

# Embedded feature selection using dual-network architecture

Abderrahim Abbassi<sup>a,\*</sup>, Arved Dörpinghaus<sup>a</sup>, Niklas Römgens<sup>a</sup>, Tanja Griebmann<sup>a</sup>,  
Raimund Rolfes<sup>a</sup>

<sup>a</sup>*Leibniz University Hannover, Institute of Structural Analysis, Appelstraße 9A, 30167 Hanover, Germany*

---

## 1 Abstract

2 Feature selection is essential for eliminating noise, reducing redundancy, simpli-  
3 fying computational complexity, and lowering data collection and processing costs.  
4 However, existing methods often face challenges due to the complexity of feature  
5 interdependencies, uncertainty regarding the exact number of relevant features, and  
6 the need for hyperparameter optimization, which increases methodological complex-  
7 ity.

8 This research proposes a novel dual-network architecture for feature selection  
9 that addresses these issues. The architecture consists of a task model and a selection  
10 model. First, redundant features are fed into the selection model, which generates  
11 a binary mask aligned with the input feature dimensions. This mask is applied to  
12 a shifted version of the original features, serving as input to the task model. The  
13 task model then uses the selected features to perform the target supervised task.  
14 Simultaneously, the selection model aims to minimize the cumulative value of the  
15 mask, thus selecting the most relevant features with minimal impact on the task  
16 model's performance.

17 The method is evaluated using benchmark and synthetic datasets across different  
18 supervised tasks. Comparative evaluation with state-of-the-art techniques demon-  
19 strates that the proposed approach exhibits superior or competitive feature selection  
20 capabilities, achieving a reduction of 90% or more in feature count. This is par-  
21 ticularly notable in the presence of non-linear feature interdependencies. The key

---

\*Corresponding author

*Email addresses:* `a.abbassi@isd.uni-hannover.de` (Abderrahim Abbassi),  
`arved.doerpinghaus@stud.uni-hannover.de` (Arved Dörpinghaus),  
`n.roemgens@isd.uni-hannover.de` (Niklas Römgens), `t.griessmann@isd.uni-hannover.de`  
(Tanja Griebmann), `r.rolfes@isd.uni-hannover.de` (Raimund Rolfes)

advantages of the proposed method are its ability to self-determine the number of relevant features needed for the supervised task and its simplicity, requiring the pre-definition of only a single hyperparameter, for which an estimation approach is suggested.

**Keywords:** Feature selection, dual-network architecture, feature redundancy, non-linear interdependencies

---

## 1. Introduction

Due to the rising versatility of machine learning applications and the increasing intricacy of the tasks they accomplish, the underlying datasets are constantly growing in size and complexity. Often, data are gathered without the specific model and job in mind, resulting in vast amounts of data being collected and stored. These abundant quantities induce demand for methods that select relevant variables for a given predictive task and filter out correlations and dependencies. These variables are referred to as *features*, and methods for selecting relevant subsets perform *feature selection* (Kira and Rendell, 1992; Li et al., 2017; Kumar and Minz, 2014). By selecting a meaningful feature subset, reduces the dataset in size, resulting in a more manageable and interpretable representation. Pruning datasets of irrelevant or redundant variables offers several benefits, like a reduction in model complexity, an improvement in computational efficiency, the mitigation of overfitting (Ying, 2019), and possible enhancements of the predictive performance by removing noisy data (Li et al., 2017). Focusing on the most informative features, which exhibit a strong relationship with the target variable, the model better captures the underlying patterns (Li et al., 2017). A smaller dataset (or model) reduces sensor and power usage, computing requirements, and storage demands, leading to cost savings, streamlined data collection processes, efficient data management, and better accessibility.

Applying sufficient feature selection has become unquestionably important as dataset sizes grow and machine learning is employed across a broadening range of applications.

## 2. Related work

The wide variety of existing feature selection methods can broadly be classified into three categories using the taxonomy described in (Guyon and Elisseeff, 2003a; Chandrashekar and Sahin, 2014). Methods are categorized as *filter*, *wrapper*, and

55 *embedded* methods depending on their working principles. The first evaluates the rel-  
 56 evance of features by looking at the intrinsic properties of the data, without involving  
 57 any learning algorithm. These methods typically use statistical techniques to rank  
 58 features based on their correlation with the output variable. Common techniques  
 59 include Pearson’s correlation coefficient, the Chi-square test, mutual information,  
 60 and various statistical tests (e.g., t-test, ANOVA). Those methods have the main  
 61 advantage of being computationally efficient, easy to compute, and can handle high-  
 62 dimensional data but suffer from a lack of specificity and overly broad assumptions,  
 63 e.g., data distribution or linearity (Bommert et al., 2020). Wrapper methods involve  
 64 a learning algorithm to evaluate the importance of feature subsets. These methods  
 65 search through the space of possible feature subsets and use a predictive model to  
 66 evaluate their performance. Techniques like recursive feature elimination (Guyon  
 67 et al., 2002a), forward selection (Guyon and Elisseeff, 2003b), backward elimination  
 68 (Koller and Sahami, 1996), and genetic algorithms (Leardi, 1994) are commonly used  
 69 in wrapper methods. While these methods tend to provide better performance than  
 70 filter methods by considering feature dependencies and interactions, they are com-  
 71 putationally expensive, especially with large datasets and high-dimensional feature  
 72 spaces (El Aboudi and Benhlila, 2016). Therefore, performance is often approxi-  
 73 mated rather than computed for each subset (El Aboudi and Benhlila, 2016; Kursu  
 74 and Rudnicki, 2010; Lundberg and Lee, 2017; Mlambo et al., 2016). Embedded  
 75 methods perform feature selection during the process of training a learning algo-  
 76 rithm. These methods integrate feature selection into the model training process,  
 77 often leading to more accurate and efficient feature selection. Examples of embed-  
 78 ded methods include decision trees and tree-based methods like random forests (Ho,  
 79 1995) and gradient boosting (Xu et al., 2014), LASSO (Least Absolute Shrinkage  
 80 and Selection Operator) (Tibshirani, 1996), and elastic net regularization (Zou and  
 81 Hastie, 2005). Embedded methods strike a balance between the computational effi-  
 82 ciency of filter methods and the accuracy of wrapper methods. By considering feature  
 83 interactions while being less computationally intensive than wrapper methods, these  
 84 methods are increasingly gaining interest. This work focuses on some of the leading  
 85 approaches for feature selection within this category.

86 The **concrete autoencoder** (Balm et al., 2019) uses an autoencoder architecture  
 87 that determines a pre-defined number of important features, which itself requires  
 88 certain knowledge about the task and dataset. It employs a concrete selector layer in  
 89 the encoder, which learns to select a specified number of features from the input while  
 90 allowing end-to-end training through gradient-based optimization. This method is  
 91 unsupervised and hence does not result in a trained model for inference. While the  
 92 method is effective, it requires prior knowledge about the number of relevant features

123 necessitating extensive experimentation to optimize.

124 **Shapley Additive Explanations** (SHAP) values (Lundberg and Lee, 2017) are  
125 used for interpreting machine learning models by quantifying the contribution of  
126 each feature to the model’s predictions. They are based on cooperative game theory  
127 and provide a consistent way to attribute the importance of features. For feature  
128 selection, SHAP values are calculated for each feature across the dataset, representing  
129 their contribution to the prediction outcomes. By aggregating these values, features  
130 can be ranked according to their importance (Samek et al., 2017). Heuristics like  
131 LinearSHAP or DeepLIFT (Lundberg and Lee, 2017; Mishra, 2021) approximate  
132 SHAP values, addressing the combinatorial explosion of subsets and making SHAP  
133 practical for feature selection (Lundberg and Lee, 2017; Marćilio and Eler, 2020).  
134 Computing SHAP values can be computationally expensive, especially for complex  
135 models and large datasets, resulting in high processing times. Additionally, the  
136 method assumes feature independence and the absence of redundancy, which might  
137 not hold in all datasets, potentially impacting the accuracy of the selection task.

138 **Feature Importance Deep Learning** (FIDL)(Wojtas and Chen, 2020) is a method  
139 for measuring feature importance in Deep Neural Networks (DNNs). It uses a dual-  
140 net strategy to find optimal feature subsets via binary masks, where large values  
141 indicate important features. FIDL employs a multiphase selection process with an  
142 exploration-exploitation strategy to generate mask candidates. The selector network  
143 aims to minimize the operator network’s loss, and new masks are chosen based on  
144 the input gradient size of the selector (exploitation) and random generation (explo-  
145 ration). Unlike other methods, FIDL operates by changing subsets of masks, making  
146 it a suitable choice for comparison.

147 **Stochastic Gate** feature selection (STG) (Yamada et al., 2020, 2018) employs a  
148 continuously relaxed Bernoulli variable for feature selection using stochastic gates.  
149 This method does not require specifying the number of features in advance. How-  
150 ever, it necessitates a large number of parameters to be trained in the first layer,  
151 leading to a risk of overfitting, particularly with deep neural networks handling high-  
152 dimensional data or when the training dataset is limited (Li et al., 2023; Singh et al.,  
153 2023).

154 The previously discussed approaches, while effective, face several limitations such  
155 as dependency on the choice of hyper-parameters, susceptibility to redundant infor-  
156 mation, and the need for prior knowledge about the number of relevant features.  
157 This research addresses these challenges by introducing a novel embedded method  
158 for feature selection that employs a dual-net architecture. The architecture con-  
159 sists of a task model and a selection model. Initially, the input data are fed into  
160 the selection model, which generates a quasi-binary mask (that is a mask that ap-

proaches a binary mask) aligned with the input feature dimensions. This mask is applied to a shifted version of the original features, which serves as input to the task model. The task model then uses the selected features to perform the target supervised task. Concurrently, the selection model aims to minimize the cumulative value of the mask, thereby selecting the most relevant features with minimal impact on the task model’s performance. The contributions and advantages of the proposed methodology are threefold: (1) Firstly, it effectively performs feature selection even in the presence of information redundancy and non-linear feature dependencies, achieving accuracy comparable to or better than the best of the compared state-of-the-art methods. (2) Secondly, while some state-of-the-art methods are limited by their complexity and the requirement for extensive hyper-parameter optimization, which increases the time and computational costs, this work emphasizes simplicity, requiring the optimization of only one hyper-parameter, for which a reliable estimation is suggested. (3) Lastly, selecting the minimal subset of features relevant to the supervised task often necessitates prior knowledge about the data and the number of relevant features. In the proposed approach, the number of relevant features is determined automatically, making the method more suitable for real-world applications and large, complex datasets.

