

Lecture 9-1: Trust, Opaque Algorithms, Transparency, Explanations

Atoosa Kasirzadeh

October 6, 2020

Administrative

Assignment 1

- ▶ Available: October 13
- ▶ Deadline: October 28
- ▶ Grades available: Around November 7 (before final exam)

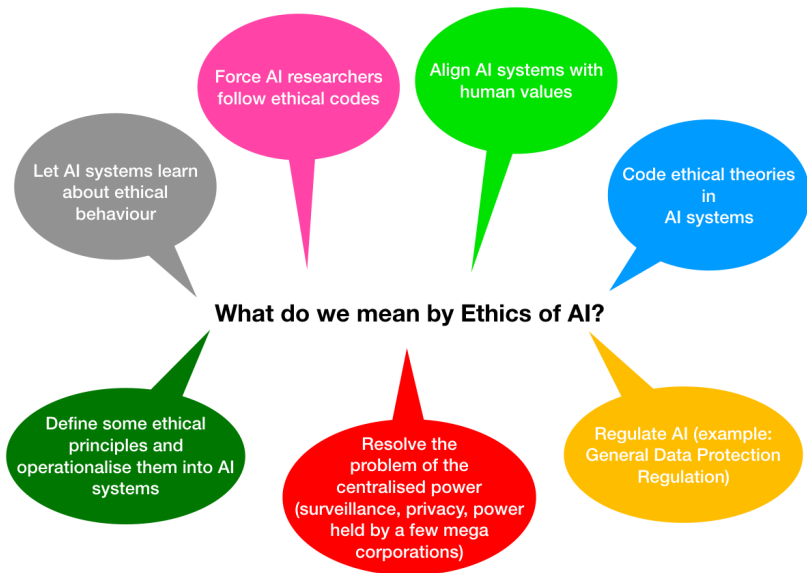
Assignment 2

- ▶ Available: November 1
- ▶ Deadline: November 14

Administrative

Tomorrow's lecture include a mini-lecture on how to write an effective paper for this course .

NO office hours tomorrow (I will be in an online conference).



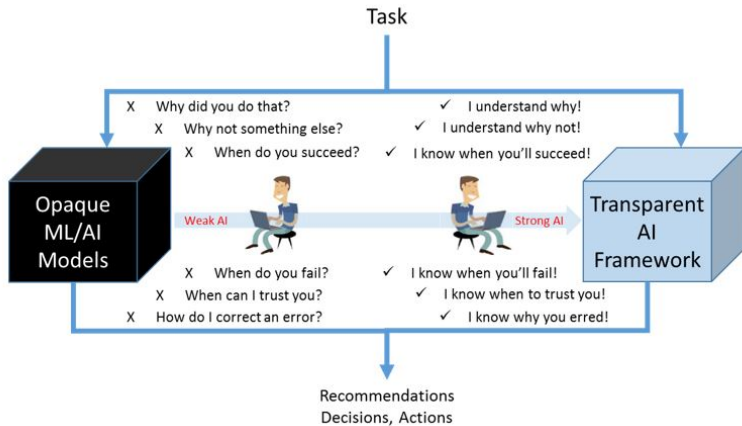
Defining some ethical principles and operationalizing them into
AI systems

This week: transparency and explainability

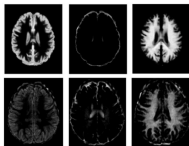
Governments and private actors are using some **truly opaque** AI algorithms to resolve critical decision-making problems.

- ▶ Hiring employees
- ▶ Assigning loans and credit scores
- ▶ Medical diagnosis
- ▶ Criminal recidivism

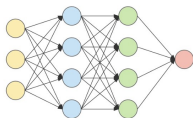
These algorithms are truly opaque: it is difficult for humans to understand why an algorithmic outcome is achieved.



Epilepsy Detection Model with Brain MRI Data



Brain MRI data



Complex ML model



Report:

Patient is
diagnosed
with **Epilepsy**
with %85
confidence.



But why?!

Can I trust this
prediction?

One popular proposal to overcome this opaqueness: require the algorithms to explain themselves (explainable AI)

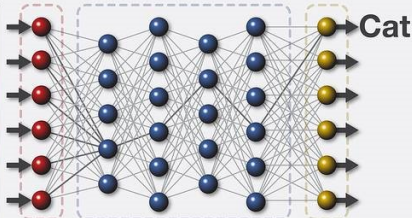
- ▶ To increase the societal acceptance of prediction-based decisions
- ▶ To establish trust in the results of these decisions
- ▶ To make these algorithms accountable to the public
- ▶ To prevent the sources of algorithmic discrimination and unfairness To legitimize the incorporation of the AI algorithms in several decision contexts
- ▶ Legal right to explanation for those affected by algorithmic decisions
- ▶ To facilitate a fruitful conversation among different stakeholders concerning the justification of using these algorithms for decision making

But what is “an explanation”

and

explainable for whom?

Machine Learning System



This is a cat.

