

Explainability & Transparency



What is “Explainable AI” ?

- **Explainable-AI** explores and investigates methods to produce or complement **AI models** to make **accessible and interpretable** the internal logic and the outcome of the algorithms, making such process **understandable by humans**.
- **Explicability**, understood as incorporating both **intelligibility** (“*how does it work?*”) for non-experts, e.g., patients or business customers, and for experts, e.g., product designers or engineers) and **accountability** (“*who is responsible for?*”).

Interpretability

- To ***interpret*** means to give or provide the meaning or to explain and present in understandable terms some concepts.
- In data mining and machine learning, interpretability is the ***ability to explain*** or to provide the meaning ***in understandable terms to a human***.



- <https://www.merriam-webster.com/>

- Finale Doshi-Velez and Been Kim. 2017. ***Towards a rigorous science of interpretable machine learning***. arXiv:1702.08608v2.

Motivating Examples

- Criminal Justice
 - People wrongly denied
 - Recidivism prediction
 - Unfair Police dispatch
- Finance:
 - Credit scoring, loan approval
 - Insurance quotes
- Healthcare
 - AI as 3rd-party actor in physician - patient relationship
 - Learning must be done with available data: cannot randomize cares given to patients!
 - Must validate models before use.

Opinion

The New York Times

OP-ED CONTRIBUTOR

When a Computer Program Keeps You in Jail

The Big Read **Artificial intelligence**

+ Add to myFT

Insurance: Robots learn the business of covering risk

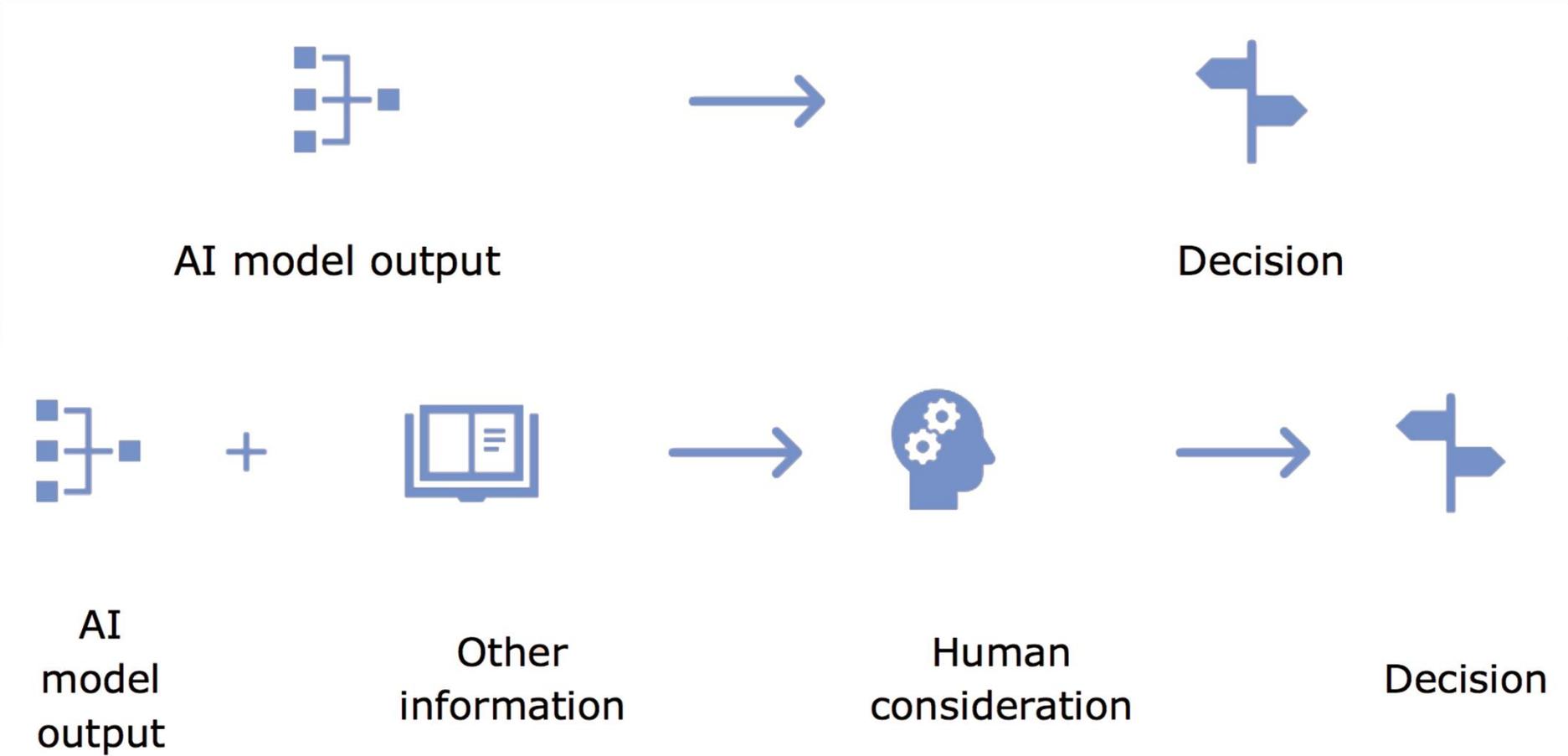
 **Stanford**
MEDICINE | News Center

[Email](#) [Tweet](#)

Researchers say use of artificial intelligence in medicine raises ethical questions

In a perspective piece, Stanford researchers discuss the ethical implications of using machine-learning tools in making health care decisions for patients.

What is AI-assisted decision making?



What is a Black Box Model?



A ***black box*** is a model, whose internals are either unknown to the observer or they are known but uninterpretable by humans.

- Guidotti, R., Monreale, A., Ruggieri, S., Turini, F., Giannotti, F., & Pedreschi, D. (2018). *A survey of methods for explaining black box models*. *ACM Computing Surveys (CSUR)*, 51(5), 93.



Needs For Interpretable Models

Right of Explanation



General Data Protection Regulation

Since 25 May 2018, GDPR establishes a right for all individuals to obtain “meaningful explanations of the logic involved” when “automated (algorithmic) individual decision-making”, including profiling, takes place.

COMPAS recidivism black bias

DYLAN FUGETT

Prior Offense
1 attempted burglary

Subsequent Offenses
3 drug possessions

LOW RISK **3**

BERNARD PARKER

Prior Offense
1 resisting arrest
without violence

Subsequent Offenses
None

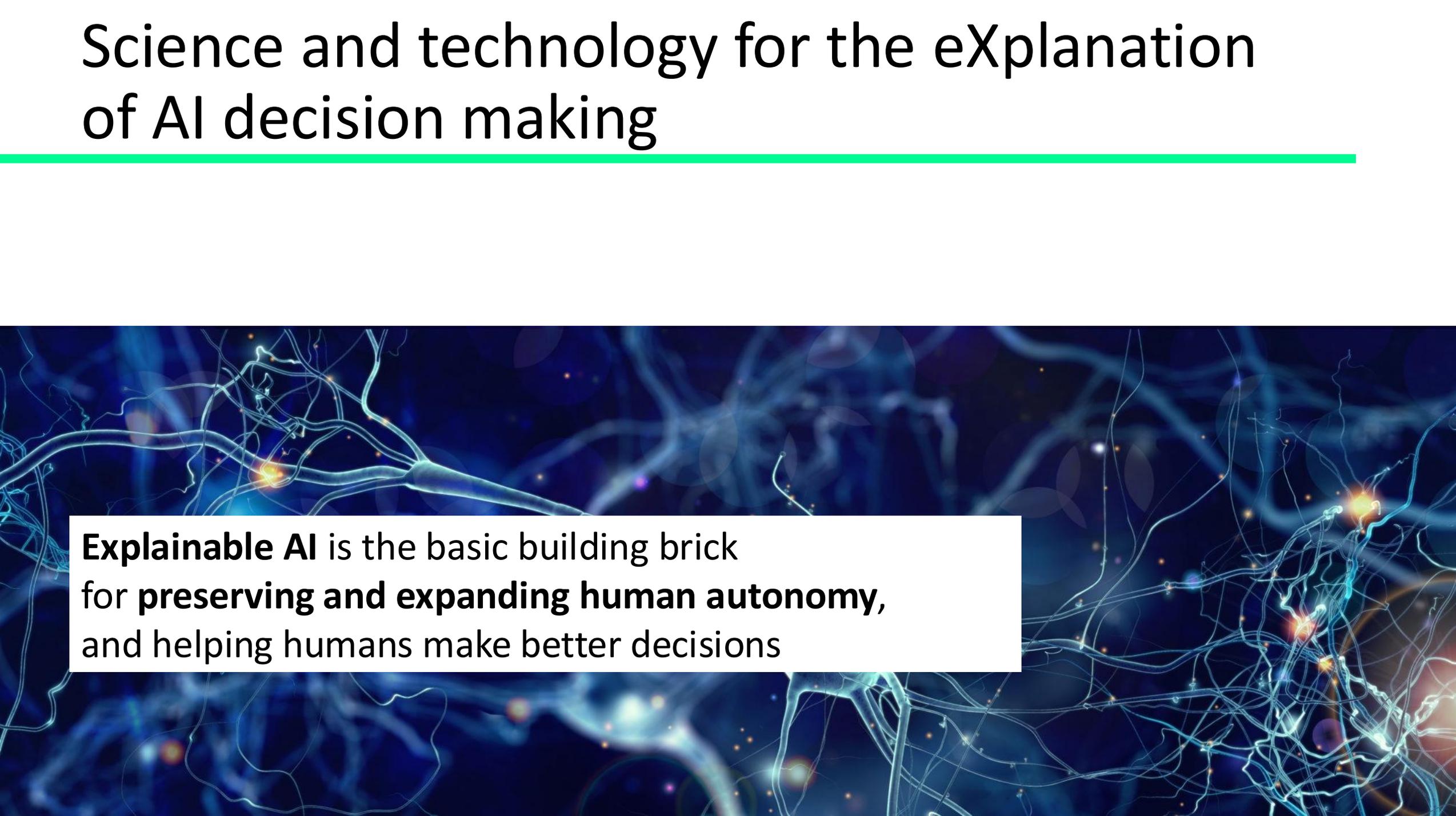
HIGH RISK **10**

Fugett was rated low risk after being arrested with cocaine and marijuana. He was arrested three times on drug charges after that.

