# FYS: AI in Healthcare
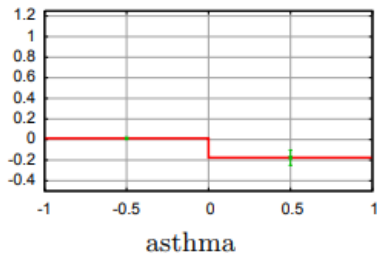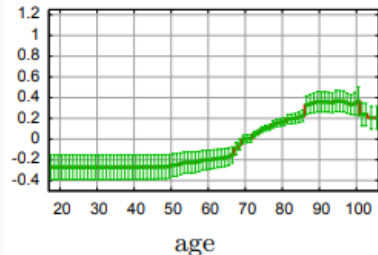
Interpretability in Machine Learning

John Lalor

October 23, 2018
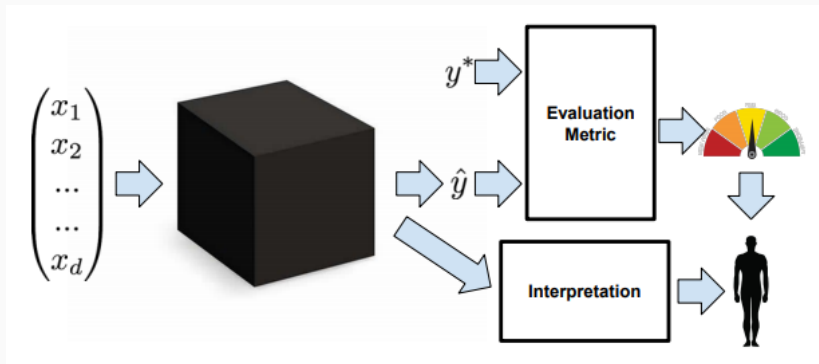
Midterm questions?

Caruana, Rich, et al. "Intelligible models for healthcare: Predicting pneumonia risk and hospital 30-day readmission." Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining. ACM, 2015.

# Evaluation-Interpretability Relationship



Lipton, Zachary C. "The Mythos of Model Interpretability." Queue 16.3 (2018): 30.

Trust

## Causality



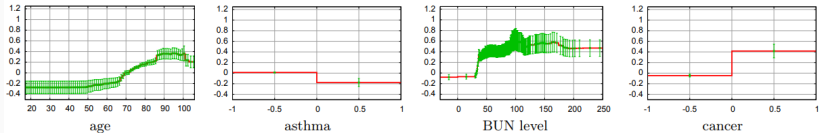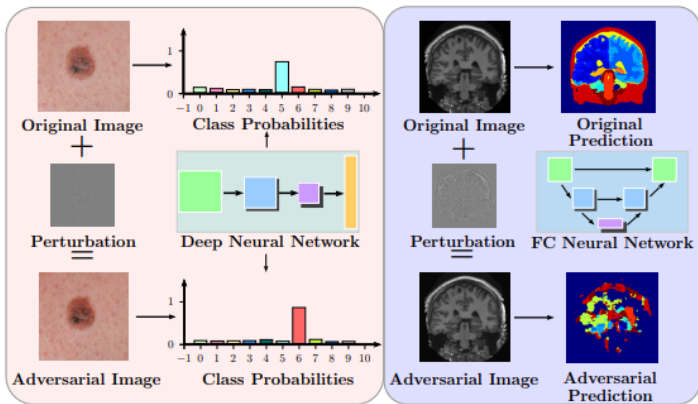Caruana, Rich, et al. "Intelligible models for healthcare: Predicting pneumonia risk and hospital 30-day readmission." Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining. ACM, 2015.
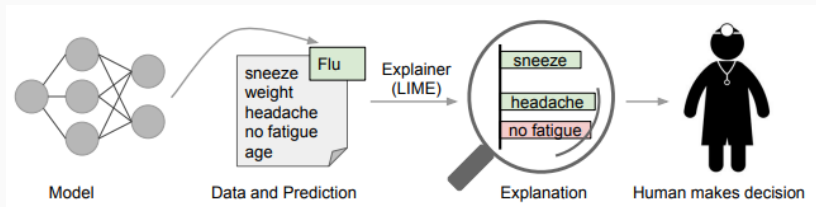
## Transferability



Paschali, Magdalini, et al. "Generalizability vs. Robustness: Adversarial Examples for Medical Imaging." arXiv:1804.00504 (2018).

Informativeness



Ribeiro, Marco Tulio, Sameer Singh, and Carlos Guestrin. "Why should i trust you?: Explaining the predictions of any classifier." Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining. ACM, 2016.