This article, after having presented an introduction to the problem and existing possible solutions will explain the working principle of the new approach, present results proving its capability to achieve the stated objectives, compare the performance to the previously described methods, and discuss the determination of the hyper-parameter and possible limitations.

### 3. Fundamentals of dual-net feature selection

This paper proposes a dual-net approach, referred to as DNFS in the rest of this paper, which involves two jointly trained DNNs, one *task* network and one *selection* network. In the following, the task and selection networks will be denoted by  $\Pi$  and  $\Gamma$ , respectively.  $\Pi_t$  and  $\Gamma_t$  refer to the specific weights, biases, and optimizer states at a time  $t$ .  $\mathcal{D} = \{\mathcal{X}, \mathcal{Y}\}$  represents the dataset, where  $\mathcal{X} \subset \mathbb{R}^n$  denotes the set of instances as real-valued feature vectors in  $\mathbb{R}^n$ , and  $\mathcal{Y} \subset \mathbb{R}$  represents the target variable. Furthermore,  $\mathcal{V}$  refers to the set of all variables (features) of a given dataset with the feature count  $n := |\mathcal{V}|$ . Obviously, we have  $|\mathcal{X}| = |\mathcal{Y}|$ .  $\mathcal{D}' = \{\mathcal{X}', \mathcal{Y}'\}$  denotes a shuffled version of the dataset and  $x' \in \mathcal{X}'$  acts as a placeholder for some random instance. To clarify the distinction, we define variables as raw data inputs or measurements obtained from the system, while features are derived or transformed representations of these variables, relevant for predictive modeling or classification

167 tasks. Despite this distinction, the proposed approach can be applied to the selection  
 168 of either variables or features, depending on the input to the task model.

The task network’s target is to accomplish the supervised learning task on a feature subset provided by the selection network. The goal of the selection network is to generate a quasi-binary mask  $\mathbf{m} \in [0, 1]^n$  where 0 indicates a redundant or unimportant variable, and 1 indicates an important variable. The modified dataset

$$\mathcal{D}_{\mathbf{m}} := \mathcal{D} \otimes \mathbf{m} := \{(x \otimes \mathbf{m}, y)\}_{x \in \mathcal{X}, y \in \mathcal{Y}}$$

169 (here  $\otimes$  refers to the Hadamard (Horn, 1990) product) ought to maximize the per-  
 170 formance of the task network while simultaneously ensuring that  $\mathbf{m}$  is as sparse as  
 171 possible. The objectives of the networks are controlled via their loss functions

$$\mathcal{L}_{task}(x, y, \mathbf{m}; \Pi) = l(\Pi(x \otimes \mathbf{m}), y), \quad (1)$$

$$\mathcal{L}_{sel}(x, y; \Pi, \Gamma) = \lambda \cdot \mathcal{L}_{task}(x, y, \mathbf{m}; \Pi) + p(\mathbf{m}), \quad (2)$$

where  $\mathbf{m} := \Gamma(x')$ .

172 The task network’s loss  $\mathcal{L}_{task}$  consists of an instance-level task-dependent loss  
 173 function  $l$ , e.g., mean squared error (MSE) for regression or cross-entropy for classi-  
 174 fication, respectively. The selection network’s loss  $\mathcal{L}_{sel}$  is a weighted sum of the task  
 175 network’s loss and a penalization function

$$p : [0, 1]^n \rightarrow \mathbb{R}_{\geq 0}$$

176 that punishes the use of too many features, encouraging the selection of a smaller,  
 177 more relevant subset. The parameter  $\lambda$  weighs the sum to achieve the desired selec-  
 178 tion aggressiveness, where a higher  $\lambda$  increases the penalty on feature usage, promot-  
 179 ing sparser feature selection, and a lower  $\lambda$  reduces the penalty, allowing for more  
 180 features to be included. The concept of weighting and summing different loss func-  
 181 tions to achieve multiple training goals is uncommon in machine learning and can be  
 182 used when multiple goals are related to the optimization task, such as regularization  
 183 terms and multi-task learning like physics-informed neural networks (Abbassi et al.,  
 184 2024). The novelty of DNFS lies in outsourcing the penalization to a separate model  
 185 so that the training of the model performing the supervised task is not constrained  
 186 by a more complex loss function. Unlike other embedded methods, such as Lasso,  
 187 where the task model simultaneously handles both prediction and feature selection,  
 188 our approach decouples these tasks, allowing the selection model to focus solely on

189 optimizing the feature subset. The working concept of the proposed approach is  
 190 illustrated in Figure 1a.

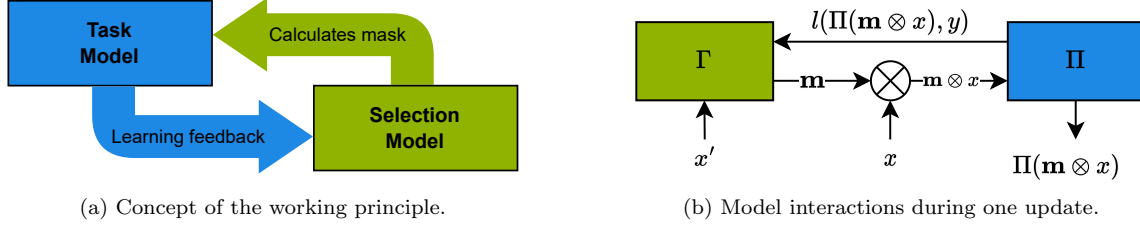


Figure 1: Illustration of the workflow of the proposed DNFS approach.

191 The selection network receives feedback through the learning performance of the  
 192 task model. It then generates a quasi-binary mask  $\mathbf{m} \in [0, 1]^n$ . A new dataset is  
 193 generated via

$$\mathcal{D}_{\mathbf{m}} = \mathcal{D} \otimes \mathbf{m}.$$

194 The task model is then trained for one iteration on  $\mathcal{D}_{\mathbf{m}}$ . Features with lower impor-  
 195 tance are eliminated by multiplication by 0, while important features are preserved  
 196 by the multiplication with non-zero values (for mask entries  $0 < \mathbf{m}_i < 1$ , the fea-  
 197 ture undergoes a linear transformation which the task model’s first layer can easily  
 198 compensate). The selection model incorporates the task model’s learning feedback  
 199 and adjusts its output correspondingly. The interactions of the networks are illus-  
 200 trated in Figure 1b. Since the selection model should output a global mask, i.e., an  
 201 input-independent mask, it is trained on a shuffled dataset denoted by  $\mathcal{D}'$ . A global  
 202 mask refers to a mask that does not depend on a specific instance  $x \in \mathcal{D}$  but rather  
 203 remains unchanged for all  $x \in \mathcal{D}$ , i.e., it contains the significant features for the  
 204 whole dataset. The process is described in detail in algorithm 1.

205 Before the training starts, both models are initialized (the specific implementa-  
 206 tion uses Glorot uniform initialization (Glorot and Bengio, 2010)). The selection  
 207 model generates a quasi-binary mask  $\mathbf{m}_0$  based on an instance  $x' \in \mathcal{D}'$  of the shuffled  
 208 dataset, which, at this point, is just a random output. The training data are mul-  
 209 tiplied by  $\mathbf{m}_0$  to generate the initial dataset  $\mathcal{D}_{\mathbf{m}_0}$  and passed on to the task model.  
 210 This updates according to the employed optimizer (cf. Appendix F). Subsequently,  
 211 the loss between the updated task model’s predicted values and the actual values is  
 212 recalculated. The loss of the selection model is computed as demonstrated in Eq.  
 213 2 and the selection model itself undergoes an update, initiating the process anew.  
 214 The procedure is detailed in Figure 2. For efficiency reasons, in practice, the model  
 215 updates happen batch-wise, unlike on the whole dataset. The training continues

---

**Algorithm 1** Dual-Net Feature Selection

---

**Input:** dataset  $\mathcal{D}$ ,  $y$  target variables, networks  $\Pi, \Gamma$ , networks parameter  $\theta_\Pi, \theta_\Gamma$  epochs  $E$ , hyperparameter  $\lambda$ ,  $p$  penalization function (cf. Section 6).

**Output:** trained task model  $\Pi$  and mask  $\mathbf{m}$

Generate randomized dataset  $\mathcal{D}' := \text{shuffle}(\mathcal{D})$

Initialize  $\mathbf{m} = \Gamma(x')$ ,  $x' \in \mathcal{D}'$ .

**repeat**

**for**  $\mathcal{B}$  (*batch*) **in**  $\mathcal{D}$  **do**

    calculate  $\mathcal{L}_{task}(\mathcal{B}, y, \mathbf{m}; \Pi) = l(\Pi(\mathcal{B} \otimes \mathbf{m}), y)$

    calculate  $\mathcal{L}_{sel}(\mathcal{B}, y; \Pi, \Gamma) = \lambda \cdot \mathcal{L}_{task}(\mathcal{B}, y, \mathbf{m}; \Pi) + p(\mathbf{m})$

    Update  $\theta_\Pi$  using  $\mathcal{L}_{task}$

    Update  $\theta_\Gamma$  using  $\mathcal{L}_{sel}$

$\mathbf{m} = \Gamma(x')$ .

**end for**

**until** stopping criterion is fulfilled.

---

216 either for a predetermined number of epochs or until a termination criterion is met.  
217 Examples of such a criterion are a certain number of features being set to 0, the  
218 mask  $\mathbf{m}$  converging to a certain subset, or a predefined early stopping criterion (cf.  
219 section 6.5). The method returns a feature subset  $\mathcal{S} \subset \mathcal{V}$  containing the features  
220 deemed relevant. The cardinality  $|\mathcal{S}|$  denotes the number of selected features.

221 An optimal feature selection method should be able to minimize the number of  
222 features while maintaining the same or a similar information density in the dataset,  
223 ensuring that a model can be trained on the remaining data without loss in per-  
224 formance. Further, it should effectively **prune highly correlated or redundant**  
225 **features**, even when the interdependencies are **non-linear** and the feature space is  
226 **high-dimensional**. Many existing feature selection methods require the target num-  
227 ber of features to be specified beforehand. This requires an initial understanding of  
228 the dataset and an estimate of the necessary percentage of features, which is often  
229 not a trivial task. The proposed method, by directly incorporating learning feedback,  
230 aims to **determine the number of relevant features** autonomously. These are  
231 the primary objectives of this approach.

