

# Interpretable Neural Networks using EAGGA

Simon Stürzebecher

SIMON.STUERZEBECHER@CAMPUS.LMU.DE

## Abstract

**Keywords:** tabular data, multi-objective optimization, deep learning

## 1 Introduction

Tabular Data - most common type of data - bad performance of deep learning models on tabular data, recent research suggests that heavy regularisation is necessary

regularisation - LM: ridge (L2) + lasso (L1) \* ridge: overcome multi-collinearity by pulling coeffs close to 0 \* lasso: variable selection by setting coeffs = 0 - NN: \* L2 regularisation implicitly included in SGD when using weight decay (pull params towards 0) - does this combat multi-collinearity as well or serve another purpose? \* L1 + key claim of DL: feature selection obsolete, network finds important ones itself + still, feature usually always has some impact (back this up), no explicit feature selection + explicit feature selection simply by not including a feature \* another classical NN regularisation technique: dropout against co-adaptation \* early stopping against overfitting when using iterative method - all of those require specifying Hyperparameters in advance

## 2 Background and Related Works

### 2.1 Interpretability

As there is no clear definition for interpretability, we will consider it as “the ability to provide *explanations* in *understandable terms* to a human”, where *explanations* are logical decision rules and *understandable terms* relate to commonly used terms in the domain of the problem, as suggested by Zhang et al. (2021, chap. 1). Further, we use the term “explainability” in an exchangeable manner with “interpretability”, as is commonly done.

(Why Interpretability is desirable?) - in many ways desirable, as per Zach (2019, pp. 3-4), e.g. \* **gaining trust**: most notable use case, e.g. medical diagnosis: predictions need to be validated by doctors (Antamis et al., 2024, p. 1) \* **model debugging** - model is optimised w.r.t. loss and judged on accuracy - might have unexpected performance drops in certain situations, makes model unreliable (Zhang et al., 2021, 1B) - having an interpretable model + model induction process can help find ways to improve accuracy + reliability \* **scientific understanding**: when only models can make sense of increasingly complex data anymore, interpretability enables extraction of learnt knowledge encoded in the model to make it accessible + reliable to humans (also mentioned here (Antamis et al., 2024, p. 1)) \* **subconscious bias**, e.g. loan approval: must ensure that decision is not discriminatory \* **regulatory** - such as EU’s ”Right to Explanation” warranted by the GDPR (Antamis et al., 2024, p. 1) or - drug approval processes in situations where a machine learning model

discovered a new drug, process needs to be transparent for the regulator to approve (Zhang et al., 2021, 1B)

Commonly, methods for model interpretation are divided into intrinsic methods, where the search space only comprises models with a structure simple enough to be considered “explainable” (such as tree-based or simple linear models), and post-hoc methods, where interpretation methods are applied after model training. Amongst post-hoc techniques, we can further divide the space into model-specific (such as analysing GLM coefficients) and model-agnostic (e.g. partial dependence plots, ALE) methods. Molnar (2022, chap. 3.2)

(Taxonomy) Zhang et al. (2021, chap. 2) extend this distinction to a 3-dimensional taxonomy, allowing for better categorisation of neural networks (NN), a model class that in its regular feedforward form is inherently non-interpretable. - **Passive vs Active Approaches** \* passive: post-hoc \* active: changing architecture or training process to make model interpretable - **Type of Explanations** \* examples: providing examples of what leads to desired output \* attribution: attributing effect on output for a specific feature \* hidden semantics: examine what type of inputs specific neurons / layers pick up on \* rules: logical rules, e.g. if-then, trees - **Local vs Global Interpretability** \* local: explanation based on individual samples \* semi-local: explanation based on set of samples, e.g. grouped together by some criterion \* global: explaining network as a whole

## 2.2 Hyperparameter Optimization (HPO)

- grid + random search - Bayesian optimisation - evolutionary algorithms \* CMA-ES \* differential evolution

