# Explaining Al Why, why not and how

Inga Strümke

# The engine: An explanation



Courtesy of Teknisk Ukeblad

## The engine: An explanation

(relax)



Courtesy of Teknisk Ukeblad

#### Today's menu

#### → Why?

3 perspectives and needs for explanations.

#### → Why not?

Mathematics, we have a problem.

#### → How?

Many directions, methods and open questions







AI = Artificial intelligence.

XAI = eXplainable AI

In the present discussion, AI means machine learning

Machine learning: A program that learns from data (i.e. from experience)



#### Machine learning: A program that learns from data (i.e. from experience)



#### "Artificial intuition"







#### Part 1: Legal requirements

In the EU (and parts of the US) customers are protected by regulations

GDPR:

"meaningful information about the logic of processing"

The end user doesn't care about your model; they care about how they can affect the outcome

 $\Rightarrow$  A convincing story about how the relevant part of the world works







# Note: The legal status isn't clear yet

GDPR:

"an explanation of the decision reached after [algorithmic] assessment"

However, it's not specified what such an explanation entails

#### THE RIGHT TO EXPLANATION, EXPLAINED

Margot E. Kaminski<sup>†</sup>

DOI: https://doi.org/10.15779/Z38TD9N83H © 2019 Margot E. Kaminski. Associate Professor of Law, University of Colorado Law School.

Problem? Yes. But: New EU proposal for AI regulation released April 2021.







#### Part 1: Legal requirements

Stakeholder / business leader must understand in order to evaluate risk and defend decisions. Risk is a central aspect in business decisions.

New EU proposal for AI regulation also has a risk based approach.

Question: Is the system trustworthy?

The traditional approach is testing...

![](_page_10_Picture_5.jpeg)

#### 84%

Of respondents think that AI based decisions need to be explainable in order to be trusted (PwC CEO Survey 2019)

![](_page_10_Picture_8.jpeg)

![](_page_10_Picture_9.jpeg)

![](_page_10_Picture_10.jpeg)

#### Regarding testing:

![](_page_11_Picture_1.jpeg)

Robust Physical-World Attacks on Deep Learning Visual Classification, Kevin Eykholt et al

(How can this happen?? Stay tuned)

#### Part 2: Researchers and developers

Wants to know

→ Which part of the world has my model understood?

![](_page_12_Picture_3.jpeg)

![](_page_12_Picture_4.jpeg)

![](_page_12_Picture_5.jpeg)

#### What, 'which part'? My model is perfect.

"any two optimization algorithms are equivalent when their performance is averaged across all possible problems"

-No free lunch theorems, Wolpert and Macready (2005)

![](_page_13_Picture_3.jpeg)

Your model is never perfect. But it can have a high accuracy where you tested it.

#### The curse of dimensionality

![](_page_14_Picture_1.jpeg)

#### The curse of dimensionality

You want to solve some problem using machine learning, and collect a data set...

Say... feature values in the range [0, 10], one data point per integer

Two features (dimensions):  $10^2 = 100$  data points

Three features:  $10^3 = 1000$  data points

13 features:  $10^{13}$  = ten thousand billion data points. Good luck.

![](_page_15_Figure_6.jpeg)

#### What to do?

Explainable AI (XAI) is a fairly young and active field of research. Nobody knows exactly what to do. But:

- 1. Accept that human intuition won't cut it. We can't understand complex models by looking at them.
- 2. Methods. New methods are developed all the time.

(et's have a look at methods

![](_page_16_Picture_5.jpeg)

#### Map of Explainability Approaches

![](_page_17_Figure_1.jpeg)

#### A mystery box fell from the sky

![](_page_18_Picture_1.jpeg)

#### A mystery box fell from the sky

![](_page_19_Picture_2.jpeg)

 $\rightarrow$  Predictions

#### A mystery box fell from the sky

![](_page_20_Picture_1.jpeg)

![](_page_20_Picture_2.jpeg)

#### What do do with the mystery box?

![](_page_21_Picture_1.jpeg)

![](_page_21_Picture_2.jpeg)

→ Poke from outside
→ Look inside

#### Poke the box like a boss

![](_page_22_Picture_1.jpeg)

![](_page_22_Picture_2.jpeg)

The Shapley decomposition: 
$$\phi_i(v) = \sum_{S \subseteq N \setminus \{i\}} \frac{|S|!(N-|S|-1)!}{N!} (v(S \cup \{i\}) - v(S))$$

# Shapley value example: Cab sharing

![](_page_23_Picture_1.jpeg)

