# Merci à nos sponsors

Et en partenariat avec le meetup Azure Nantes

clever cloud

# XAI : eXplainable Artificial Intelligence

clever cloud

**Victor Ballu**

Global AI Nights Nantes

# Explainability and AI – Definitions

*Artificial Intelligence :*

*"[Artificial Intelligence] refers to a programme whose ambitious objective is to understand and reproduce human cognition; creating cognitive processes comparable to those found in human beings"* *for a meaningful artificial intelligence towards a french and european strategy,* Villani report, 2018
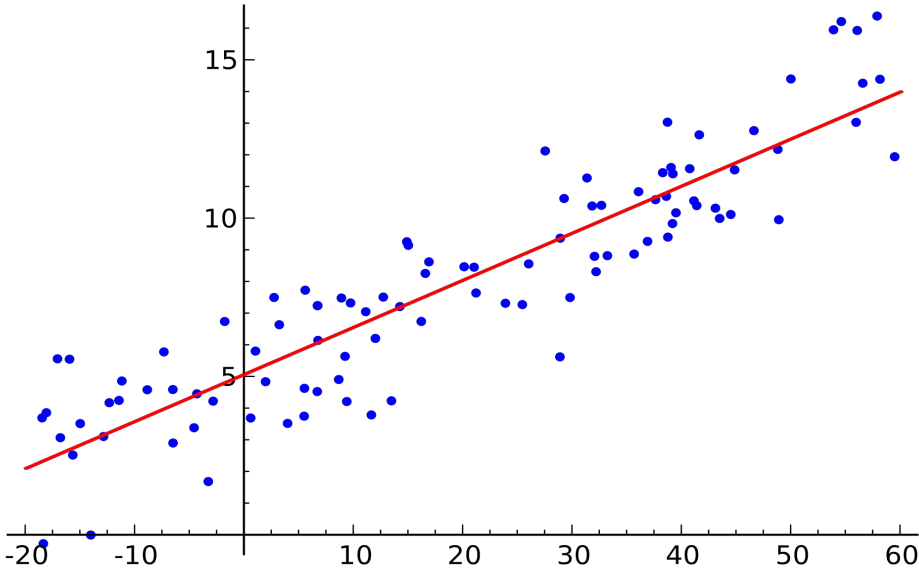
*Explainability :*

*Make the algorithms inner state understable by humans*
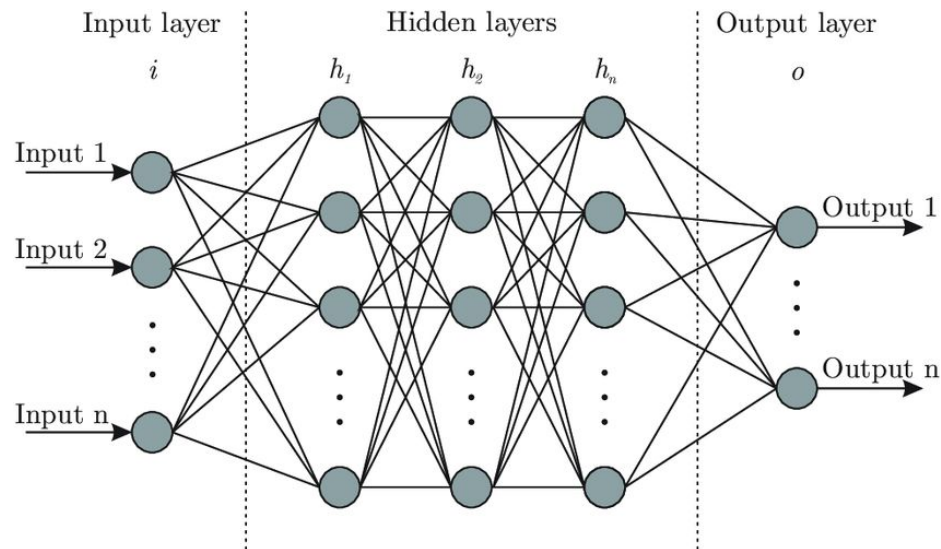
## *Linear regression*