## 2.3 Neural Architecture Search (NAS)

give background on NAS (e.g. one-shot, etc.), but for this project wanted to have a “minimal viable product” to test whether NNs can yield competitive performance on tabular ML + be interpretable using EAGGA at all before optimising the inner networks to the max

## 2.4 Multi-Objective Optimization

- problem: all aforementioned approaches are a-priori approaches -> we choose best HPs (weight decay, features used, dropout) w.r.t. performance - intro MOO

### 2.4.1 PARETO-OPTIMALITY

### 2.4.2 A-PRIORI

### 2.4.3 A-POSTERIORI

- NSGA-II (MOEA) \* non-dominated sorting \* crowding distance -> give intuition why non-dom-sort (asc) + crowding dist (desc) makes sense for selecting individuals

## 3 EAGGA

- for our problem (interpretability) we cannot use regular EA algorithms as we are also interested in reducing pairwise interactions + non-monotone feature effects - this creates

high-dim search space as shown in EAGGA paper - they provide a framework to efficiently traverse through a space that is reduced to only sensible HP configuration (e.g. interaction of A,B with A monot incr and B monot decr - couldn't be guaranteed - don't look at it in first place) - NSGA-II inspired evolutionary / genetic algorithm - introduces group structure for more efficient optimization over the entire search space

## 4 Extension to neural networks

- why extend it to NNs? \* NNs notoriously uninterpretable due to complex transformation of the feature space \* due to (lin alg) non-linear activation functions, interactions can be modelled \* no monotonicity guarantees - our approach (general algorithm) \* EAGGA algorithm implemented just as described in the paper \* HP init - NN total layers  $\in \{3, \dots, 10\}$ , i.e. hidden layers  $\in \{1, \dots, 8\}$ , init from trunc geom with  $p=0.5$  - NN nodes per hidden layer (for each group)  $\in \{3, \dots, 20\}$ , init from trunc geom with  $p=0.5$  - NN dropout %  $\in [0, 1]$ , init from trunc gamma with shape=2, scale=0.15 \* group structure init - feature detector \* as described in paper \* but instead of fitting 10 trees + taking rel. # of features used as  $p$  for trunc geom, we simply use 0.5 (sklearn dectree examination not straightforward) \* also in preliminary experiments we found the sampled # of features to be used occasionally to be # of non-0 values in normalised information gain filter (e.g. if former is = total # of features and one feature is indep. from target) in these edge cases we only use the features with non-0 filter values - interaction detector \* as described in paper \* same as for feat detect: use  $p=0.5$  to sample # of interactions used instead of 10 dec trees \* also for FAST algorithm don't use RSS (lin reg metric) but mean accuracy, as we fit log reg - monotonicity detector \* as described in paper \* use default HPs (max depth 30, minsplie = 20) of mlr3 classification tree implementation, as this is the library the paper uses - both interaction + monotonicity detectors fit their models on 80% of train split from holdout + eval on remaining 20% (implementation details) \* group structures + datasets - group structures implemented as described in the paper with additional list encoding sign of monotonicity of the individual features as detected by the monotonicity detector - included features are passed to dataset, dataset outputs included features, multiplied by the features' individual monotonicity (-1 / 1) - entire group structure is passed to NN, where architecture is built accordingly \* neural network - "instead of XGB as in the original paper, we apply the method to NNs to examine whether this type of regularisation can achieve interpretability on NNs while outperforming EAGGA on XGB (original paper), ..., this requires special architecture, etc. pp." - NN \* comprises of "sub-NNs", one for each equivalence relation - basically non-fully connected NN \* hidden layers use ReLU activation and dropout afterwards, then a shared output layer with concatenated sub-NNs' outputs as input and sigmoid activation (implicit in the loss function for better numerical stability, NN itself outputs "logits", which refers to pre-activation output in pytorch) - feature sparsity thus achieved by only training on included feature groups - goes somewhat against of deep learning where model is supposed to decide itself, which feature to "use" / put importance on - ELABORATE - feature interaction achieved by grouping, max-operation in ReLU induces interaction effect (different kind of interaction than e.g. in LM with multiplication) - equation why the interacting features need to be grouped together when using ReLU, cf. photos - monotonicity constraint achieved by clipping weights to  $[0, \infty)$  for restricted

