space: $\max_{z \in Z} \log p(w_c | g(z)) - \lambda \|z\|^2$
 $x \rightarrow g(z)$

Inter-Loop: choose interpretable model and train it

Post-Loop: choose working model and make interpretable afterwards

\hookrightarrow possible since 2015

Layer-wise relevance propagation (LRP)

- deep Taylor decomposition
- applicable to all NNs with monotonous activation functions
- shows which pixels contribute to classification