Fairness (e.g. *the right to an explanation*)

$X \longrightarrow$ Black box $\longrightarrow \hat{Y}$

## Transparency



Model transparency (simulatability)

Parameter transparency (decomposability)

Training transparency (algorithmic transparency)

Given the input data and model parameters:

Could a human produce an output?
In some reasonable amount of time?

## Model transparency

Given the input data and model parameters:

Could a human produce an output?
    In some reasonable amount of time?
Example: BOW logistic regression

## Model transparency

Given the input data and model parameters:

Could a human produce an output?
    In some reasonable amount of time?

Example: BOW logistic regression

Example: fully-connected DNN with hidden layer size 10

Each part of the model is intuitive
    Inputs
    Parameters
    Calculations

Each part of the model is intuitive

Inputs

Parameters

Calculations

Ex.: Descriptive decision tree nodes

## Parameter transparency

Each part of the model is intuitive

    Inputs

    Parameters

    Calculations

Ex.: Descriptive decision tree nodes

Ex.: Linear model parameters

Each part of the model is intuitive

Inputs

Parameters

Calculations

Ex.: Descriptive decision tree nodes

Ex.: Linear model parameters

Caveat: Can be fragile depending on pre-processing

## Algorithmic transparency

Insight into the decision-making process

## Algorithmic transparency

Insight into the decision-making process

    Linear models?

## Algorithmic transparency

Insight into the decision-making process

Linear models?

DNNs?

## Algorithmic transparency

Insight into the decision-making process

Linear models?

DNNs?

Humans?

Human interpretability
After the fact interpretation
    Not part of model training

# Text

Figure 3: Examples of extracted rationales indicating the sentiments of various aspects. The extracted texts for appearance, smell and palate are shown in red, blue and green color respectively. The last example is shortened for space.

Lei, Tao, Regina Barzilay, and Tommi Jaakkola. "Rationalizing Neural Predictions." EMNLP 2016.

# Visualizations



Maaten, Laurens van der, and Geoffrey Hinton. "Visualizing data using t-SNE." Journal of machine learning research 9.Nov (2008): 2579-2605.
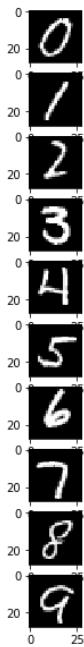
# Local explanations

Simonyan, Karen, Andrea Vedaldi, and Andrew Zisserman. "Deep inside convolutional networks: Visualising image classification models and saliency maps." arXiv:1312.6034 (2013).
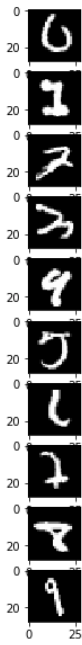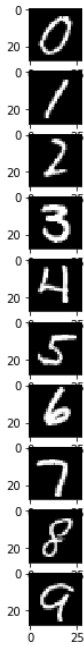
# Explanations by example

| Newspapers | | | |
|---|---|---|---|
| New York | New York Times | Baltimore | Baltimore Sun |
| San Jose | San Jose Mercury News | Cincinnati | Cincinnati Enquirer |
| NHL Teams | | | |
| Boston | Boston Bruins | Montreal | Montreal Canadiens |
| Phoenix | Phoenix Coyotes | Nashville | Nashville Predators |
| NBA Teams | | | |
| Detroit | Detroit Pistons | Toronto | Toronto Raptors |
| Oakland | Golden State Warriors | Memphis | Memphis Grizzlies |
| Airlines | | | |
| Austria | Austrian Airlines | Spain | Spainair |
| Belgium | Brussels Airlines | Greece | Aegean Airlines |
| Company executives | | | |
| Steve Ballmer | Microsoft | Larry Page | Google |
| Samuel J. Palmisano | IBM | Werner Vogels | Amazon |

Mikolov, Tomas, et al. "Distributed representations of words and phrases and their compositionality." Advances in neural information processing systems. 2013.
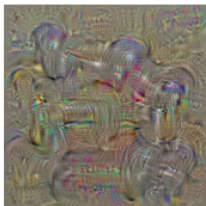
# Explanations by example

# Explanations by example

# Interpretability Examples
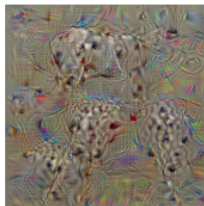
# Saliency maps



$$\arg \max_I S_c(I) - \lambda ||I||_2^2$$

Simonyan, Karen, Andrea Vedaldi, and Andrew Zisserman. "Deep inside convolutional networks: Visualising image classification models and saliency maps." arXiv:1312.6034 (2013).
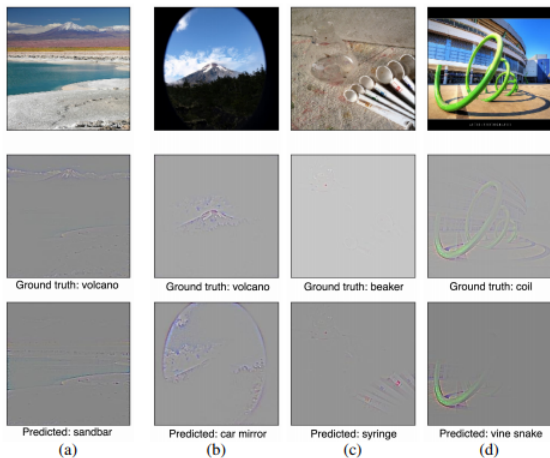
# Saliency maps