equivalence relations (monotonic decrease achieved via dataset object multiplying features with their individual signs), bias clipping not necessary (constant additive term) -  $\downarrow$  equation how monotonicity is enforced with this \* evaluation, holdout, cv, early stopping - evaluation via dominated hypervolume along AUC-ROC, NF, NI, NNM as defined in original paper \* NF simply rel. # of included features \* NI = sum of all possible pairwise interactions in each group over all possible pairwise interactions among all features =  $\frac{\sum_g \binom{p_g}{2}}{\binom{p}{2}}$  \* NNM = rel. # unconstrained features - outer holdout split: 2/3 train, 1/3 test (as in paper) \* run EAGGA on holdout train split, i.e. train + select best  $\mu = 100$  individuals based on non-dom-sorting (ascending) + crowding distance (descending) as tie breaker - inner CV split on outer train portion: 5-fold (as in paper) \* in each fold fit model on 80% of CV-train portion, use remaining 20% of CV-train for early stopping - early stopping criterion \* for each fold, always train for min 200 epochs, keep track of model with lowest loss \* after that, use patience of 100: if current model's loss is  $\downarrow$  mean of last 100 epochs' losses \* if no early stopping, train each fold for max 10 minutes \* then stop training, return model with lowest loss -  $\downarrow$  add graphic of entire dataset split - then evaluate on last fold from CV + keep best  $\mu$  individuals + generate  $\lambda = 10$  offspring for new generation + repeat \* hardware: training on Sagemaker Notebook instance - initial experiments on ml.g4dn.xlarge instance (2 vCPUs, 16GiB RAM, 1 NVIDIA T4, cf. <https://aws.amazon.com/de/ec2/instance-types/g4/>) using cuda not much faster than on ml.t3.medium (2 Intel Xeon 8000 vCPUs, 4GiB RAM, cf. <https://aws.amazon.com/de/ec2/instance-types/t3/>) - thus decided for more economic + ressourcen-schonend t3.medium \* known bugs - in rare cases (anecdotally every 2 datasets), gga\_mutate seems to be generating -1 as monotonicity attribute, despite np.random.randint(low=0, high=2, size=1) \* loop crashes at group\_structure creation, for these cases subtract previous runtime from 8hrs + load from last generation in output via load\_population \* so far only happened for madeline in after gen-8.json after 5hrs 45mins (at start of gen-9, which wasn't exported yet) -  $\downarrow$  loaded this and ran for another 2hrs 15mins

## 5 Experimental Results and Discussion

## 6 Conclusion and Future Outlook

- MO BO on group structure space possible? - here is a citation Schneider et al. (2023).

## Appendix A. Software used

for implementation we used openml Vanschoren et al. (2014), Feurer et al. (2021), numpy Harris et al. (2020), pandas pandas development team (2023), Wes McKinney (2010), pytorch Ansel et al. (2024), scikit-learn Pedregosa et al. (2011), scipy Virtanen et al. (2020), pymoo Blank and Deb (2020), and tqdm da Costa-Luis et al. (2024)

## Appendix B.

In this appendix we prove the following theorem from Section 6.2:

**Theorem** *Let  $u, v, w$  be discrete variables such that  $v, w$  do not co-occur with  $u$  (i.e.,  $u \neq 0 \Rightarrow v = w = 0$  in a given dataset  $\mathcal{D}$ ). Let  $N_{v0}, N_{w0}$  be the number of data points for which  $v = 0, w = 0$  respectively, and let  $I_{uv}, I_{uw}$  be the respective empirical mutual information values based on the sample  $\mathcal{D}$ . Then*

$$N_{v0} > N_{w0} \Rightarrow I_{uv} \leq I_{uw}$$

with equality only if  $u$  is identically 0. ■

## Appendix C.

**Proof.** We use the notation:

$$P_v(i) = \frac{N_v^i}{N}, \quad i \neq 0; \quad P_{v0} \equiv P_v(0) = 1 - \sum_{i \neq 0} P_v(i).$$

These values represent the (empirical) probabilities of  $v$  taking value  $i \neq 0$  and 0 respectively. Entropies will be denoted by  $H$ . We aim to show that  $\frac{\partial I_{uv}}{\partial P_{v0}} < 0 \dots$

*Remainder omitted in this sample. See <http://www.jmlr.org/papers/> for full paper.*

## References

Jason Ansel, Edward Yang, Horace He, Natalia Gimelshein, Animesh Jain, Michael Voznesensky, Bin Bao, Peter Bell, David Berard, Evgeni Burovski, Geeta Chauhan, Anjali Chourdia, Will Constable, Alban Desmaison, Zachary DeVito, Elias Ellison, Will Feng, Jiong Gong, Michael Gschwind, Brian Hirsh, Sherlock Huang, Kshiteej Kalambarkar, Laurent Kirsch, Michael Lazos, Mario Lezcano, Yanbo Liang, Jason Liang, Yinghai Lu, C. K. Luk, Bert Maher, Yunjie Pan, Christian Puhersch, Matthias Reso, Mark Saroufim, Marcos Yukio Siraichi, Helen Suk, Shunting Zhang, Michael Suo, Phil Tillet, Xu Zhao, Eikan Wang, Keren Zhou, Richard Zou, Xiaodong Wang, Ajit Mathews, William Wen, Gregory Chanan, Peng Wu, and Soumith Chintala. Pytorch 2: Faster machine learning through dynamic python bytecode transformation and graph compilation. In *Proceedings of the 29th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 2*, ASPLOS '24, page 929–947, New York, NY, USA, 2024. Association for Computing Machinery. ISBN 9798400703850. doi: 10.1145/3620665.3640366. URL <https://doi.org/10.1145/3620665.3640366>.