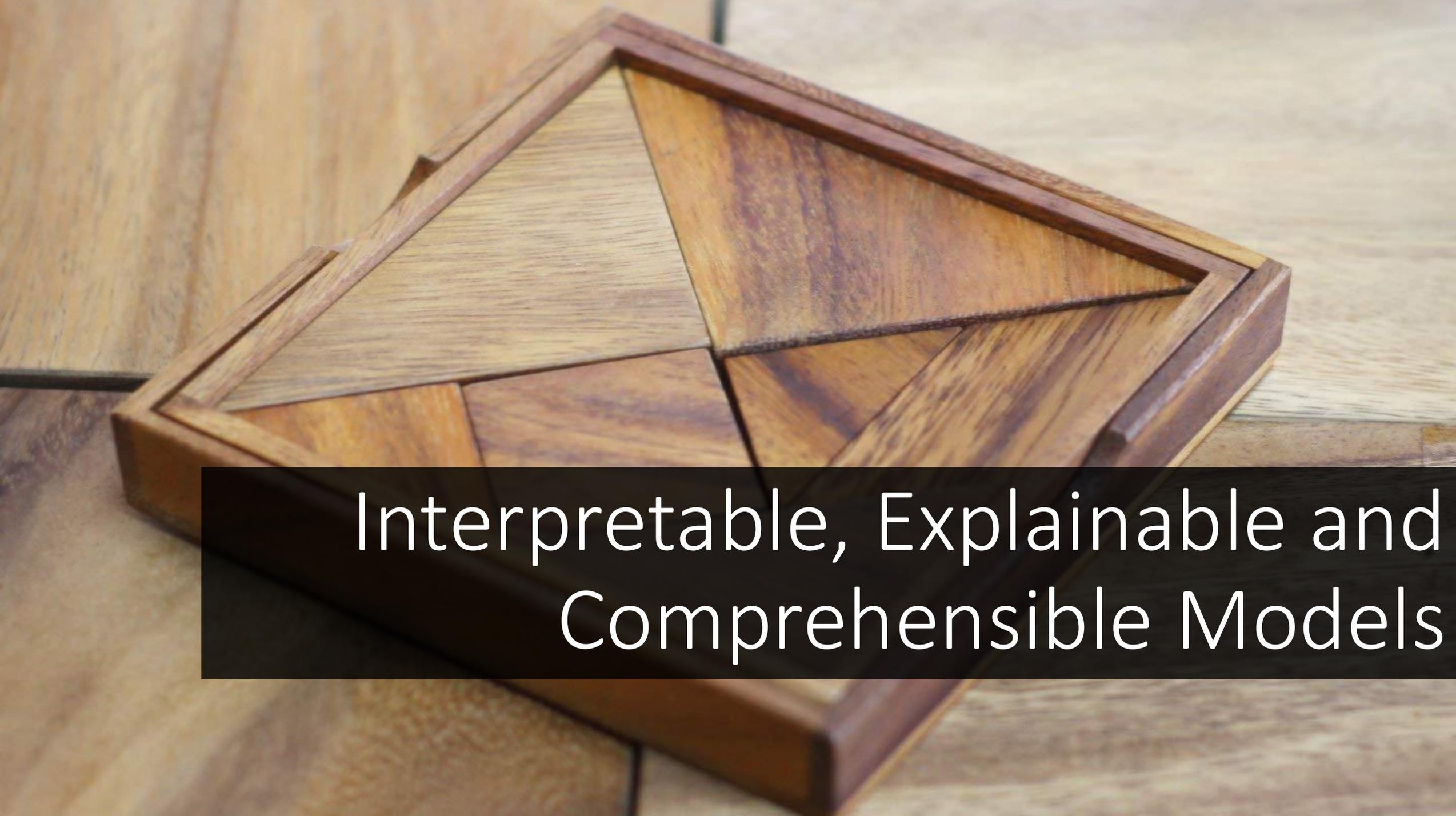
Military tank classification depends on the background



Science and technology for the eXplanation of AI decision making



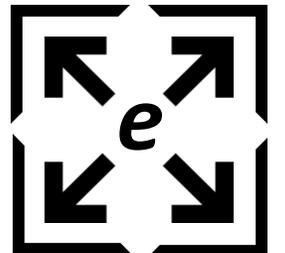
Explainable AI is the basic building brick for **preserving and expanding human autonomy**, and helping humans make better decisions

A close-up photograph of a wooden tray with a geometric, triangular pattern. The tray is made of light-colored wood with a darker wood inlay forming a large triangle. The tray is resting on a wooden surface. A black rectangular box is overlaid on the bottom right of the image, containing white text.

Interpretable, Explainable and
Comprehensible Models

Dimensions of Interpretability

- ***Global and Local Interpretability:***
 - *Global:* understanding the whole logic of a model
 - *Local:* understanding only the reasons for a specific decision
- ***Time Limitation:*** the time that the user can spend for understanding an explanation.
- ***Nature of User Expertise:*** users of a predictive model may have different background knowledge and experience in the task. The nature of the user expertise is a key aspect for interpretability of a model.



Desiderata of an Interpretable Model

- ***Interpretability*** (or comprehensibility): to which extent the model and/or its predictions are human understandable. Is measured with the *complexity* of the model.
- ***Fidelity***: to which extent the model imitate a black-box predictor.
- ***Accuracy***: to which extent the model predicts unseen instances.

- Alex A. Freitas. 2014. *Comprehensible classification models: A position paper*. ACM SIGKDD Explor. Newslett.



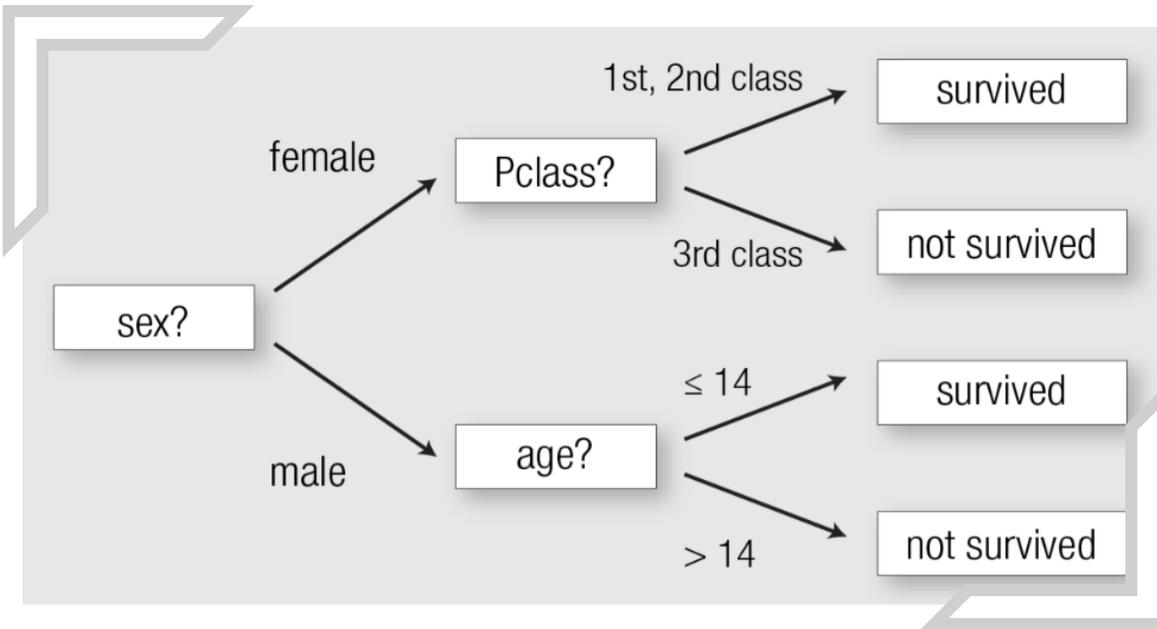
Desiderata of an Interpretable Model

- ***Fairness***: the model guarantees the protection of groups against discrimination.
- ***Privacy***: the model does not reveal sensitive information about people.
- ***Respect Monotonicity***: the increase of the values of an attribute either increase or decrease in a monotonic way the probability of a record of being member of a class.
- ***Usability***: an interactive and queryable explanation is more usable than a textual and fixed explanation.