Two correlated axes



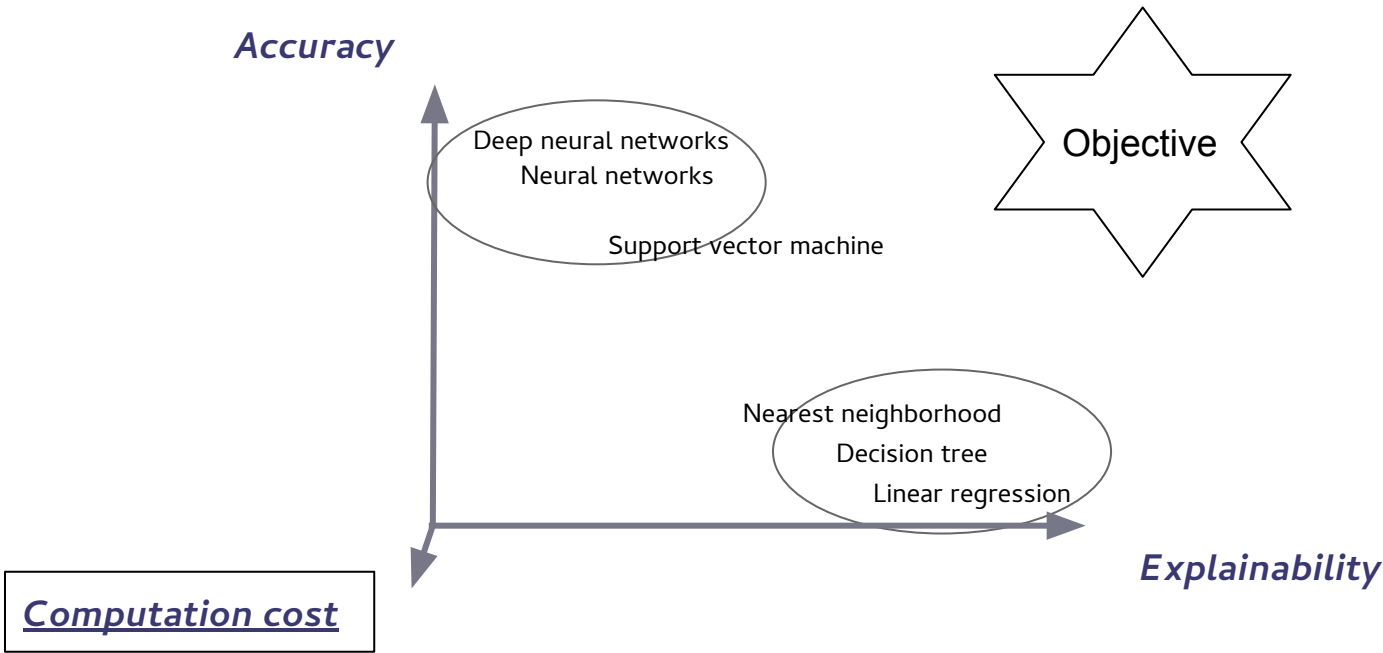*Wikipedia: Linear regression article*
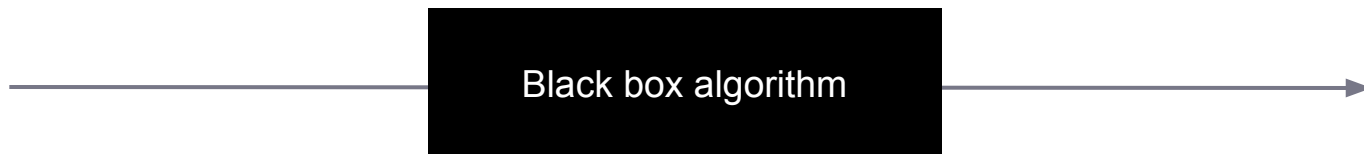
## *Artificial neural network*

### Many fields of low value, very entangled



*Bre, Facundo & Gimenez, Juan & D. Fachinotti, Víctor. (2017). Prediction of wind pressure coefficients on building surfaces using Artificial Neural Networks. Energy and Buildings.*

clever cloud

Victor Ballu

Global AI Nights Nantes

**Accuracy**

Deep neural networks

Neural networks

Support vector machine

Objective

Nearest neighborhood

Decision tree

Linear regression

**Explainability**

**Computation cost**

# Why is understanding the black box important ?

Black box algorithm

- General Data Protection Regulation (GDPR)

