

# **Quality Assurance in Explainable Artificial Intelligence Engineering**

Alex Cummaudo

BSc *Swinburne*, BIT(Hons)



Applied Artificial Intelligence Institute  
Deakin University  
Melbourne, Australia

August 9, 2018



# Contents

<b>Contents</b>	<b>ii</b>
<b>List of Figures</b>	<b>iv</b>
<b>List of Tables</b>	<b>v</b>
<b>List of Listings</b>	<b>vii</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Background . . . . .	1
1.1.1 Problem . . . . .	2
1.2 Motivation . . . . .	2
1.3 Research Goals . . . . .	2
1.4 Strategy and Roadmap . . . . .	2
1.5 Contributions . . . . .	2
<b>2 Literature Review</b>	<b>3</b>
<b>3 Methodology</b>	<b>5</b>
3.1 Data Collection and Ethics . . . . .	5
3.2 Approach . . . . .	5
3.3 Evaluation Methods . . . . .	5
3.4 Threats to Validity . . . . .	5
3.4.1 Internal Threats . . . . .	5
3.4.2 External Threats . . . . .	5
3.4.3 Construct . . . . .	5

<b>4</b>	<b>Project Status</b>	<b>7</b>
4.1	Completed Work . . . . .	7
4.2	Impact . . . . .	7
4.3	Timeline . . . . .	7
<b>5</b>	<b>Conclusion</b>	<b>9</b>
	<b>References</b>	<b>18</b>

# List of Figures



# List of Tables





# List of Listings



# Chapter 1

## Introduction

### 1.1 Background

As the rise of AI increases, the need for engineering model interpretability around AI becomes paramount. Model interpretability has been stressed since early machine learning research in the late 1980s and 1990s (such as Quinlan [56] and Michie [50]), and while there has since been a significant body of work in the area [2, 3, 6, 7, 10, 15, 22, 28, 38, 46–48, 53, 57, 60, 64, 71, 72], it is evident that ‘accuracy’ or model ‘confidence’ is still used as a primary criterion for AI evaluation [33, 35, 65]. Indeed, much research into NN or SVM development stresses that ‘good’ models are those with high accuracy. However, is accuracy enough to justify a model’s quality?

To answer this, we revisit what it means for a model to be accurate. Accuracy is an indicator for estimating how well a model’s algorithm will work with future or unforeseen data. It is quantified in the AI testing stage, whereby the algorithm is tested against cases known by humans to have ground truth but such cases are unknown by the algorithm. In production, however, all cases are unknown by both the algorithm *and* the humans behind it, and therefore a single value of quality is “not reliable if the future dataset has a probability distribution significantly different from past data” [24], a problem commonly referred to as the *datashift* problem [67]. Analogously, Freitas [24] provides the following description of the problem:

*The military trained [a NN] to classify images of tanks into enemy and friendly tanks. However, when the [NN] was deployed in the field (corresponding to “future data”), it had a very poor accuracy rate. Later, users noted that all photos of friendly (enemy) tanks were taken on a sunny (overcast) day. I.e., the [NN]*

*learned to discriminate between the colors of the sky in sunny vs. overcast days!*  
*If the [NN] had output a comprehensible model (explaining that it was discriminating between colors at the top of the images), such a trivial mistake would immediately be noted. [24]*

So, why must we interpret models? While the trade-off between accuracy and interpretability is critical [18, 23, 29, 37, 40, 76], a quantifiable value cannot satisfy more subjective needs of end-users: trust, anomaly detection, legality and insight. The ever-growing domains of applied ML in areas such as medicine [5, 21, 38, 44, 54, 58, 71, 75, 78], bioinformatics [17, 25, 36, 39, 69], finance [3, 16, 34] and customer analytics [46, 72] engages application-end needs....

How do we interpret models? Methods for developing interpretation models include: decision trees [9, 14, 30, 55, 59], decision tables [4, 46] and decision sets [43, 51]; input gradients, gradient vectors or sensitivity analysis [3, 45, 57, 60, 62]; exemplars [26, 41]; generalised additive models [11]; classification (*if-then*) rules [8, 13, 52, 70, 74] and falling rule lists [64]; nearest neighbours [48, 63, 68, 73, 77] and Naïve Bayes analysis [5, 12, 27, 32, 42, 44, 49, 78]. Several cross-domain studies have assessed the interpretability of these techniques against end-users, measuring response time, accuracy in model response and user confidence [1, 25, 31, 34, 48, 61, 66, 72], although it is generally agreed that decision rules and decision tables provide the most interpretation in non-linear models such as SVM's or NN's [25, 48, 72]. For an extensive survey of the benefits and fallbacks of these techniques, we refer to Freitas [24] and Doshi-Velez and Kim [19].

