

University of Groningen

Intrinsically Interpretable Machine Learning In Computer Aided Diagnosis

S. Ghosh, Sreejita

DOI:
[10.33612/diss.175627883](https://doi.org/10.33612/diss.175627883)

IMPORTANT NOTE: You are advised to consult the publisher's version (publisher's PDF) if you wish to cite from it. Please check the document version below.

Document Version
Publisher's PDF, also known as Version of record

Publication date:
2021

[Link to publication in University of Groningen/UMCG research database](#)

Citation for published version (APA):
S. Ghosh, S. (2021). *Intrinsically Interpretable Machine Learning In Computer Aided Diagnosis*. University of Groningen. <https://doi.org/10.33612/diss.175627883>

Copyright

Other than for strictly personal use, it is not permitted to download or to forward/distribute the text or part of it without the consent of the author(s) and/or copyright holder(s), unless the work is under an open content license (like Creative Commons).

The publication may also be distributed here under the terms of Article 25fa of the Dutch Copyright Act, indicated by the "Taverne" license. More information can be found on the University of Groningen website: <https://www.rug.nl/library/open-access/self-archiving-pure/taverne-amendment>.

Take-down policy

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

Downloaded from the University of Groningen/UMCG research database (Pure): <http://www.rug.nl/research/portal>. For technical reasons the number of authors shown on this cover page is limited to 10 maximum.



rijksuniversiteit
groningen

Intrinsically Interpretable Machine Learning In Computer Aided Diagnosis

Proefschrift

ter verkrijging van de graad van doctor aan de
Rijksuniversiteit Groningen
op gezag van de
rector magnificus Prof. Dr. C. Wijmenga
en volgens besluit van het College voor Promoties.

De openbare verdediging zal plaatsvinden op
vrijdag 10 September 2021 om 14.30 uur

door

Sreejita Ghosh

geboren op 31 maart 1991
te Kolkata, India

Promotores:

Prof. dr. M. Biehl
Prof. dr. N. Petkov

Copromotor:

Dr. K. Bunte

Beoordelingscommissie:

Prof. dr. D. Karastoyanova
Prof. dr. B. Hammer
Prof. dr. J. A. Lee

Contents

Acknowledgements	v
List of figures	ix
List of tables	xi
Preface	1
1 Introduction	3
1.1 Scope of this thesis	4
1.2 Outline	5
I Interpretability, Missingness, and Datasets explored	7
2 Interpretability in Machine Learning	9
2.1 Introduction	9
2.2 Taxonomy	10
2.3 Properties of an XAI	11
2.4 Types of explainability and interpretability and transparency	12
2.4.1 Types of interpretability	13
2.4.2 Post-hoc explainability	13
2.4.3 Levels of transparency	15
2.5 Some transparent models	15
2.5.1 Decision trees	15
2.5.2 K Nearest Neighbours	16
2.5.3 Learning Vector Quantization	16

2.6	Models needing post-hoc explainability	17
2.6.1	Random forest	17
2.7	Trade-off: interpretability vs performance	17
2.8	Conclusion	18
3	Missing values and Imbalanced classes	19
3.1	Introduction	19
3.2	Categories of missingness	19
3.2.1	Missing completely at random (MCAR)	20
3.2.2	Missing at random (MAR)	20
3.2.3	Missing not at random (MNAR)	20
3.3	Missingness handling strategies	21
3.3.1	Case deletion	21
3.3.2	Generative modelling of data distribution	22
3.3.3	Imputation	22
3.3.4	Techniques which could be applied after imputation	24
3.3.5	Techniques which can intrinsically deal with missing values	25
3.3.6	Prototype-based machine learning methods	27
3.4	Imbalanced classes	28
3.4.1	Sampling	29
3.4.2	Bagging	29
3.4.3	Boosting	29
3.4.4	Geodesic SMOTE	30
3.4.5	Incorporation of expert-knowledge: modified boosting	31
3.5	Conclusion	31
4	Datasets Explored	33
4.1	Introduction	33
4.2	Synthetic dataset	33
4.3	Real-life medical dataset	35
II	New variants of Learning Vector Quantization	39
5	Angle Learning Vector Quantization	41
5.1	Introduction	41
5.2	Learning Vector Quantization (LVQ)	42
5.2.1	Angle Learning Vector Quantization (Angle LVQ)	43
5.3	Experiments on synthetic data	45
5.3.1	Preliminary findings from experiments on synthetic data	46

