

# Gradient-based Optimization for Compact and Explainable Fuzzy Rule-based Classification

Javier Fumanal-Idocin<sup>\*1</sup>, Raquel Fernandez-Peralta<sup>2</sup>, and Javier Andreu-Perez<sup>1</sup>

<sup>1</sup>School of Computer Science and Electronic Engineering, University of Essex

<sup>2</sup>Mathematical Institute, Slovak Academy of Sciences, Bratislava, Slovakia

j.fumanal-idocin@essex.ac.uk, raquel.fernandez@mat.savba.sk, j.andreu-perez@essex.ac.uk

## Abstract

Rule-based models are valued in high-stakes decision-making for their transparency, but their discrete nature limits optimization and scalability. We propose the Fuzzy Rule-based Reasoner (FRR), a gradient-based rule learner that enforces user-defined complexity constraints while maintaining strong performance. FRR combines interpretable fuzzy logic partitions with sufficient (single-rule) decision-making, avoiding the combinatorial growth of additive ensembles. Across 40 datasets, FRR outperforms traditional rule-based methods (by about 5% over RIPPER), matches the accuracy of tree-based models like CART with rule bases 90% smaller, and achieves 96% of the accuracy of additive rule-based models while using only 3% of their rule base size.

## 1 Introduction

Deep neural networks excel in handling large volumes of unstructured data, such as images and videos [1], but their lack of transparency limits their use in high-stakes domains like medicine and finance [2]. Rule-based algorithms, by contrast, are inherently interpretable, as they reveal explicit decision patterns and their relevance, offering more trustworthy explanations than many post-hoc XAI methods [3, 4]. Studying learned rules allows practitioners to validate or challenge model insights and even use them to approximate complex models like deep networks [5, 6]. However, rule-based classifiers often face a trade-off between interpretability and accuracy: larger rule sets improve performance but reduce transparency [7]. While fuzzy logic and genetic optimization have been used to balance this trade-off [8], they struggle to scale, and gradient-based approaches [9–13] often produce overly complex rule bases. To address these issues, we propose the Fuzzy Rule-Based Reasoner (FRR), a fully differentiable rule-based classifier that integrates user-defined complexity constraints (maximum rules and conditions per rule) with interpretable fuzzy partitions, achieving a balance between performance, simplicity, and transparency.

<sup>\*</sup>Corresponding Author.

## 2 Fuzzy Rule-Based Reasoner

The FRR is a hierarchical model composed of four layers of matrix operations that mimic fuzzy logical reasoning (Figure 1). It receives an input vector  $X_i$  and produces a class prediction through a sequence of fuzzification, inference, and decision steps.

- Fuzzification layer:** The fuzzification layer transforms input features into interpretable fuzzy membership values. Continuous variables are mapped to linguistic terms such as “low,” “medium,” and “high” using trapezoidal fuzzy sets defined according to the data’s quantile distribution, while categorical variables are encoded using a one-hot representation.
- Logic inference layers:** The logic inference process in the FRR has two main steps. First, for each feature, the model selects the most relevant fuzzy label, the one with the highest weight, representing which linguistic partition is activated. Next, it determines which features form the rule’s antecedent, keeping a fixed number of conditions per rule. The truth value of a rule is then computed as the product of the selected condition contributions, ensuring the result remains within the  $[0, 1]$  range.
- Decision layer:** Implements sufficient-rule prediction by selecting the consequent of the highest-scoring activated rule.

Beyond these layers, the FRR also incorporates a parsimony mechanism that automatically prunes irrelevant rule conditions during training. Through a competitive cancellation process, conditions that contribute little to the prediction are deactivated and penalized via a regularization term, ensuring compact and efficient rule bases without sacrificing accuracy.