- Thanasis Antamis, Anastasis Drosou, Thanasis Vafeiadis, Alexandros Nizamis, Dimosthenis Ioannidis, and Dimitrios Tzovaras. Interpretability of deep neural networks: A review of methods, classification and hardware. *Neurocomputing*, 601:128204, 2024. ISSN 0925-2312. doi: <https://doi.org/10.1016/j.neucom.2024.128204>. URL <https://www.sciencedirect.com/science/article/pii/S0925231224009755>.
- J. Blank and K. Deb. pymoo: Multi-objective optimization in python. *IEEE Access*, 8: 89497–89509, 2020.
- Casper da Costa-Luis, Stephen Karl Larroque, Kyle Altendorf, Hadrien Mary, richard-sheridan, Mikhail Korobov, Noam Yorav-Raphael, Ivan Ivanov, Marcel Bargull, Nishant Rodrigues, Shawn, Mikhail Dektyarev, Michał Górny, mjstevens777, Matthew D. Pagel, Martin Zugnoni, JC, CrazyPython, Charles Newey, Antony Lee, pgajdos, Todd, Staffan Malmgren, redbug312, Orivej Desh, Nikolay Nechaev, Mike Boyle, Max Nordlund, MapleCCC, and Jack McCracken. tqdm: A fast, extensible progress bar for python and cli, November 2024. URL <https://doi.org/10.5281/zenodo.14231923>.
- Matthias Feurer, Jan N. van Rijn, Arlind Kadra, Pieter Gijsbers, Neeratyoy Mallik, Sahithya Ravi, Andreas Müller, Joaquin Vanschoren, and Frank Hutter. Openml-python: an extensible python api for openml. *Journal of Machine Learning Research*, 22(100):1–5, 2021. URL <http://jmlr.org/papers/v22/19-920.html>.
- Charles R. Harris, K. Jarrod Millman, Stéfan J. van der Walt, Ralf Gommers, Pauli Virtanen, David Cournapeau, Eric Wieser, Julian Taylor, Sebastian Berg, Nathaniel J. Smith, Robert Kern, Matti Picus, Stephan Hoyer, Marten H. van Kerkwijk, Matthew Brett, Allan Haldane, Jaime Fernández del Río, Mark Wiebe, Pearu Peterson, Pierre Gérard-Marchant, Kevin Sheppard, Tyler Reddy, Warren Weckesser, Hameer Abbasi, Christoph Gohlke, and Travis E. Oliphant. Array programming with NumPy. *Nature*, 585(7825):357–362, September 2020. doi: 10.1038/s41586-020-2649-2. URL <https://doi.org/10.1038/s41586-020-2649-2>.
- Christoph Molnar. *Interpretable Machine Learning*. 2 edition, 2022. URL <https://christophm.github.io/interpretable-ml-book>.
- The pandas development team. pandas-dev/pandas: Pandas, January 2023. URL <https://doi.org/10.5281/zenodo.7549438>.
- F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- Lennart Schneider, Bernd Bischl, and Janek Thomas. Multi-objective optimization of performance and interpretability of tabular supervised machine learning models. In *Proceedings of the Genetic and Evolutionary Computation Conference, GECCO '23*, page 538–547, New York, NY, USA, 2023. Association for Computing Machinery. ISBN 9798400701191. doi: 10.1145/3583131.3590380. URL <https://doi.org/10.1145/3583131.3590380>.

- Joaquin Vanschoren, Jan N. van Rijn, Bernd Bischl, and Luis Torgo. Openml: networked science in machine learning. *SIGKDD Explor. Newsl.*, 15(2):49–60, June 2014. ISSN 1931-0145. doi: 10.1145/2641190.2641198. URL <https://doi.org/10.1145/2641190.2641198>.
- Pauli Virtanen, Ralf Gommers, Travis E. Oliphant, Matt Haberland, Tyler Reddy, David Cournapeau, Evgeni Burovski, Pearu Peterson, Warren Weckesser, Jonathan Bright, Stéfan J. van der Walt, Matthew Brett, Joshua Wilson, K. Jarrod Millman, Nikolay May-  
orov, Andrew R. J. Nelson, Eric Jones, Robert Kern, Eric Larson, C J Carey, İlhan Polat, Yu Feng, Eric W. Moore, Jake VanderPlas, Denis Laxalde, Josef Perktold, Robert Cimr-  
man, Ian Henriksen, E. A. Quintero, Charles R. Harris, Anne M. Archibald, Antônio H. Ribeiro, Fabian Pedregosa, Paul van Mulbregt, and SciPy 1.0 Contributors. SciPy 1.0: Fundamental Algorithms for Scientific Computing in Python. *Nature Methods*, 17:261–272, 2020. doi: 10.1038/s41592-019-0686-2.
- Wes McKinney. Data Structures for Statistical Computing in Python. In Stéfan van der Walt and Jarrod Millman, editors, *Proceedings of the 9th Python in Science Conference*, pages 56 – 61, 2010. doi: 10.25080/Majora-92bf1922-00a.
- J. Zach. Interpretability of deep neural networks. 2019. URL <https://api.semanticscholar.org/CorpusID:198979458>.
- Yu Zhang, Peter Tiño, Aleš Leonardis, and Ke Tang. A survey on neural network inter-  
pretability. *IEEE Transactions on Emerging Topics in Computational Intelligence*, 5(5):  
726–742, 2021. doi: 10.1109/TETCI.2021.3100641.