Assessment of a Trustworthy AI System in Healthcare: Smartphone-based Monitoring of Bipolar Disorder

Katarzyna Kaczmarek-Majer, Ph.D., D.Sc.

Computational Intelligence for Healthcare, Trustworthy AI Lab for Heatlhcare Systems Research Institute Polish Academy of Sciences, Warsaw, Poland

E-mail: k.kaczmarek@ibspan.waw.pl

Workshop 'Social Impact of AI' Madrid, 30-31 May 2024

Assessment of a Trustworthy AI System in Healthcare: Smartphone-based

Monitoring of Bipolar Disorder

- Trustworthiness
- ② Explainability

Assessment of a Trustworthy AI System in Healthcare: Smartphone-based

Monitoring of Bipolar Disorder

- Trustworthiness
- ② Explainability

Trustworthy Artificial Intelligence¹:

- robust (technical and social perspective)
- 2 lawful
- ethical
 - Respect for human autonomy
 - Prevention of harm,
 - Fairness
 - Explicability



This illustration of artificial intelligence has in fact been generated by Al

Policy-makers pledged to develop a 'human-centric' approach to Al²

¹High-Level Expert Group on Artificial Intelligence, "Ethics guidelines for trustworthy AI," European Commission, Text

² "EU AI Act: first regulation on artificial intelligence". European Parliament News. Archived from the original on 10 January 2024. Retrieved 22 January 2024

Use case: Smartphone-based Monitoring of Bipolar Disorder



Participants of BDMon study⁶ (100 patients diagnosed with BD) received a dedicated

mobile application, called BDMon, able to collected acoustic data.

⁵A. Z. Antosik-Wojcinska, M. Dominiak, M. Chojnacka, K. Kaczmarek-Majer, K. R. Opara, W. Radziszewska, A. Olwert, and L. Swiecicki, Smartphone as a monitoring tool for bipolar disorder: a systematic review including data analysis, machine learning algorithms and predictive modelling, Int J Med Inform, vol. 138:104131, 2020.

⁶Study was conduced within the CHAD project entitled "Smartphone-based diagnostics of phase changes in the course of bipolar disorder" (RPMA.01.02.00-14-5706/16-00) that was financed from EU funds (Regional Operational Program for Mazovia) in 2017-2018 • The fundament of the

Z-Inspection process is the

identification and discussion

of ethical issues and tensions

through the elaboration of

socio-technical scenarios. This

can only be achieved with

interdisciplinary team

• Assessment of trustworthy AI systems with

the Z-Inspection **process** (https://z-inspection.org/).



⁷Zicari, R. V., Brodersen, J., Brusseau, J., Dudder, B., Eichhorn, T., Ivanov, T., et al. (2021). Z-inspection: A Process to Assess Trustworthy Al. IEEE Trans. Technol. Soc.

Illustrative example. Socio-Technical Scenario #1

Doctor using the BDMon system based on Model's Predictions and Alarms

Overview of the BDMon pipeline can be summarized as follows:

- measurements and data collection (data acquisition, data annotation with ground truth, and signal processing)
- ② feature extraction and feature selection
- Itraining of the classifier using the annotated examples
- once the model is trained, actions are taken for new data (alarm or no alarm), based on the model's prediction and interpreted by the doctor and discussed with the patient
 - (1) The BDMon system predicts an alarm for a person
 - (2) The BDMon system predicts no alarm for a person

Illustrative example.

Question to the doctors/team members to describe problems they see—or expect to see—arising from the BDMon system's use

- false positives, false negatives, true positives
- extensive trust into alarms linked to the perception issues of BD patients
- Iow trust in alarms (even though approx. 90% sensitivity of the system!)
- occasional lack of communication with the patient
- various groups of variables result differently efficient for various types of patient (acoustic, locomotion sensor, armband, etc.)
- delays in data transmissions

٢

- database disconnected (technical watchdogs added)
- doctors overloaded with IT systems (need for extremely simple interface)
- no answer when the voice bot is calling (goodwill of patient needed)



Consulting the requirements listed in the Ethics Guidelines for Trustworthy AI

for individual perspectives for Illustrative Socio-Technical Scenario #1
Ethical Issue #1: Description: Using a black-box algorithm might impair the

trust of the psychiatrists in the BDMon app, especially if the app has not been verified by multiple prospective studies. This complicates the accountability. Neither the reasoning behind alarms not the responsibility of alarms are not clearly defined for clinicians, healthcare institutions, etc.

