

Simple interpretability methods for black-box machine learning systems

Adriano Koshiyama

Table of Contents



- Key Components of Algorithmic Impact Assessment
- Importance of Algorithmic Explainability
- Types and Methods of Algorithmic Explainability

Key Components of Al Assessment



- What do we mean by Algorithmic Impact (AI) Assessment?
- Assessment vs by Design in Al Assessment

What do we mean by Al Assessment



- Algorithmic Impact Assessment focus on evaluating an Automated Decision-making system mainly from a Robustness, Fairness and Explainability point of view
- The goals of Al Assessment are
 - Set the boundary, usage and shelf-life of a system
 - Build trust between the stakeholders of a system
 - ❖ Be the entry point to hold the system's creators accountable of the results of its decision-making
- We should also mention other areas of Al Assessment, such as Transparency, Accountability, etc.

What do we mean by Al Assessment



In a nutshell

- Robustness: systems should be safe and secure, not vulnerable to tampering or compromising of the data they are trained on.
- Fairness: systems should use training data and models that are free of bias, to avoid unfair treatment of certain groups.
- Explainability: systems should provide decisions or suggestions that can be understood by their users and developers.

To avoid these cases

In the news



Microsoft deletes 'teen girl' Al after it became a Hitler-loving sex robot within 24 hours

Telegraph.co.uk - 5 hours ago

To chat with Tay, you can tweet or DM her by finding @tayandyou on Twitter, or add her as a ...

Microsoft Releases Al Twitter Bot That Immediately Learns How To Be Racist Kotaku - 3 hours ago

Microsoft Created a Twitter Bot to Learn From Users. It Quickly Became a Racist Jerk. New York Times - 3 hours ago



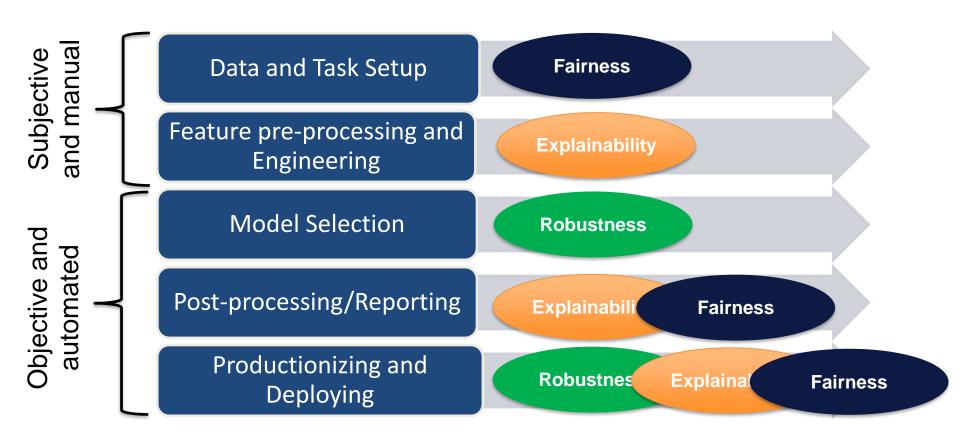


Assessment vs Modelling



From an Assessment point of view

Areas in the modelling pipeline where a AI Assessment Analyst (AI²) should analyse using the different criteria



Assessment vs Modelling



From a by Design point of view

It is possible to have systems that by design are able to increase or fulfil the stakeholders demand for Fairness, Robustness and Explainability