Contents

5.3.2	Comparison of power of multiple learners	47
5.4	Geodesic average model	48
5.5	Experiments on GCMS datasets	49
5.5.1	Preliminary findings on the GCMS dataset	52
5.6	Knowledge extraction from averaged model	53
5.7	Conclusion	56
6	Angle Learning Vector Quantization for Non-Linear Decision Boundaries	57
6.1	Introduction	57
6.2	Datasets	57
6.2.1	Synthetic non-linear dataset (Football)	58
6.2.2	Heart disease dataset from UCI	58
6.3	Methodology	60
6.3.1	Local relevance matrix	60
6.3.2	Two-matrix decomposition for visualisation	61
6.4	Experiments	62
6.4.1	Synthetic data	62
6.4.2	Heart disease data	63
6.4.3	GCMS dataset	64
6.5	Results	64
6.5.1	Synthetic football dataset results	64
6.5.2	Heart disease (binary class problem) results	64
6.5.3	Heart disease (5-class problem) results	69
6.5.4	GCMS dataset: results and interpretation	71
6.6	Conclusion and future Work	75
7	Probabilistic Learning Vector Quantization	77
7.1	Introduction	77
7.2	Methodology	78
7.2.1	Exclusive KL divergence	78
7.2.2	Inclusive KL divergence	79
7.2.3	Probabilistic GMLVQ	79
7.2.4	Probabilistic Angle LVQ	83
7.3	Experiments on synthetic dataset and findings	87
7.3.1	Findings from the synthetic dataset	87
7.3.2	Findings from the GCMS dataset in terms of performance	89
7.3.3	Findings from the GCMS dataset in terms of interpretability	90
7.4	Conclusion and Future Work	94

8 Conclusion	95
8.1 Outlook	95
8.2 Future Work	97
Bibliography	99
Summary	105
Samenvatting	107

Acknowledgements

When I had come to Groningen in August 2014 for my Masters in Biomedical Engineering I was absolutely certain that I am not going to put myself through PhD. How could I make such a *mistake* after having watched the *PhDmovies*? But then Life happened: I got to know a friendly, funny and caring human named Prof. Dr. Michael Biehl, who taught the course *Modelling and Simulation*.

When I approached Michael about possible projects to learn Machine Learning (ML) he told me of a project which a former PhD student of his, Dr. Gert-Jan de Vries had for a Masters student, at Philips Research in Eindhoven. My master thesis project at Philips Research was when I first encountered Learning Vector Quantization and made a nerdy friend from Computing Science, Rick van Veen. It will sound cheesy but I fell for LVQ during this time. A few months after that I started my PhD journey with another of Michael's former PhD student, Dr. Kerstin Bunte as my daily supervisor, and Prof. Dr. Nicolai Petkov, along with Michael, as my promoters. Kerstin, you have not just been my daily supervisor in research, you have been an incredible emotional support at some of my worst moments during the last few years, and an inspiration.

During the very first week of PhD I got to know some of the people of the Intelligent Systems and the Distributed Systems, who influenced my life and personality in different ways and made me the current version of myself. They were Nicola, Estefania, Ahmad, Laura Fiorini, Xiaoxuan, Godliver, Jiapan and Astone, and Mohammad Babai. Also during that week I met Dr. Michael Wilkinson. When he told me that he made an extremely spicy sambal, I did not take him seriously. Being an Indian I used to take significant pride in my ability to *enjoy* spicy food and sauces. So I had thought to myself, *What does this European guy know about spicy food?* I took half a table spoon of his sambal and gulped it down. The rest was an episode of raging fire racing through my GI tract while I ran to get yoghurt. But then on the brighter

side, it was my first week of PhD and I was already on fire, thanks to Michael (W)'s sambal Satan. Michael (W) also introduced me to the Discworld series, and some other materials with dark humour.