### **1.1.1 Problem**

As it stands, AI presents an issue with. (For a detailed discussion, see Doshi-Velez et al. [20].

## **1.2 Motivation**

## **1.3 Research Goals**

## **1.4 Strategy and Roadmap**

## **1.5 Contributions**

## **Chapter 2**

### **Literature Review**



# **Chapter 3**

## **Methodology**

### **3.1 Data Collection and Ethics**

### **3.2 Approach**

### **3.3 Evaluation Methods**

### **3.4 Threats to Validity**

#### **3.4.1 Internal Threats**

#### **3.4.2 External Threats**

#### **3.4.3 Construct**





# **Chapter 4**

## **Project Status**

### **4.1 Completed Work**

### **4.2 Impact**

### **4.3 Timeline**



## **Chapter 5**

## **Conclusion**



# References

- [1] H. Allahyari and N. Lavesson, “User-oriented assessment of classification model understandability,” in *11th scandinavian conference on Artificial intelligence*. IOS Press, 2011.
- [2] M. G. Augasta and T. Kathirvalavakumar, “Reverse engineering the neural networks for rule extraction in classification problems,” *Neural processing letters*, vol. 35, no. 2, pp. 131–150, 2012.
- [3] D. Baehrens, T. Schroeter, S. Harmeling, M. Kawanabe, K. Hansen, and K.-R. MÅžller, “How to explain individual classification decisions,” *Journal of Machine Learning Research*, vol. 11, no. Jun, pp. 1803–1831, 2010.
- [4] B. Baesens, C. Mues, M. De Backer, J. Vanthienen, and R. Setiono, “Building intelligent credit scoring systems using decision tables,” in *Enterprise Information Systems V*. Springer, 2004, pp. 131–137.
- [5] R. Bellazzi and B. Zupan, “Predictive data mining in clinical medicine: current issues and guidelines,” *International journal of medical informatics*, vol. 77, no. 2, pp. 81–97, 2008.
- [6] A. Ben-David, “Monotonicity maintenance in information-theoretic machine learning algorithms,” *Machine Learning*, vol. 19, no. 1, pp. 29–43, 1995.
- [7] O. Boz, “Extracting decision trees from trained neural networks,” in *Proceedings of the eighth ACM SIGKDD international conference on Knowledge discovery and data mining*. ACM, 2002, pp. 456–461.
- [8] M. Bramer, *Principles of data mining*. Springer, 2007, vol. 180.
- [9] L. Breiman, J. H. Friedman, R. A. Olshen, and C. J. Stone, *Classification and Regression Trees*. CRC press, 1984.

- [10] A. Bussone, S. Stumpf, and D. O’Sullivan, “The Role of Explanations on Trust and Reliance in Clinical Decision Support Systems.” *ICHI*, 2015.
- [11] R. Caruana, Y. Lou, J. Gehrke, P. Koch, M. Sturm, and N. Elhadad, “Intelligible Models for HealthCare - Predicting Pneumonia Risk and Hospital 30-day Readmission.” *KDD*, pp. 1721–1730, 2015.
- [12] J. Cheng and R. Greiner, “Learning bayesian belief network classifiers: Algorithms and system,” in *Conference of the Canadian Society for Computational Studies of Intelligence*. Springer, 2001, pp. 141–151.
- [13] P. Clark and R. Boswell, “Rule induction with CN2: Some recent improvements,” in *European Working Session on Learning*. Springer, 1991, pp. 151–163.
- [14] M. Craven and J. W. Shavlik, “Extracting Tree-Structured Representations of Trained Networks.” *NIPS*, 1995.
- [15] K. Dejaeger, F. Goethals, A. Giangreco, L. Mola, and B. Baesens, “Gaining insight into student satisfaction using comprehensible data mining techniques,” *European Journal of Operational Research*, vol. 218, no. 2, pp. 548–562, 2012.
- [16] V. Dhar, D. Chou, and F. Provost, “Discovering Interesting Patterns for Investment Decision Making with GLOWER—A Genetic Learner Overlaid with Entropy Reduction,” *Data Mining and Knowledge Discovery*, vol. 4, no. 4, pp. 251–280, 2000.
- [17] M. Doderer, K. Yoon, J. Salinas, and S. Kwek, “Protein subcellular localization prediction using a hybrid of similarity search and error-correcting output code techniques that produces interpretable results,” *In silico biology*, vol. 6, no. 5, pp. 419–433, 2006.
- [18] P. Domingos, “Occam’s two razors: The sharp and the blunt,” in *KDD*, 1998, pp. 37–43.
- [19] F. Doshi-Velez and B. Kim, “Towards a rigorous science of interpretable machine learning,” 2017.
- [20] F. Doshi-Velez, M. Kortz, R. Budish, C. Bavitz, S. Gershman, D. O’Brien, S. Schieber, J. Waldo, D. Weinberger, and A. Wood, “Accountability of AI Under the Law: The Role of Explanation,” *arXiv.org*, Nov. 2017.