232 While the dual-network structure has been employed in other methods, their  
233 working principles differ significantly from DNFS. For instance, FIDL uses an oper-  
234 ator net and a selector net, trained alternately. The operator net performs super-  
235 vised learning, while the selector net refines feature subsets based on performance  
236 feedback. FIDL focuses on ranking feature importance and incorporates exploration-



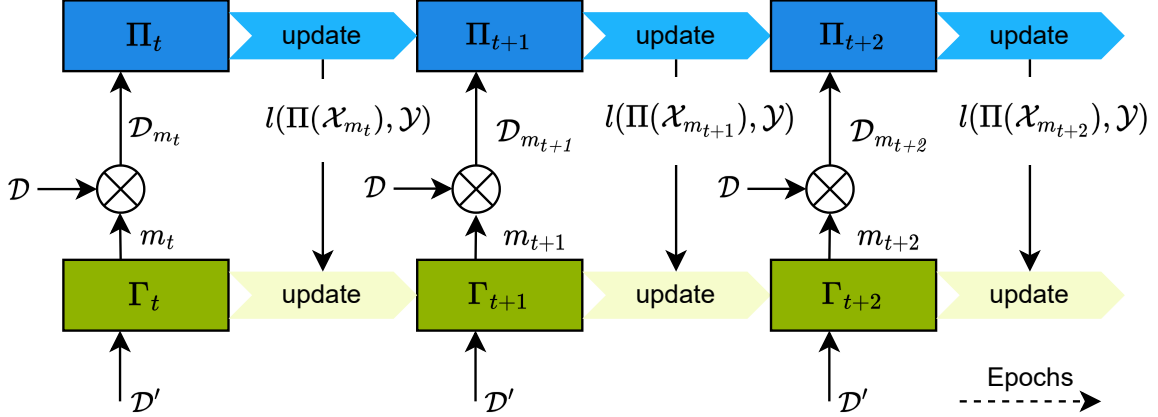


Figure 2: Model interactions over time. The concrete implementation performs batchwise instead of epoch-wise updates (cf. Algorithm 1).

237 exploitation with perturbation to avoid local optima. In contrast, DNFS trains its  
 238 task network and selection network jointly. The selection network generates a sparse,  
 239 quasi-binary mask by combining task performance loss with a sparsity penalty, en-  
 240 abling autonomous feature selection. This makes DNFS highly efficient for high-  
 241 dimensional datasets, as it emphasizes sparse feature selection and uses global masks.  
 242 Compared to CA, which employs a concrete selector layer to sample feature subsets  
 243 and focuses on reconstructing input data, DNFS prioritizes task performance and  
 244 sparsity. While CA excels in scenarios requiring feature interpretability and re-  
 245 construction accuracy, DNFS avoids the computational overhead of reconstruction,  
 246 making it more efficient for large-scale tasks. Similarly, STG introduces a proba-  
 247 bilistic approach, using stochastic binary gates to mask features. STG encourages  
 248 sparsity by penalizing the probability of gates being active, making it suitable for  
 249 tasks requiring highly sparse feature subsets. However, unlike DNFS, which uses  
 250 a deterministic quasi-binary mask, STG relies on probabilistic gates. This makes  
 251 DNFS more interpretable, as mask values directly indicate feature importance, while  
 252 STG’s probabilistic nature is advantageous for exploring uncertain feature interac-  
 253 tions. DNFS stands out for its efficiency, interpretability, and autonomous feature  
 254 selection, making it ideal for high-dimensional datasets. While FIDL ranks feature  
 255 importance, CA focuses on reconstruction, and STG offers probabilistic sparsity,  
 256 DNFS uniquely balances sparsity, performance, and scalability.

257 In the following sections, we present investigation results that demonstrate the  
 258 method’s ability to achieve these targets and compare its performance to other state-

259 of-the-art feature selection techniques.

## 260 4. Empirical evaluation

261 To assess the capabilities of DNFS, the method was tested on various datasets  
262 and tasks.

### 263 4.1. Datasets

264 To evaluate DNFS’s ability to eliminate redundancies in feature sets, synthetic  
265 datasets containing redundant features were engineered. The underlying synthetic  
266 dataset contains ten features, where the vector

$$(X_0, X_1, X_2, X_3, X_4, X_9) \sim \mathcal{N}(0, I_6)$$

267 follows a standard normal distribution. The remaining feature set  $\{X_5, X_6, X_7, X_8\}$   
268 is populated via linear and non-linear (n.l.) dependencies, resulting in two distinct  
269 datasets. The dependencies for the linear case are introduced as follows:

$$X_5 = 6X_2, \quad X_6 = -\frac{1}{2}X_3, \quad X_7 = 5X_0, \quad X_8 = -2X_1 \quad (3)$$

270 For the non-linear case, the  $X_5$  to  $X_8$  were populated via:

$$X_5 = X_2^3, \quad X_6 = -\frac{1}{2} \cdot \frac{X_3^5}{1 + \exp(X_0)}, \quad X_7 = 5 \tan^{-1}(X_0), \quad X_8 = -2 \exp(X_1^3) \quad (4)$$

271 For each dataset, a regression (Wojtas and Chen, 2020; Chen et al., 2018) and a  
272 classification (Wojtas and Chen, 2020; J. Friedman) task were defined with targets  
273  $Y_{reg}, Y_{cl}$  set as follows:

$$Y_{reg} = -2 \sin(X_0) + \max(X_1, 0) + X_2 + \exp(-X_3) + \varepsilon \quad (5)$$

274

$$Y_{cl} = \begin{cases} 0 & \text{if } 4 \leq \sum_{i=0}^3 X_i^2 \leq 16, \\ 1 & \text{else.} \end{cases} \quad (6)$$

275 where  $\varepsilon \sim \mathcal{N}(0, 0.1)$  represents noise added to the target. The condition in Eq. 6  
276 leads to 60%-40% distribution of ones and zeros, making for a balanced classification  
277 task.

278 For the benchmark tasks, the well-known MNIST (Lecun et al., 1998), Basehock,  
279 and Ames housing (Cock, 2011) datasets were chosen. The first involves an image  
280 classification task, the second a quasi-binary classification task, and the third a

281 regression task predicting housing prices. All these datasets include redundancies  
 282 in their features and can be significantly reduced in size without degradation of the  
 283 information density. These benchmarks may validate the achievement of the stated  
 284 objectives by the proposed method,, even when applied to more complex, real-world  
 285 tasks. For a detailed overview of the datasets and their properties, refer to Table  
 286 H.10 in Section Appendix H.

#### 287 4.2. Evaluation setup

288 Analyzing the results of the synthetic case is a comparatively simple task since  
 289 these can be classified deterministically. Because the underlying dependencies and  
 290 redundancies are known, it is possible to precisely quantify the quality of selected  
 291 subsets. To simplify the analysis of the results obtained using the synthetic datasets,  
 292 we introduce the following definitions for a subset of the set of all variables  $\mathcal{S} \subseteq \mathcal{V}$ :

293 **Definition 1.**  *$\mathcal{S}$  is called valid if it contains all information about  $Y$ , i.e.,  $X_0$  to  $X_3$   
 294 can be isolated via equivalence transformations and simultaneously does not contain  
 295 random noise, i.e., features.*

296 **Definition 2.**  *$\mathcal{S}$  is called minimal if it has size four, i.e.,  $N := |\mathcal{S}| = 4$ . This  
 297 originates from the fact that there are four unique variables and only bijective trans-  
 298 formations.*

299 **Definition 3.**  *$\mathcal{S}$  is called optimal if it is valid and minimal. A selection is called  
 300 optimal when the selected subset is optimal.*

301 A valid subset is defined as one which contains all information about the target.  
 302 For both datasets, the linear and non-linear sets, the minimal size of a valid subset  
 303 is 4. If the method has selected valid subsets of minimal size (within a margin), the  
 304 result is deemed optimal. Note here that if a method requires a target feature count  
 305 the percentages of valid and optimal subsets are equal. Although the minimal size is  
 306 well defined, the specific subsets that qualify as valid are not unique. For instance,  
 307 both  $\{X_0, X_1, X_2, X_3\}$  and  $\{X_5, X_6, X_7, X_8\}$  are optimal subsets, yet they share  
 308 no common variables. Each dataset presents 16 potential valid subsets. The concrete  
 309 subset varies depending on the dataset, the model weight initialization, and several  
 310 random states in the optimizer. To account for this variability, each selection process  
 311 was repeated multiple times, and the results were averaged to provide a robust  
 312 measure of performance.

313 Evaluating the performance of subsets using benchmark datasets poses a chal-  
 314 lenge due to the inherent ambiguity in both their sizes and the specific features  
 315 they contain. As a result, performance assessment was conducted by comparing the

model’s performance before and after subset selection. To establish a baseline, a machine learning model with an architecture identical to that of the task model was trained on the entire dataset. The validation performance of this full model was then recorded as the baseline performance for the dataset. Following the subset selection process, the model’s performance was re-evaluated using the reduced subset of features, employing the same training and validation procedures. This comparative approach makes a consistent evaluation of the impact of feature selection on model performance possible.

## 5. Results

We performed feature selection using DNFS on the previously presented datasets employing the proposed methodology. The results of the selection of the synthetic and benchmark datasets are summarized in Tables 1 and 2.

Table 1: Results of DNFS on synthetic datasets.

DATASET	VALID	N	OPTIMAL
SYNTH. REG.	100%	4	100%
SYNTH. CL.	100%	4	100%
SYNTH. REG. N.L.	100%	$4.132 \pm 0.4$	95%
SYNTH. CL. N.L.	100%	$4.325 \pm 0.52$	93%

For all four synthetic datasets, DNFS consistently produced optimal results, with a standard deviation in subset size of less than 10%. Remarkably, the method yielded valid subsets almost 100% of the time, underscoring that the performance of the task model is prioritized over the removal of features. The best performance, indicated by the highest percentage of optimal subsets, was observed in the linear regression task, where DNFS universally selected optimal subsets. The greatest deviation in subset size was seen with the non-linear classification task, with a standard deviation of 0.52.

An illustrative outcome of  $\Gamma_{200}$  applied to synthetic regression with linear dependencies is exemplified by the following vector:

$$\mathbf{m}_{200} \approx (0.1, 0.2, 0.1, 0, 0, 0, 0.4, 0, 0, 0)^T$$

Interpreting the non-zero entries as selected features translates to the subset

$$\mathcal{S} = \{X_0, X_1, X_2, X_6\}, \quad \text{with} \quad |\mathcal{S}| = 4$$

being selected in epoch 200.

Table 2: Results of DNFS on benchmark datasets.