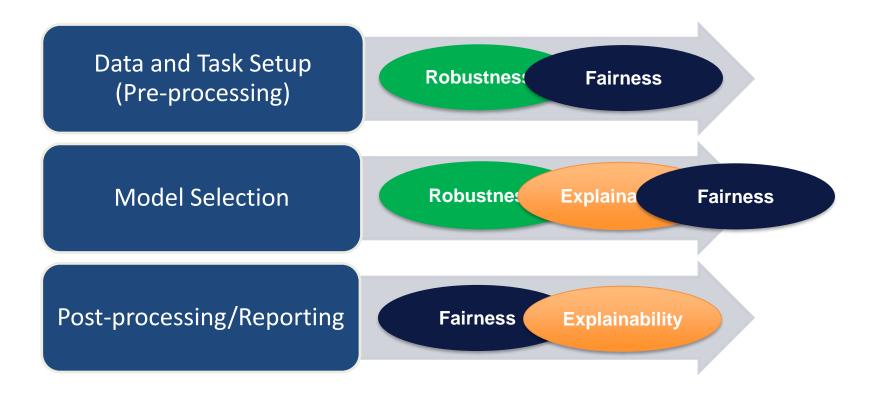


Table of Contents



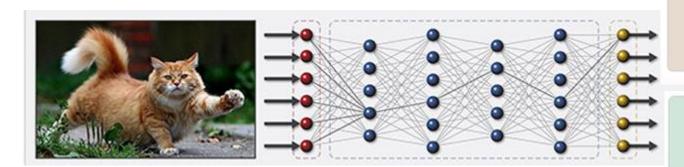
- Key Components of Algorithmic Impact Assessment
- Importance of Algorithmic Explainability
- Types and Methods of Algorithmic Explainability



- What do we mean by an explainable decision
- Why and what type of Explainability
- Legal basis for Explainability
- Different types of Explainability
- Technological solutions for Explainability
- Explainability: an Al Assessment checklist
- Further reading

What do we mean by an explainable decision UCL

Object recognition



This is a cat.

This is a cat:

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





Healthcare



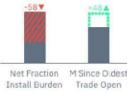
Finance

Soi If i

Sorry, your loan application has been rejected.

If instead you had the following values, your application would have been approved:

- MSinceOldestTradeOpen: 161
- · NumSatisfactoryTrades: 36
- NetFractionInstallBurden: 38
- NumRevolvingTradesWBalance: 4
- NumBank2NatlTradesWHighUtilization: 2



Input Value

(b) Counterfactual explanation

Why and what type of Explainability



- Explicability is crucial for building and maintaining users' and designers' trust in Al-based decisions
 - Users: contest decisions, learning
 - Creators: knowledge discovery, debugging systems, uncover unfair decisions
- Hence, the capabilities and purpose of AI systems should be
 - openly communicated
 - decisions explainable to those directly and indirectly affected
 - timely and adapted to the expertise of the stakeholder concerned (e.g. layperson, regulator or researcher)

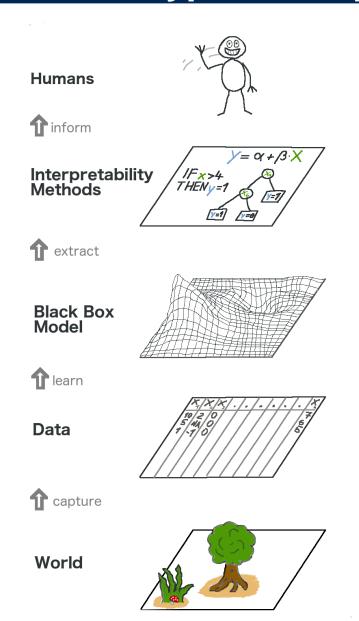
Legal basis for Explainability

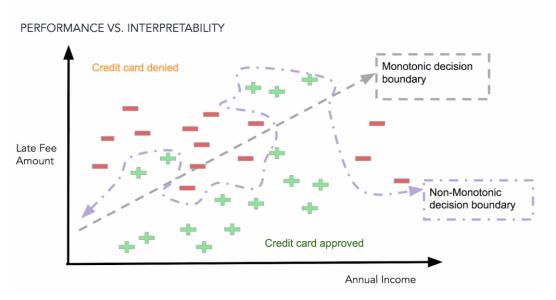


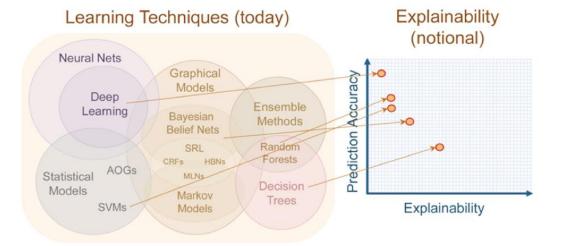
- Credit Scoring in the US have a well-established right to explanation
 - The Equal Credit Opportunity Act (1974)
- Credit agencies and data analysis firms such as FICO comply with this regulation by providing a list of reasons (generally at most 4, per interpretation of regulations)
- ❖ From an Al standpoint, there are new regulations that gives the system's user the right (?) to know why a certain automated decision was taken in a certain form
 - ❖ Right to an Explanation EU General Data Protection Regulation (2018)

Different types of Explainability









Different types of Explainability

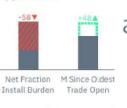




Sorry, your loan application has been rejected.