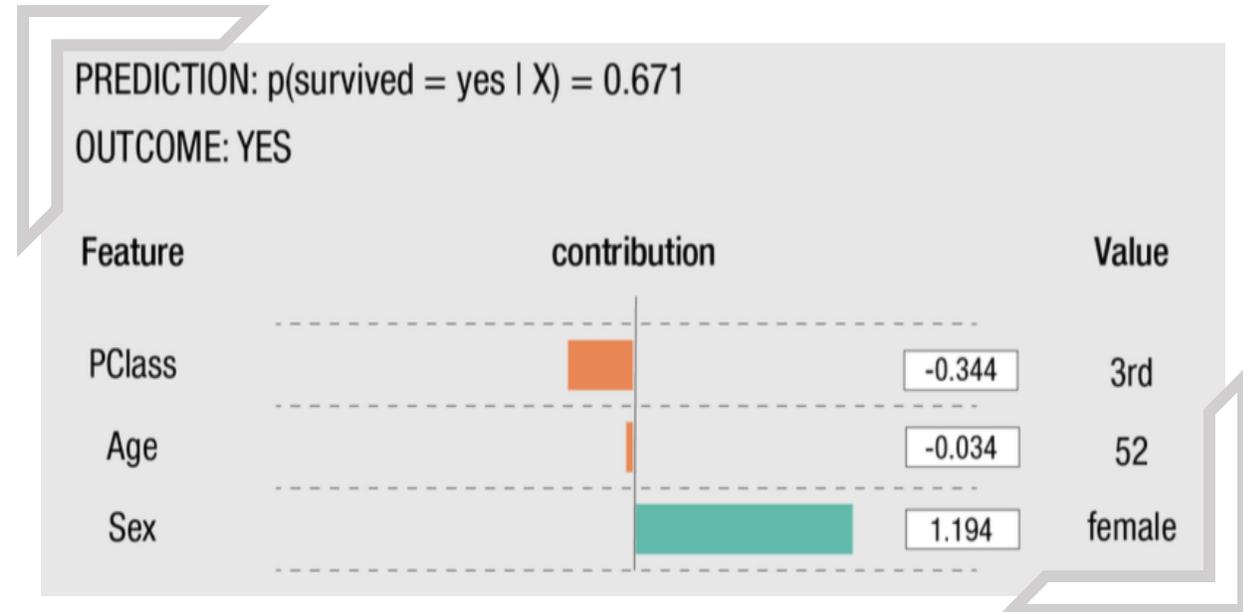
- Andrea Romei and Salvatore Ruggieri. 2014. ***A multidisciplinary survey on discrimination analysis***. Knowl. Eng.
- Yousra Abdul Alsaheb S. Aldeen, Mazleena Salleh, and Mohammad Abdur Razzaque. 2015. ***A comprehensive review on privacy preserving data mining***. SpringerPlus .
- Alex A. Freitas. 2014. ***Comprehensible classification models: A position paper***. ACM SIGKDD Explor. Newslett.



Recognized Interpretable Models



Decision Tree



Linear Model

if condition₁ \wedge condition₂ \wedge condition₃ then outcome

Rules

There are several kinds of explanations



Sorry, your loan application has been rejected.

Our analysis:

The following features were too high:

PercentInstallTrad...

NetFractionRevolv...

NetFractionInstall...

NumRevolvingTra...

NumBank2NatTra...

PercentTradesWB...

The following features were too low:

MSinceOldestTrad...

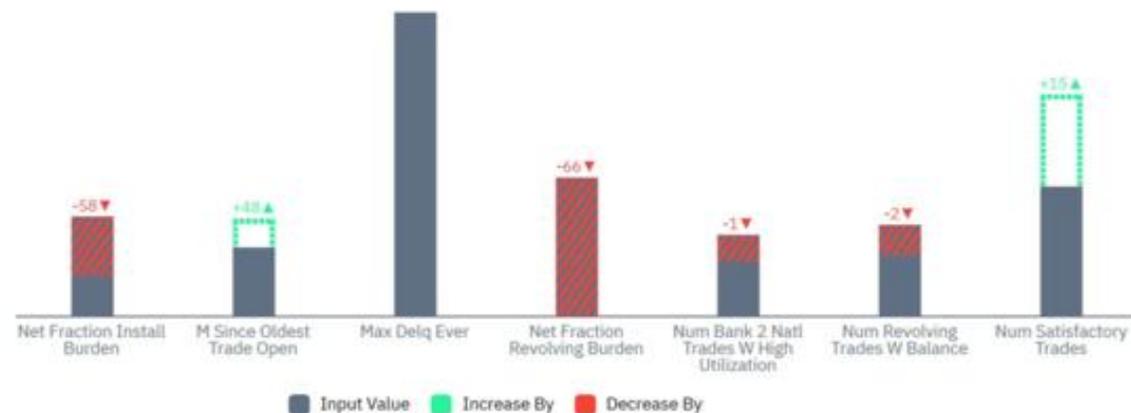
AverageMInFile

NumTotalTrades

The following features require changes:

MaxDelq2PublicR...

MaxDelqEver

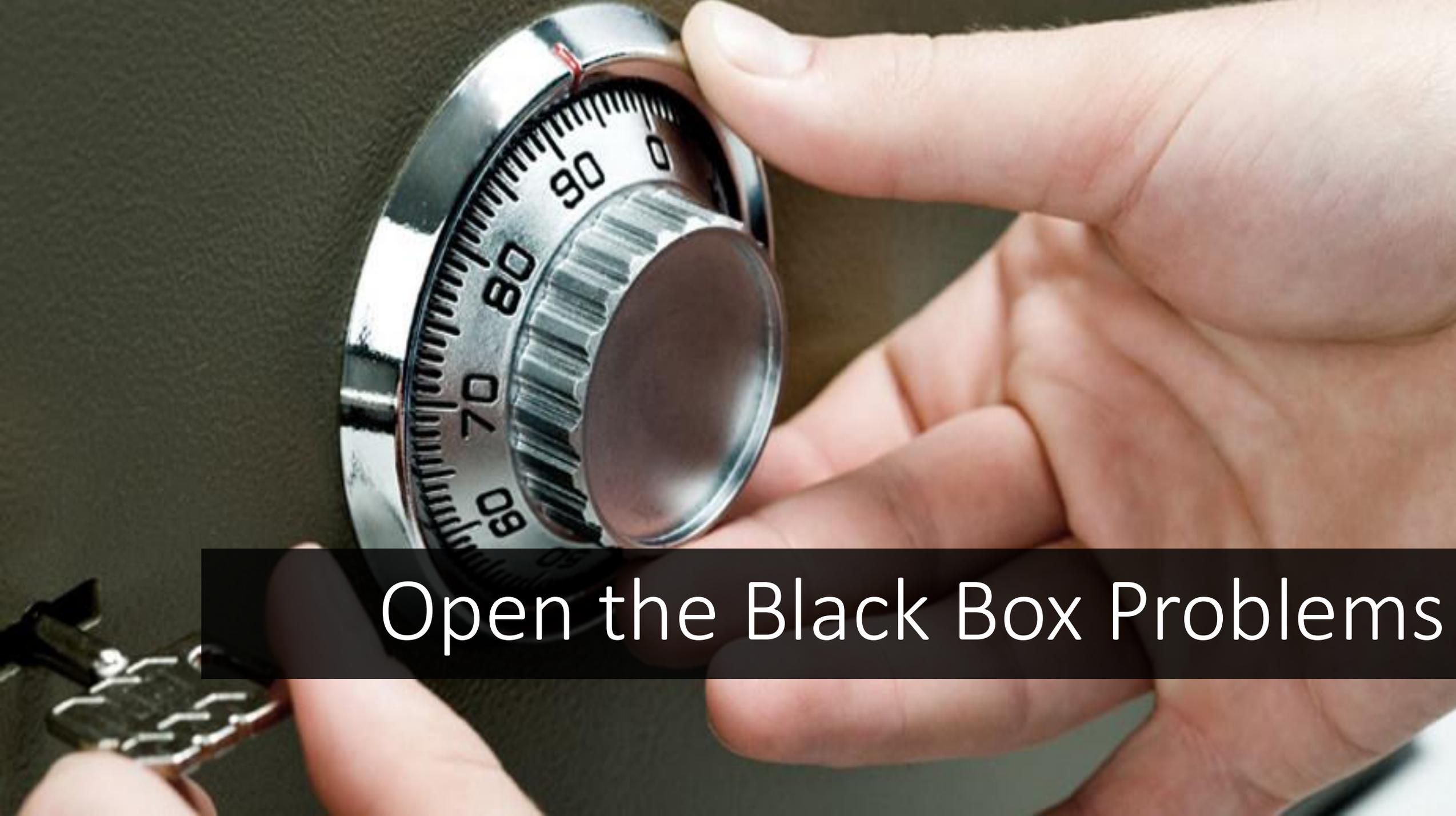


Counterfactuals suggest where to increase (green, dashed) or decrease (red, striped) each feature.

Complexity

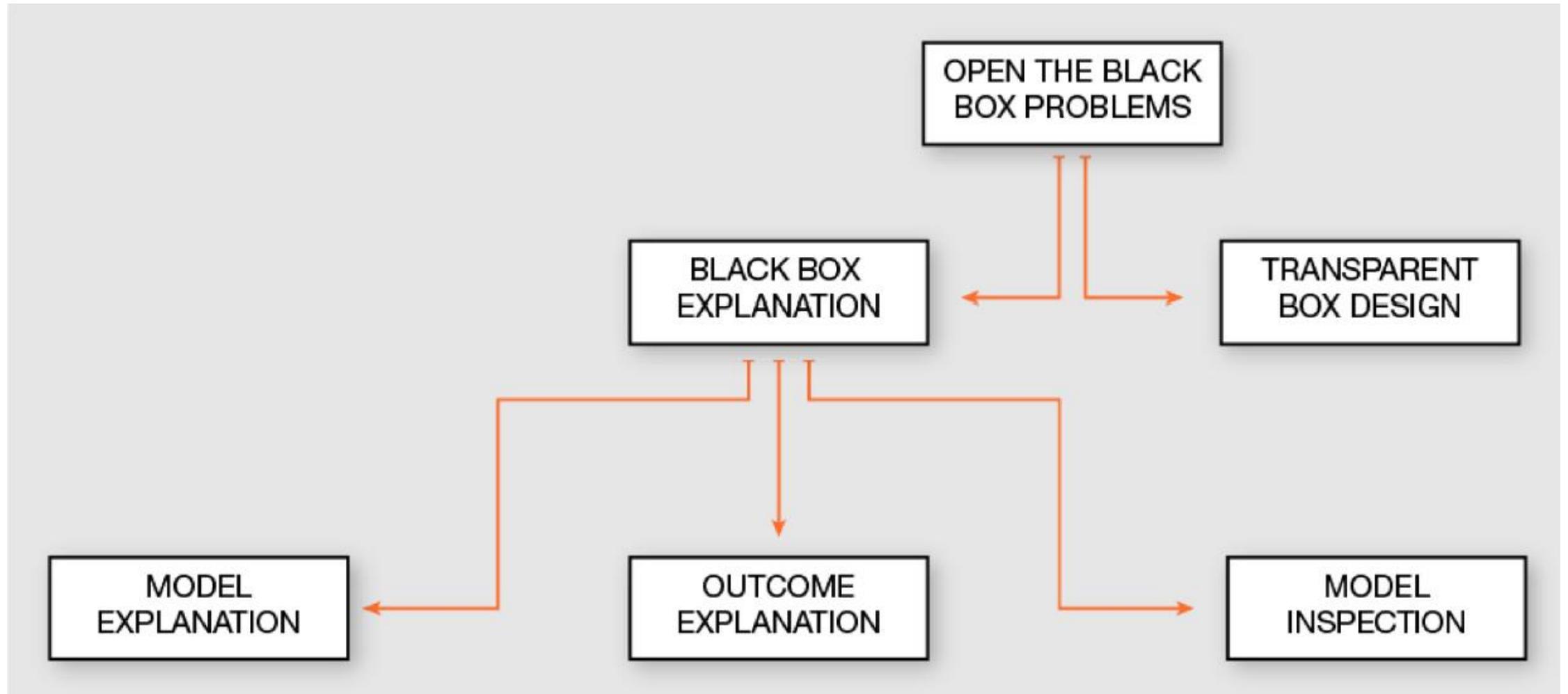
- Opposed to *interpretability*.
- Is only related to the model and not to the training data that is unknown.
- Generally estimated with a rough approximation related to the **size** of the interpretable model.
- Linear Model: number of non zero weights in the model.
- Rule: number of attribute-value pairs in condition.
- Decision Tree: estimating the complexity of a tree can be hard.

- Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. 2016. *Why should i trust you?: Explaining the predictions of any classifier*. KDD.
- Houtao Deng. 2014. *Interpreting tree ensembles with intrees*. arXiv preprint arXiv:1408.5456.
- Alex A. Freitas. 2014. *Comprehensible classification models: A position paper*. ACM SIGKDD Explor. Newslett.

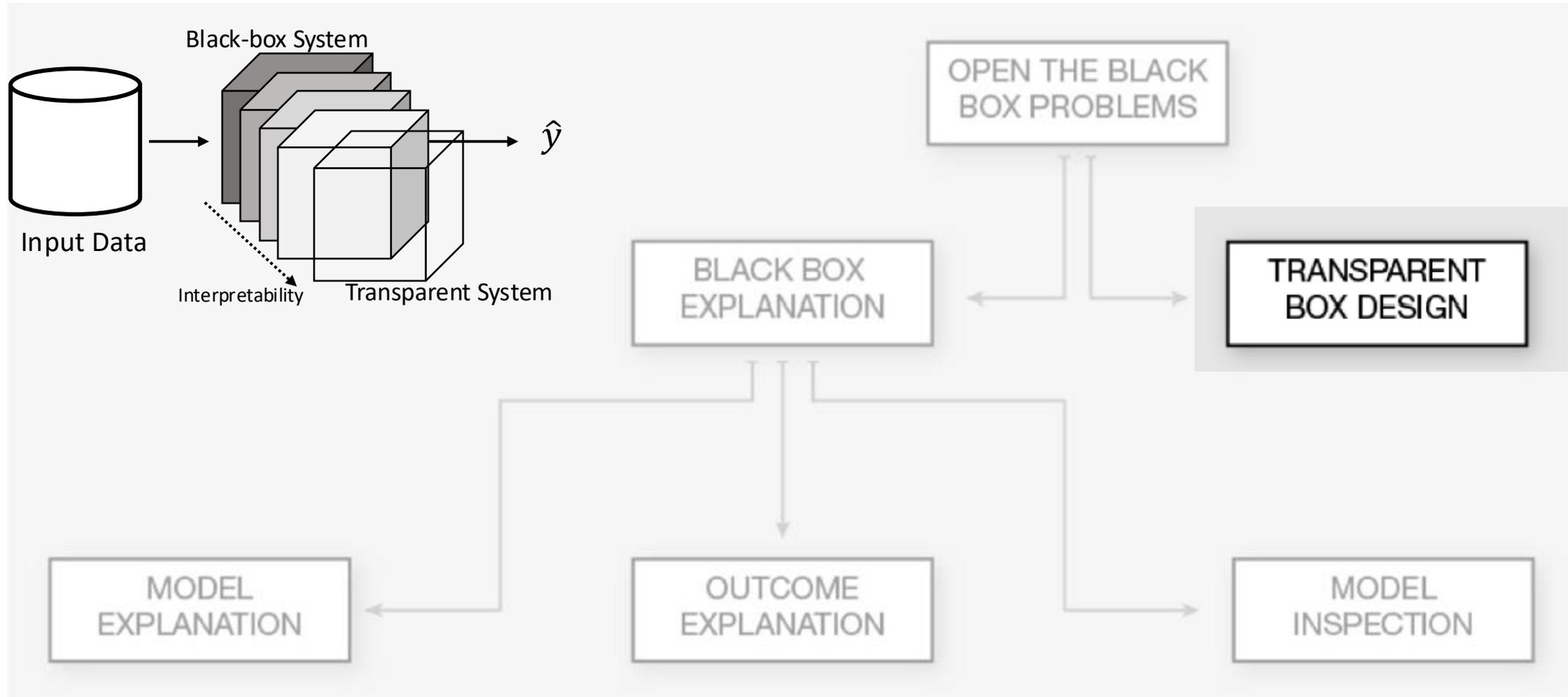
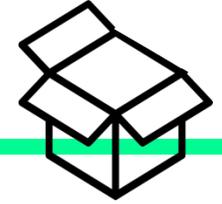