*"the data subject should have the right [...] to obtain an explanation of the decision reached" – GDPR, Recital 71*
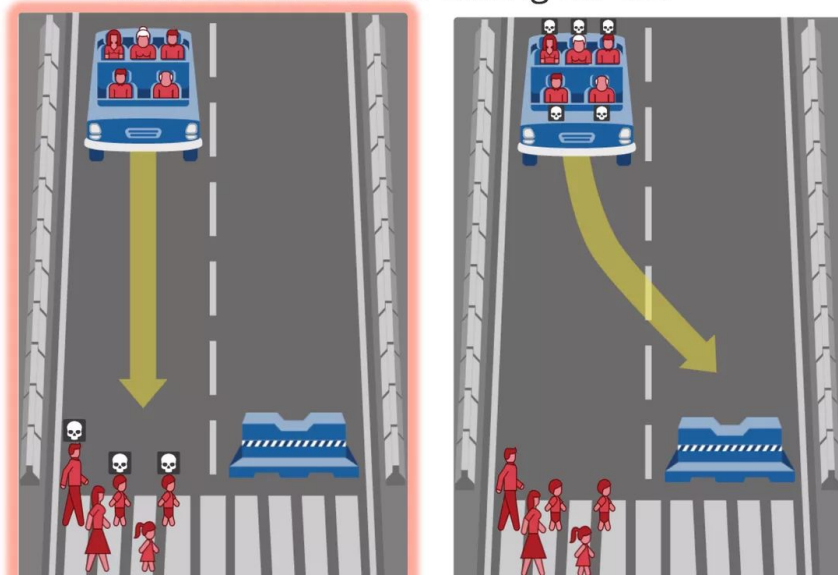
- GDPR

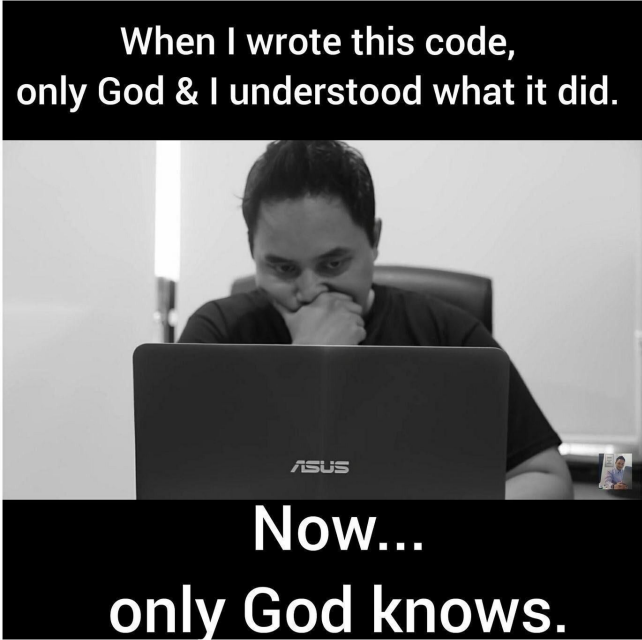## What should the self-driving car do?

- Ethics

*Awad, Edmond & Dsouza, Sohan & Kim, Richard & Schulz, Jonathan & Henrich, Joseph & Shariff, Azim & Bonnefon, Jean-François & Rahwan, Iyad. (2018). The Moral Machine Experiment. Nature. http://moralmachine.mit.edu/*

- GDPR    - Ethics

- Users confidence



When I wrote this code,
only God & I understood what it did.

Now...
only God knows.