 $\begin{array}{l} v(\{\})=0 \quad (\text{no passengers costs nothing}) \\ v(\{1\})=3, \quad v(\{2\})=7, \quad v(\{3\})=10 \\ v(\{1,2\})=7, \quad v(\{1,3\})=10, \quad v(\{2,3\})=10 \\ v(\{1,2,3\})=10 \end{array}$ 

![](_page_24_Picture_2.jpeg)

Characteristic function values -

 $\begin{array}{l} v(\{\})=0 \quad (\text{no passengers costs nothing}) \\ v(\{1\})=3, \quad v(\{2\})=7, \quad v(\{3\})=10 \\ v(\{1,2\})=7, \quad v(\{1,3\})=10, \quad v(\{2,3\})=10 \\ v(\{1,2,3\})=10 \end{array}$ 

![](_page_25_Figure_2.jpeg)

$$\phi_i(v) = \sum_{S \subseteq N \setminus \{i\}} \frac{|S|!(N - |S| - 1)!}{N!} \left( v(S \cup \{i\}) - v(S) \right), \quad i = 1, \dots, N$$

 $\begin{array}{l} v(\{\})=0 \quad (\text{no passengers costs nothing}) \\ v(\{1\})=3, \quad v(\{2\})=7, \quad v(\{3\})=10 \\ v(\{1,2\})=7, \quad v(\{1,3\})=10, \quad v(\{2,3\})=10 \\ v(\{1,2,3\})=10 \end{array}$ 

$$\frac{1}{2} \frac{2}{3}$$

$$\phi_i(v) = \sum_{S \subseteq N \setminus \{i\}} \frac{|S|!(N - |S| - 1)!}{N!} \left( v(S \cup \{i\}) - v(S) \right), \quad i = 1, \dots, N$$

N = total passengers

 $\begin{array}{l} v(\{\})=0 \quad (\text{no passengers costs nothing}) \\ v(\{1\})=3, \quad v(\{2\})=7, \quad v(\{3\})=10 \\ v(\{1,2\})=7, \quad v(\{1,3\})=10, \quad v(\{2,3\})=10 \\ v(\{1,2,3\})=10 \end{array}$ 

$$1 2 3$$

$$3 3$$

$$3 3$$

$$3 4$$

$$3 4$$

$$3 4$$

$$3 4$$

$$3 4$$

$$\phi_i(v) = \sum_{S \subseteq N \setminus \{i\}} \frac{|S|!(N - |S| - 1)!}{N!} \left( v(S \cup \{i\}) - v(S) \right), \quad i = 1, \dots, N$$

 $\begin{array}{l} v(\{\})=0 \quad (\text{no passengers costs nothing}) \\ v(\{1\})=3, \quad v(\{2\})=7, \quad v(\{3\})=10 \\ v(\{1,2\})=7, \quad v(\{1,3\})=10, \quad v(\{2,3\})=10 \\ v(\{1,2,3\})=10 \end{array}$ 

1

 $\begin{array}{l} v(\{\})=0 \quad (\text{no passengers costs nothing}) \\ v(\{1\})=3, \quad v(\{2\})=7, \quad v(\{3\})=10 \\ v(\{1,2\})=7, \quad v(\{1,3\})=10, \quad v(\{2,3\})=10 \\ v(\{1,2,3\})=10 \end{array}$ 

$$\frac{1}{2}$$

$$\phi_i(v) = \sum_{S \subseteq N \setminus \{i\}} \frac{|S|!(N - |S| - 1)!}{N!} \left( v(S \cup \{i\}) - v(S) \right), \quad i = 1, \dots, N$$

Sets excluding passenger 1:  $\{\}, \{2\}, \{3\}, \{2, 3\}$ 

 $\begin{array}{l} v(\{\})=0 \quad (\text{no passengers costs nothing}) \\ v(\{1\})=3, \quad v(\{2\})=7, \quad v(\{3\})=10 \\ v(\{1,2\})=7, \quad v(\{1,3\})=10, \quad v(\{2,3\})=10 \\ v(\{1,2,3\})=10 \end{array}$ 

![](_page_30_Figure_2.jpeg)

$$\phi_i(v) = \sum_{S \subseteq N \setminus \{i\}} \frac{|S|!(N - |S| - 1)!}{N!} \left( v(S \cup \{i\}) - v(S) \right), \quad i = 1, \dots, N$$

Sets excluding passenger 1:  $\{\}, \{2\}, \{3\}, \{2, 3\}$ 

 $\phi_1 = \underbrace{\frac{1(3-0-1)!}{3!} \left( v(\{1\}) - v(\{\}) \right) + \underbrace{\frac{1(3-1-1)!}{3!} \left( v(\{1,2\}) - v(\{2\}) \right) + \underbrace{\frac{1(3-1-1)!}{3!} \left( v(\{1,3\}) - v(\{3\}) \right) + \underbrace{\frac{1(3-2-1)!}{3!} \left( v(\{1,2,3\} - v(\{2,3\}) \right) = 1}_{3!} \right)}_{3!}$ 

True story: The fair and unique way to distribute the cost of the journey (at a price of 1NOK per km), is when passenger 1 pays 1, passenger 2 pays 3 and passenger 3 pays 6.