DATASET	BASELINE	SEL. PERF.	N
MNIST	$96,9 \pm 0,01\%$	$95,5\% \pm 0,3\%$	$50 \pm 2$
AMES H.	$17.000 \pm 300\$$	$15.000\$ \pm 500\$$	$20 \pm 3$
BASEHOCK	96%	95%	$53 \pm 4$

Table 2 shows the results for the MNIST, Ames, and Basehock datasets. The baseline was established by training a model with the same architecture as the task model on the whole dataset and measuring its validation performance. The intention behind using the same architecture is to ensure a fair and controlled evaluation, where the impact of the feature selection process on performance can be isolated from architectural or hyperparameter influences. For the MNIST dataset, DNFS achieved  $\approx 94\%$  feature reduction while maintaining a high validation accuracy of 95.5%. On the Basehock dataset, the feature count was reduced by more than 99% with a stable accuracy of 95%. For the Ames housing dataset, the reduction was around 85%, accompanied by a significant improvement in validation MAE, with a decrease of approximately 12%. To ensure that the observed performance improvements are not dependent on the architecture or hyperparameter settings of the task model, we additionally evaluated the selected features using a support vector machine (SVM). The SVM was trained on both the full feature set and the subset selected by DNFS. The results, provided in Table C.6, show similar performance trends to those reported in Table 2, thereby validating the effectiveness of the selected features in a model-independent manner.

Overall, DNFS successfully identified a significant number of relevant features across various datasets and achieved reductions in size surpassing 90% depending on the specific task. Notably, DNFS demonstrated consistent results with minimal variation, and the validation performance of the task model either remained similar or even improved.

## 6. Discussion

To contextualize the presented results, the same feature selection tasks were performed by the methods discussed in Section 2. This section will compare the results and provide a detailed analysis of the properties, advantages, and disadvantages of DNFS.

### 6.1. Comparative results on synthetic datasets

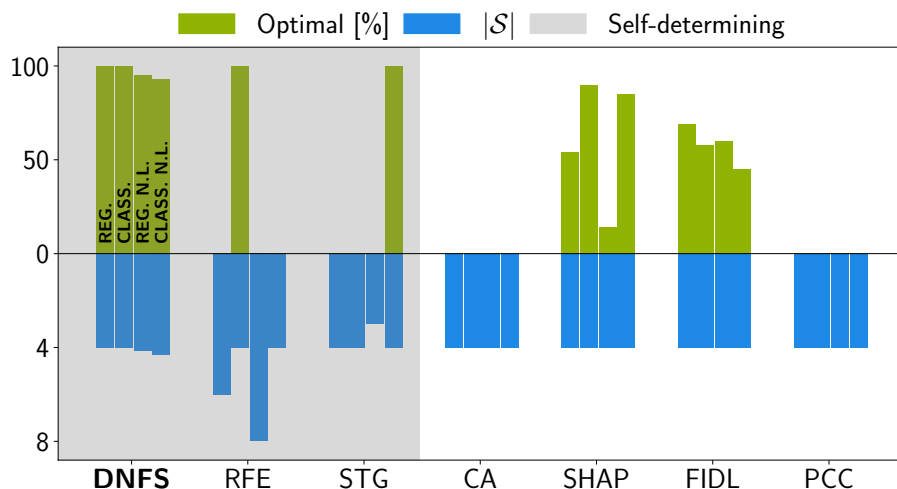


Figure 3: DNFS, RFE, STG, CA, SHAP, FIDL, and PCC results. The datasets are (from left to right): linear regression, linear classification, non-linear regression, non-linear classification. CA, SHAP, and FIDL return a fixed number, i.e., 4 features. For exact results see [Appendix A](#).

Figure 3 shows the performance of the seven different methods on the synthetic datasets. In addition to state-of-the-art embedded methods, **P**earson **C**orrelation **C**oefficient (PCC) and **R**ecursive **F**eature **E**limination (RFE) ([Guyon et al., 2002b](#)) are included in the comparison. PCC identifies relevant features by calculating absolute correlation values, discarding those with little influence on the outcome. RFE, on the other hand, is a feature selection technique that iteratively removes the least important features to improve model performance. It works by fitting a model, ranking features by importance, and recursively eliminating the least significant ones until a desired number of features remains. RFE is commonly used with algorithms like support vector machines (SVMs) and decision trees, making it an effective wrapper method for selecting relevant features while maintaining accuracy.

379 The definitions for valid, minimal, and optimal results (subsets) are consistent  
 380 with the criteria outlined in Section 4. The results highlight the superiority of DNFS  
 381 in the feature selection task compared to all other implemented methods. DNFS  
 382 successfully identifies optimal feature subsets for 100 % of the linear tasks and 94 % of  
 383 the non-linear regression and classification tasks, demonstrating clear dominance over  
 384 the compared methods (detailed results are provided in Table A.4). An interesting  
 385 observation is that while RFE is theoretically expected to deliver optimal results,  
 386 it often falls short in practice, selecting redundant features for nearly every dataset  
 387 except the linear classification task. This limitation stems from the absence of deep  
 388 learning models in the RFE process, which restricts its ability to capture complex,  
 389 non-linear feature interactions. The use of synthetic datasets was intended to test  
 390 the methods' ability to identify and prune even complex non-linear dependencies.  
 391 While the other methods may also result in subsets delivering good performance  
 392 in the linear tasks, they often struggle with either selecting the minimal number  
 393 of features or, if so, the correct combinations for optimal performance that don't  
 394 contain redundant information. Specifically, CA selects redundant features, and  
 395 STG doesn't succeed in selecting all relevant features for the fulfillment of the task or  
 396 selects redundant ones as well. FIDL and SHAP produce mediocre results, frequently  
 397 failing on multiple datasets to meet one of the objectives by selecting redundant  
 398 features in many trials. Although these methods can achieve optimal results, they  
 399 do not consistently reproduce them. PCC fails to remove redundancies and instead  
 400 consistently selects the redundant features with the highest correlation to the target  
 401 in all cases.

## 402 6.2. Comparison of Embedded Methods on Benchmark Datasets

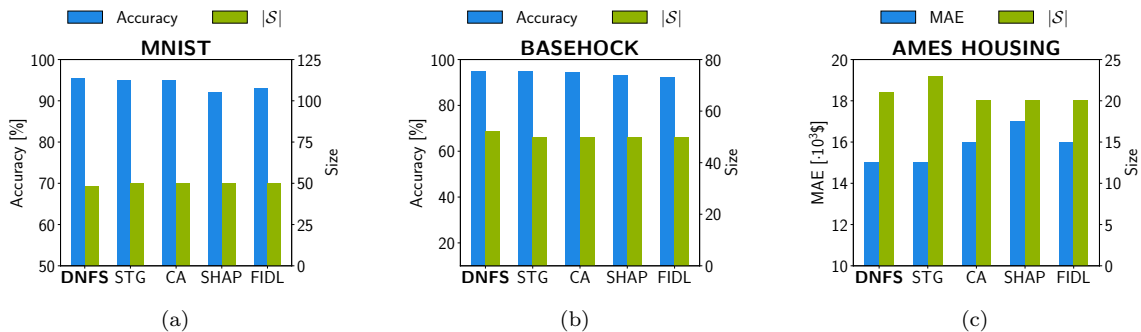


Figure 4: Comparative results on the benchmark datasets.

403 After analyzing a multitude of different methods, this section focuses on compar-  
 404 ing DNFS to other embedded methods on benchmark datasets. Since DNFS is itself

an embedded method, we limit our analysis to this category to ensure a fair and meaningful comparison. On MNIST, DNFS achieves the best results with the best validation performance and the highest validation accuracy among all methods for the given feature count. The validation performance with DNFS is nearly identical to the baseline, showing no significant deviation and remaining well within the margin of error, which is markedly better than the declines observed with the other methods. The targeted feature count was 50 features, as this was the average number selected by DNFS, and further shrinking the target did not result in improved performance of the other methods.

On the Basehock dataset, DNFS demonstrates the capability to perform competitively in a high-dimensional setting, matching the performance of methods such as STG, CA, FIDL, and SHAP. Although DNFS selected more features than the other methods, it also achieved the highest accuracy among them. This may be due to the elimination of relevant features by the other methods, which, as observed with STG in non-linear tasks, results in smaller selected feature subsets but at the cost of performance.

Regarding the Ames Housing dataset, DNFS successfully identified the optimal subset of features, leading to superior performance in the regression task compared to subsets selected by other methods. Notably, when compared to STG—the only other method capable of autonomously determining the number of relevant features—DNFS not only selected fewer features but also achieved better results. This dataset is particularly challenging due to its high proportion of categorical features, which introduces significant sparsity, especially after one-hot encoding.

Table 3: Comparison of study contents with other works. \*While the authors outline how to apply their method in a supervised setting (Balm et al., 2019), no implementation is provided, therefore it hasn’t been tested.

	DNFS	RFE	STG	CA	SHAP	FIDL	PCC
SELF-DETERMINING	✓	✓	✓	✗	✗	✗	✓
REDUNDANCY-ELIMINATION	✓	✓	—	✗	✗	✗	✗
FEW HYPERPARAMETERS	✓	✓	✗	✓	✓	✗	✓
SUPERVISED	✓	✓	✓	—*	✓	✓	✓
SIMPLICITY	✓	✗	✗	✓	✓	✗	✓

427



428 The comparison results demonstrated the superiority of DNFS in selecting the  
 429 optimal subset of features, particularly in the presence of non-linear dependencies,  
 430 as observed in the synthetic data, and in high-dimensional datasets like MNIST  
 431 and Ames Housing. DNFS excels in selecting all relevant features for the target,  
 432 avoiding redundant features, and maintaining or improving task model performance.  
 433 A key advantage of DNFS is its ability to determine the optimal number of fea-  
 434 tures without requiring any prior specification, setting it apart from methods such  
 435 as CA, FIDL, SHAP, and other state-of-the-art feature selection techniques. In  
 436 practical applications, identifying the optimal feature count often requires exten-  
 437 sive optimization. However, DNFS inherently addresses this challenge without the  
 438 need for predefined targets. Moreover, DNFS offers simplicity, requiring minimal  
 439 hyperparameter tuning. This simplicity is further exemplified by the use of a single  
 440 hyperparameter,  $\lambda$ , whose optimal value is suggested and discussed in section 6.3.  
 441 The training duration of DNFS increases by approximately  $1.2\times$  to  $2.6\times$ , depending  
 442 on the dataset characteristics and whether the task is regression or classification,  
 443 compared to training the baseline model (task model without the selection task).  
 444 This increase in training time can be attributed to the additional computational  
 445 complexity introduced by the selection process, which requires the model to evaluate  
 446 feature relevance during training. Despite this increase, the enhanced interpretability  
 447 and improved performance obtained through feature selection justify the additional  
 448 computational effort. Further details, including runtime comparisons and hardware  
 449 specifications, are provided in [Appendix E.4](#) and [Table E.9](#). STG is particularly  
 450 sensitive to a range of hyperparameters, which can result in substantially different  
 451 outcomes. For instance, modifications to the optimizer or number of epochs lead to  
 452 significant variations in performance. To the best of our knowledge, the authors have  
 453 not explicitly addressed this issue. While other methods exhibit similar limitations,  
 454 they are generally less pronounced compared to STG. A high number of sensitive  
 455 hyperparameters diminishes the practicality of a method, as it necessitates exten-  
 456 sive computational resources to identify optimal parameters through grid search or  
 457 comparable algorithms, thereby significantly increasing runtime and complexity.