- [21] W. Elazmeh, W. Matwin, D. O’Sullivan, W. Michalowski, and W. Farion, “Insights from predicting pediatric asthma exacerbations from retrospective clinical data,” in *Evaluation Methods for Machine Learning II—Papers from 2007 AAAI Workshop*, 2007, pp. 10–15.
- [22] A. J. Feelders, “Prior knowledge in economic applications of data mining,” in *European Conference on Principles of Data Mining and Knowledge Discovery*. Springer, 2000, pp. 395–400.
- [23] A. A. Freitas, “A critical review of multi-objective optimization in data mining: a position paper,” *ACM SIGKDD Explorations Newsletter*, vol. 6, no. 2, pp. 77–86, 2004.
- [24] —, “Comprehensible classification models,” *ACM SIGKDD Explorations Newsletter*, vol. 15, no. 1, pp. 1–10, Mar. 2014.
- [25] A. A. Freitas, D. C. Wieser, and R. Apweiler, “On the importance of comprehensible classification models for protein function prediction,” *IEEE/ACM Transactions on Computational Biology and Bioinformatics (TCBB)*, vol. 7, no. 1, pp. 172–182, 2010.
- [26] B. J. Frey and D. Dueck, “Clustering by Passing Messages Between Data Points,” *Science*, vol. 315, no. 5814, pp. 972–976, Feb. 2007.
- [27] N. Friedman, D. Geiger, and M. Goldszmidt, “Bayesian network classifiers,” *Machine Learning*, vol. 29, no. 2-3, pp. 131–163, 1997.
- [28] G. Fung, S. Sandilya, and R. B. Rao, “Rule extraction from linear support vector machines,” in *Proceedings of the eleventh ACM SIGKDD international conference on Knowledge discovery in data mining*. ACM, 2005, pp. 32–40.
- [29] P. D. Grünwald, *The minimum description length principle*. MIT press, 2007.
- [30] T. Hastie, R. Tibshirani, and J. H. Friedman, *The Elements of Statistical Learning*, ser. Data Mining, Inference, and Prediction. Springer Science & Business Media, Jan. 2001.
- [31] B. Hayete and J. R. Bienkowska, “Gotrees: predicting go associations from protein domain composition using decision trees,” in *Biocomputing 2005*. World Scientific, 2005, pp. 127–138.