 $\phi_{1} = \frac{1}{3} \left( v(\{1,2,3\} - v(\{2,3\})) + \frac{1}{6} \left( v(\{1,2\}) - v(\{2\}) \right) + \frac{1}{6} \left( v(\{1,3\}) - v(\{3\}) \right) + \frac{1}{3} \left( v(\{1\}) - v(\emptyset) \right) = 1$   $\phi_{2} = \frac{1}{3} \left( v(\{1,2,3\} - v(\{1,3\})) + \frac{1}{6} \left( v(\{1,2\}) - v(\{1\}) \right) + \frac{1}{6} \left( v(\{2,3\}) - v(\{3\}) \right) + \frac{1}{3} \left( v(\{2\}) - v(\emptyset) \right) = 3$  $\phi_{3} = \frac{1}{3} \left( v(\{1,2,3\} - v(\{1,2\})) + \frac{1}{6} \left( v(\{1,3\}) - v(\{1\}) \right) + \frac{1}{6} \left( v(\{2,3\}) - v(\{2\}) \right) + \frac{1}{3} \left( v(\{3\}) - v(\emptyset) \right) = 6$ 

#### Shapley values for machine learning...

Shapley values do dependence attribution. From game theory to ML

Nall players --> features

i player --> feature

Scoalition of players --> set of features

vgame --> model

The Shapley value takes as input a set function  $v: 2^N \rightarrow R$  and produces attributions  $\varphi_i$  for each player  $i \in N$  that add up to v(N)

![](_page_32_Picture_7.jpeg)

![](_page_32_Picture_8.jpeg)

#### Shapley values for machine learning...

#### Various libraries: SHAP, SAGE, ...

Model prediction drivers Model loss drivers

				and the second		Location	Trysil
						Sam	130 m <sup>2</sup>
-offer ; a	Will Strategies					Elevation	898 m
HIT T	The second second	a service	To prove the st	Same Black Brown .		Built	1976
-	Manarter			Chigan (the i	A state of the second of the	Distance to ski	6.2 km
ET S		the -	and the may	one of the Sol America		Distance to road	18 meters
	THE STATE		A DE			Distance to sea	175 km
and the second second		at a	THE REAL PROPERTY OF			Distance to lake	10 km
Andrew		V. and				Neighbours	776 w/1000 m
No de la		E.				Near neighbour	31 m
	in the	State State		HERE A	- Anti- C	Population	76 w/5 km
		- 3			5. Ask	0 MNOK 5.2 M KING PRICE PREDIC	NOK
		2	MNOK	3 MNOK	4 MNOK	6 MNOK	
			$\rangle \rangle \rangle \rangle \rangle \rangle$	$\rightarrow$	$\rangle$		
1							
PwC - R	Distance to ski 6.2 km responsible AI - Inga Stri	Latitude 61.327 ümke	Elevation 898 meters	Neighbours w/1000m 776 neighbours	Square meters 130 m <sup>2</sup> Driving price up	Built Population w 1976 76 inhabitat Driving price down	/5km nts

#### Intrinsic explanations

#### What does the model focus on?

**INPUT IMAGE ACTIVATIONS** of neuron groups

#### Intrinsic explanations

What do the <u>different parts</u> of the model focus on?

![](_page_35_Picture_2.jpeg)

![](_page_35_Picture_3.jpeg)

![](_page_35_Picture_4.jpeg)

#### Intrinsic explanations

A convincing story and the truth are not necessarily the same thing

(this holds for all aspects in life)

	Test Image	Evidence for Animal Being a Siberian Husky	Evidence for Animal Being a Transverse Flute
Explanations Using Attention Maps			

#### The story these days

How doe	es a model ork?	What is driving decisions?	Can I trust the model?	
		Key stakeholders		
Data Scientist	Business Owner	Model Risk	Regulator	Consumer
<ul> <li>Understand the model</li> <li>De-bug it</li> <li>Improve its performance</li> </ul>	<ul> <li>Understand the model</li> <li>Evaluate fit for purpose</li> <li>Agree to use</li> </ul>	<ul> <li>Challenge the model</li> <li>Ensure its robustness</li> <li>Approve it</li> </ul>	<ul> <li>Check its impact on consumers</li> <li>Verify reliability</li> </ul>	<ul> <li>"What is the Impact on me?"</li> <li>"What actions can I take?"</li> </ul>

#### ... can you spot what's missing?

#### Part 3: Ethics

I cannot stress enough how important ethical development is in AI.

