FRANK PASQUALE

# THE BLACK BOX SOCIETY

The Secret Algorithms
That Control Money
and Information

(1)
Data

# Meaning of Explainability
# High Level Group on AI

Explicability is crucial for building and maintaining users' trust in AI systems.

➤ **Processes need to be transparent**

➤ **Capabilities & purpose of AI systems need to be openly communicated**

➤ **Decisions need to be – to the extent possible – explainable** to those directly and indirectly affected.

Without such information, a decision cannot be duly contested.

cerre

# Meaning of Explainability
# High Level Group on AI

Explanation as to **why** a model has generated a particular output or decision (and **what** combination of input factors contributed) is not always possible.

**'BLACK BOX' ALGORITHMS REQUIRE SPECIAL ATTENTION**

**Other explicability measures** may be required, provided that the system as a whole respects fundamental rights:

- **traceability**

- **auditability**

- **transparent communication on system capabilities**

# Meaning of Explainability
# High Level Group on AI

The degree to which explicability is needed is

**highly dependent on the**

**context & severity of the consequences**

if that output is erroneous or otherwise inaccurate.

# Meaning of Explainability
# High Level Group on AI

**EXPLAINABILITY**: ability to explain the **technical processes** of an AI system & the related human **decisions**.

**Technical explainability** requires that the decisions made by an AI system can be **understood and traced** by human beings.

**Trade-offs** might have to be made:

- enhancing a system's explainability - may reduce accuracy
- increasing its accuracy - at the cost of explainability

cerre

# Meaning of Explainability
# High Level Group on AI

If AI system has **significant impact on people's lives** -> possibility to request suitable explanation of the system's decision-making process.

**Explanation should be timely and adapted** to the expertise of the stakeholder concerned (e.g. layperson, regulator, researcher).

Explanations should be available on the degree to which an AI system influences and shapes:

- the organisational decision-making process

- design choices of the system

- the rationale for deploying it

→ business model transparency

# Meaning of Explainability
# ICO Draft Guidance

cerre

| | | |
|---|---|---|
| Rationale explanation | Responsibility explanation | Data explanation |
| Fairness explanation | Safety & performance explanation | Impact explanation |

# Meaning of Explainability
# Computer science

- Ability for an abstract **mathematical model** to be **understood** by its users

- **Interpretable models**
  - Simple mathematical expression (e.g. linear models)
  - Representation allows users to understand their mathematical expression (e.g. decision trees)

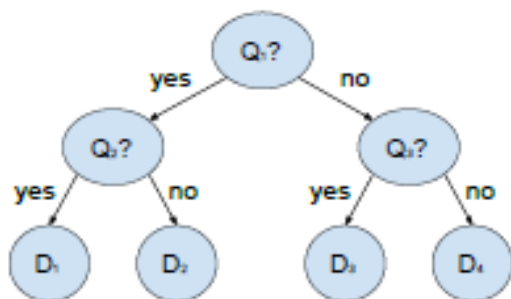- **Black-box models**
  - Neural networks

# Legal Obligations on XAI

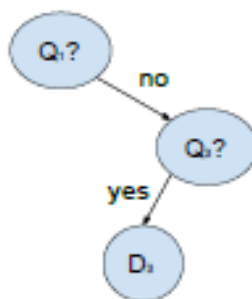| | Horizontal legislation | Sector-specific legislation |
|---|---|---|
| **AI-specific obligations** | - Personal data protection: GDPR, Convention 108+<br>- Consumer acquis<br>- P2B Regulation | - Finance<br>- Health<br>- Automotive … |
| **General obligations** | Consumer acquis | |

# Legal Obligations on XAI

| | |
|---|---|
| **Main features** | Directive 2011/83 on Consumer Rights, art. 6(a):<br><br>*obligation to provide "**the main parameters**" and "**the relative importance of those parameters**"*<br><br>Regulation 2019/1150 on promoting fairness and transparency of online intermediation services, art. 5:<br><br>*obligation to provide "the **main parameters**" and "the **relative importance** of those parameters"* |
| **All features** | GDPR, art.22 and Guidelines on Automated individual decision-making & Profiling<br><br>*obligation to provide "the **criteria relied on in reaching the decision**"* |
| **Combination of features** | GDPR, art.22 and Guidelines on Automated individual decision-making and Profiling<br><br>*obligation to provide "the **rationale behind the decision**"* |
| **Whole model** | Directive 2014/65 on Markets in Financial Instruments, art. 17<br><br>obligation to provide "**information [...] about its algorithmic trading and the systems used for that trading**" |

1) Decision Tree

2) Particular Decision of the Decision Tree

3) Features Involved in the Decision

# Implementation in ML models

| | |
|---|---|
| **Main features** | Well developed in ML<br><br>Linear models: weight to features, strongly and weakly relevant features<br><br>Black-box models: features sorted by importance |
| **All features** | Possible for all ML models<br><br>Some models show trade-off between accuracy and complexity |
| **Combination of features** | Require the use of transparent models: decision tree, linear<br><br>Create new ones to explain black box models: local explanation |
| **Whole model** | Possible only for some models |

# Issues

## 50 SHADES OF TRANSPARENCY

Explanation – audit
Types of explanations
Timing: ex ante, ex post

## MANY RULES EXIST ALREADY

Non AI specific and AI specific
Horizontal and sector specific

## RISK-BASED AND PROPORTIONALITY

## TRADE-OFFS

Accuracy & explainability
Rights of users
Rights of AI developers/owners

## cerre

### Centre on Regulation in Europe

📍 Avenue Louise, 475 (box 10)
   1050 Brussels, Belgium

📱 +32 2 230 83 60

✉ info@cerre.eu

🖱 cerre.eu

🐦 @CERRE_ThinkTank

**cerre.eu**

Improving network and digital industries regulation