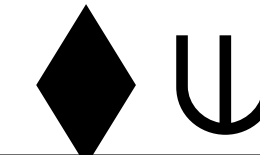# Backgrounds

## Black-Box Models

- High accuracy & low error
- Lacks transparency in decisions

## Instance-wise Feature Selection (IFS)

▶ Select instance-specific set of important features, by assigning an *<u>importance score</u>* to each feature

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

• Instances: $x_i = (x_i^1, x_i^2, \ldots, x_i^d)$

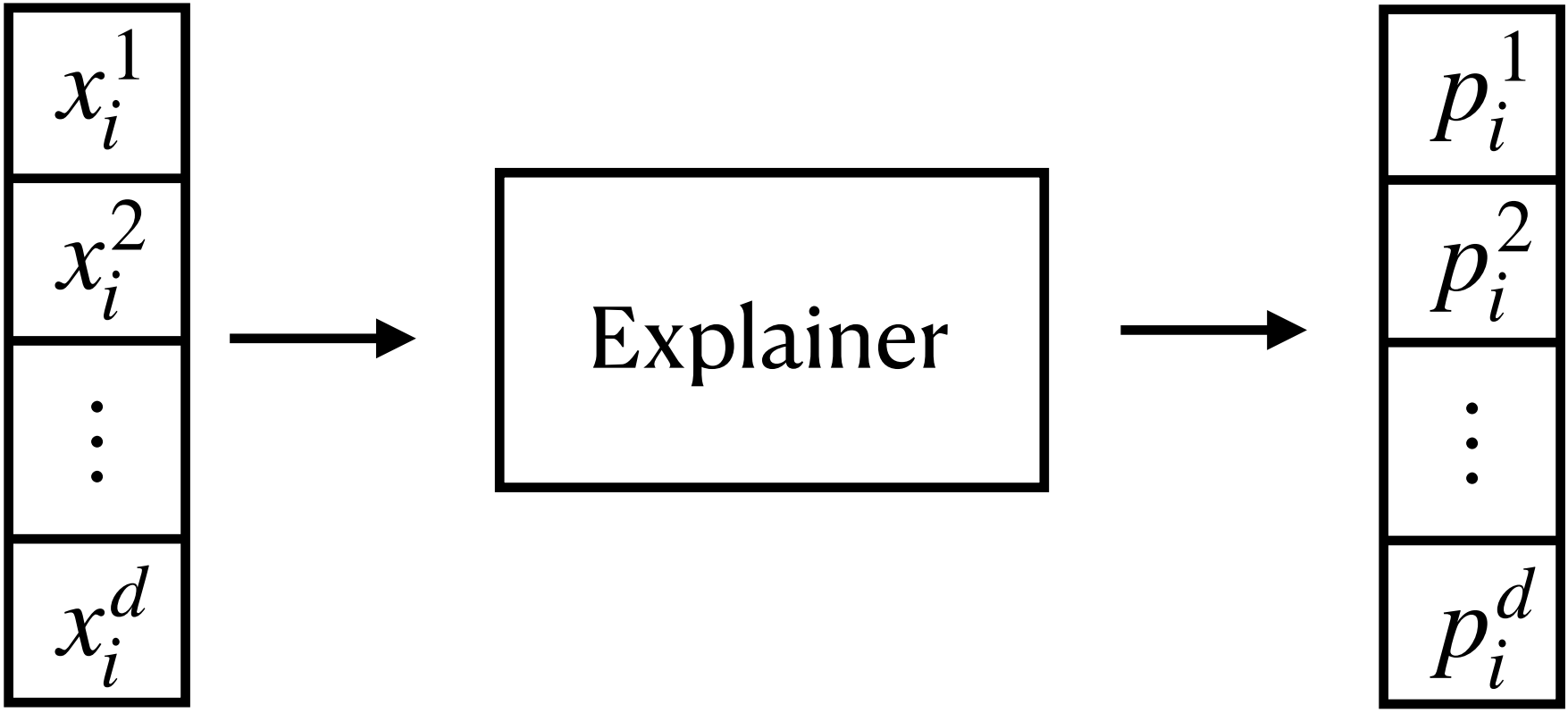• Black-box model to explain: $M$

➡ Response variable: $y_i = M(x_i)$

➡ Dataset: $D = (x_i, y_i)_{i=1}^n$

## Explainable AI

- **Techniques** that make decisions of AI models *understandable* and *interpretable* to humans