Nicola, you've always given me some of the best advice, and taught me how to teach. Estefania, from the *girl who calls me for free food* you became my neighbour and thesis-support. You gave me the proper motivation and push for starting my thesis seriously, and after I shipped off my boxes to Denmark, you and Maik gave me a roof over my head and food and ensured that I could submit the thesis before leaving NL. Thanks you again Estefania for helping me prepare for interviews for jobs in NL when Denmark did not work out. I loved and will tremendously miss our impromptu dinners at your place. Maik and you have been my loveliest neighbours. Maik thank you for the Dutch translation of my thesis summary. Fiorini, you inspire me to read more and I will miss our discussions in office. Jiapan, Astone, and Cheng, you are the most adorable humans. The day I moved into Planetenlaan I would have been so lost without your warm hospitality. When I was preparing to leave Planetenlaan you both helped me again. And you both gave my painted tables a home as well and made me feel like a real artist.

During the next few years I met some other interesting and wonderful fellow PhD students and postdocs, such as Heerko, Laura from Spain, Mohammad*(1+i), Abolfazl, Michiel, Aleke, Caroline, Jiwoo, Simon, Maria, Hyoyin, Swaroop, Arash, Samira, Xueyi, Robert (Bwana), Henry, and Htet Htet from our wing of fifth floor of Bernoulliborg. Caroline, you are my best Friday dinner companion and my best (and only) frisbee mate. Abol and Michiel, I enjoyed your foosball fights. Mohammad, thank you so much for all the Mathematics trouble-shooting and discussion. Malihe, I will miss your chicken. Aleke and Jiwoo, you two are the kindest intelligent people I have ever met and you both, in your own ways, make me want to be a better human being. Heerko, thank you for introducing me to miniature painting and to South Park, and for our dark discussions about humans in general. Rick, thanks for teaching me \LaTeX back in 2016 and being there for me when I accidentally burnt my hand. Maria, thanks for your help during my move to Planetenlaan. It was great knowing you, Abel and Bilbo. Hyoyin, the Christmas trip to Slovenia was so much fun because of you.

Ahmad I will miss your dad jokes, bouldering with you, and sharing food with you. You give the best food compliments. Thanks for sharing your compositions with me and for introducing me to Buckethead. I am still waiting for our fusion musical piece together. Liz, my medical counterpart, this journey would have been extremely lame without you (from all perspectives). Without you, my prime data provider, I would not have had such an interesting project to work on. Working with you has been so much fun and that is one of the main reasons why I am sticking

around in academia a bit longer after PhD. You've been my academic significant other, despite the physical distance, and the voice of sanity and strength every time I gave up on myself and PhD. Ahmad and Liz, you both have been my emotional support humans throughout this journey. During some of my bad days it was chats with you both which helped me get out of bed, and face the day.

The people I want to thank now are the finance team, especially Martin Sanders, the Peregrine team, and the LWP team, especially Remco, Chris, Heiko, and Jurjen. Remco and Heiko I wish some day I will have your patience to tolerate someone's stupidity while simultaneously troubleshooting their technical issues, just like you both did with my mails. Without your support working from home would not have been possible. The lovely secretaries of Bernoulli Institute, you never failed to make us, the lost PhD students, smile even on a sad day. Ineke I will miss our chats. And now for my friends and supports from outside academia and immediate research group: Yoshita and Viraj, Laura, Ricardo, Akhil, Aniket, Kushagradhi, Suvam, Pratik kaku, Dipayan and Swati, Margherita, Pallab da, Rumki di, and my Googly-Mowgli, Pranjal, and Renger. Yoshita, without you and Laura, Masters would have been boring and life-less. You have been by my side through some of my dark days, be it a horribly burnt hand or a terribly broken heart. Kushagradhi, though we never properly spoke to each other in school you are one of the very few people I can now share my issues with despite the ocean between us. Thanks for being so understanding and non-judgmental despite being such a good listener. Suvam, thanks for being the annoying little brother VIT gave me. Pallab da and Rumki di thank you for driving me to Groningen from Schiphol back in August 2014 when I arrived in NL for Masters, and then welcoming me so warmly to your family. Google, you've been such a sweet nephew and student when I taught you origami and the falling techniques of jiu-jitsu. Renger, thanks for teaching me to stand up for myself while teaching me how to break a fall, and for making me love sports in the first place. Pranjal, whoever I am today I owe it partly to the wide range of interactions we've had over the years; thanks man, especially for being my comrade and defying "power" during Diwali 2014, the dinners at Winscho, and all the free legal advice. But I still hate your European Maya Sarabhai *chachi's* shop on Poelestraat. Margherita, you were the best next-door neighbours I have ever had, and you inspired me to pursue my non-academic interests.