- Map to Ethical Pillars and Requirements: Explicability > Transparency > Traceability, Explainability
- Ethical Tensions: Accuracy versus Transparency/Explainability





https://www.healthine.com/health-news/does-insulin-resistance-cause-fibromysigia

Assessment of a Trustworthy AI System in Healthcare: Smartphone-based Monitoring of Bipolar Disorder:

Trustworthiness

② Explainability

As extracted from the Cambridge Dictionary of English Language³:

an explanation - the details or reasons that someone gives to make something clear or easy to understand.

Definition

Given a certain audience, explainability refers to the details and reasons a model gives

to make its functioning clear or easy to understand⁴.

³E. Walter, Cambridge advanced learner's dictionary, Cambridge University Press, 2008

⁴A.B. Arrieta et al, Explainable artificial intelligence (xai): Concepts, taxonomies, opportunities and challenges toward responsible ai, Inform. Fusion 58 (2020) 82–115.

- Interpretability: level of understanding how the underlying technology works⁸
- Explainability: level of understanding how the AI-based system came up with a given result⁸

Explanation is textual or visual artifact that provide a qualitative understanding of the relationship between the instance's components (e.g. words in text, patches in an image) and the model's prediction¹⁰

⁸ ISO/IEC TR 29119-11:2020, Software and systems engineering, Software testing, Part 11: Guidelines on the testing of Al-based systems. ISO. 2020. Retrieved 25 January 2024

⁹D. A. Broniatowski, "Psychological Foundations of Explainability and Interpretability in Artificial Intelligence," National Institute of Standards and Technology, Tech. Rep., Apr. 2021.

¹⁰MT Ribeiro, S Singh, C Guestrin (2016) "Why should I trust you?"Explaining the predictions of any classifier Proceedings of the 22nd ACM Knowledge Discovery and Data Mining (ACM KDD)



(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)

Predictions of Google's Pre-trained Inception neural network¹

explained by a post-hoc method ²

¹C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich. Going deeper with convolutions. In Computer Vision and Pattern Recognition (CVPR), 2015.

 $^{^2}$ M.T. Ribeiro, S. Singh, C. Guestrin, why should I trust you? explaining the predictions of any classifier, in: Proceedings of the 22nd ACM

SIGKDD International Conference on Knowledge Discovery and Data Mining, 2016, pp. 1135-1144.





 Symptoms in the common rating scales used in psychiatry, such as Hamilton Depression Rating Scale and Young Mania Rating Scale



Fuzzy linguistic summarization

• Fuzzy linguistic summaries (LSs) are statements in natural language that describe numerical datasets¹¹. I Ss have been confirmed as human-consistent information granules with applications in various domains.

Most young people are tall. Few young people are tall. Most young people are short

Most calls with high loudness in mania have low spectrum

The jump height achieved is lower since phase 1 is extended in time. The jump height achieved is lower since the first maximum is much greater than the second one in phase 3. It represents an excessive lowering of the centerof gravity.

¹¹J. Kacprzyk, R. R. Yager, and J. M. Merigo (2019) Towards human-centric aggregation via ordered weighted aggregation operators and linguistic data summaries: A new perspective on zadeh's inspirations," IEEE Computational Intelligence Magazine, vol. 14, no. 1, pp. 16–30

¹² J. Moreno-Garcia, J. Abian - Vicen, L. Jimenez-Linares, L. Rodriguez-Benitez, Description of multivariate time series by means of trends characterization in the fuzzy domain, Fuzzy Sets and Systems 285 (2016) 118 - 139.

• Frequent patterns e.g., If A happens before B and in the meantime we do not

observe C, then it is a failure of class X^{14} .

• Fuzzy Association Rules e.g., IF Strength of Seasonality is Small AND

Coefficient of Variation is Roughly Small THEN Weight of the j-th method is Big¹⁵.

 with Intermediate Quantifiers¹⁶ e.g., <u>If tone of central banks news is medium</u> and central banks are medium experienced in inflation targeting then customers expect small inflation support 0.24 and confidence 0.78.¹⁷.

¹⁴F. Höppner, S. Peter, and M. Berthold, Enriching multivariate temporal patterns with context information to support

classification, Computational Intelligence in Intelligent Data Analysis. vol. 445, pp. 195-206, 2013.