### 458 6.3. The impact of $\lambda$

459 The hyperparameter  $\lambda$ , which weights the sum in the selection network’s ob-  
 460 jective function (cf. [Eq. 2](#)), is the only hyperparameter which directly influences  
 461 the selection. Consequently, variations in  $\lambda$  can potentially significantly impact the  
 462 outcomes. Other parameters are set indirectly, like the model architectures, the em-  
 463 ployed optimizer, and the method of weight initialization. Compared to  $\lambda$ , we have  
 464 found none of these to be nearly as impactful, or impactful at all.

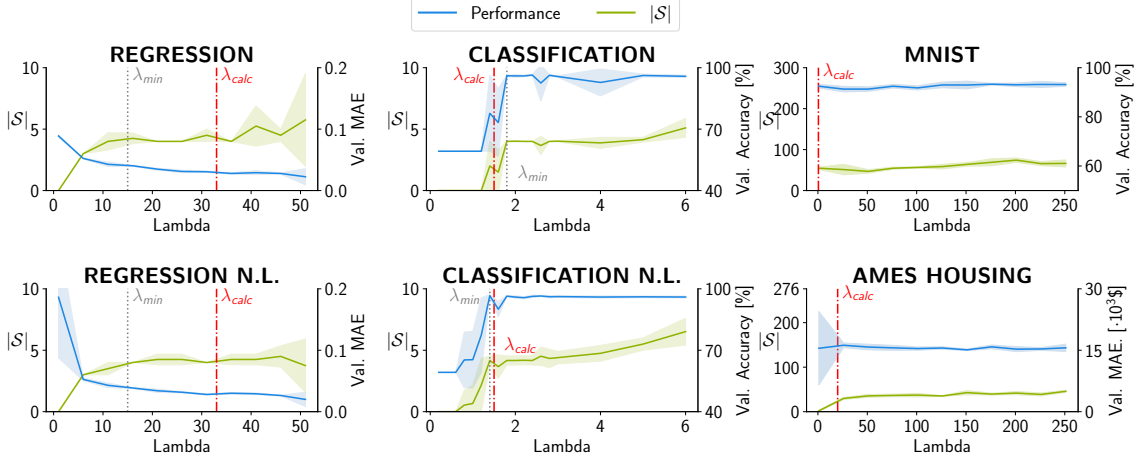


Figure 5: Selection behavior of DNFS for the synthetic datasets for increasing values of  $\lambda$ . Depending on  $\lambda$ 's magnitude, the summands of (2) are weighted differently, resulting in an in- or decreased number of selected features.  $\lambda_{min}$  marks the first values at which DNFS converges reliably to the optimal solution.  $\lambda_{calc}$  represents the values calculated via the approximations presented in section 6.3. The regressions were performed on normalized target (cf. section 6.3) and, for Ames, scaled back to the original scale.

The aforementioned dependency on  $\lambda$  makes it the only interesting parameter to analyze. To better visualize its impact, Figure 5 shows the size  $|\mathcal{S}|$  of the outputted subsets, their yielded performance, and the standard deviation of those metrics for all presented datasets. For all synthetic datasets, there exists a threshold value  $\lambda_{min}$  beyond which DNFS exhibits a stable convergence to optimal subsets. On the benchmark datasets, the intervals of convergence are even larger, spanning over one order of magnitude in length. It is important to note that no method converges at  $\lambda = 0$ , even if some of the graphs might suggest otherwise. This is because, for some datasets, DNFS still converges for small values.

Choosing  $\lambda$  is crucial for a successful selection. It depends on the dataset, the loss function  $\mathcal{L}_{task}$ , and the penalization function  $p$ . We employed categorical cross-entropy (resp. MSE) for classification (resp. regression) as loss function  $l$  in Eq. 1. Both losses can theoretically take values in whole  $\mathbb{R}_{\geq 0}$  but can be approximated if one makes certain reasonable assumptions. The first is that the task model is reasonably well suited for the predictive task and, in the case of a regression, that the target  $Y$  is normalized, which can always be achieved via min-max transformation:

$$Y_{norm} = \frac{Y - \min(Y)}{\max(Y) - \min(Y)}$$

481 For classification, it is then reasonable to approximate  $l$  as

$$l \leq -\log\left(\frac{1}{|\text{Classes}|}\right) = \log(|\text{Classes}|),$$

482 which is the loss resulting from the task model uniformly guessing ( $|\text{Classes}|$  denotes  
 483 the number of classes of a given classification task). In the case of regression, we  
 484 assume that the mse deviates approximately 10% around the mean  $\bar{Y}_{norm}$  which,  
 485 especially for later epochs, is fairly generous. Choosing  $\lambda$  such that  $\lambda \cdot \mathcal{L}_T$  and  $p$  are  
 486 of the same order of magnitude (for  $p(\mathbf{m}) = \bar{\mathbf{m}}$  we have  $0 \leq p \leq 1$ ) leads to:

$$\lambda_{class} \geq \frac{1}{\log(|\text{Classes}|)}, \quad \lambda_{reg} \geq \frac{1}{0.1 \cdot \bar{Y}_{norm}} \quad (7)$$

487 These are not strict boundaries but result from, in this context, justified generaliza-  
 488 tions. The calculated values  $\lambda_{calc}$  are shown in Figure 5 (cf. [Appendix E.2 Table](#)  
 489 [E.7](#)). All lie in the area of convergence except for the linear classification dataset, for  
 490 which the calculation is slightly too small. This further emphasizes that, for optimal  
 491 selection of  $\lambda$ , an initial value can be derived based on the Eq. 7, with further possible  
 492 improvement of results through trials with other  $\lambda$  values. Additionally, this suggests  
 493 that the model’s uniform guessing approximation might be overly pessimistic, which  
 494 is likely the case.

#### 495 6.4. Convergence behavior

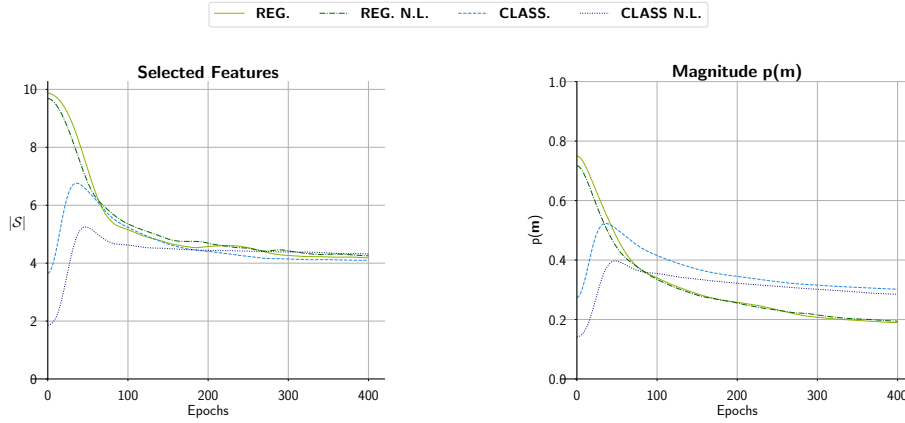


Figure 6: Convergence behavior of the subset size DNFS selects (left) and the magnitude of the penalization function  $p$  (right) on the synthetic datasets. The values were averaged over multiple trials.

DNFS exhibited the desired performance on all datasets and converged to optimal solutions. Figure 6 examines further the convergence behavior of DNFS on the synthetic datasets for a specific  $\lambda$  (cf. Appendix E Table E.7). The left graphs in Figure 6 show the number of features being selected in every epoch, i.e., the number of non-zero entries in  $\mathbf{m}$ , while the right show the magnitude of the penalization function  $p$ . In the concrete implementation, we chose  $p(\mathbf{m}) = \overline{\mathbf{m}}$  as the arithmetic mean of  $\mathbf{m}$ . Comparing the two plots reveals that while  $|\mathcal{S}|$  converges,  $p$  seems not to. The cause lies in the continuity of  $p$ . Since  $p$  can be reduced by linearly decreasing all non-zero values, its magnitude decreases while simultaneously the number of selected features stays constant. This holds true for every continuous function  $p$ .

**Proposition 1.** *For every non-constant continuous function  $f \in \mathcal{C}(\mathbb{R}^n)$  and  $\varepsilon > 0$  with  $f(x) = 0$  for some  $x \in \mathbb{R}^n$  there exist  $y \in \mathbb{R}^n$  such that  $f(y) \neq 0$  and  $|f(y)| < \varepsilon$ .*

*Proof.* This is an immediate deduction from the intermediate values theorem for multivariate functions (Shimamoto, 2019) (cf. Appendix B).  $\square$

The result shows that the previously described problem can arise with any non-constant continuous penalization function. The following two masks are an example of the phenomenon on the synthetic datasets. After 200 epochs of training the selection network’s output is:

$$\mathbf{m}_{200} \approx (0.2, 0, 0.2, 0.5, 0, 0, 0, 0, 0.2, 0)^T$$

After 400 epochs  $\Gamma_{400}$  generated

$$\mathbf{m}_{400} \approx (0.1, 0, 0.1, 0.3, 0, 0, 0, 0, 0.1, 0)^T \approx \frac{1}{2} \mathbf{m}_{200}$$

which corresponds to a halving of  $\mathbf{m}_{200}$ . Notably, no new feature has been set to zero, insinuating that the feature count has stabilized. Therefore, convergence criteria that solely depend on  $p$  are unfeasible.

### 6.5. Classifying selections and stopping criteria

Determining which features should be classified as selected and when to stop the selection is a non-trivial task in and of itself. In practice,  $\mathbf{m}$  approaches gradually zero at indices whose corresponding features are being deemed unimportant. Nevertheless, this convergence can be slow and tends to decrease in pace as the neuron values inch toward 0. To address this, a classification threshold  $\delta > 0$  is introduced to determine whether an entry should be classified as selected or not. Moreover,

the dropout layers (cf. Table E.8) help prevent the method from converging prematurely to a suboptimal subset but also introduce volatility into  $\mathbf{m}$ , causing mask entries to occasionally spike even after extended training if the feature is deemed unimportant. This problem can be mitigated using different strategies. A possible way is to calculate the moving average of the last  $t_{av}$  epochs and perform the classification on this averaged mask. This method alleviates the volatility problem but is still prone to the threshold  $\delta$  beyond which a feature is classified as selected. The current implementation employs a more sophisticated statistical approach. Starting at epoch  $t_{start}$ , which gives the process time to stabilize, it records whether a feature is selected in each subsequent epoch. We refer to the interval  $[t_{start}, t_{end}]$  as the selection interval. Post training, features deemed important are those selected in at least  $p\%$  of the epochs. This generally is a stronger condition than the moving average method and yields better results. Determining the selected features in each epoch  $t$  requires setting a threshold  $\delta$ . The hyperparameters  $p$  and  $\delta$  are robust and generally independent of the dataset (cf. Appendix D Figures D.8 and D.7).