Kerstin and Sebastian, Michael (W) and Meiny, Nicolai and Michael (B) you have taught me how to appreciate good food, whisky and wine. Meiny thank you so much for making me the Halloween skirt. I will miss the barbeque at your place. Sebastian, Kerstin, and Michael (B) thank you for teaching me to appreciate good art. And thanks to this PhD I met two other amazing humans, Prof. Dr. Peter Tino and Prof. Dr. Wiebke Arlt, my external supervisors from Birmingham. Peter, when

I see or talk to you I see the sweet Neville Longbottom from Harry Potter. You and Barbara (Hammer) explain complicated ideas so smoothly that it inspired me to try and become a science communicator. Your paintings and music inspire me to carry on my hobbies outside of academia. Wiebke, you are a power house of energy and enthusiasm. Thanks for being there for me whenever I reached out to you, for academic and non-academic reasons. Having strong female supervisors such as you and Kerstin, and such sweet male supervisors such as Michael and Peter shielded me from what I have heard some other female PhD students have had to tolerate. I aspire to become a human/scientist you will be proud to call your student. Kerstin, thank you for taking me for mud-walking back in 2017. It was the perfect analogy to what to expect from PhD. Also, thank you for teaching me kayaking at Mittweida, for being there for me when I needed support, and for pushing me to try harder. The four of you have made me a much better scientist, though I know there is still so much more to learn. Nicolai, thank you for arranging the Allersmaborg research retreats, the group dinners, and inviting us to bike trips around Groningen. Michael (B) you are the reason I got introduced to LVQ, Gert-Jan, Kerstin, Peter, Wiebke, Liz, the Intelligent Systems group, Villy and team, Barbara and team, and so many other wonderful people and started this PhD journey in the first place. I wish some day I can change the minds of young people towards research, the way you changed mine. Other amazing people I am grateful to have known over the last few years and will miss are Prof. Dr. Dimka Karastoyanova, Dr. Vasilios Andrikopoulos, Prof. Alexander Lazovik, Dr. George Azzopardi, and Dr. Charmaine Borg. Dimka, I miss our random chats in the corridor, and Vasilios I will miss your thundering laughter during post-exam times.

Now to mention those few people without whose disagreement during 2009-2010 I would not have been on this journey: My parents. I can only imagine how much you both have to put up against, for me and my career/life choices. You have inspired and supported me and provided me the science capital I needed to come this far. I still do not have an exact answer to when I can return to you physically and I will not dare to *thank* you for your patience and trust in me. Maa thanks for being tough on me and for being my best bitching partner and social filter. During my bleakest times in Denmark, it was you and Baba, Gurumaharaj from Prem Bihar, Mohammad*(1+i), Estefania, Suvam, Michael (W), Michael (B), Kerstin, Peter, Pratik kaku, Thamma, Ahmad, and Liz who provided me the support to keep fighting and trying. Dida, Dadu, and Dada, wherever you are you know that I think of you everyday. I would not have had the little golden windows of childhood if it were not for the three of you.