¹⁵M. Burda, M. Stepnicka, and L. Stepnicková, Fuzzy rule-based ensamble for time series prediction: Progresses with associations mining, in Strength. Links Between Data Analysis and Soft Computing, vol. 315, pp. 261-271, Springer, 2014

¹⁶Novák V. and Murinová P.: A formal model of the intermediate quantifiers "A few, Several, A little" 2019

¹⁷Association Rules on Data using Intermediate Quantifiers, Murinová Petra, Karel Fiala, Katarzyna Kaczmarek-Majer and Aleksandra Rutkowska, FSTA 2024)

How to efficiently communicate with the psychiatrist about an alarming situation basing on the data collected from smartphone ?



Most outgoing calls of patient P are long [T=0.8]Depressive episode may have started. • We construct fuzzy linguistic summaries using evolving membership functions based on prototypes from semi-supervised learning following the idea of **evolving** fuzzy systems ¹⁸.

¹⁸P. Angelov, D. P. Filev, N. Kasabov, Evolving Takagi-Sugeno Fuzzy Systems from Streaming Data (eTS+), 2010, pp. 21 - 50

¹⁹Hryniewicz, O. Kaczmarek-Majer, K. Opara, K. (2019) Control Charts Based on Fuzzy Costs for Monitoring Short Autocorrelated Time Series. International Journal of Approximate Reasoning, p 166-181, 10.1016/j.ijar.2019.08.013

LS-FC: Explaining partially-labeled data in natural language²⁰



Rysunek: Overview of the proposed *LS-FC* approach that explains data stream by means of summarization and online learning

²⁰K. Kaczmarek-Majer, G. Casalino, G. Castellano, O. Hryniewicz, M. Dominiak, Explaining smartphone-based acoustic data in bipolar disorder: Semi-supervised fuzzy clustering and relative linguistic summaries, Information Sciences 588 (2022) 174–195.

²¹Kaczmarek-Majer, K. and Hryniewicz, O. (2019) Application of linguistic summarization methods in time series forecasting. Information Sciences Vol 478, p 580 – 594

Table 6

Relative linguistic summaries based on short protoforms for mania and hypomania episodes (LS with T = 1.0) and extended protoforms for mania and hypomania episodes (LS with T > 0.5).

Relative LS based on short protoform	Т
Most calls in the state of mania have low spectrum compared to the state of euthymia.	1.0
Most calls in the state of mania have low quality compared to the state of euthymia.	1.0
Most calls in the state of hypomania have low spectrum compared to the state of euthymia.	1.0
Most calls in the state of hypomania have low loudness compared to the state of euthymia.	1.0
Most calls in the state of hypomania have low qualty compared to the state of euthymia.	1.0
Relative LS based on extended protoform - HYPOMANIA	Т
Most calls with low loudness in hypomania have low spectrum compared to the state of euthymia.	1.0
Most calls with low loudness in hypomania have low quality compared to the state of euthymia.	1.0
Most calls with high loudness in hypomania have high spectrum compared to the state of euthymia.	1.0
Most calls with high loudness in hypomania have high quality compared to the state of euthymia.	1.0
Most calls with low pitch in hypomania have low spectrum compared to the state of euthymia.	1.0
Most calls with low pitch in hypomania have low loudness compared to the state of euthymia.	1.0
Most calls with low pitch in hypomania have low quality compared to the state of euthymia.	1.0
Most calls with low spectrum in hypomania have low loudness compared to the state of euthymia.	1.0
Most calls with low spectrum in hypomania have low quality compared to the state of euthymia.	1.0
Most calls with high spectrum in hypomania have high loudness compared to the state of euthymia.	1.0
Most calls with high spectrum in hypomania have high quality compared to the state of euthymia.	1.0
Most calls with low quality in hypomania have low loudness compared to the state of euthymia.	1.0
Most calls with low quality in hypomania have low spectrum compared to the state of euthymia.	1.0
Most calls with high quality in hypomania have high loudness compared to the state of euthymia.	1.0
Most calls with high quality in hypomania have high spectrum compared to the state of euthymia.	1.0

PLENARY: Explaining black-box models in natural language using fuzzy linguistic summaries²⁴

LS1: Among records that contribute positively to predicting depression class,

most of them have energy-related features at low level.



²⁴Katarzyna Kaczmarek-Majer, Gabriella Casalino, Giovanna Castellano, Monika Dominiak, Olgierd Hryniewicz, Olga Kamińska, Gennaro Vessio, Natalia Díaz-Rodríguez, PLENARY: Explaining black-box models in natural language through fuzzy linguistic summaries, Information Sciences, Volume 614, 2022, Pages 374-399.