Open the Black Box Problems

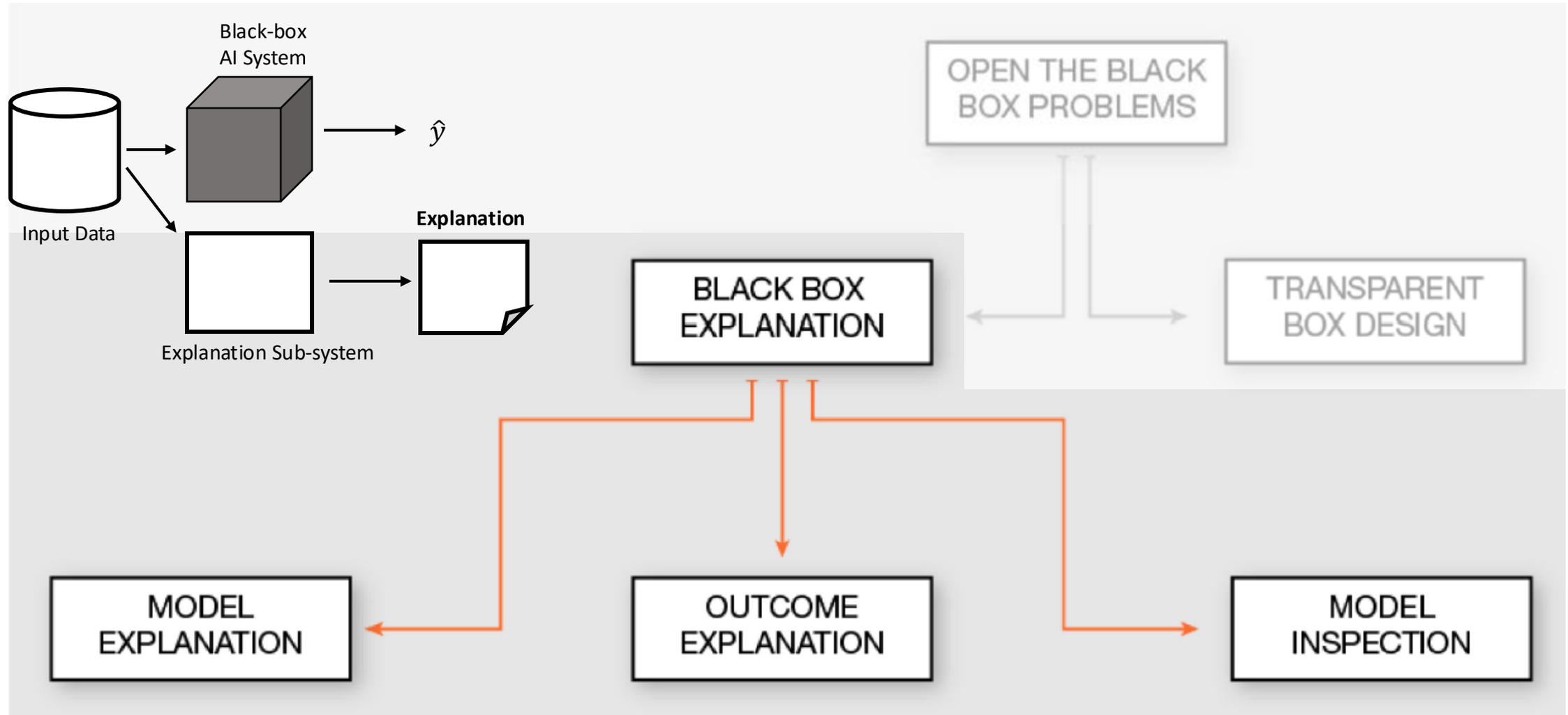
Problems Taxonomy



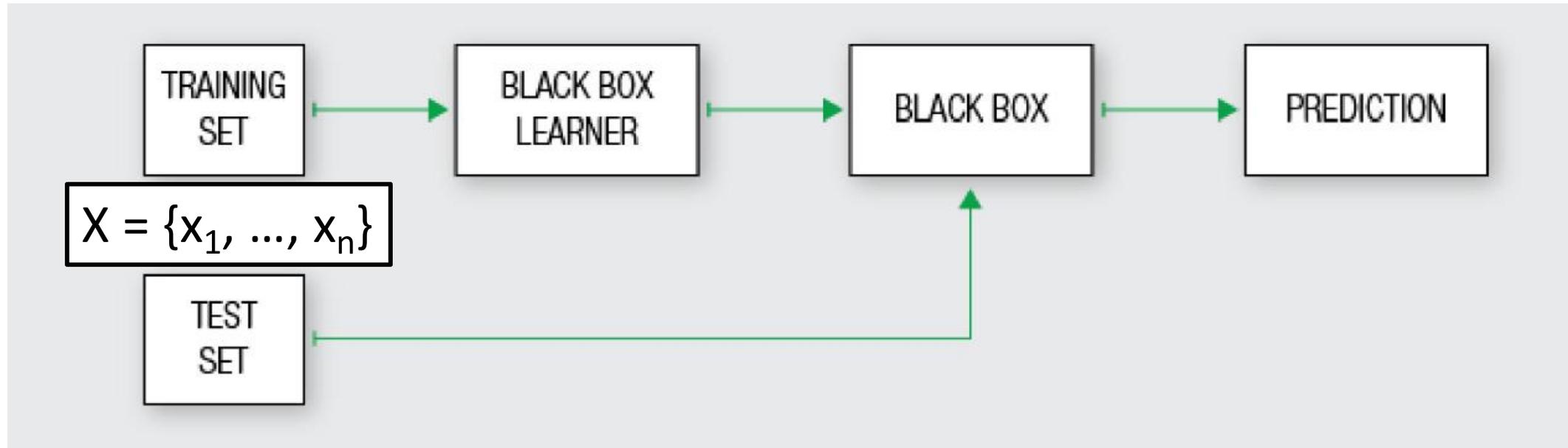
XbD – eXplanation by Design



BBX - Black Box eXplanation



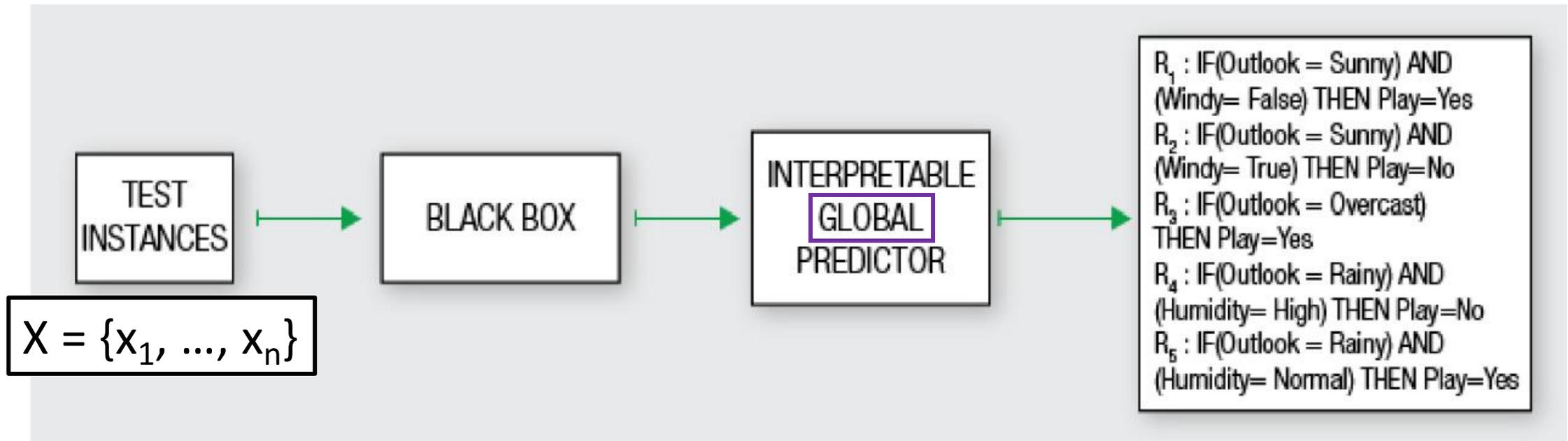
Classification Problem



Model Explanation Problem



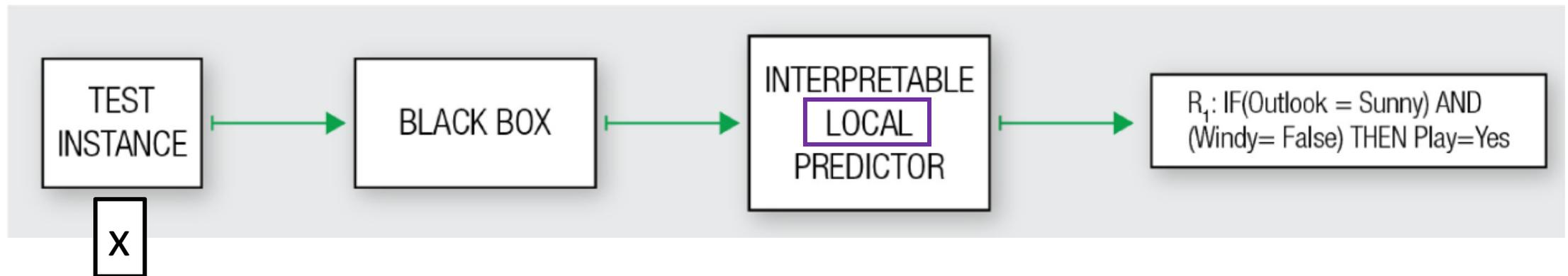
Provide an interpretable model able to mimic the *overall logic/behavior* of the black box and to explain its logic.



Outcome Explanation Problem



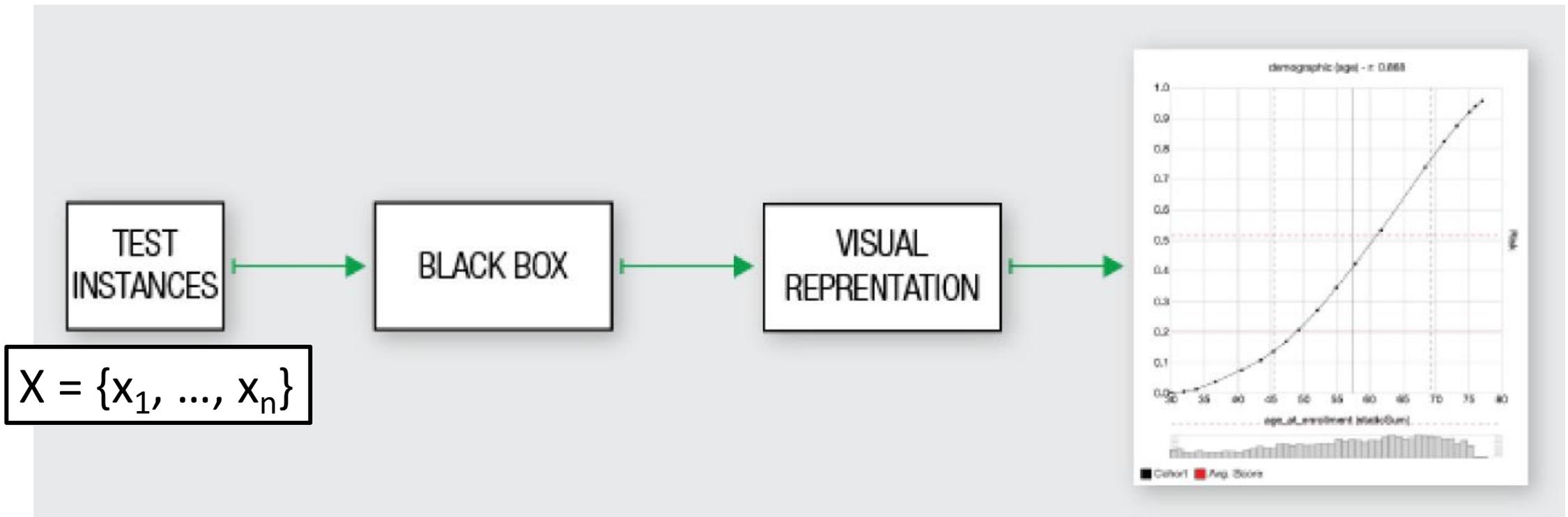
Provide an interpretable outcome, i.e., an *explanation* for the outcome of the black box for a *single instance*.



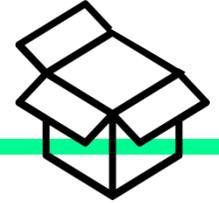
Model Inspection Problem



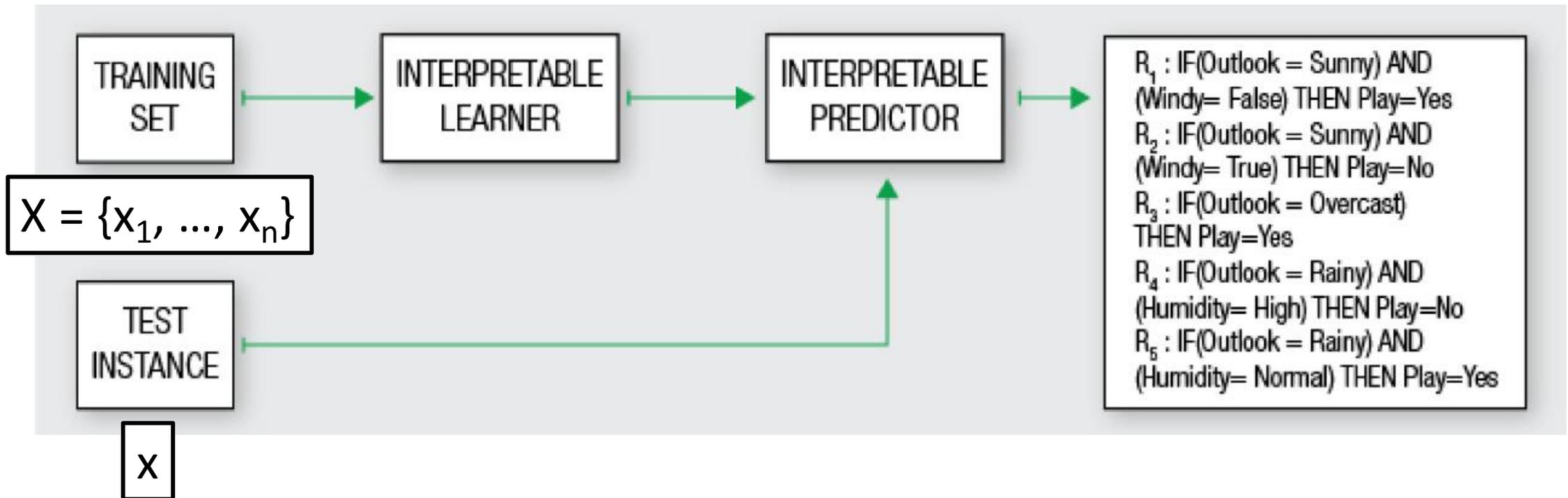
Provide a representation (visual or textual) for understanding either how the black box model works or why the black box returns certain predictions more likely than others.