- [32] D. Heckerman, D. M. Chickering, C. Meek, R. Rounthwaite, and C. Kadie, “Dependency networks for inference, collaborative filtering, and data visualization,” *Journal of Machine Learning Research*, vol. 1, no. Oct, pp. 49–75, 2000.
- [33] J. Huang and C. X. Ling, “Using AUC and accuracy in evaluating learning algorithms,” *IEEE Transactions on knowledge and Data Engineering*, vol. 17, no. 3, pp. 299–310, 2005.
- [34] J. Huysmans, K. Dejaeger, C. Mues, J. Vanthienen, and B. Baesens, “An empirical evaluation of the comprehensibility of decision table, tree and rule based predictive models,” *Decision Support Systems*, vol. 51, no. 1, pp. 141–154, Apr. 2011.
- [35] N. Japkowicz and M. Shah, *Evaluating learning algorithms: a classification perspective*. Cambridge University Press, 2011.
- [36] T. Jiang and A. E. Keating, “AVID: an integrative framework for discovering functional relationships among proteins,” *BMC bioinformatics*, vol. 6, no. 1, p. 136, 2005.
- [37] Y. Jin, *Multi-objective machine learning*. Springer Science & Business Media, 2006, vol. 16.
- [38] U. Johansson and L. Niklasson, “Evolving decision trees using oracle guides,” in *IEEE Symposium on Computational Intelligence and Data Mining (CIDM)*. IEEE, 2009.
- [39] A. Karwath and R. D. King, “Homology induction: the use of machine learning to improve sequence similarity searches,” *BMC bioinformatics*, vol. 3, no. 1, p. 11, 2002.
- [40] K. A. Kaufman and R. S. Michalski, “Learning from inconsistent and noisy data: the AQ18 approach,” in *International Symposium on Methodologies for Intelligent Systems*. Springer, 1999, pp. 411–419.
- [41] B. Kim, C. Rudin, and J. A. Shah, “The Bayesian Case Model - A Generative Approach for Case-Based Reasoning and Prototype Classification.” *NIPS*, 2014.
- [42] I. Kononenko, “Inductive and Bayesian learning in medical diagnosis,” *Applied Artificial Intelligence an International Journal*, vol. 7, no. 4, pp. 317–337, 1993.
- [43] H. Lakkaraju, S. H. Bach, and J. Leskovec, “Interpretable Decision Sets - A Joint Framework for Description and Prediction.” *KDD*, pp. 1675–1684, 2016.



- [44] N. Lavrač, “Selected techniques for data mining in medicine,” *Artificial intelligence in medicine*, vol. 16, no. 1, pp. 3–23, 1999.
- [45] T. Lei, R. Barzilay, and T. Jaakkola, “Rationalizing Neural Predictions,” *arXiv.org*, Jun. 2016.
- [46] E. Lima, C. Mues, and B. Baesens, “Domain knowledge integration in data mining using decision tables: Case studies in churn prediction,” *Journal of the Operational Research Society*, vol. 60, no. 8, pp. 1096–1106, 2009.
- [47] Z. C. Lipton, “The Mythos of Model Interpretability.” *CoRR*, 2016.
- [48] D. Martens, J. Vanthienen, W. Verbeke, and B. Baesens, “Performance of classification models from a user perspective,” *Decision Support Systems*, vol. 51, no. 4, pp. 782–793, 2011.
- [49] D. Michie, D. J. Spiegelhalter, and C. C. Taylor, “Machine Learning and Statistical Classification of Artificial intelligence,” 1994.
- [50] D. Michie, “Machine Learning in the Next Five Years.” *EWSL*, 1988.
- [51] M. Narayanan, E. Chen, J. He, B. Kim, S. Gershman, and F. Doshi-Velez, “How do Humans Understand Explanations from Machine Learning Systems? An Evaluation of the Human-Interpretability of Explanation.” *CoRR*, 2018.
- [52] F. E. Otero and A. A. Freitas, “Improving the interpretability of classification rules discovered by an ant colony algorithm,” in *Proceedings of the 15th annual conference on Genetic and evolutionary computation*. ACM, 2013, pp. 73–80.
- [53] M. Pazzani, “Comprehensible knowledge discovery: gaining insight from data,” in *First Federal Data Mining Conference and Exposition*, 1997, pp. 73–82.
- [54] M. J. Pazzani, S. Mani, and W. R. Shankle, “Acceptance of rules generated by machine learning among medical experts,” *Methods of information in medicine*, vol. 40, no. 05, pp. 380–385, 2001.
- [55] J. R. Quinlan, “C4. 5: Programming for machine learning,” *Morgan Kauffmann*, vol. 38, p. 48, 1993.