Think

AI impact  $\infty$  medicine + atomic bomb

![](_page_38_Picture_4.jpeg)

![](_page_38_Picture_5.jpeg)

![](_page_38_Picture_6.jpeg)

# Facebook language predicts depression in medical records

#### Johannes C. Eichstaedt<sup>a,1,2</sup>, Robert J. Smith<sup>b,1</sup>, Raina M. Merchant<sup>b,c</sup>, Lyle H. Ungar<sup>a,b</sup>, Patrick Crutchley<sup>a,b</sup>, Daniel Preoțiuc-Pietro<sup>a</sup>, David A. Asch<sup>b,d</sup>, and H. Andrew Schwartz<sup>e</sup>

<sup>a</sup>Positive Psychology Center, University of Pennsylvania, Philadelphia, PA 19104; <sup>b</sup>Penn Medicine Center for Digital Health, University of Pennsylvania, Philadelphia, PA 19104; <sup>c</sup>Department of Emergency Medicine, Perelman School of Medicine, University of Pennsylvania, PA 19104; <sup>d</sup>The Center for Health Equity Research and Promotion, Philadelphia Veterans Affairs Medical Center, Philadelphia, PA 19104; and <sup>e</sup>Computer Science Department, Stony Brook University, Stony Brook, NY 11794

Edited by Susan T. Fiske, Princeton University, Princeton, NJ, and approved September 11, 2018 (received for review February 26, 2018)

Depression, the most prevalent mental illness, is underdiagnosed and undertreated, highlighting the need to extend the scope of current screening methods. Here, we use language from Facebook posts of consenting individuals to predict depression recorded in electronic medical records. We accessed the history of Facebook statuses posted by 683 patients visiting a large urban academic emergency department, 114 of whom had a diagnosis of depression in their medical records. Using only the language preceding their first documentation of a diagnosis of depression, we could identify depressed patients with fair accuracy [area under the curve (AUC) = 0.69, approximately matching the accuracy of screening surveys benchmarked against medical records. Restricting Facebook data to only the 6 months immediately preceding the first documented diagnosis of depression yielded a higher prediction accuracy (AUC = 0.72) for those users who had sufficient Facebook data. Significant prediction of future depression status was possible as far as 3 months before its first documentation. We found that language predictors of depression include emotional (sadness), interpersonal (loneliness, hostility), and cognitive (preoccupation with the self, rumination) processes. Unobtrusive depression assessment through social media of consenting individuals may become feasible as a scalable complement to existing screening and monitoring procedures.

the diagnosis of depression, which prior research has shown is feasible with moderate accuracy (15). Of the patients enrolled in the study, 114 had a diagnosis of depression in their medical records. For these patients, we determined the date at which the first documentation of a diagnosis of depression was recorded in the EMR of the hospital system. We analyzed the Facebook data generated by each user before this date. We sought to simulate a realistic screening scenario, and so, for each of these 114 patients, we identified 5 random control patients without a diagnosis of depression in the EMR, examining only the Facebook data they created before the corresponding depressed patient's first date of a recorded diagnosis of depression. This allowed us to compare depressed and control patients' data across the same time span and to model the prevalence of depression in the larger population ( $\sim$ 16.7%).

#### Results

**Prediction of Depression.** To predict the future diagnosis of depression in the medical record, we built a prediction model by using the textual content of the Facebook posts, post length, frequency of posting, temporal posting patterns, and demographics (*Materials and Methods*). We then evaluated the performance of this model by comparing the probability of depression estimated by our algorithm

SANG

#### The full picture?

We don't even know what we expect from an explanation

Explanation for human-beings: Not just the model but its impact (depends on environment and context)

![](_page_40_Picture_3.jpeg)

![](_page_40_Picture_4.jpeg)

## Thank you!

inga.strumke@ntnu.no

#### Discussion points:

- 1. I said that the discussion on XAI is about machine learning. Does it have to be?
  - a. Can non-learning approaches to AI require explanations?
- 2. What are heat maps?
  - a. What are their weaknesses?
  - b. Do they match the way humans would explain a visual decision?
- 3. What is a feature importance ranking?
  - a. How could a bank use feature ranking to tell you why you got a loan?
  - b. Would you accept such an explanation, and why (not)?
  - c. What weaknesses does this have when used as an explanation?