Transparent Box Design Problem

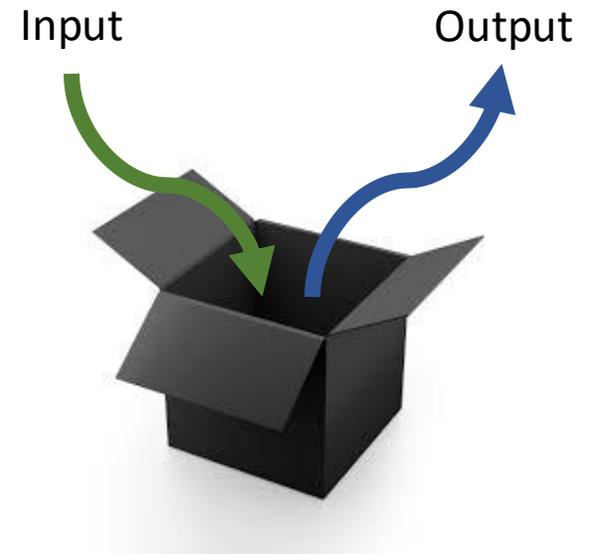


Provide a model which is locally or globally interpretable on its own.

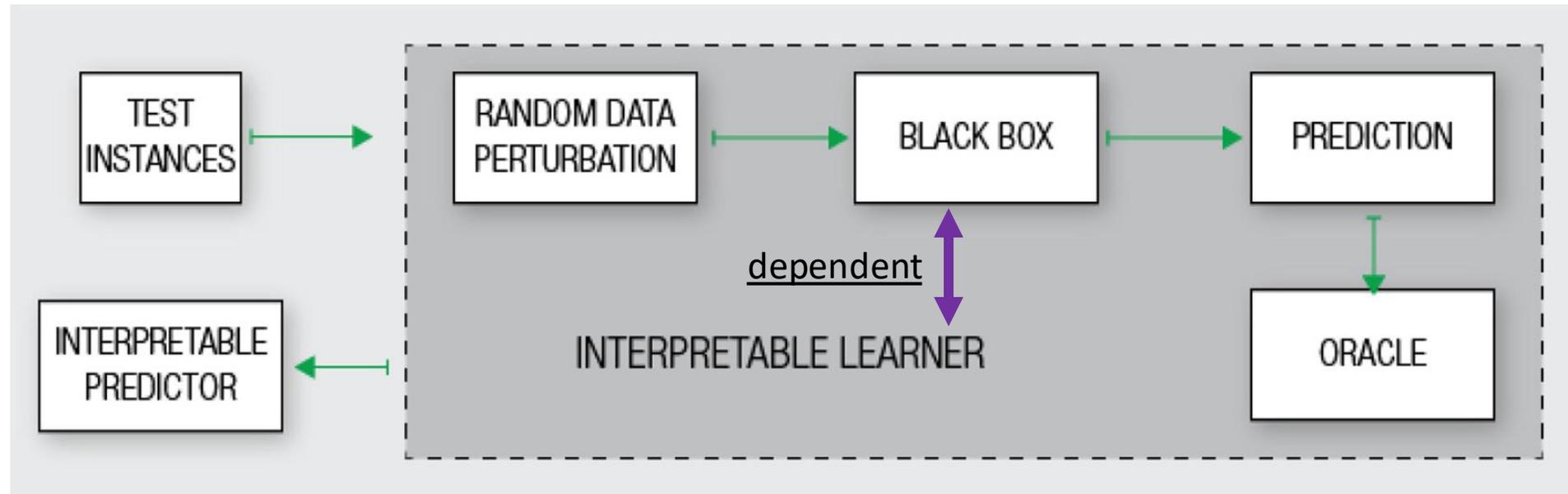
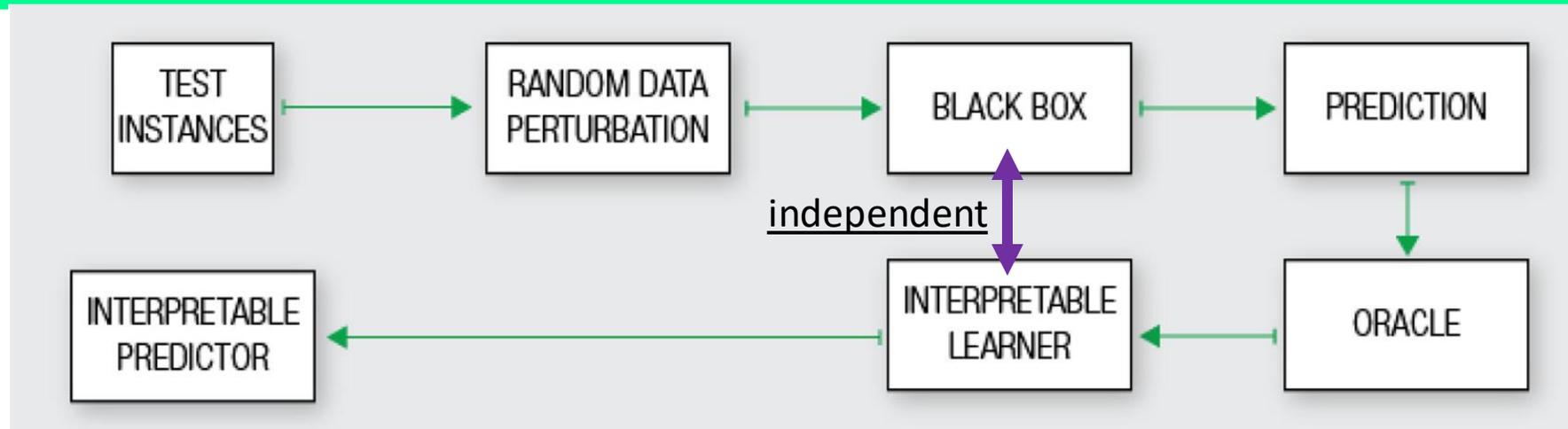


Reverse Engineering

- The name comes from the fact that we can only **observe** the **input** and **output** of the black box.
- Possible actions are:
 - **choice** of a particular comprehensible predictor
 - querying/auditing the black box with input records created in a controlled way using **random perturbations** w.r.t. a certain prior knowledge (e.g. train or test)
- It can be **generalizable or not**:
 - Model-Agnostic
 - Model-Specific



Model-Agnostic vs Model-Specific



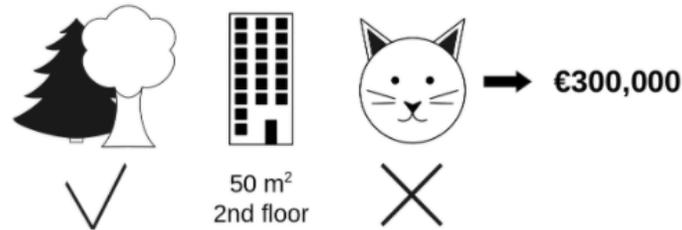
<i>Name</i>	<i>Ref.</i>	<i>Authors</i>	<i>Year</i>	<i>Explanator</i>	<i>Black Box</i>	<i>Data Type</i>	<i>General</i>	<i>Random</i>	<i>Examples</i>	<i>Code</i>	<i>Dataset</i>
–	[134]	Xu et al.	2015	SM	DNN	IMG			✓	✓	✓
–	[30]	Fong et al.	2017	SM	DNN	IMG			✓		
CAM	[139]	Zhou et al.	2016	SM	DNN	IMG			✓	✓	✓
Grad-CAM	[106]	Selvaraju et al.	2016	SM	DNN	IMG			✓	✓	✓
–	[109]	Simonian et al.	2013	SM	DNN	IMG			✓		✓
PWD	[7]	Bach et al.	2015	SM	DNN	IMG			✓		✓
–	[113]	Sturm et al.	2016	SM	DNN	IMG			✓		✓
DTD	[78]	Montavon et al.	2017	SM	DNN	IMG			✓		✓
DeapLIFT	[107]	Shrikumar et al.	2017	FI	DNN	ANY			✓	✓	
CP	[64]	Landecker et al.	2013	SM	NN	IMG			✓		
–	[143]	Zintgraf et al.	2017	SM	DNN	IMG			✓	✓	✓
VBP	[11]	Bojarski et al.	2016	SM	DNN	IMG			✓	✓	✓
–	[65]	Lei et al.	2016	SM	DNN	TXT			✓		✓
ExplainD	[89]	Poulin et al.	2006	FI	SVM	TAB		✓	✓		
–	[29]	Strumbelj et al.	2010	FI	AGN	TAB	✓	✓	✓		✓

Solving The Outcome Explanation Problem

SHAP

A prediction can be explained by assuming that **each feature value of the instance is a "player"** in a game where the **prediction is the payout**. **Shapley values** tells us how to fairly distribute the "payout" among the features.