Current Explanation

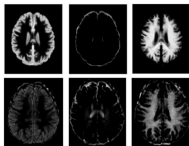
This is a cat:

- It has fur, whiskers, and claws.
- It has this feature:

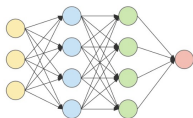


XAI Explanation

Epilepsy Detection Model with Brain MRI Data



Brain MRI data



Complex ML model



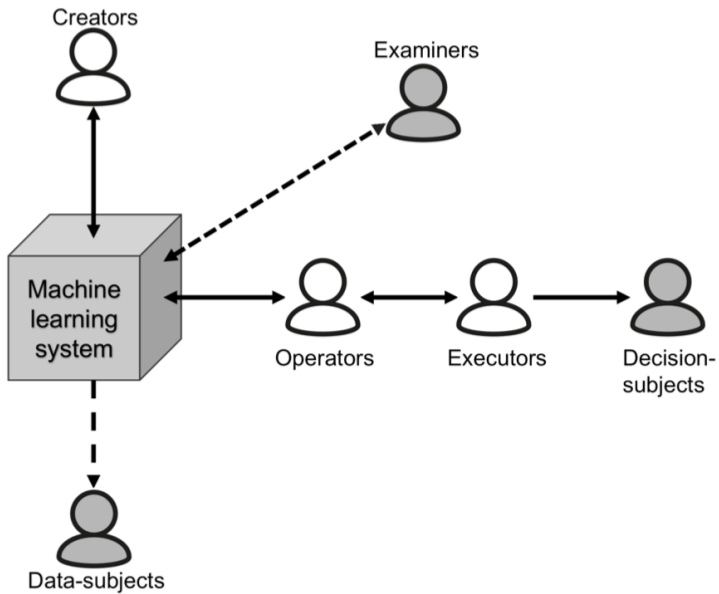
Report:

Patient is
diagnosed
with **Epilepsy**
with %85
confidence.



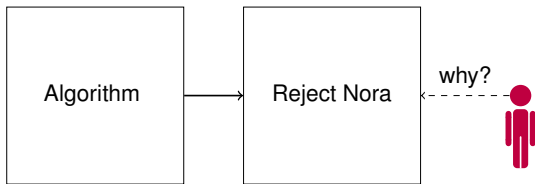
But why?!

Can I trust this
prediction?

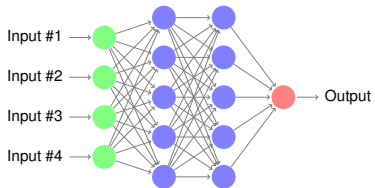
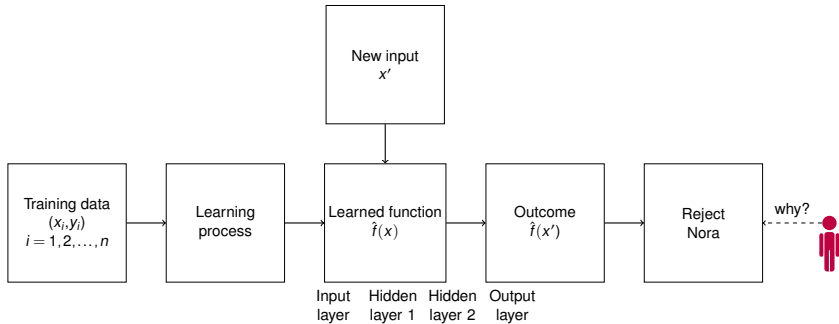


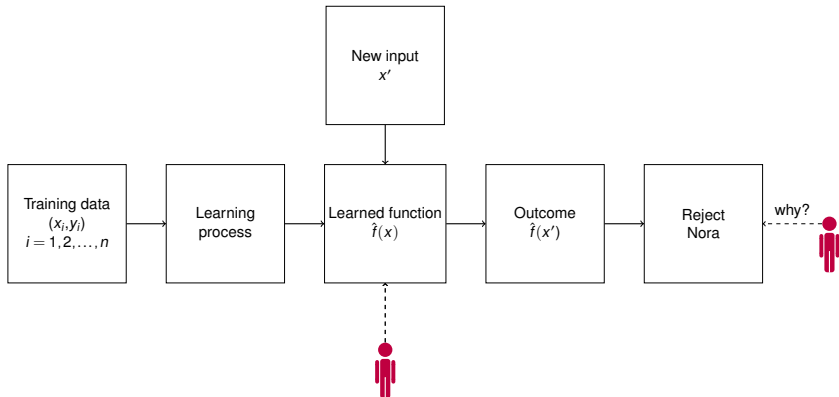
A machine-learning algorithm sifts through several job applications to recommend a hire for company X.

Nora, a competent candidate, applies for the job. Her application gets rejected by the algorithmic decision. Nora wants to know why she is rejected: she searches for an explanation.



How nice would it be if the algorithm were able to explain its decision?





What kinds of explanations do you seek?

