

Explainable Artificial Intelligence (XAI) - Concepts, taxonomies, opportunities and challenges toward responsible AI

HYU AI LAB 세미나

2021.06.21

조환희

목차

1. XAI?
2. XAI in Deep Learning
3. Challenges of XAI

1. XAI?

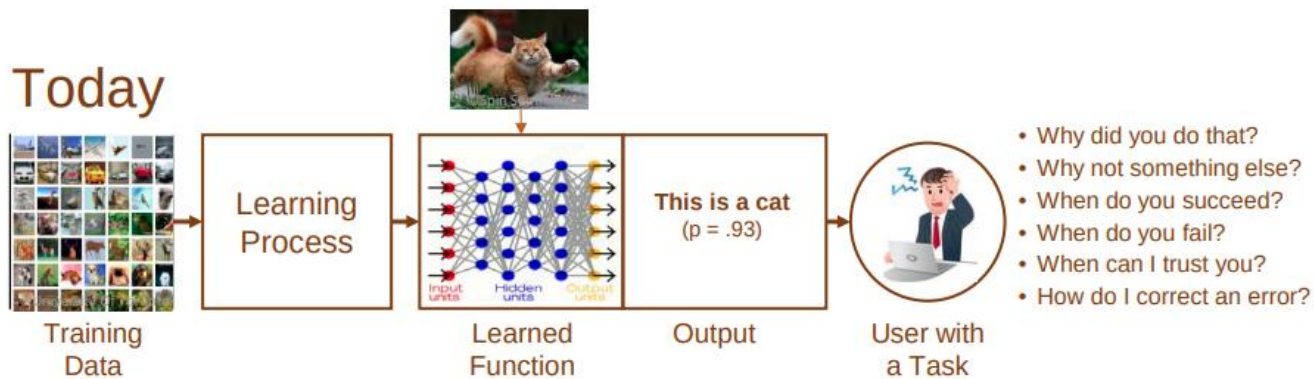
XAI

=eXplainable Artificial Intelligence

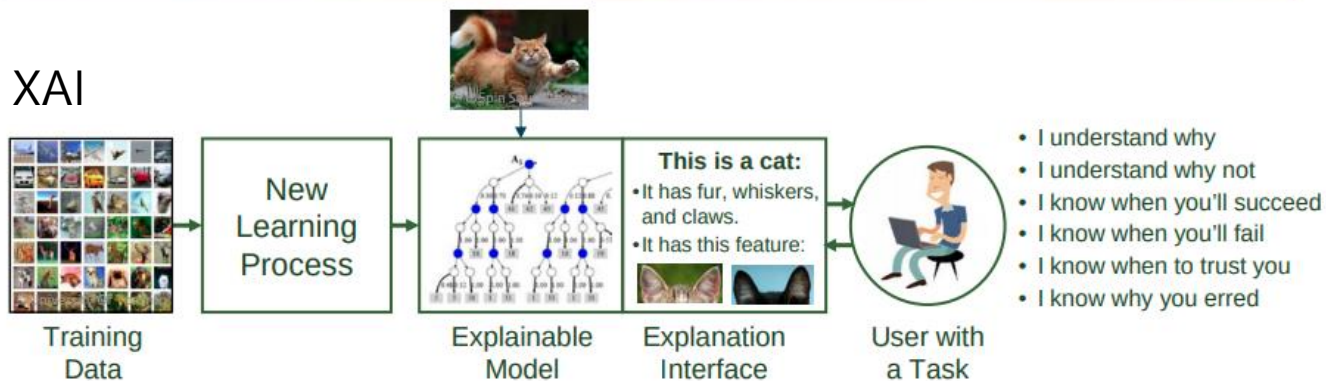
: **artificial intelligence** (AI) in which the results of the solution can be understood by humans

1. XAI?

Today



XAI



1. XAI?

Deep Learning

→ black-box

Increasingly being employed to make important predictions in critical contexts

Decisions are made not justifiable, legitimate, that do not allow obtaining detailed explanations

→ XAI가 필요하다!!