# Why is understanding the black box important ?

- GDPR
- Ethics
- Users confidence

- Colleagues confidence

- GDPR
- Ethics
- Users confidence
- Colleagues confidence

- Self confidence

- GDPR
- Ethics
- Users confidence
- Colleagues confidence
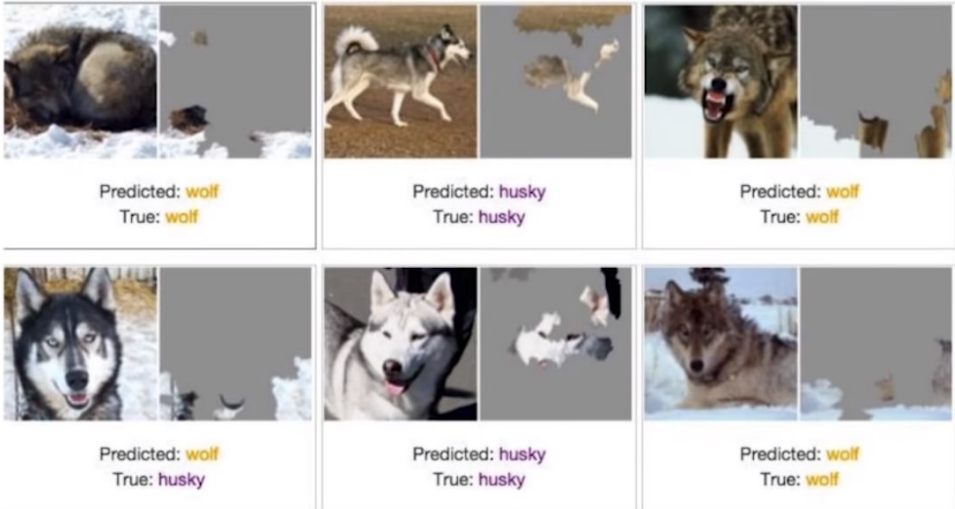- Self confidence

# CONFIDENCE

- GDPR
- Ethics
- Users confidence
- Colleagues confidence
- Self confidence

# CONFIDENCE

## By cross validation !!!

# Why is understanding the black box important ?

Predicted: wolf
True: wolf

Predicted: husky
True: husky

Predicted: wolf
True: wolf

Predicted: wolf
True: husky

Predicted: husky
True: husky

Predicted: wolf
True: wolf

*Tulio Ribeiro, Marco & Singh, Sameer & Guestrin, Carlos. (2016). "Why Should I Trust You?": Explaining the Predictions of Any Classifier*

clever cloud    **Victor Ballu**
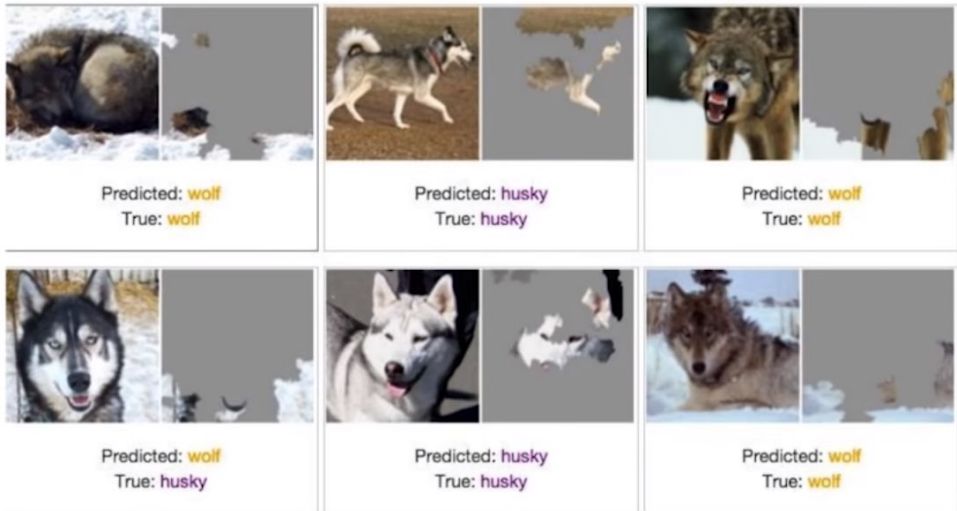
**Global AI Nights Nantes**

# Why is understanding the black box important ?

*Tulio Ribeiro, Marco & Singh, Sameer & Guestrin, Carlos. (2016). "Why Should I Trust You?": Explaining the Predictions of Any Classifier*