Example



Prediction: You have trained a machine learning model to predict apartment prices. For a certain apartment it predicts €300,000 and you need to explain this prediction.

The apartment has an area of 50 m², is located on the 2nd floor, has a park nearby and cats are banned.

The average prediction is €310,000.

How much has each feature value contributed to the prediction compared to the average prediction?

SHAP

The average prediction is **€310,000** while the prediction is **€300,000**

How much has each feature value contributed to the prediction compared to the average prediction?

The answer is simple for linear regression models. **The effect of each feature is the weight of the feature times the feature value.** This only works because of **the linearity of the model.**

For more complex models, we need a different solution!!!!

GOAL: explain the difference between the actual prediction (€300,000) and the average prediction (€310,000): a difference of -€10,000.

SHAP

GOAL: explain the difference between the actual prediction (€300,000) and the average prediction (€310,000): a difference of -€10,000.

Game theory:

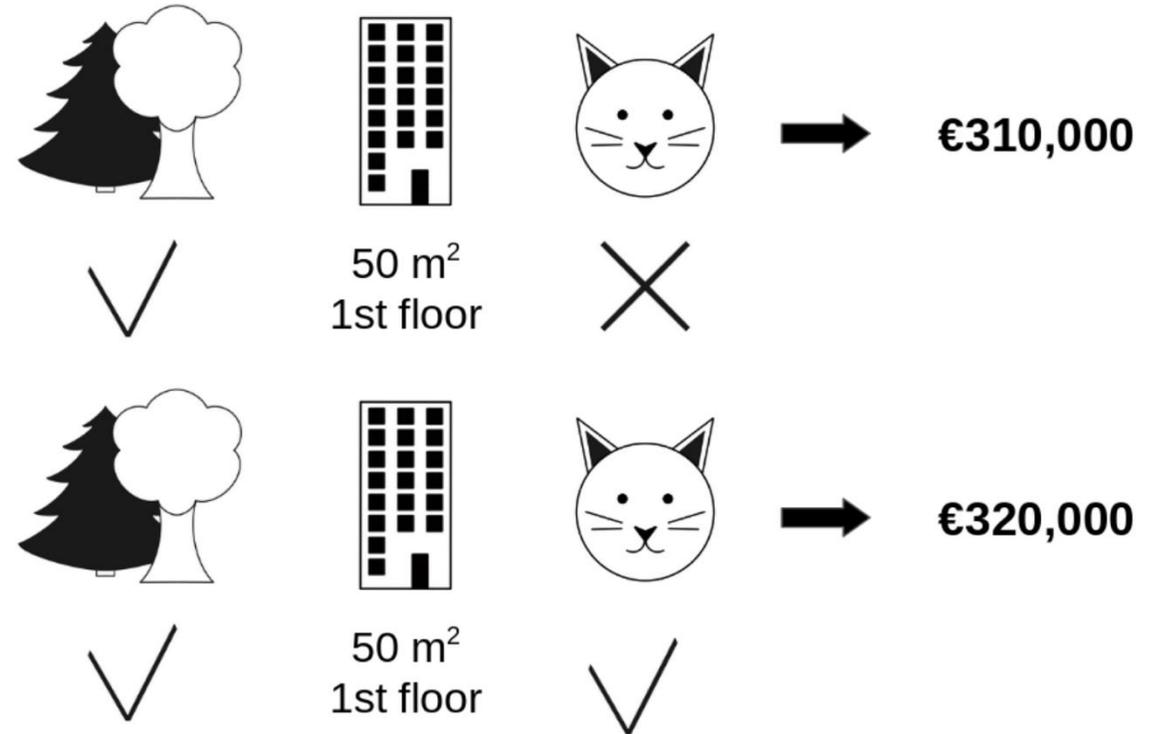
- The "game" is the prediction task for a single instance of the dataset.
- The "gain" is the **actual prediction** for this instance **minus** the **average prediction for all instances**.

The "players" are the feature values of the instance that collaborate to receive the gain (= predict a certain value).

The Shapley value is the average marginal contribution of a feature value **across all possible coalitions**.

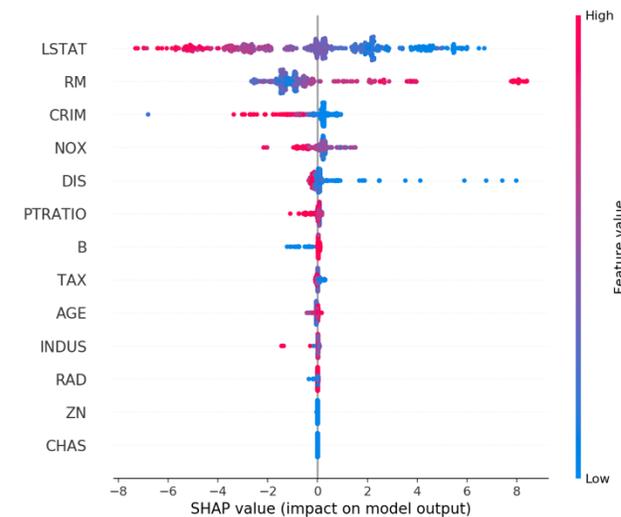
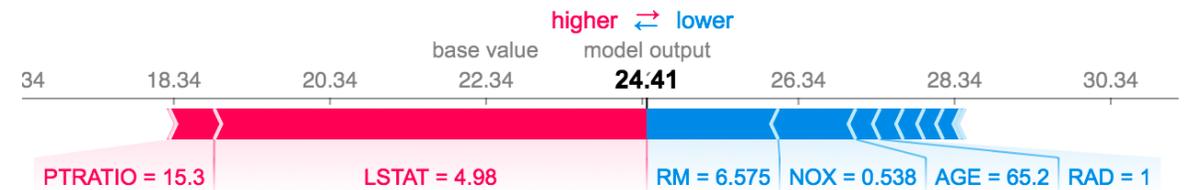
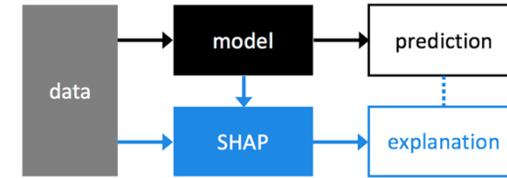
Shapely Values

One sample repetition to estimate the contribution of **cat-banned** to the prediction when added to the coalition of *park-nearby* and *area-50*.



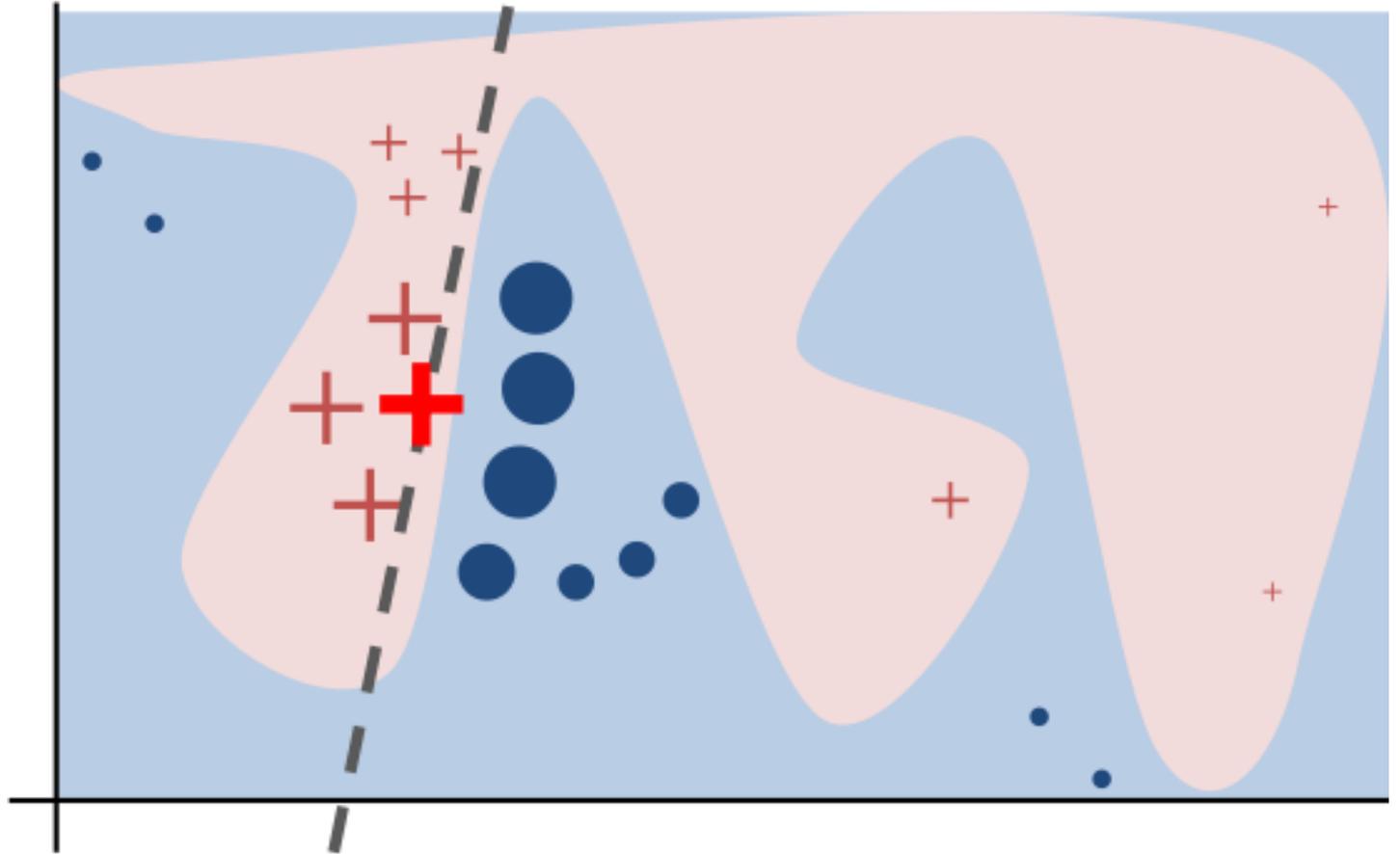
SHAP

- SHAP (SHapley Additive exPlanations) assigns each feature an importance value for a particular prediction by means of an additive feature attribution method.
- It assigns an importance value to each feature that represents the effect on the model prediction of including that feature