Deciding when to stop the selection can be achieved either by training for a fixed number of epochs or by using a classic early stopping criterion like validation loss monitoring of the task model’s performance. The latter has been implemented in the current version.

## 6.6. Limitations

While DNFS introduces an innovative approach and demonstrates promising results in various scenarios, it is important to acknowledge its potential limitations. Like all approaches, it encounters certain constraints that may impact its performance and applicability in specific contexts. It necessitates a certain level of understanding, if not of the relationship between the features and targets, then of the task itself. First and foremost, the current implementation and theoretical foundation are designed for multilayer perceptrons (MLPs) and were not tested on different model architectures. Therefore, we cannot make any assertions about performance beyond classical MLPs. The specific architectures vary from task to task and need to be determined by a prior investigation. The models used by us are stated in Table E.8. Furthermore, requiring the training of deep neural networks DNFS is inherently more computing intensive than statistical methods or methods that employ simpler networks, e.g., linear models. The increase in the necessary time per epoch varies from 40% for the synthetic datasets to 160% for MNIST (for used hardware and precise runtimes see Appendix E.4).

To guarantee a comparable and fair evaluation, we trained a new network on the selected subsets. In praxis, though, it is recommendable to use the already

562 trained task network to perform the predictive task. This, however, is perhaps not  
563 always feasible since task model performance tends to lag behind a newly trained  
564 model. Convergence issues may arise with DNFS when variables in a dataset exhibit  
565 extremely uneven information density. In such cases, DNFS tends to select too  
566 few variables, failing to include those that, despite offering less information about  
567 the target compared to the already selected variables, still contribute to increases  
568 in performance, even if small. This problem may also arise if a lot of categorical  
569 features are present.

## 570 **7. Conclusion**

571 This study presents a new dual-network architecture for embedded feature se-  
572 lection. The main innovation is its ability to identify a subset of features that are  
573 essential for achieving the target variable, effectively addressing issues of redundancy  
574 and irrelevance. This method also tackles the challenge of not having prior knowledge  
575 about the number of features. Through evaluation and comparison with state-of-the-  
576 art methods using various synthetic and benchmark datasets in different supervised  
577 tasks, this approach demonstrates its potential to achieve high accuracy for the in-  
578 tended tasks while using a minimal feature set. The effectiveness of the selection  
579 process depends on the choice of the hyperparameter  $\lambda$ , which balances accuracy  
580 and the number of selected features.

## 581 **Acknowledgments**

582 The authors gratefully acknowledge the financial support provided by the Federal  
583 Ministry for Economic Affairs and Climate Action of the Federal Republic of Ger-  
584 many through the project “SONYA - Increasing the reliability of segmented rotor  
585 blades through hybrid condition monitoring” (FKZ 03EE3026B).

## 586 **Declaration of conflicting interests**

587 The authors declared no potential conflicts of interest with respect to the research,  
588 authorship, and/or publication of this article.

## 589 **Declaration of generative AI and AI-assisted technologies in the writing** 590 **process**

591 During the preparation of this work, the authors used ChatGPT in order to  
592 improve the language. After using this tool, the authors reviewed and edited the  
593 content as needed and take full responsibility for the content of the publication.

## 594 References

- 595 Abbassi, A., Römgers, N., Griebmann, T., Rolfes, R., 2024. Data augmentation with  
596 data-based simulation in structural health monitoring using knowledge embedded  
597 generative adversarial network. Available at SSRN 4687651 .
- 598 Baln, M.F., Abid, A., Zou, J., 2019. Concrete autoencoders: Differentiable feature  
599 selection and reconstruction, in: Chaudhuri, K., Salakhutdinov, R. (Eds.), Pro-  
600 ceedings of the 36th International Conference on Machine Learning, PMLR. pp.  
601 444–453. URL: <https://proceedings.mlr.press/v97/balin19a.html>.
- 602 Basodi, S., Ji, C., Zhang, H., Pan, Y., 2020. Gradient amplification: An efficient  
603 way to train deep neural networks. Big Data Mining and Analytics 3, 196–207.  
604 doi:[10.26599/BDMA.2020.9020004](https://doi.org/10.26599/BDMA.2020.9020004).
- 605 Bommert, A., Sun, X., Bischl, B., Rahnenführer, J., Lang, M., 2020. Benchmark  
606 for filter methods for feature selection in high-dimensional classification data.  
607 Computational Statistics & Data Analysis 143, 106839. URL: <https://www.sciencedirect.com/science/article/pii/S016794731930194X>, doi:<https://doi.org/10.1016/j.csda.2019.106839>.
- 610 Chandrashekar, G., Sahin, F., 2014. A survey on feature selection methods. Com-  
611 puters & electrical engineering 40, 16–28.
- 612 Chen, J., Stern, M., Wainwright, M.J., Jordan, M.I., 2018. Kernel feature selection  
613 via conditional covariance minimization. [arXiv:1707.01164](https://arxiv.org/abs/1707.01164).
- 614 Cock, D.D., 2011. Ames, iowa: Alternative to the boston housing data as an end of  
615 semester regression project. Journal of Statistics Education URL: <https://jse.amstat.org/v19n3/decock.pdf>.
- 617 El Aboudi, N., Benhlima, L., 2016. Review on wrapper feature selection approaches,  
618 in: 2016 International Conference on Engineering & MIS (ICEMIS), pp. 1–5.  
619 doi:[10.1109/ICEMIS.2016.7745366](https://doi.org/10.1109/ICEMIS.2016.7745366).
- 620 Glorot, X., Bengio, Y., 2010. Understanding the difficulty of training deep feed-  
621 forward neural networks, in: Teh, Y.W., Titterton, M. (Eds.), Proceedings  
622 of the Thirteenth International Conference on Artificial Intelligence and Statis-  
623 tics, PMLR, Chia Laguna Resort, Sardinia, Italy. pp. 249–256. URL: <https://proceedings.mlr.press/v9/glorot10a.html>.
- 624

- 625 Guyon, I., Elisseeff, A., 2003a. An introduction to variable and feature  
626 selection. CoRR URL: [https://www.jmlr.org/papers/volume3/guyon03a/](https://www.jmlr.org/papers/volume3/guyon03a/guyon03a.pdf)  
627 [guyon03a.pdf](https://www.jmlr.org/papers/volume3/guyon03a/guyon03a.pdf).
- 628 Guyon, I., Elisseeff, A., 2003b. An introduction to variable and feature selection.  
629 Journal of machine learning research 3, 1157–1182.
- 630 Guyon, I., Weston, J., Barnhill, S., Vapnik, V., 2002a. Gene selection for cancer  
631 classification using support vector machines. Machine learning 46, 389–422.
- 632 Guyon, I., Weston, J., Barnhill, S., Vapnik, V., 2002b. Gene selection for cancer  
633 classification using support vector machines. Machine learning 46, 389–422.
- 634 Ho, T.K., 1995. Random decision forests, in: Proceedings of 3rd international con-  
635 ference on document analysis and recognition, IEEE. pp. 278–282.
- 636 Horn, R.A., 1990. The hadamard product, in: Proc. Symp. Appl. Math, pp. 87–169.
- 637 J. Friedman, T. Hastie, R.T., . The Element of Statistical Learning. Springer, Berlin,  
638 2001.
- 639 Kingma, D.P., Ba, J., 2017. Adam: A method for stochastic optimization.  
640 [arXiv:1412.6980](https://arxiv.org/abs/1412.6980).
- 641 Kira, K., Rendell, L.A., 1992. A practical approach to feature selection, in: Slee-  
642 man, D., Edwards, P. (Eds.), Machine Learning Proceedings 1992. Morgan Kauf-  
643 mann, San Francisco (CA), pp. 249–256. URL: [https://www.sciencedirect.](https://www.sciencedirect.com/science/article/pii/B9781558602472500371)  
644 [com/science/article/pii/B9781558602472500371](https://www.sciencedirect.com/science/article/pii/B9781558602472500371), doi:[https://doi.org/10.](https://doi.org/10.1016/B978-1-55860-247-2.50037-1)  
645 [1016/B978-1-55860-247-2.50037-1](https://doi.org/10.1016/B978-1-55860-247-2.50037-1).
- 646 Koller, D., Sahami, M., 1996. Toward optimal feature selection. Technical Report.  
647 Stanford InfoLab.
- 648 Kumar, V., Minz, S., 2014. Feature selection: a literature review. SmartCR 4,  
649 211–229.
- 650 Kursa, M.B., Rudnicki, W.R., 2010. Feature selection with the boruta package.  
651 Journal of Statistical Software 36, 1–13. URL: [https://www.jstatsoft.org/](https://www.jstatsoft.org/index.php/jss/article/view/v036i11)  
652 [index.php/jss/article/view/v036i11](https://www.jstatsoft.org/index.php/jss/article/view/v036i11), doi:[10.18637/jss.v036.i11](https://doi.org/10.18637/jss.v036.i11).
- 653 Leardi, R., 1994. Application of a genetic algorithm to feature selection under full  
654 validation conditions and to outlier detection. Journal of Chemometrics 8, 65–79.