## 3 Training

Training the FRR addresses three main challenges: non-differentiability, gradient sparsity, and vanishing gradients. The non-differentiable arg max used for rule and feature selection is approximated with the

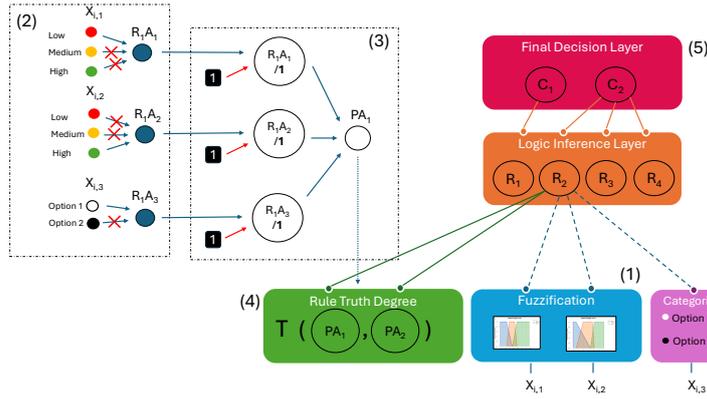


Figure 1. FRR scheme, using an input  $X_i$  with three features, four rules, and two target classes.

Table 1. 5-fold results for all the datasets considered.

Method	Sufficient Rule-based					Tree-based			Additive Rule-based		LR	GB
	FRR	FGA	DRNet	DINA	RIPPER	CART	C4.5	GOT	SIRUS	RRL		
Accuracy	79.51	70.46	56.08	54.99	75.22	81.06	79.99	76.91	82.17	81.99	82.12	86.04
Number of Rules	13.77	7.12	24.04	18.48	16.04	39.75	131.92	5.23	286.71	99.35	-	-
Conditions/Rule	1.94	2.23	6.37	4.59	1.96	5.75	8.10	2.27	2.90	8.85	-	-
Rule base Size	26.71	15.87	153.13	84.82	31.43	228.56	1068.55	11.87	831.45	879.24	-	-
Unique Conditions	10.78	10.52	16.26	9.18	21.30	34.72	68.56	11.00	357.05	125.16	-	-

084 Straight-Through Estimator (STE) [14, 15], enabling  
 085 gradient flow through discrete choices. To reduce  
 086 gradient sparsity, a relaxed indicator function with  
 087 parameter  $\beta$  allows partial updates to non-selected  
 088 features;  $\beta$  gradually decreases during training to  
 089 recover discrete behavior. Finally, to combat van-  
 090 ishing gradients from multiplicative rule inference,  
 091 two strategies are applied: root-normalized activa-  
 092 tion, which rescales small truth values, and residual  
 093 connections [16], which preserve gradient flow early  
 094 in training and fade out over time. Together, these  
 095 mechanisms ensure stable and efficient optimization  
 096 of interpretable fuzzy rules.

## 097 4 Experiments

098 We evaluated the proposed Fuzzy Rule-based Rea-  
 099 soner (FRR) across 40 widely used classification  
 100 datasets, ranging from 80 to 19,020 samples and 2  
 101 to 85 features, using 5-fold cross-validation and ac-  
 102 curacy as the primary metric. Statistical differences  
 103 between classifiers were assessed using the Friedman  
 104 Test and Post-hoc Nemenyi procedure [17]. The  
 105 FRR was compared against three groups of base-  
 106 lines: (i) rule-based methods, including the Fuzzy  
 107 Genetic Algorithm classifier (FGA) [18], RIPPER  
 108 [19], SIRUS [20], DRNet [9], DINA [21], and Rule-  
 109 based Representation Learning (RRL) [22]; (ii) tree-  
 110 based models, such as CART [23], C4.5 [24], and the  
 111 Generalized Optimal Sparse Decision Tree (GOT)  
 112 [25]; and (iii) non-rule-based classifiers, namely Lo-  
 113 gistic Regression (LR) and Gradient Boosting (GB)  
 114 [26]. While GB achieved the highest average accu-