2. XAI in Deep Learning

How to implement XAI in Deep Learning?

: Transparent models and Post-hoc explainability

2. XAI in Deep Learning

1) Transparent models : convey interpretability by themselves

Linear/Logistic Regression

Decision Trees

K-Nearest Neighbors

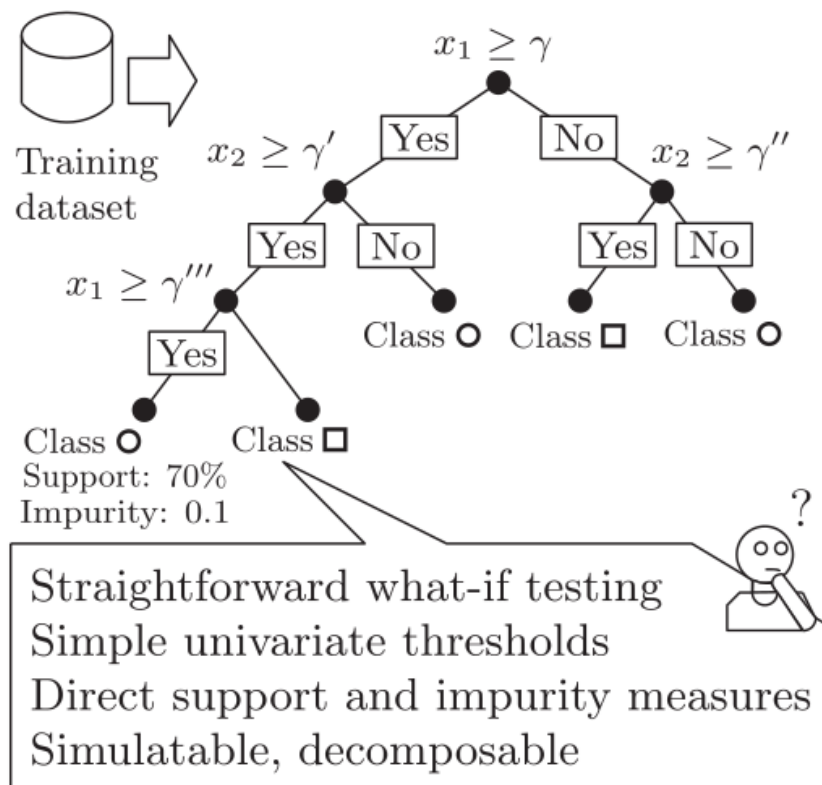
2. XAI in Deep Learning

ex) Decision Trees

features : x_1, x_2

thresholds : r, r', r'', \dots

→ results can be understood by humans



2. XAI in Deep Learning

2) Post-hoc explainability : models that are not readily interpretable by design by resorting to diverse means to enhance their interpretability

text explanations

visual explanations

local explanations

explanations by example

explanations by simplifications

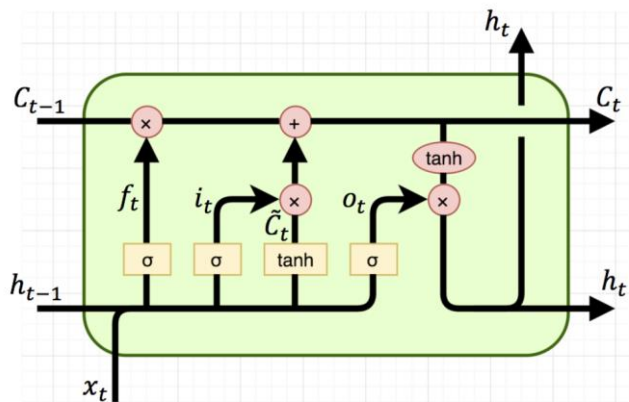
feature relevance explanations

2. XAI in Deep Learning

Visual Explanations : interpretable, long-range **LSTM cells**

LSTMs can use its **memory cells** to remember long-range information.