- 655 Lecun, Y., Bottou, L., Bengio, Y., Haffner, P., 1998. Gradient-based learning applied  
656 to document recognition. *Proceedings of the IEEE* 86, 2278–2324. doi:[10.1109/  
657 5.726791](https://doi.org/10.1109/5.726791).
- 658 Li, A., Peng, H., Zhang, L., Huang, J., Guo, Q., Yu, H., Liu, Y., 2023. Fedsdg-  
659 fs: Efficient and secure feature selection for vertical federated learning, in: *IEEE*  
660 *INFOCOM 2023-IEEE Conference on Computer Communications*, IEEE. pp. 1–  
661 10.
- 662 Li, J., Cheng, K., Wang, S., Morstatter, F., Trevino, R.P., Tang, J., Liu, H., 2017.  
663 Feature selection: A data perspective. *ACM Comput. Surv.* 50. URL: [https:  
664 //doi.org/10.1145/3136625](https://doi.org/10.1145/3136625), doi:[10.1145/3136625](https://doi.org/10.1145/3136625).
- 665 Lundberg, S.M., Lee, S., 2017. A unified approach to interpreting model pre-  
666 dictions. *CoRR* abs/1705.07874. URL: <http://arxiv.org/abs/1705.07874>,  
667 [arXiv:1705.07874](https://arxiv.org/abs/1705.07874).
- 668 Marçílio, W.E., Eler, D.M., 2020. From explanations to feature selection: assessing  
669 shap values as feature selection mechanism, in: *2020 33rd SIBGRAPI Confer-*  
670 *ence on Graphics, Patterns and Images (SIBGRAPI)*, pp. 340–347. doi:[10.1109/  
671 SIBGRAPI51738.2020.00053](https://doi.org/10.1109/SIBGRAPI51738.2020.00053).
- 672 Mishra, P., 2021. Explainability for deep learning models, in: *Practical Explainable*  
673 *AI Using Python: Artificial Intelligence Model Explanations Using Python-based*  
674 *Libraries, Extensions, and Frameworks*. Springer, pp. 243–263.
- 675 Mlambo, N., Cheruiyot, W.K., Kimwele, M.W., 2016. A survey and comparative  
676 study of filter and wrapper feature selection techniques. *International Journal of*  
677 *Engineering and Science (IJES)* 5, 57–67.
- 678 Samek, W., Wiegand, T., Müller, K., 2017. Explainable artificial intelligence:  
679 Understanding, visualizing and interpreting deep learning models. *CoRR*  
680 abs/1708.08296. URL: <http://arxiv.org/abs/1708.08296>, [arXiv:1708.08296](https://arxiv.org/abs/1708.08296).
- 681 Shimamoto, D., 2019. *Multivariable Calculus*. Don Shimamoto.
- 682 Singh, D., Climente-González, H., Petrovich, M., Kawakami, E., Yamada, M., 2023.  
683 Fsnets: Feature selection network on high-dimensional biological data, in: *2023*  
684 *International Joint Conference on Neural Networks (IJCNN)*, IEEE. pp. 1–9.
- 685 Tibshirani, R., 1996. Regression shrinkage and selection via the lasso. *Journal of the*  
686 *Royal Statistical Society Series B: Statistical Methodology* 58, 267–288.

- 687 Wojtas, M., Chen, K., 2020. Feature importance ranking for deep learn-  
 688 ing. CoRR abs/2010.08973. URL: <https://arxiv.org/abs/2010.08973>,  
 689 [arXiv:2010.08973](https://arxiv.org/abs/2010.08973).
- 690 Xu, Z., Huang, G., Weinberger, K.Q., Zheng, A.X., 2014. Gradient boosted feature  
 691 selection, in: Proceedings of the 20th ACM SIGKDD international conference on  
 692 Knowledge discovery and data mining, pp. 522–531.
- 693 Yamada, M., Tang, J., Lugo-Martinez, J., Hodzic, E., Shrestha, R., Saha, A.,  
 694 Ouyang, H., Yin, D., Mamitsuka, H., Sahinalp, C., Radivojac, P., Menczer, F.,  
 695 Chang, Y., 2018. Ultra high-dimensional nonlinear feature selection for big biolog-  
 696 ical data. IEEE Transactions on Knowledge and Data Engineering 30, 1352–1365.  
 697 doi:[10.1109/TKDE.2018.2789451](https://doi.org/10.1109/TKDE.2018.2789451).
- 698 Yamada, Y., Lindenbaum, O., Negahban, S., Kluger, Y., 2020. Feature selection  
 699 using stochastic gates, in: III, H.D., Singh, A. (Eds.), Proceedings of the 37th  
 700 International Conference on Machine Learning, PMLR. pp. 10648–10659. URL:  
 701 <https://proceedings.mlr.press/v119/yamada20a.html>.
- 702 Ying, X., 2019. An overview of overfitting and its solutions. Journal of Physics: Con-  
 703 ference Series 1168, 022022. URL: [https://dx.doi.org/10.1088/1742-6596/](https://dx.doi.org/10.1088/1742-6596/1168/2/022022)  
 704 [1168/2/022022](https://dx.doi.org/10.1088/1742-6596/1168/2/022022), doi:[10.1088/1742-6596/1168/2/022022](https://doi.org/10.1088/1742-6596/1168/2/022022).
- 705 Zou, H., Hastie, T., 2005. Regularization and variable selection via the elastic net.  
 706 Journal of the Royal Statistical Society Series B: Statistical Methodology 67, 301–  
 707 320.

708 **Appendix A. Results of synthetic Datasets 1**

Table A.4: Averaged results of DNFS compared to FIDL, STG, CA, and SHAP on the synthetic datasets. Green highlights the method with the most optimal selections.

		Linear			Non-Linear		
	Method	Val.	N	Opt.	Val.	N	Opt.
Regression	<b>DNFS</b>	<b>100</b>	<b>4</b>	<b>100</b>	<b>100</b>	<b>4.132</b>	<b>95</b>
	<i>FIDL</i>	54	4	54	90	4	90
	<i>STG</i>	100	4	0	100	3	0
	<i>CA</i>	100	4	0	100	4	0
	<i>SHAP</i>	69	4	69	60	4	60
	<i>PCC</i>	100	8	0	100	7	0
Classification	<b>DNFS</b>	<b>100</b>	<b>4</b>	<b>100</b>	<b>100</b>	<b>4.325</b>	<b>93</b>
	<i>FIDL</i>	14	4	14	85	4	85
	<i>STG</i>	100	4	0	100	4	100
	<i>CA</i>	100	4	0	100	4	0
	<i>SHAP</i>	58	4	58	45	4	45
	<i>PCC</i>	0	0	0	0	0	0

709 **Appendix B. Proof of proposition 1**

710 Here we provide a proof of proposition 1, which in our case tells us that every  
711 possible penalization function can further be reduced to zero.

*Proof.* Since  $f$  is non-constant there exists  $z \in \mathbb{R}^n$  such that  $f(z) \neq 0$ . Since  $f$  is continuous, we can use the intermediate value theorem for multivariate (Shimamoto, 2019) functions and deduce that for all

$$\xi \in [f(x), f(z)] = [0, f(z)] \subset \mathbb{R}$$

712 there exists an  $y \in [x, y]$  such that  $f(y) = \xi$ . Expecially, for all  $\varepsilon > 0$  we can find an  
713  $y$  such that  $|f(y)| < \varepsilon$ .

Table A.5: Results of DNFS compared to FIDL, CA, STG, and SHAP on the benchmark datasets. Green highlights the method with the most optimal selections.

	MNIST			Ames Housing			Basehock		
Method	Val.	Acc. [%]	N	Val.	MAE [ $\cdot 10^3 \$$ ]	N	Val.	Acc [%]	N
<b>DNFS</b>	<b>97</b>	<b>48</b>		<b>15</b>	20		<b>95</b>	<b>53</b>	
<i>FIDL</i>	93	50		16	20		91	50	
<i>STG</i>	96	47		15	23		93	51	
<i>CA</i>	97	50		17	20		92	50	
<i>SHAP</i>	92	50		16	20		92	50	

714 Note that this is a more general result than necessary in our case. The result also  
715 covers  $z$  and (or)  $f(z)$  being negative, which is impossible for a sensible penalization  
716 function  $p$ . Also we have the special case of  $f(x) = f(0) = 0$  which is also covered  
717 by this result.  $\square$

## 718 Appendix C. Performance validation

719 The goal of the feature reduction task is to retain strong performance in the  
720 downstream supervised task, even after significantly reducing the number of input  
721 features. In Section 5, the performance of the task model is compared when trained  
722 on the full feature set and on the reduced set selected by DNFS, to ensure that the  
723 reduction does not negatively impact the task outcome. To further validate the effec-  
724 tiveness and generalizability of the selected features, we additionally trained an SVM  
725 or, in the case of regression tasks, a support vector regression (SVR) model with an  
726 RBF kernel, regularization parameter  $C = 1.0$ , and  $\epsilon = 0.01$ . These models were  
727 evaluated on the benchmark datasets using both the full and the reduced feature sets.  
728 The results, presented in Table C.6, are consistent with those in Table 2, confirming  
729 that DNFS reliably selects a compact yet effective subset of features without degrad-  
730 ing task performance. For the Basehock dataset, only the classes "sci.space" and  
731 "rec.autos" were used in a binary classification setting, due to computational con-  
732 straints that prevented successful training of the SVM on the full multiclass dataset.  
733 However, this does not affect the validity of the performance evaluation.

Table C.6: Validation performance of SVM/SVR models trained on the full and DNFS-reduced feature sets across benchmark datasets. \*For the Basehock dataset, only the classes "sci.space" and "rec.autos" were used in the evaluation due to computational constraints.

DATASET	BASELINE	SEL. PERF.	N
MNIST	95, 2%	93, 5% $\pm$ 0, 2%	50 $\pm$ 2
AMES H.	18.191	15.000\$ $\pm$ 500\$	20 $\pm$ 3
BASEHOCK*	99%	98.8% $\pm$ 0.2 %	53 $\pm$ 4

## Appendix D. Stable statistical selection

Deciding when to stop the selection can be achieved either by training for a fixed number of epochs or by classic early stopping criteria like validation loss monitoring of the task model's performance. The latter has been implemented in the current version and leads to a severe decrease in runtime (cf. Appendix E.4) while maintaining high performance and a good selection.

Figure D.8 shows how often a given feature was classified as selected in the selection interval for different thresholds  $\delta$ . Most of the features are selected either in nearly 100% or 0% of epochs, leaving just a small portion of features not decisively classified. The heatmap in Figure D.7 shows that changes in  $\delta$  and  $p$  have minimal impact on the number of selected features, which is the desired behavior. Therefore, it suffices to use the default values (cf. Appendix E) when performing selection with DNFS. Convergence tends to get stabler for larger datasets with more features. The most unstable convergence was yielded on the synthetic classification. Here, changes in  $\delta$  have the largest impact.

## Appendix E. Experimental setup

### Appendix E.1. Comparative methods

The experimental setups are similar but vary slightly from task to task. For the regularization the built in TensorFlow<sup>1</sup> implementation was used. For SHAP<sup>2</sup> and FIDL<sup>3</sup> the authors' implementations were used. For comparison, the following

<sup>1</sup>TensorFlow: <https://github.com/tensorflow/tensorflow>

<sup>2</sup>SHAP: <https://github.com/slundberg/shap>

<sup>3</sup>FeatureImportanceDeepLearning: <https://github.com/maksym33/FeatureImportanceDL>

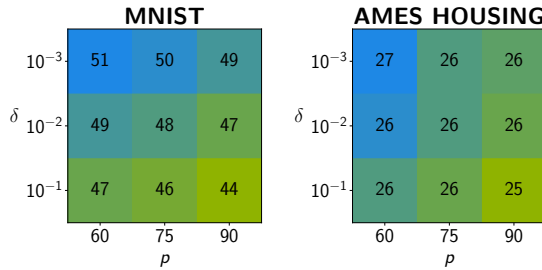


Figure D.7: The number of selected features for different  $p$  and  $\delta$  on the benchmark datasets.