Simonyan, Karen, Andrea Vedaldi, and Andrew Zisserman. "Deep inside convolutional networks: Visualising image classification models and saliency maps." arXiv:1312.6034 (2013).

Analyzing failures



Selvaraju, Ramprasaath R., et al. "Grad-CAM: Visual Explanations from Deep Networks via Gradient-Based Localization." ICCV. 2017.

Handling adversarial noise



(a) Original image — Boxer: 0.40 Tiger Cat: 0.18
(b) Adversarial image — Airliner: 0.9999
(c) Grad-CAM "Dog" — Boxer: 1.1e-20
(d) Grad-CAM "Cat" — Tiger Cat: 6.5e-17

Selvaraju, Ramprasaath R., et al. "Grad-CAM: Visual Explanations from Deep Networks via Gradient-Based Localization." ICCV. 2017.

Counterfactuals



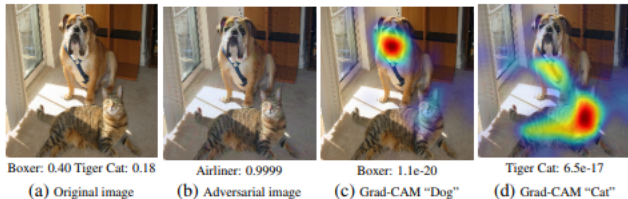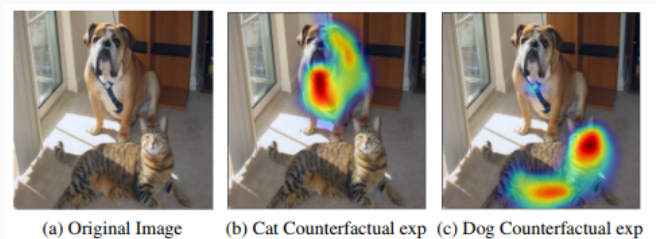(a) Original Image    (b) Cat Counterfactual exp    (c) Dog Counterfactual exp

Selvaraju, Ramprasaath R., et al. "Grad-CAM: Visual Explanations from Deep Networks via Gradient-Based Localization." ICCV. 2017.

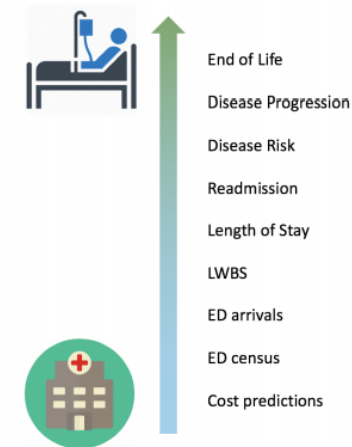Fig. 1: Prediction Use Cases vs. Need for Interpretability (LWBS: left without being seen)

# Activity: Interpretability

# Is it interpretable?

- Food

- Food

  Grocery store

## Is it interpretable?

- Food
    - Grocery store
    - Restaurant

## Is it interpretable?

- Food
    - Grocery store
    - Restaurant
- Travel

## Is it interpretable?

- Food
    - Grocery store
    - Restaurant
- Travel
    - Airlines

## Is it interpretable?

- Food
    - Grocery store
    - Restaurant
- Travel
    - Airlines
    - Google maps

## Is it interpretable?

- Food
    - Grocery store
    - Restaurant
- Travel
    - Airlines
    - Google maps

- Learning

## Is it interpretable?

- Food
  - Grocery store
  - Restaurant
- Travel
  - Airlines
  - Google maps

- Learning
  - In the classroom

## Is it interpretable?

- Food
  - Grocery store
  - Restaurant
- Travel
  - Airlines
  - Google maps

- Learning
  - In the classroom
  - Learning by doing

## Is it interpretable?

- Food
    - Grocery store
    - Restaurant
- Travel
    - Airlines
    - Google maps

- Learning
    - In the classroom
    - Learning by doing
- Law

## Is it interpretable?

- Food
  - Grocery store
  - Restaurant
- Travel
  - Airlines
  - Google maps

- Learning
  - In the classroom
  - Learning by doing
- Law
- Taxes