If instead you had the following values, your application would have been approved:

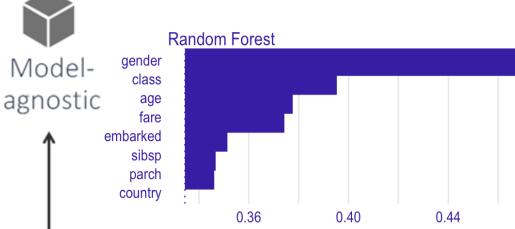
- MSinceOldestTradeOpen: 161
- · NumSatisfactoryTrades: 36
- NetFractionInstallBurden: 38 NumRevolvingTradesWBalance: 4
- NumBank2NatlTradesWHighUtilization: 2



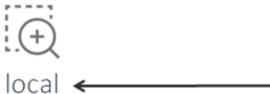


Model-specific









global

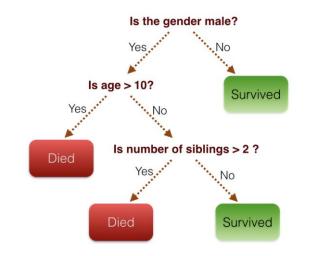
amazon.com **Recommended** for You Amazon.com has new recommendations for you based on items you purchased or told us you own.

Deciphered: Compute in Administrator Guide: A the Cloud to Streamline Your Desktop

Workspace



Ultimate Google Resource (3rd Edition)



Technical solutions for Explainability



global

Local Interpretable Model-Agnostic explanations (LIME)



- Shapley values (SHAP)
- Counterfactual explanations

- Partial Dependence
- **❖** Feature Importance

.()

local

- Linear model
- Decision tree
- Rule-based system

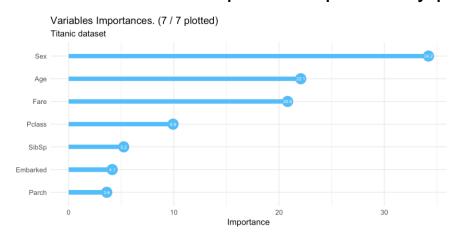
- Linear model
- Decision tree
- Rule-based system

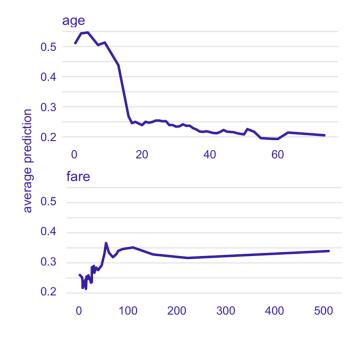
Model-specific

Technical solutions for Explainability



- Some examples of how it works
 - Feature and partial dependency plots





Local Interpretable Model-Agnostic explanations



(a) Original Image



(b) Explaining Electric guitar



(c) Explaining Acoustic guitar



(d) Explaining Labrador

Explainability: an Al assessment checklist 🚊 📗 🤇



Explainability

- ❖ Did you assess to what extent the decisions and hence the outcome made by the AI system can be understood?
- ❖ Did you ensure an explanation as to why the system took a certain choice resulting in a certain outcome that all users can understand?
- ❖Did you design the AI system with interpretability in mind from the start?
- ❖ Did you research and try to use the simplest and most interpretable model possible for the application in question?
- ❖ Did you assess whether you can examine interpretability after the model's training and development, or whether you have access to the internal workflow of the model?

Further reading



Legislation

- The Equal Credit Opportunity Act: https://www.justice.gov/crt/equal-credit-opportunity-act-3
- GDPR: https://ico.org.uk/for-organisations/guide-to-data-protection/guide-to-the-general-data-protection-regulation-gdpr/individual-rights/rights-related-to-automated-decision-making-including-profiling/

Papers and books

- Ribeiro, Marco Tulio, Sameer Singh, and Carlos Guestrin. "Model-agnostic interpretability of machine learning." arXiv preprint arXiv:1606.05386 (2016).
- Hall, Patrick. "On the Art and Science of Machine Learning Explanations." arXiv preprint arXiv:1810.02909 (2018).
- Hall, Patrick, and Navdeep Gill. Introduction to Machine Learning Interpretability. O'Reilly Media, Incorporated, 2018.
- Wachter, Sandra, Brent Mittelstadt, and Chris Russell. "Counterfactual explanations without opening the black box: Automated decisions and the GDPR." Harvard Journal of Law & Technology 31, no. 2 (2017): 2018.