Demonstrate that LSTM learn powerful, and often interpretable long-range interactions on real-world data.



using $\tanh(c)$ for visualization

(a) Long Short-Term Memory

2. XAI in Deep Learning

Cell that turns on inside quotes:

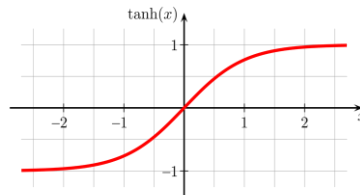
"You mean to imply that I have nothing to eat out of.... On the contrary, I can supply you with everything even if you want to give dinner parties," warmly replied Chichagov, who tried by every word he spoke to prove his own rectitude and therefore imagined Kutuzov to be animated by the same desire.

Kutuzov, shrugging his shoulders, replied with his subtle penetrating smile: "I meant merely to say what I said."

Cell that robustly activates inside if statements:

```
static int __dequeue_signal(struct sigpending *pending, sigset_t *mask,
    siginfo_t *info)
{
    int sig = next_signal(pending, mask);
    if (sig) {
        if (current->notifier) {
            if (sigismember(current->notifier_mask, sig)) {
                if (!(current->notifier)(current->notifier_data)) {
                    clear_thread_flag(TIF_SIGPENDING);
                    return 0;
                }
            }
        }
        collect_signal(sig, pending, info);
    }
    return sig;
}
```

$\tanh(c)$: -1 red
 $\tanh(c)$: +1 blue



2. XAI in Deep Learning

feature relevance explanations : **Layer-Wise Relevance Backpropagation**

decompose the network classification decision into **contributions** of its **input elements**. → **Relevance Score**



Relevance score : score that how much pixel p contributes to explaining the classification decision, $R_p(x)$

2. XAI in Deep Learning

Consider each **neuron** as an object that can **be decomposed and expanded**.

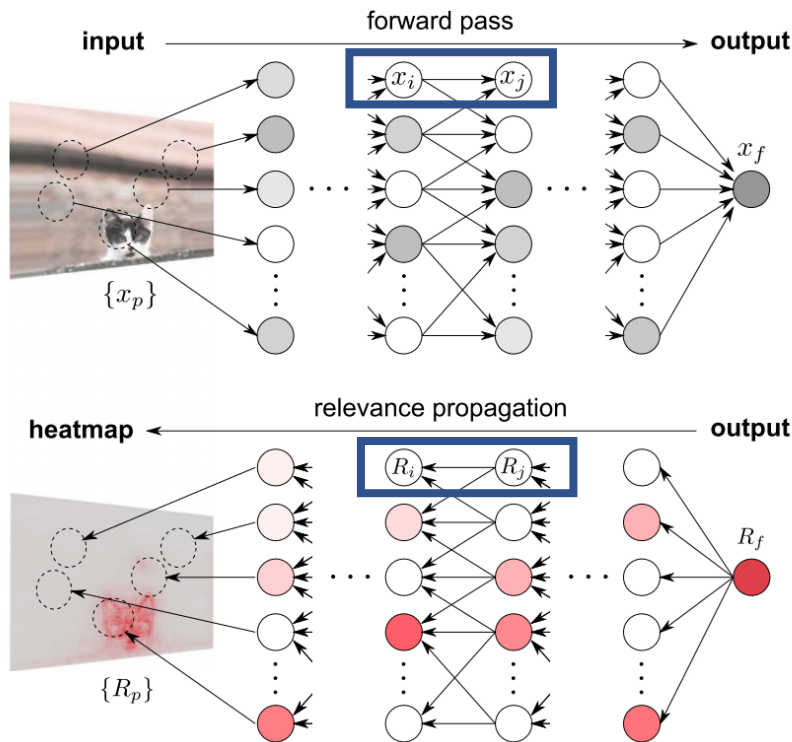
Taylor Decomposition

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

well-chosen *root point* : $f(\tilde{\mathbf{x}}) = 0$.

the relevances $R_p(\mathbf{x})$ assigned to pixels in the image.