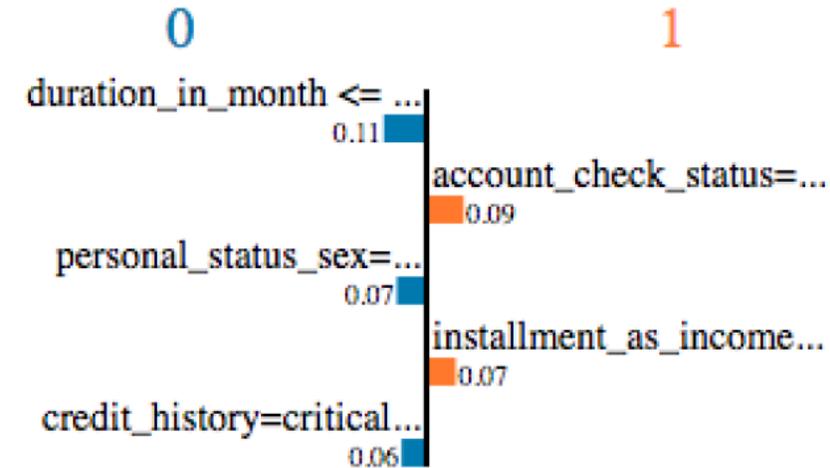
Local Explanation

- The overall decision boundary is complex
- In the neighborhood of a single decision, the boundary is simple
- A single decision can be explained by auditing the black box around the given instance and learning a *local* decision.



LIME

```
01  Z = {}
02  x instance to explain
03  x' = real2interpretable(x)
04  for i in {1, 2, ..., N}
05      zi = sample_around(x')
06      z = interpretabel2real(z')
07      Z = Z U {<zi, b(zi), d(x, z)>}
08  w = solve_Lasso(Z, k)
09  return w
```



LIME

- LIME *turns* an image x to a vector x' of interpretable superpixels expressing presence/absence.
- It *generates* a synthetic neighborhood Z by randomly perturbing x' and labels them with the black box.
- It *trains* a linear regression model (interpretable and locally faithful) and assigns a weight to each superpixel.



LIME – tab data

- LIME does not really generate images with different information: it randomly removes some superpixels, i.e. it suppresses the presence of an information rather than modifying it.
- On tabular data LIME generates the neighborhood by changing the feature values with other values of the domain.

$x = \{\text{age}=24, \text{sex}=\text{male}, \text{income}=1000\}$ ($x = x'$)

$z = \{\text{age}=30, \text{sex}=\text{male}, \text{income}=800\}$ ($z = z'$)

LORE – DR, AGN, TAB

01 x instance to explain

02 $Z_{=}$ = `geneticNeighborhood`(x , $fitness_{=}$, $N/2$)

03 Z_{\neq} = `geneticNeighborhood`(x , $fitness_{\neq}$, $N/2$)

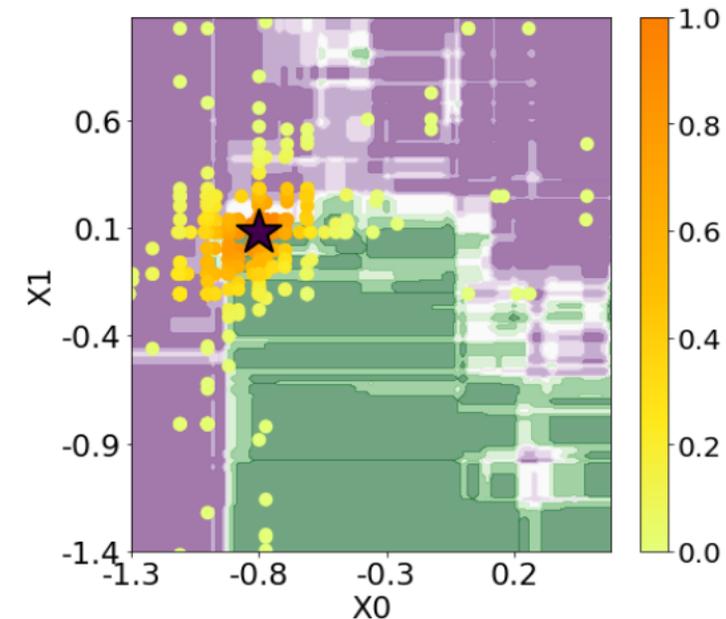
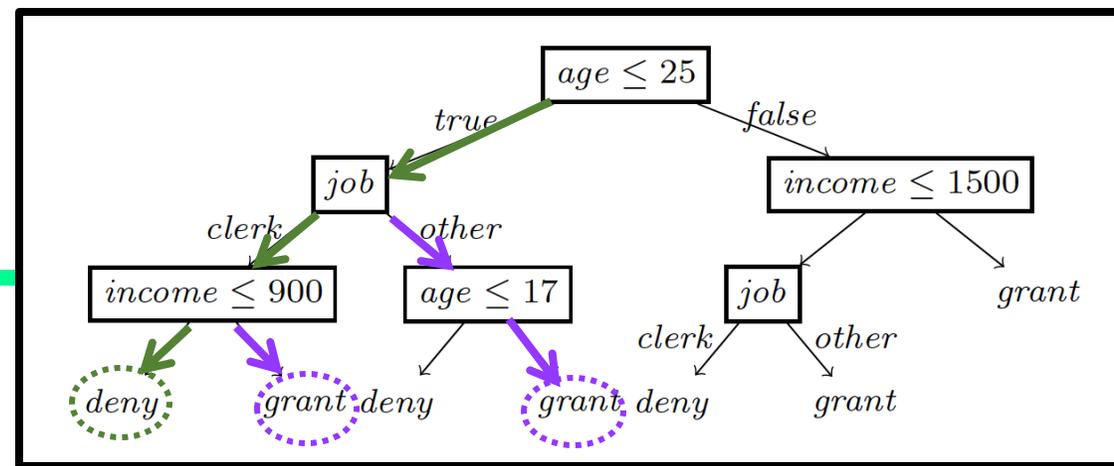
04 Z = $Z_{=}$ \cup Z_{\neq}

05 c = `buildTree`(Z , $b(Z)$) *black box auditing*

06 r = ($p \rightarrow y$) = `extractRule`(c , x)

07 ϕ = `extractCounterfactual`(c , r , x)

08 **return** $e = \langle r, \phi \rangle$



$r = \{age \leq 25, job = clerk, income \leq 900\} \rightarrow deny$

$\Phi = \{(\{income > 900\} \rightarrow grant),$
 $(\{17 \leq age < 25, job = other\} \rightarrow grant)\}$

Pedreschi, Franco Turini,
of black box decision

Adversarial Black box Explainer generating Latent Exemplars

- Explaining image classification
- Solving the drawback of LIME
- Exploit adversarial autoencoders
- Providing explanations based on exemplars and counter exemplars

References

- Guidotti, R., Monreale, A., Ruggieri, S., Turini, F., Giannotti, F., & Pedreschi, D. (2018). ***A survey of methods for explaining black box models***. *ACM Computing Surveys (CSUR)*, 51(5), 93
- Finale Doshi-Velez and Been Kim. 2017. ***Towards a rigorous science of interpretable machine learning***. arXiv:1702.08608v2
- Alex A. Freitas. 2014. ***Comprehensible classification models: A position paper***. ACM SIGKDD Explor. Newslett.
- Andrea Romei and Salvatore Ruggieri. 2014. ***A multidisciplinary survey on discrimination analysis***. Knowl. Eng.
- Yousra Abdul Alsaheb S. Aldeen, Mazleena Salleh, and Mohammad Abdur Razzaque. 2015. ***A comprehensive review on privacy preserving data mining***. SpringerPlus
- Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. 2016. ***Why should i trust you?: Explaining the predictions of any classifier***. KDD.
- Houtao Deng. 2014. ***Interpreting tree ensembles with intrees***. arXiv preprint arXiv:1408.5456.
- Mark Craven and JudeW. Shavlik. 1996. ***Extracting tree-structured representations of trained networks***. NIPS.

References

- M. Gethsiyal Augasta and T. Kathirvalavakumar. 2012. ***Reverse engineering the neural networks for rule extraction in classification problems***. NPL
- Riccardo Guidotti, Anna Monreale, Salvatore Ruggieri, Dino Pedreschi, Franco Turini, and Fosca Giannotti. 2018. ***Local rule-based explanations of black box decision systems***. arXiv preprint arXiv:1805.10820
- Ruth Fong and Andrea Vedaldi. 2017. ***Interpretable explanations of black boxes by meaningful perturbation***. arXiv:1704.03296 (2017).
- Paulo Cortez and Mark J. Embrechts. 2011. ***Opening black box data mining models using sensitivity analysis***. CIDM.
- Ruth Fong and Andrea Vedaldi. 2017. ***Interpretable explanations of black boxes by meaningful perturbation***. arXiv:1704.03296 (2017).
- Xiaoxin Yin and Jiawei Han. 2003. ***CPAR: Classification based on predictive association rules***. SIAM, 331–335
- Angelino, E., Larus-Stone, N., Alabi, D., Seltzer, M., & Rudin, C. 2017. ***Learning certifiably optimal rule lists***. KDD.