Sreejita Ghosh
Groningen

List of Figures

2.1	A colour-coded summary of the types of explainability and under which category each XAI technique fall. It is inspired from figure 6 of XAI:Concepts,Taxonomies, Opportunities and Challenges towards Responsible AI.	14
3.1	The three broad categories of missingness created on the exact same dataset of 20 dimensions and 900 instances.	21
3.2	Oversampling using SMOTE on a geodesic sphere.	30
4.1	Visualization of the 3 informative dimensions of the synthetic dataset (left panel) and the heatmaps for each class showing a subject-wise average of 60% missingness of type MNAR and MCAR per sample (right panel).	34
4.2	The heat-maps show the presence of systematic missing in each of the conditions of the GCMS dataset.	37
5.1	Overview of the classification error plots of 6 methods on the hold-out test set for missingness of type MCAR (top) and MNAR (bottom). The letters F and R mark the full and to 20% of the original size of the training set, while the number 1-3 indicates 0%, 30% and 60% missingness respectively.	48
5.2	Cumulative variance shows that the 57 and 100 dimensions together explain 95% and 99% of the variance of the dataset respectively. . . .	50
5.3	Influence of number of trees bagged for RF, on its performance on the GCMS dataset.	51

5.4	The blue parts of bars indicate how often a biomarker had appeared in the denominator for a condition, and the red parts indicate how often they were in the numerator.	54
5.5	Numbers 1-32 indicate the 32 biomarkers extracted from urine for investigation of the inborn disorders of steroidogenesis.	55
6.1	Football: 3 different views of a non-linearly separable synthetic dataset.	58
6.2	Three different perspectives of the nonlinear decision boundaries in the spherical classification space of LVQ^{2MA} trained on the Football dataset.	65
6.3	Eigenvalues of Λ across 5 folds and 5 runs. The box and whiskers show the mean and the 25th and 75th percentiles of amount of variance explained by each intrinsic dimension.	66
6.4	Feature relevances (top panel), as well as healthy and disease prototypes (bottom row) obtained by LVQ^{A1} on the binary HD classification.	68
6.5	Summary of the feature importance determined by RF over 5 folds and 5 runs, trained for the binary class problem. The boxes and whiskers denote the mean and the 25th and 75th percentiles.	68
6.6	Three example perspectives of the classification sphere depicting the decision boundaries as determined by LVQ^{2MA} on the 5-class HD problem.	71
6.7	Local relevances (left) and prototypes (right) of Healthy, and HD-patients of types 1-4, from LVQ^{LA} over 5 folds, for the 5-class problem. The box and whiskers plot denote the mean and 25th and 75th percentiles.	72
6.8	Global relevance, and prototypes (w) of Healthy and the 4 types of HD patients, from LVQ^{2MA^5} , over 5 folds, trained for the 5-class problem. The box and whiskers plot denote the mean and 25th and 75th percentiles.	73
6.9	The blue parts of bars indicate how often a biomarker was in the denominator for a condition, and the red parts indicate how often it was in the numerator.	74
6.10	LVQ^{LA} indicates which biomarkers are fingerprints for each condition.	75
7.1	Mollweide projection of 3D dataset for investigation of the effect of the value of Θ	86
7.2	Effect of the value of Θ in classification uncertainty. The \mathbf{o} represents the prototype of the class highlighted in each row.	86

7.3	Generalization abilities of $P_1LVQ^A, P_ELVQ^A, P_1LVQ^E$, LDA, KNN and RF.	89
7.4	The blue parts of bars indicate how often a biomarker had appeared in the denominator for a condition, and the red parts indicate how often they were in the numerator.	91
7.5	Numbers 1-32 indicate the 32 biomarkers extracted from urine for investigation of the inborn disorders of steroidogenesis.	92
7.6	Mollweide projection of the decision boundaries of the classifier, and where the data-points lie wrt these decision boundaries.	93