115 racy (86.04%), FRR maintained competitive perfor-  
 116 mance (79.51%) with significantly lower complexity.  
 117 Indeed, its rule base was only 3% the size of SIRUS  
 118 and 11% that of CART. FRR was statistically supe-  
 119 rior to other sufficient rule-based classifiers, particu-  
 120 larly RIPPER, and demonstrated scalability across  
 121 datasets of varying dimensionality. Gradient-based  
 122 rule learners (DRNet and DINA) performed worse,  
 123 likely due to their sensitivity to hyperparameters  
 124 and smaller dataset sizes. Overall, FRR achieved  
 125 a strong balance between interpretability and accu-  
 126 racy, offering an efficient and explainable alternative  
 127 to traditional and gradient-based rule learning meth-  
 128 ods.

## 129 5 Conclusion

130 We introduced the Fuzzy Rule-based Reasoner  
 131 (FRR), an explainable classifier that learns inter-  
 132 pretable rules through gradient-based optimization  
 133 while allowing users to control model complexity  
 134 by setting limits on rule count and length. This  
 135 design maintains interpretability without sacrificing  
 136 performance, achieving a strong balance between  
 137 accuracy and simplicity. FRR outperforms other  
 138 gradient-based rule learners and reaches accuracy  
 139 comparable to tree-based models with far lower com-  
 140 plexity. Future work will integrate FRR into deep  
 141 learning frameworks to provide rule-based expla-  
 142 nations within gradient flows and explore how its  
 143 learned rules relate to epistemic and aleatoric uncer-  
 144 tainty.