Further reading

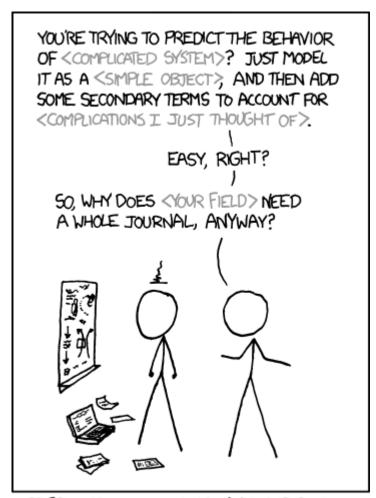


Tools

- https://pair-code.github.io/what-if-tool/
- https://github.com/marcotcr/lime
- https://github.com/microsoft/interpret
- https://github.com/slundberg/shap

Other good online resources

- https://christophm.github.io/interpretableml-book/
- https://distill.pub/2018/building-blocks/
- https://www.oreilly.com/ideas/ideas-on-interpreting-machine-learning



LIBERAL-ARTS MAJORS MAY BE ANNOYING SOMETIMES, BUT THERE'S NOTHING MORE OBNOXIOUS THAN A PHYSICIST FIRST ENCOUNTERING A NEW SUBJECT.

Table of Contents



- Key Components of Algorithmic Impact Assessment
- Algorithmic Explainability
- Types and Methods of Algorithmic Explainability

Algorithmic Explainability



- Model-specific explanations: Decision Trees
 - Global explanations: tree structure
 - Local explanations: decision rule
- Model-agnostic explanations: Neural Networks
 - Global explanations: feature importance and partial dependence plots
 - Local explanations: Local Interpretable Modelagnostic Explanations (LIME)

Model-specific: Global



Decision Tree: main Idea

Split the feature space into a set of rectangles, and inside each subspace fit a simple predictive model

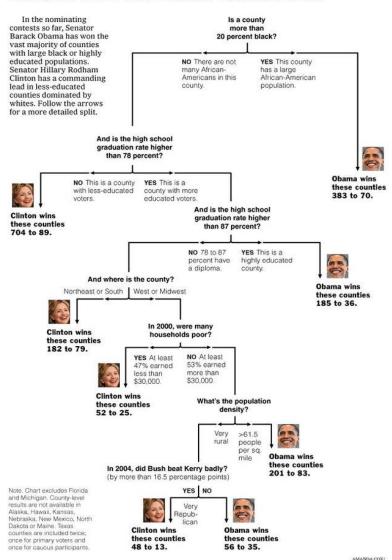
Some interesting characteristics

- Recursive algorithm
- Greedy construction
- Automatic feature selection
- A good way to start data exploration and analysis

Interpretability

- The deeper the tree, the harder is to analyse and reason
- The variables near the root are the most 'important'
- Aim for depth = 3 for a good balance of performance and explanation to a large audience

Decision Tree: The Obama-Clinton Divide

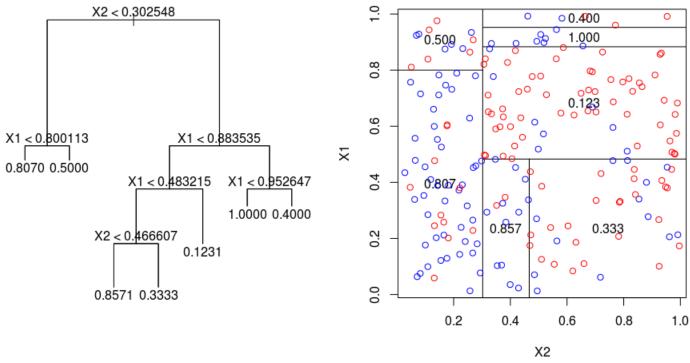


Model-specific: Local



Decision propagation or rule unwrapping

- ❖ Input and output space example, Regression Tree
- ❖ Two input variables, X1 and X2, and a single output variable $Y \in [0, 1]$



❖ If Subject 1 has X1 = 0.6 and X2 = 0.2, then ?