The program code and running examples of are available at the following link: https://github.com/ PLENARY ITPsychiatry/plenary

Table 6

Evaluation of linguistic summaries from PLENARY for the prediction of BD classes with the sequential and compositional MLP model. Degree of truth, degree of support, degree of focus, and expert-based degree of usefulness are applied as criteria. Post-processing criteria: DoT > 0.1 and DoF > 0.05. Summaries that contribute positively to predicting a class are presented in bold. The font colors of the LS description indicate the high-level semantic groups of acoustic features. LS related to: the energy-related features are marked in black; the spectral-related features are marked in olive; the pitch-related features in orange; and the quality-related features are marked in purple.

Id	LS description	DoT	DoS	DoF	DoU
001	Among records that contribute around zero to predicting euthymia, most of them have energy-related features at low level.	0.58	0.17	0.06	1
002	Among records that contribute positively to predicting euthymia, most of them have energy-related features at low level.	0.24	0.17	0.21	5
003	Among records that contribute against predicting euthymia, most of them have spectral-related features at high level.	0.19	0.54	0.63	2
004	Among records that contribute around zero to predicting euthymia, most of them have spectral-related features at low level.	0.53	0.17	0.06	1
005	Among records that contribute positively to predicting euthymia, most of them have spectral-related features at low level.	1.00	0.30	0.21	4
006	Among records that contribute against predicting euthymia, most of them have quality-related features at high level.	0.26	0.70	0.63	3
007	Among records that contribute positively to predicting euthymia, most of them have quality-related features at low level.	0.23	0.19	0.21	4
101	Among records that contribute around zero to predicting depression, most of them have energy-related features at high level.	0.12	0.17	0.06	1
102	Among records that contribute positively to predicting depression, most of them have spectral-related features at high level.	1.00	0.29	0.31	5
103	Among records that contribute against predicting depression, most of them have quality-related features at low level.	0.51	0.61	0.76	4
104	Among records that contribute positively to predicting depression, most of them have quality-related features at low level.	1.00	0.18	0.31	5
201	Among records that contribute against predicting mania, most of them have energy-related features at low level.	0.33	0.68	0.73	4
202	Among records that contribute around zero to predicting mania, most of them have energy-related features at low level.	1.00	0.19	0.03	1
203	Among records that contribute against predicting mania, most of them have pitch-related features at low level.	0.25	0.45	0.73	4
204	Among records that contribute around zero to predicting mania, most of them have pitch-related features at low level.	1.00	0.05	0.03	1
205	Among records that contribute positively to predicting mania, most of them have pitch-related features at high level.	0.59	0.39	0.44	5
206	Among records that contribute positively to predicting mania, most of them have spectral-related features at low level.	1.00	0.27	0.44	5
301	Among records that contribute positively to predicting mixed state, most of them have energy-related features at high level.	0.11	0.16	0.31	5
302	Among records that contribute positively to predicting mixed state, most of them have pitch-related features at low level.	0.45	0.34	0.31	5
303	Among records that contribute against predicting mixed state, most of them have spectral-related features at low level.	0.11	0.50	0.63	3
304	Among records that contribute positively to predicting mixed state, most of them have spectral-related features at high level.	1.00	0.27	0.31	5
305	Among records that contribute against predicting mixed state, most of them have quality-related features at low level.	0.75	0.66	0.63	3

²⁴ Katarzyna Kaczmarek-Majer, Gabriella Casalino, Giovanna Castellano, Monika Dominiak, Olgierd Hryniewicz, Olga Kamińska, Gennaro Vessio, Natalia Díaz-Rodríguez, **PLENARY: Explaining black-box models in natural language through fuzzy linguistic summaries**, Information Sciences, Volume 614, 2022, Pages 374-399.

The program code and running examples of are available at the following link: https://github.com/ PLENARY ITPsychiatry/plenary



- Z-Inspection as a holistic process to assess Trustworthy AI.
- Fuzzy linguistic summaries as explanations / human-consistent information granules with applications in various domains. Individual sentences have not always proven sufficient without exposing a clarification for the group of sentences to explain the rationale of the intelligent system's behaviour. Semi-Supervised Fuzzy-C Means as an explainable-by-design method for partially-labeled data.

Thank to all my collaborators!

Thank you for your attention!

Katarzyna Kaczmarek-Majer

k.kaczmarek@ibspan.waw.pl

https://www.ibspan.waw.pl/ kaczmar/

http://z-inspection.org/ http://bipolar.ibspan.waw.pl https://github.com/ITPsychiatry/ssfclust https://github.com/ITPsychiatry/bipolar https://github.com/ITPsychiatry/plenary