145 **References**

- 146 [1] Y. LeCun, Y. Bengio, and G. Hinton. “Deep  
147 learning”. In: *Nature* 521.7553 (2015), pp. 436–  
148 444. DOI: [10.1038/nature14539](https://doi.org/10.1038/nature14539).
- 149 [2] A. B. Arrieta, N. Díaz-Rodríguez, J. Del Ser,  
150 A. Bennetot, S. Tabik, A. Barbado, S. García,  
151 S. Gil-López, D. Molina, R. Benjamins, et  
152 al. “Explainable Artificial Intelligence (XAI):  
153 Concepts, taxonomies, opportunities and chal-  
154 lenges toward responsible AI”. In: *Information*  
155 *fusion* 58 (2020), pp. 82–115. DOI: [10.1016/  
156 j.inffus.2019.12.012](https://doi.org/10.1016/j.inffus.2019.12.012).
- 157 [3] R. Tomsett, D. Harborne, S. Chakraborty, P.  
158 Gurram, and A. Preece. “Sanity checks for  
159 saliency metrics”. In: *Proceedings of the AAAI*  
160 *conference on artificial intelligence*. Vol. 34. 04.  
161 2020, pp. 6021–6029. DOI: [10.5555/3327546.  
162 3327621](https://doi.org/10.5555/3327546.3327621).
- 163 [4] C. Rudin, C. Chen, Z. Chen, H. Huang, L.  
164 Semenova, and C. Zhong. “Interpretable ma-  
165 chine learning: Fundamental principles and  
166 10 grand challenges”. In: *Statistic Surveys* 16  
167 (2022), pp. 1–85. DOI: [10.1214/21-SS133](https://doi.org/10.1214/21-SS133).
- 168 [5] J. Fumanal-Idocin, J. Andreu-Perez, O. Cord,  
169 H. Hagra, H. Bustince, et al. “Artxai: Explain-  
170 able artificial intelligence curates deep repre-  
171 sentation learning for artistic images using  
172 fuzzy techniques”. In: *IEEE Transactions on*  
173 *Fuzzy Systems* (2023). DOI: [10.1109/TFUZZ.  
174 2023.3337878](https://doi.org/10.1109/TFUZZ.2023.3337878).
- 175 [6] R. Li, Q. Li, Y. Zhang, D. Zhao, Y. Jiang, and  
176 Y. Yang. “Interpreting unsupervised anomaly  
177 detection in security via rule extraction”. In:  
178 *Advances in Neural Information Processing*  
179 *Systems* 36 (2024). DOI: [10.5555/3666122.  
180 3668840](https://doi.org/10.5555/3666122.3668840).
- 181 [7] L. Breiman. “Random forests”. In: *Machine*  
182 *learning* 45 (2001), pp. 5–32. DOI: [10.1023/A:  
183 1010950718922](https://doi.org/10.1023/A:1010950718922).
- 184 [8] J. Alcalá-Fdez, R. Alcalá, and F. Herrera.  
185 “A fuzzy association rule-based classification  
186 model for high-dimensional problems with ge-  
187 netic rule selection and lateral tuning”. In:  
188 *IEEE Transactions on Fuzzy systems* 19.5  
189 (2011), pp. 857–872. DOI: [10.1109/TFUZZ.  
190 2011.2147794](https://doi.org/10.1109/TFUZZ.2011.2147794).
- 191 [9] L. Qiao, W. Wang, and B. Lin. “Learning  
192 accurate and interpretable decision rule sets  
193 from neural networks”. In: *Proceedings of the*  
194 *AAAI Conference on Artificial Intelligence*.  
195 Vol. 35. 5. 2021, pp. 4303–4311. DOI: [10.5555/  
196 AAI27997821](https://doi.org/10.5555/AAI27997821).
- [10] W. Zhang, Y. Liu, Z. Wang, and J. Wang. 197  
“Learning to Binarize Continuous Features 198  
for Neuro-Rule Networks”. In: *IJCAI*. 2023, 199  
pp. 4584–4592. DOI: [10.24963/ijcai.2023/  
200 510](https://doi.org/10.24963/ijcai.2023/510). 201
- [11] Z. Wang, W. Zhang, N. Liu, and J. Wang. 202  
“Learning interpretable rules for scalable data 203  
representation and classification”. In: *IEEE*  
204 *Transactions on Pattern Analysis and Machine*  
205 *Intelligence* (2024). DOI: [10.1109/TPAMI.  
206 2023.3328881](https://doi.org/10.1109/TPAMI.2023.3328881). 207
- [12] S. Xu, N. P. Walter, J. Kalofolias, and J. 208  
Vreeken. “Learning Exceptional Subgroups by 209  
End-to-End Maximizing KL-Divergence”. In:  
210 *International Conference on Machine Learn-*  
211 *ing*. PMLR. 2024, pp. 55267–55285. DOI: [10.  
212 5555/3692070.3694348](https://doi.org/10.5555/3692070.3694348). 213
- [13] Y. Yang, W. Ren, and S. Li. “Hyperlogic: 214  
Enhancing diversity and accuracy in rule 215  
learning with hypernets”. In: *Advances in*  
216 *Neural Information Processing Systems* 37  
217 (2024), pp. 3564–3587. DOI: [10.5555/3737916.  
218 3738034](https://doi.org/10.5555/3737916.3738034). 219
- [14] P. Yin, J. Lyu, S. Zhang, S. J. Osher, Y. Qi, 220  
and J. Xin. “Understanding Straight-Through 221  
Estimator in Training Activation Quantized 222  
Neural Nets”. In: *International Conference on*  
223 *Learning Representations*. 2019. 224
- [15] M. Schoenbauer, D. Moro, L. Lew, and A. 225  
Howard. *Custom Gradient Estimators are*  
226 *Straight-Through Estimators in Disguise*. 2024.  
227 DOI: [10.48550/arXiv.2405.05171](https://doi.org/10.48550/arXiv.2405.05171). arXiv:  
228 2405.05171 [cs.LG]. 229
- [16] K. He, X. Zhang, S. Ren, and J. Sun. “Deep 230  
residual learning for image recognition”. In:  
231 *Proceedings of the IEEE conference on com-*  
232 *puter vision and pattern recognition*. 2016,  
233 pp. 770–778. DOI: [10.1109/CVPR.2016.90](https://doi.org/10.1109/CVPR.2016.90). 234
- [17] J. Demšar. “Statistical Comparisons of Classi- 235  
fiers over Multiple Data Sets”. In: *Journal of*  
236 *Machine Learning Research* 7 (2006), pp. 1–  
237 30. ISSN: 1532-4435. DOI: [10.5555/1248547.  
238 1248548](https://doi.org/10.5555/1248547.1248548). 239
- [18] J. Fumanal-Idocin and J. Andreu-Perez. “Ex- 240  
Fuzzy: A library for symbolic explainable AI 241  
through fuzzy logic programming”. In: *Neuro-*  
242 *computing* (2024), p. 128048. ISSN: 0925-2312.  
243 DOI: [10.1016/j.neucom.2024.128048](https://doi.org/10.1016/j.neucom.2024.128048). 244
- [19] W. W. Cohen et al. “Fast effective rule induc- 245  
tion”. In: *Proceedings of the twelfth interna-*  
246 *tional conference on machine learning*. 1995,  
247 pp. 115–123. DOI: [10.1016/B978-1-55860-  
248 377-6.50023-2](https://doi.org/10.1016/B978-1-55860-377-6.50023-2). 249