Model-agnostic: Global

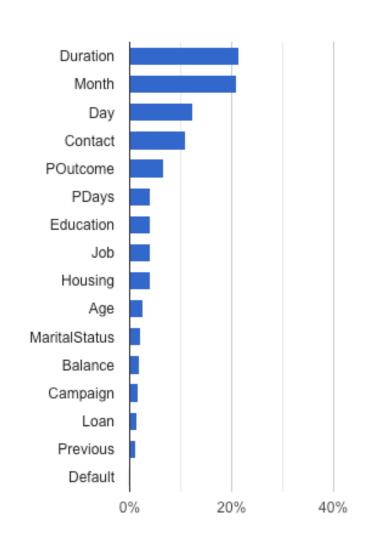


Feature Importance

- Often, the input predictor variables are seldom equally relevant. Often, only a few of them have substantial relevance
- Feature Importance metrics is a way in which "black-box" models become transparent
 - Weight of each feature
 - Positive/negative relationship

Some ways to calculate

- Weighted node impurity reduction
- Out-of-bag error degradation via Row permutation
- Caveat: will underestimate "importance" in the presence of correlated variables



Model-agnostic: Global



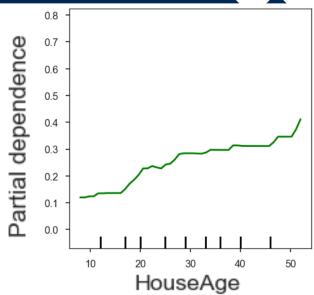
Partial Dependence

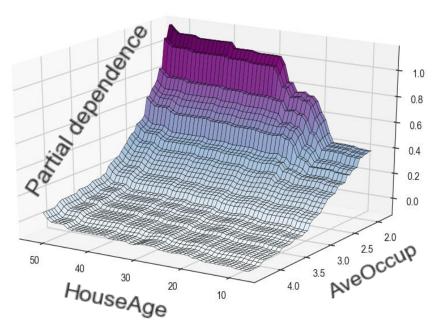
- Partial dependence functions can be used to interpret the results of any "black box" learning method.
- Given a feature set S (feature A, say) and its complement C, partial dependence can be estimated by:

$$\hat{y}_S(\mathbf{x}_S) = \frac{1}{n} \sum_{i=1}^n F(\mathbf{x}_S, \mathbf{x}_{iC})$$

for every value x_S in the range

- For a given model, it involves a heavy amount of computation
- However, for decision trees it is rapidly computed from the tree itself without reference to the data



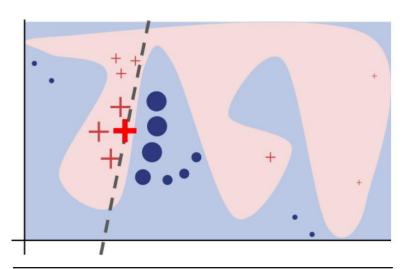


Model-agnostic: Local



LIME

- Local Interpretable Model-Agnostic Explanations (LIME)
 - Important research (2016), influencing many methods that came afterwards.
 - Provided some foundational thought, particularly about fidelity-interpretability trade-off
- What is to be Local for LIME
 - Giving more importance to samples at a vicinity of a data point x that we aim to provide a local explanation
- What is to be Interpretable for LIME
 - Belonging to a class of potentially interpretable models like decision trees, linear models, etc.
- What is to be Model-agnostic for LIME
 - Without knowledge of the underlying machine learning model, that is, if it is a Neural Network, SVM, etc.



Algorithm 1 Sparse Linear Explanations using LIME

Require: Classifier f, Number of samples N

Require: Instance x, and its interpretable version x'

Require: Similarity kernel π_x , Length of explanation K

$$\mathcal{Z} \leftarrow \{\}$$
 for $i \in \{1, 2, 3, ..., N\}$ do

$$z_i' \leftarrow sample_around(x')$$

$$\mathcal{Z} \leftarrow \mathcal{Z} \cup \langle z_i', f(z_i), \pi_x(z_i) \rangle$$

end for

 $w \leftarrow \text{K-Lasso}(\mathcal{Z}, K) \triangleright \text{with } z_i' \text{ as features, } f(z) \text{ as target }$ return w