(a) Husky classified as wolf          (b) Explanation

clever cloud          **Victor Ballu**          **Global AI Nights Nantes**

- GDPR
- Ethics
- Users confidence
- Colleagues confidence
- Self confidence

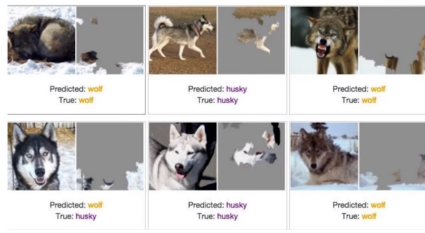- **Improve algoritms**
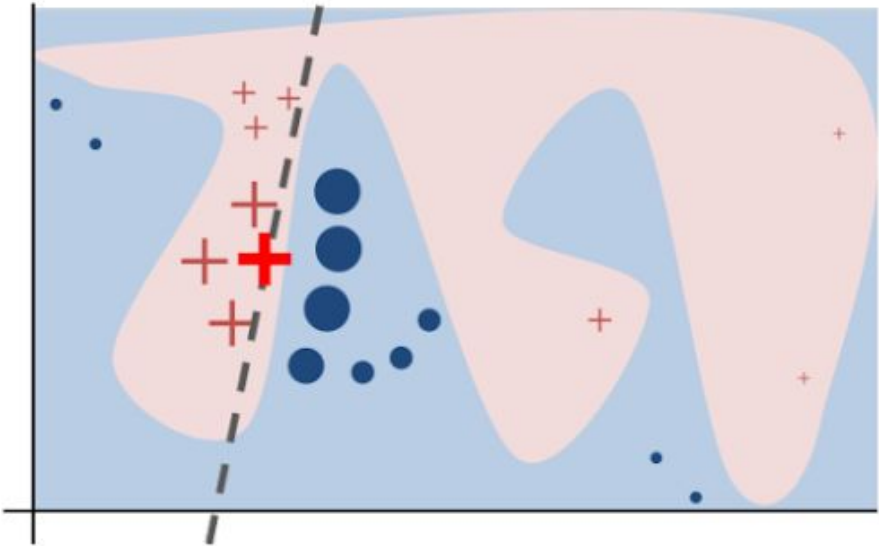
# Current approaches

- *Local interpretability*

- *Global interpretability*

- ## *Local interpretability:*

    *Explain one specific result*

    ➢ *LIME – Local Interpretable Model–Agnostic Explanations*

➤ *LIME – Local Interpretable Model-Agnostic Explanations*



(a) Original Image    (b) Explaining *Electric guitar*    (c) Explaining *Acoustic guitar*    (d) Explaining *Labrador*
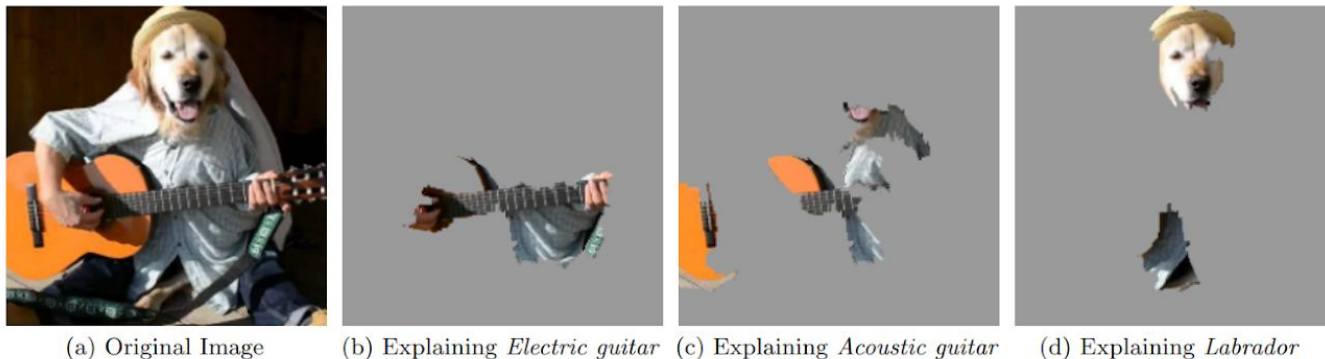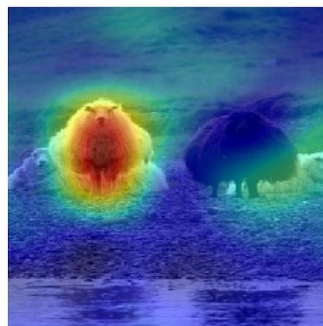
Figure 4: Explaining an image classification prediction made by Google's Inception neural network. The top 3 classes predicted are "Electric Guitar" ($p = 0.32$), "Acoustic guitar" ($p = 0.24$) and "Labrador" ($p = 0.21$)