- Understand whys begind black-box models

← Instance-specific feature inprotance

# Backgrounds

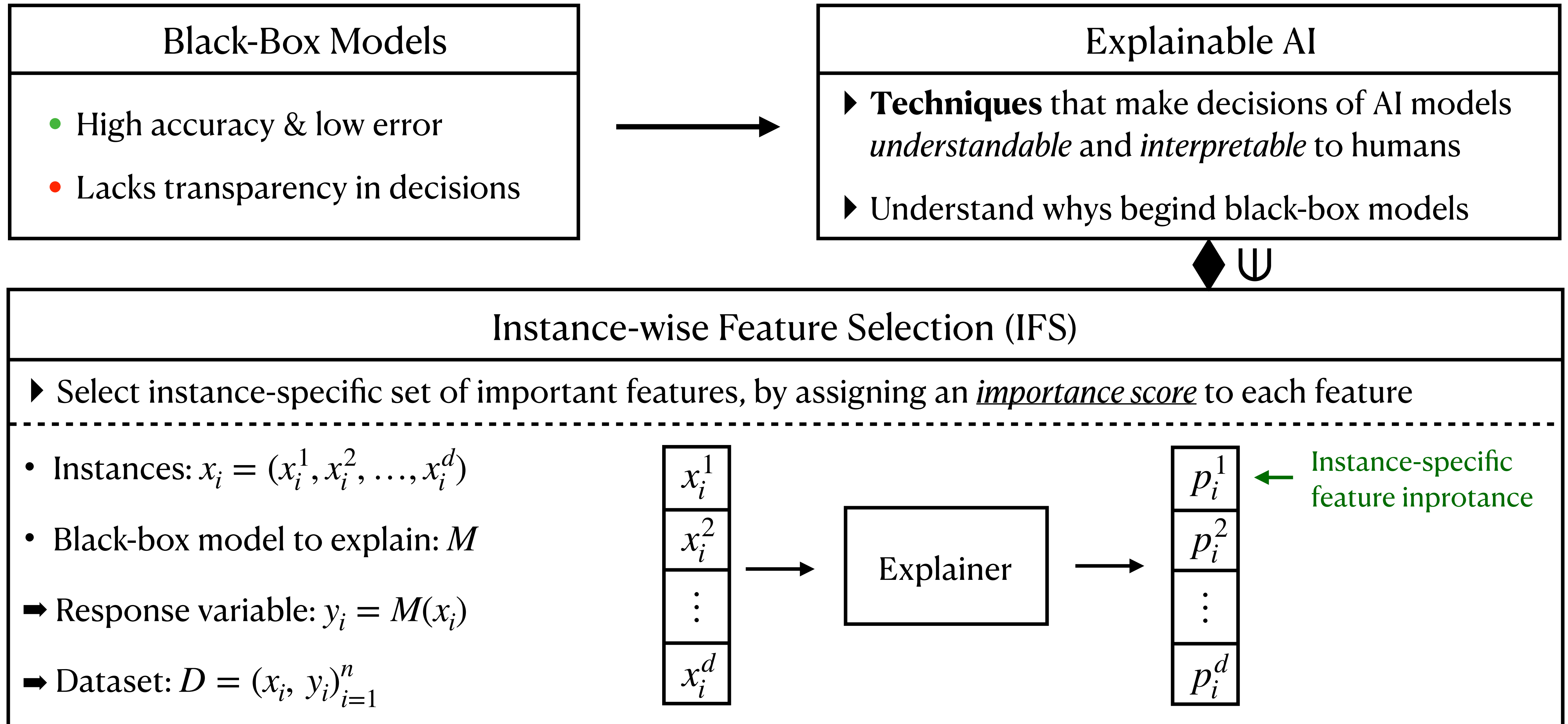| Black-Box Models | Explainable AI |
|---|---|
| • High accuracy & low error <br><br> • Lacks transparency in decisions | ▸ **Techniques** that make decisions of AI models *understandable* and *interpretable* to humans <br><br> ▸ Understand whys begind black-box models |

## Instance-wise Feature Selection (IFS)

▸ Select instance-specific set of important features, by assigning an <u>*importance score*</u> to each feature

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

• Instances: $x_i = (x_i^1, x_i^2, \ldots, x_i^d)$

• Black-box model to explain: $M$

➡ Response variable: $y_i = M(x_i)$

➡ Dataset: $D = (x_i, y_i)_{i=1}^n$

$$\begin{array}{c} x_i^1 \\ x_i^2 \\ \vdots \\ x_i^d \end{array} \longrightarrow \boxed{\text{Explainer}} \longrightarrow \begin{array}{c} p_i^1 \\ p_i^2 \\ \vdots \\ p_i^d \end{array}$$

← Instance-specific feature inprotance

# Related Work

IFS Methods $\left\{\begin{array}{l} \text{Model-specific} \\ \\ \text{Model-agnostic} \end{array}\right.$