2. XAI in Deep Learning



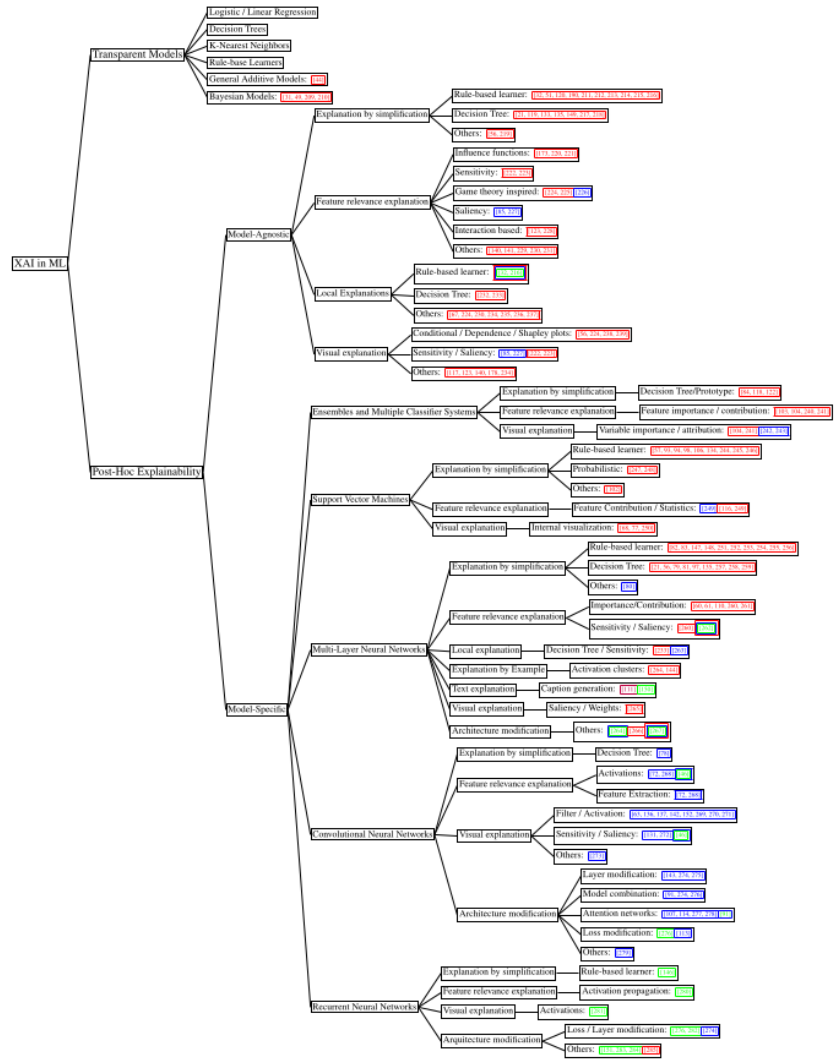
Decompose R_j on the set of lower layer neurons x_i which x_j is connected

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

$$R_i = \sum_j R_{ij}.$$

2. XAI in Deep Learning

다양한 taxonomies



3. Challenges of XAI

1. Lack of agreement on the vocabulary and definitions.
2. Trade-off between interpretability and performance
3. Providing explanations that are accessible for society, policy makers and the law.

Thank you

References

A. Karpathy, J. Johnson, L. Fei-Fei, Visualizing and understanding recurrent networks, 2015, <https://arxiv.org/pdf/1506.02078.pdf>

S. Bach , A. Binder , G. Montavon , F. Klauschen , K.-R. Müller , W. Samek , On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation, PloS one 10 (7) (2015) e0130140, <https://www.sciencedirect.com/science/article/pii/S0031320316303582>