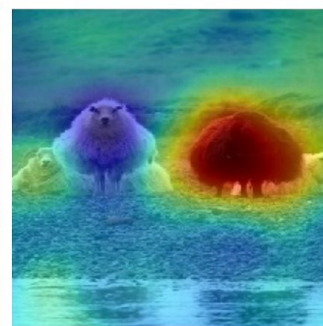➢ *RISE: Randomized Input Sampling for Explanation of Black-box Models[1]*



(a) Sheep - 26%, Cow - 17%     (b) Importance map of 'sheep'     (c) Importance map of 'cow'

(d) Bird - 100%, Person - 39%     (e) Importance map of 'bird'     (f) Importance map of 'person'

(1) Petsiuk, Vitali & Das, Abir & Saenko, Kate. (2018). RISE: Randomized Input Sampling for Explanation of Black-box Models. Boston University

- *Global interpretability:*
  *Explain the whole model*

  ➢ *LIME suggest a kind of integration of the local interpretability*

  ➢ *SHAP – A global model approximation based on linear combination of local approximations from different models*[1]

- *Hybrid approach*

- *Others methods*

  ➢ *Autoencoder*

**Promising approaches to explainability**

| CP | Performer | Explainable Model |
|---|---|---|
| Both | UC Berkeley | Deep Learning |
| | Charles River | Causal Modeling |
| | UCLA | Stochastic And-Or-Graphs |
| Autonomy | Oregon State | Deep Adaptive Programs |
| | PARC | Cognitive Modeling |
| | CMU | Explainable RL (XRL) |
| Analytics | SRI International | Deep Learning |
| | Raytheon BBN | Deep Learning |
| | UT Dallas | Probabilistic Logic |
| | Texas A&M | Mimic Learning |
| | Rutgers | Explanation by Example |

source : DARPA AI COLLOQUIUM

DARPA AI COLLOQUIUM

clever cloud   Victor Ballu

# Bibliography

# ARTICLES :

- Tulio Ribeiro, Marco & Singh, Sameer & Guestrin, Carlos. (2016). **"Why Should I Trust You?": Explaining the Predictions of Any Classifier**
- Petsiuk, Vitali & Das, Abir & Saenko, Kate. (2018). **RISE: Randomized Input Sampling for Explanation of Black-box Models**
- **GDPR - Recital 71**
- **For a meaningful artificial intelligence towards a french and european strategy,** Villani report, 2018
- Awad, Edmond & Dsouza, Sohan & Kim, Richard & Schulz, Jonathan & Henrich, Joseph & Shariff, Azim & Bonnefon, Jean-François & Rahwan, Iyad. (2018). **The Moral Machine Experiment**. Nature
- Lundberg, Scott & Lee, Su-In. (2017). **A Unified Approach to Interpreting Model Predictions**
- Zhou, Bolei & Sun, Yiyou & Bau, David & Torralba, Antonio. (2018). **Interpretable Basis Decomposition for Visual Explanation**

# Other resources :

- **DARPA** - XAI - Literature Review
- **Fairness, Accountability, and Transparency in Machine Learning** (FAT) - https://www.fatml.org/

**Clever Cloud Paris**
137 rue vieille du temple 75003 Paris


**Clever Cloud Nantes**
3 rue de l'allier 44000 Nantes
02 85 52 07 69


https://www.clever-cloud.com

# CONTACT

victor.ballu@clever-cloud.com