List of Tables

5.1	Selected average training T_{fN}^z and hold-out test HO_{fN}^z errors for fraction f number of training samples N and on average $z\%$ of missingness of type MNAR.	46
5.2	Selected average training T_{fN}^z and hold-out test HO_{fN}^z errors for fraction f number of training samples N and on average $z\%$ of missingness of type MCAR.	47
5.3	Experiments on the GCMS dataset.	51
5.4	Validation set performances of different classifiers on the GCMS dataset. The fourth and fifth blocks show the performances of the LVQ^A classifiers ensembled over $\eta = 5$ and $\eta = 100$ models/fold. The last block presents the performances of these classifiers geodesically averaged over $\eta = 100$ models/fold.	53
6.1	Experiments performed on the heart disease dataset.	63
6.2	Football comparison: mean performance (standard deviation)	65
6.3	Binary HD: mean performance (std) of global full rank ALVQ	66
6.4	Binary HD: mean performance (std) comparison.	67
6.5	5-class HD: comparison of ALVQ variants and RF.	70
6.6	Validation set performance of different classifiers on the GCMS dataset.	73
6.7	Validation set performance of ensembling and geodesic averaging over 100 models of LVQ^{LA} per fold, for fairer comparison with RF, on the GCMS dataset.	74

7.1	Selected average training T_{fN}^z and hold-out test HO_{fN}^z errors for fraction f number of training samples N and on average $z\%$ of missingness of type MNAR.	88
7.2	Selected average training T_{fN}^z and hold-out test HO_{fN}^z errors for fraction f number of training samples N and on average $z\%$ of missingness of type MCAR.	88
7.3	Validation set performance of different classifiers on the GCMS dataset.	90
7.4	Validation set performance of P_{ELVQ^A} as multiple learner. Upper block presents the performance of ensemble over $\eta = 100$ models per fold, and lower block presents performance of geodesic average of these models.	91

Preface

A significant fraction of us, the so-called well-read intellectuals like to question the age-old beliefs, social, cultural, and theological *norms*, and try to find the logic behind not just rituals, but also natural phenomena. Since the paleolithic age, or even before maybe, pre-historic hominids observed natural phenomena, questioned the *hows*, *whens*, and most importantly the *whys* of them, even with their “simpler” minds, and survived, consequently leading to us, the modern hominids-*homo sapiens*. Now I will fast-forward to a slightly nearer historic time, from the part of the Indian subcontinent which is my birth place. Before the nineteenth century, education was withheld from women so as to prevent them from becoming widows⁴⁰⁴, because that was what the *unquestionable* high priests of my parents and immediate ancestors’ religion dictated. But things changed because social reformers such as Raja Rammohan Roy and Ishwar Chandra Vidyasagar questioned those social norm and fought to change it. Therefore here I am, a woman of Indian origin, having the privilege of boring you with this thesis of more than hundred pages. It seems very obvious to us *now* that those social norms had to be questioned. Back then it was just a handful of people though who thought to, and then dared to question. But as a Machine Learning PhD student why am I rambling on about questioning norms? In recent times of such technological advancements some of us might be trusting *the algorithm* just like our distant ancestors (or some parts of our family and/or friends we consider lesser intellectuals than us) trust(ed) the so-called divine machinations and high priests. Why should we then consider ourselves *better* when we too are trusting *that which is inscrutable*? While Googling “Ramification of trusting what you cannot understand” I was trying to find some examples of adverse

⁴⁰⁴Why was that such thing of concern you ask? A woman who outlived her husband was supposed to be unlucky and sometimes even a demon, back in those times. That was the norm until it was questioned and fought against.

effects of using some black-box methods. Guess what or who I found in the results! Really, try Googling that and you will understand why I took you on this journey. If you do not want to, then this is the main message: Are we not hypocrites if we are ready to question the logic behind just the non-technology related norms and rituals, and not what *Deep Thought* or its simpler cousins tell us? How is an algorithm whose working logic we do not understand, predicting our probable financial future different from believing in psychics? A bit far-fetched? Four winters ago a tech supermarket's *algorithm* deemed me unsuitable for buying a telephoto lens with the equated monthly installment (EMI) scheme. I did not question the decision, because an algorithm *computed* it in a fancy manner. So yes, I am a hypocrite as well, but I am trying to change.