- [56] ———, “Some elements of machine learning,” in *International Conference on Inductive Logic Programming*. Springer, 1999, pp. 15–18.
- [57] M. T. Ribeiro, S. Singh, and C. Guestrin, ““Why Should I Trust You?”,” in *the 22nd ACM SIGKDD International Conference*. New York, New York, USA: ACM Press, 2016, pp. 1135–1144.
- [58] G. Richards, V. J. Rayward-Smith, P. H. Sönksen, S. Carey, and C. Weng, “Data mining for indicators of early mortality in a database of clinical records,” *Artificial intelligence in medicine*, vol. 22, no. 3, pp. 215–231, 2001.
- [59] L. Rokach and O. Z. Maimon, *Data mining with decision trees: theory and applications*. World scientific, 2008, vol. 69.
- [60] A. S. Ross, M. C. Hughes, and F. Doshi-Velez, “Right for the Right Reasons: Training Differentiable Models by Constraining their Explanations,” *arXiv.org*, Mar. 2017.
- [61] M. Schwabacher, P. Langley, and P. Norvig, “Discovering communicable scientific knowledge from spatio-temporal data,” *ICML*, 2001.
- [62] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, and D. Batra, “Grad-CAM: Visual Explanations from Deep Networks via Gradient-Based Localization,” in *2017 IEEE International Conference on Computer Vision (ICCV)*. IEEE, 2017, pp. 618–626.
- [63] S. Sen and L. Knight, “A genetic prototype learner,” in *IJCAI*. Citeseer, 1995, pp. 725–733.
- [64] S. Singh, M. T. Ribeiro, and C. Guestrin, “Programs as Black-Box Explanations,” *arXiv.org*, Nov. 2016.
- [65] M. Sokolova and G. Lapalme, “A systematic analysis of performance measures for classification tasks,” *Information Processing & Management*, vol. 45, no. 4, pp. 427–437, 2009.
- [66] G. H. Subramanian, J. Nosek, S. P. Raghunathan, and S. S. Kanitkar, “A comparison of the decision table and tree,” *Communications of the ACM*, vol. 35, no. 1, pp. 89–94, 1992.

- [67] M. Sugiyama, N. D. Lawrence, and A. Schwaighofer, *Dataset shift in machine learning*. The MIT Press, 2017.
- [68] N. R. Suri, V. S. Srinivas, and M. N. Murty, “A cooperative game theoretic approach to prototype selection,” in *European Conference on Principles of Data Mining and Knowledge Discovery*. Springer, 2007, pp. 556–564.
- [69] D. Szafron, P. Lu, R. Greiner, D. S. Wishart, B. Poulin, R. Eisner, Z. Lu, J. Anvik, C. Macdonell, and A. Fyshe, “Proteome Analyst: custom predictions with explanations in a web-based tool for high-throughput proteome annotations,” *Nucleic acids research*, vol. 32, no. suppl<sub>2</sub>, pp. W365 – –W371, 2004.
- [70] S. Thrun, “Is Learning The n-th Thing Any Easier Than Learning The First?” p. 7, 1996.
- [71] A. Van Assche and H. Blockeel, “Seeing the forest through the trees: Learning a comprehensible model from an ensemble,” in *European conference on machine learning*. Springer, 2007, pp. 418–429.
- [72] W. Verbeke, D. Martens, C. Mues, and B. Baesens, “Building comprehensible customer churn prediction models with advanced rule induction techniques,” *Expert Systems with Applications*, vol. 38, no. 3, pp. 2354–2364, 2011.
- [73] D. Wettschereck, D. W. Aha, and T. Mohri, “A review and empirical evaluation of feature weighting methods for a class of lazy learning algorithms,” *Artificial Intelligence Review*, vol. 11, no. 1-5, pp. 273–314, 1997.
- [74] I. H. Witten, E. Frank, M. A. Hall, and C. J. Pal, *Data Mining: Practical machine learning tools and techniques*. Morgan Kaufmann, 2016.
- [75] M. L. Wong and K. S. Leung, *Data mining using grammar based genetic programming and applications*. Springer Science & Business Media, 2006, vol. 3.
- [76] J. Zahálka and F. Železný, “An experimental test of Occam’s razor in classification,” *Machine Learning*, vol. 82, no. 3, pp. 475–481, 2011.
- [77] J. Zhang, “Selecting typical instances in instance-based learning,” in *Machine Learning Proceedings 1992*. Elsevier, 1992, pp. 470–479.

- [78] B. Zupan, J. Demšar, M. W. Kattan, J. R. Beck, and I. Bratko, “Machine learning for survival analysis: a case study on recurrence of prostate cancer,” *Artificial intelligence in medicine*, vol. 20, no. 1, pp. 59–75, 2000.