- 250 [20] C. B enard, G. Biau, S. Da Veiga, and E. Scornet. “Interpretable random forests via rule ex-  
251 traction”. In: *International conference on arti-*  
252 *ficial intelligence and statistics*. PMLR. 2021,  
253 pp. 937–945.  
254
- 255 [21] N. P. Walter, J. Fischer, and J. Vreeken.  
256 “Finding interpretable class-specific patterns  
257 through efficient neural search”. In: *Proceed-*  
258 *ings of the AAAI Conference on Artificial In-*  
259 *telligence*. Vol. 38. 8. 2024, pp. 9062–9070. DOI:  
260 [10.1609/aaai.v38i8.28756](https://doi.org/10.1609/aaai.v38i8.28756).
- 261 [22] Z. Wang, W. Zhang, N. Liu, and J. Wang.  
262 “Scalable rule-based representation learning  
263 for interpretable classification”. In: *Advances*  
264 *in Neural Information Processing Systems*  
265 *34* (2021), pp. 30479–30491. DOI: [10.5555/](https://doi.org/10.5555/3540261.3542593)  
266 [3540261.3542593](https://doi.org/10.5555/3540261.3542593).
- 267 [23] R. Timofeev. “Classification and regression  
268 trees (CART) theory and applications”. In:  
269 *Humboldt University, Berlin* 54 (2004).
- 270 [24] J. R. Quinlan. *C4.5: Programs for Machine*  
271 *Learning*. Morgan Kaufmann, 1993. DOI: [10.](https://doi.org/10.1007/BF00993309)  
272 [1007/BF00993309](https://doi.org/10.1007/BF00993309).
- 273 [25] H. McTavish, C. Zhong, R. Achermann, I.  
274 Karimalis, J. Chen, C. Rudin, and M. Seltzer.  
275 “Fast sparse decision tree optimization via refer-  
276 ence ensembles”. In: *Proceedings of the AAAI*  
277 *conference on artificial intelligence*. Vol. 36.  
278 9. 2022, pp. 9604–9613. DOI: [10.1609/aaai.](https://doi.org/10.1609/aaai.v36i9.21194)  
279 [v36i9.21194](https://doi.org/10.1609/aaai.v36i9.21194).
- 280 [26] J. H. Friedman. “Greedy function approxima-  
281 tion: a gradient boosting machine”. In: *Annals of statistics* (2001), pp. 1189–1232. DOI:  
282 [10.1214/aos/1013203451](https://doi.org/10.1214/aos/1013203451).  
283