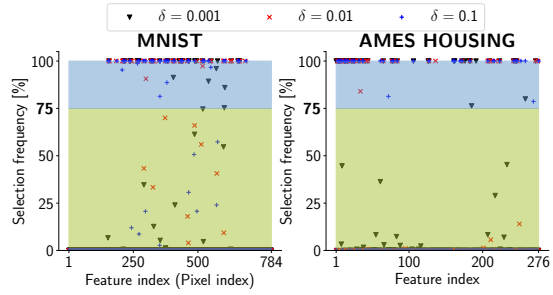


Figure D.8: Distribution of the statistically selected features for different thresholds  $\delta$  on the benchmark datasets.

methods were chosen, all of which were described detailed in section 2. The required hyperparameters were set as follows:

The **CA** parameters were set according to the guidelines and examples given by the authors [Balm et al. \(2019\)](#). The model architectures were similar to the ones chosen for implementing DNFS.

For **FIDL** the same networks and setups as described in ([Wojtas and Chen, 2020](#)) were used for the synthetic datasets and the MNIST dataset. For the Ames housing task, the same task (operator) and selection (selector) networks as in our method (cf. Table E.8) were used.

For calculating **SHAP** values, there are no hyperparameters. To speed up the calculations for the MNIST dataset, a randomly generated subset of 10% the original size was used.

## Appendix E.2. Setup

For all tasks, the penalization function  $p$  was employed as the arithmetic mean over all mask entries of  $\mathbf{m}$  since this yielded stable results. The arithmetic mean has some unsuitable properties, which can be adverted by rethinking the convergence of the method as described in section 6.4.

## Appendix E.3. Networks

Both networks are MLPs, with the selection network incorporating dropout layers between the fully connected ones. To mitigate the impact of the dropout layers on the final mask it's calculated not as the last output of the selection model, but rather as described in 6.5. To mitigate the vanishing gradient problem([Basodi et al., 2020](#)) the rectified linear unit (ReLU) activation function was selected for the hidden

Table E.7: Hyperparameters and their values for the new method.  $\lambda$  denotes the value chosen for the implementation and  $\lambda_{calc}$  the value approximated as presented in section 6.3. For  $\lambda_{calc}$  to deliver the expected results, the target has to be normalized.

DATASET	EPOCHS	$\lambda$	$\approx \lambda_{calc}$	BATCH SIZE	P	$\delta$
SYNTHETIC REG.	400	2	33	32	75%	0.01
SYNTHETIC CLASS.	400	3	1,5	32	75%	0.01
SYNTHETIC REG. N.L.	400	2	33	32	75%	0.01
SYNTHETIC CLASS. N.L.	400	3	1,5	32	75%	0.01
MNIST	300	10	0,5	64	75%	0.01
AMES HOUSING	600	$10^{-9}$	20	32	75%	0.01
BASEHOCK	200	0,5	0,3	64	75%	0.01

layers in the task and selection model, while a linear (ReLU if target was min-max-transformed) activation function (resp. a softmax activation function) was chosen for the output layer in the case of a regression (resp. classification) task.

#### Appendix E.4. Runtime and hardware

Table E.9 shows the runtime requirements of DNFS on the various datasets as well as the network complexity. As described in section 6.6, DNFS is more computationally intensive than its competitors SHAP and L1 regularization but is also the only method that delivers reliably good performance. Further, it outperforms FIDL, the only approach comparable to DNFS concept-wise. It is highly probable that further code optimization could reduce the computing requirements significantly, for example, via swapping the order in which the task and selection model update. Greater utilization of built-in keras methods also has the potential of reductions in runtime since those are highly optimized. Furthermore, it is likely that for a successful selection just a fraction of the dataset is sufficient, thereby reducing time per epoch. First, tests for the MNIST dataset have delivered promising results. Extensive code optimization is beyond the scope of this work, and therefore, we focused on demonstrating the feasibility and effectiveness of DNFS, leaving performance tuning and advanced optimizations for future research.

The time measurements have been conducted on a computing node with 16 CPU cores and 32GB of memory. The CPUs were Intel Cascade Lake Xeon Gold 6230N

Table E.8: Task and selection networks that were used for the different datasets. The dropout layers in the selection network are not shown. In the case of the synthetic datasets the network architecture remained the same, whether the dependencies were linear or non-linear.

DATASET	TASK NETWORK	SELECTION NETWORK
SYNTH. REG.	$\rightarrow 20 \rightarrow 100 \rightarrow 50 \rightarrow 50 \rightarrow 25 \rightarrow 1$	$\xrightarrow{d=0.5} 100 \xrightarrow{d=0.5} 75 \xrightarrow{d=0.5} 50 \xrightarrow{d=0.5} 10$
SYNTH. CLASS.	$\rightarrow 20 \rightarrow 60 \rightarrow 30 \rightarrow 20 \rightarrow 10 \rightarrow 2$	$\xrightarrow{d=0.5} 100 \xrightarrow{d=0.5} 75 \xrightarrow{d=0.5} 50 \xrightarrow{d=0.5} 10$
MNIST	$\rightarrow 242 \rightarrow 100 \rightarrow 10$	$\xrightarrow{d=0.5} 256 \xrightarrow{d=0.5} 784$
AMES	$\rightarrow 128 \rightarrow 64 \rightarrow 32 \rightarrow 1$	$\xrightarrow{d=0.5} 256 \xrightarrow{d=0.5} 276$
BASEHOCK	$\rightarrow 128 \rightarrow 512 \rightarrow 256 \rightarrow 2$	$\xrightarrow{d=0.5} 512 \xrightarrow{d=0.5} 256 \xrightarrow{d=0.5} 512 \xrightarrow{d=0.5} 7862$

797 (20-core, 2.3GHz, 30MB Cache, 125W)

## 798 Appendix F. Regarding code and implementation

799 The whole method is available in a GitHub repository that can be accessed via  
800 [this link](#). There, too, is described what steps to take if one desires to reproduce the  
801 results. The concrete implementation can be obtained from the code if of interest, but  
802 some key points will be discussed in the following.

803 DNFS is implemented using TensorFlow and utilizing the Keras API<sup>4</sup> for ease of  
804 use and implementation. We employed Adam(Kingma and Ba, 2017) as the optimizer  
805 solving the optimization problem. For simplicity, the generated metrics and masks  
806 are averaged over all batches of an epoch since the changes are minimal, and there  
807 would have been no point in storing every metric for every update.

808 Further, the implementation has just been tested on Linux systems. Theoretically,  
809 there should be nothing preventing you from testing the implementation on  
810 Windows, except that you perhaps have to change the optimizer to the non-legacy  
811 one. Moreover, we cannot make any statements about macOS and macOS on ARM,  
812 but due to the common Unix roots, the implementation will likely work seamlessly  
813 if all packages are installed.

---

<sup>4</sup>Keras: <https://keras.io/>



Table E.9: Overview of training runtime (ms per epoch) and network size (number of learnable parameters) of DNFS across various datasets.

DATASET	$\approx t_{base}$	$\approx t_{DNFS}$	$\approx$ INCREASE	$N_{task}$	$N_{sel.}$
SYNTHETIC REG.	51 ms	71 ms	1,4x	11.221	12.985
SYNTHETIC CLASS.	53 ms	69 ms	1,3x	3.972	12.985
MNIST	300 ms	800 ms	2,6x	20.6172	402.448
BASEHOCK	530 ms	860 ms	1,6x	2.621.698	5.246.974
AMES HOUSING	96 ms	136 ms	1,41x	45.825	141.844

## Appendix G. More results on benchmark datasets

Figure G.9 visualizes the selections of DNFS on a subset of the MNIST dataset that contains just the digits three and eight. Figure G.9a shows that DNFS predominantly selects pixels at which the two digits are distinguishable.

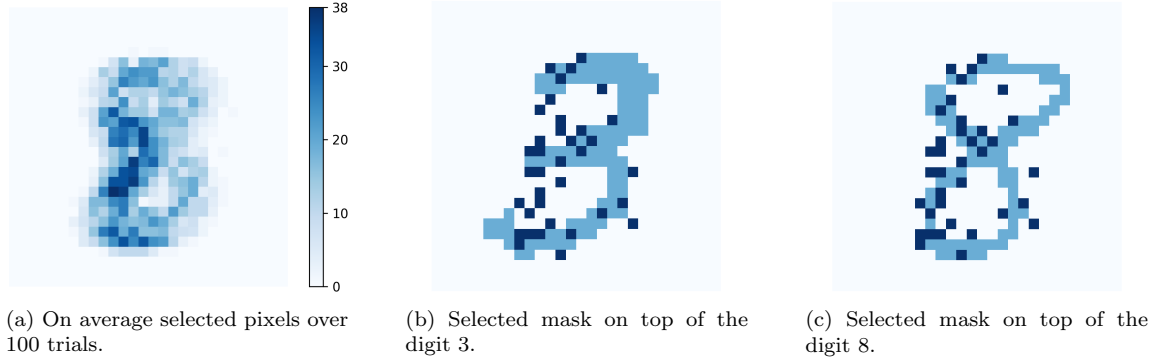


Figure G.9: Visualization of the selected masks on the MNIST-Subset dataset with just the classes 3 and 8. The mask displayed in Figures G.9b and G.9c contains 36 pixels and yields a validation accuracy of 99%.

## Appendix H. Dataset overview

Table H.10 offers a succinct summary of essential dataset metrics. The table below presents key statistics, including the number of classes, features, and overall dataset size.

Table H.10: Summary of key properties of synthetic and benchmark datasets.

DATASET	$ \mathcal{V} $	$ \mathcal{D} $	TRAIN-VAL. RATIO	CLASSES	DATA
SYNTH. REG.	10	1500	66%	REG.	NUMERIC
SYNTH. CLASS.	10	1500	66%	2	NUMERIC
SYNTH. REG. N.L.	10	1500	66%	REG.	NUMERIC
SYNTH. CLASS. N.L.	10	1500	66%	2	NUMERIC
MNIST	784	10.000	80%	10	IMAGE
BASEHOCK	7862	1992	80%	20	TEXT
AMES HOUSING	276	2930	80%	REG.	NUMERIC