003

004

006

007

008 009 010

011

013

014

015

016

017

018

019

021

025

026

027

028 029 030

031

# LCEN: A NOVEL FEATURE SELECTION ALGORITHM FOR NONLINEAR, INTERPRETABLE MACHINE LEARN-ING MODELS

Anonymous authors

Paper under double-blind review

## Abstract

Interpretable models can have advantages over black-box models, and interpretability is essential for the application of machine learning in critical settings, such as aviation or medicine. LASSO and elastic net, the most commonly used interpretable methods, are limited to linear predictions and have poor feature selection capabilities. Other important interpretable methods, such as tree-based or generalized additive models, are nonlinear but have limited performance in some scenarios. In this work, we introduce the LASSO-Clip-EN (LCEN) algorithm for the construction of nonlinear, interpretable machine learning models. LCEN is tested on a wide variety of artificial and empirical datasets, frequently creating more accurate, sparser models than other models, including those for building sparse, nonlinear models. LCEN is robust against many issues typically present in datasets and modeling, including noise, multicollinearity, data scarcity, and hyperparameter variance. LCEN is also able to rediscover multiple physical laws from empirical data and, for processes with no known physical laws, LCEN achieves better results than many other dense and sparse methods – including using 10.8fold fewer features than dense methods and 8.1-fold fewer features than EN on one dataset, and is comparable to or better than ANNs on multiple datasets.

### 1 INTRODUCTION

Models are powerful tools to explain, predict, or describe natural phenomena (Shmueli, 2010). They connect independent variables (also called "inputs" or "features") to dependent variables (also called "outputs" or "labels"). Many modeling methods and algorithms exist, including linear, ensemblebased, and deep learning models. Complex models are claimed to have greater capability to model phenomena due to their lower bias, but their intricate and numerous mathematical transformations prevent humans from understanding how an output was predicted by a model, or the relative or absolute importance of the inputs. Moreover, a lack of transparency may prevent the model from being trusted in critical or sensitive applications (Hong et al., 2020).

In a modeling context, interpretability can be defined as "how an output y = f(X) was predicted 040 for a given input X – that is, provide  $f(\cdot)$  in a form readily understandable to humans so that the 041 model's outputs may be explainable". There are two main methods to increase interpretability: the 042 use of model-agnostic algorithms, which extract interpretable explanations a posteriori and work 043 for any model, or the direct use of interpretable models (Ribeiro et al., 2016). Interpretable models 044 include "decision trees, rules, additive models, attention-based networks, and sparse linear models" 045 (Ribeiro et al., 2016). It should be noted that nonlinear models may also be made sparse, and even 046 interpretable, as described later in this section and the rest of this work. Interpretable models can 047 have many advantages over black-box or *a posteriori* explanations, including the ability to assist 048 researchers in refining the model and data, or better highlighting scenarios in which the model fails 049 or lacks robustness (Rudin, 2019). Special attention should be given to sparse models, which identify 050 the most important features, can make the model more robust to variations in the input data, and can 051 significantly improve the model's interpretability if an interpretable model is used (Rudin, 2019). A sparse model may be defined as "a model that uses few input features, particularly relative to the 052 total number of features available". However, even a linear model or decision tree/rules can become unwieldy and challenging to interpret if hundreds or thousands of coefficients or rules are present.

054 Feature selection is the process of selecting the most important features in a model to increase its robustness, interpretability, or sparsity. Many criteria for feature selection exist (Heinze et al., 2018), 056 including significance based on p-values (using a univariate, iterative/stepwise, or global method), 057 using information criteria (such as the AIC (Akaike, 1974) and BIC (Schwarz, 1978)), using penal-058 ties (such as in LASSO (Santosa & Symes, 1986; Tibshirani, 1996) and elastic net (EN) (Zou & Hastie, 2005)), criteria based on changes in estimates, and expert knowledge. More broadly, these 059 methods can be classified as filter, wrapper, or integrated methods. While no method is superior 060 for all problems, different works have evaluated and criticized these criteria. For example, stepwise 061 regression is one of the most commonly used methods in many fields thanks to its computational 062 simplicity and ease of understanding (Whittingham et al., 2006; Heinze et al., 2018; Smith, 2018). 063 However, stepwise regression is prone to ignoring features with causal effects, including irrelevant 064 features, generating excessively small confidence intervals, and producing incorrect/biased parame-065 ters (Whittingham et al., 2006; Smith, 2018). LASSO is simple and computationally cheap, and has 066 performed well for some problems (Hebiri & Lederer, 2013; Tian et al., 2015; Pavlou et al., 2016), 067 but can overselect irrelevant variables, tends to select only as many features as there are samples, 068 and does not handle multicollinear data well (Heinze et al., 2018; Zou & Hastie, 2005).<sup>1</sup>

069 Originally, most feature selection methods applied only in linear contexts (or have been applied 070 primarily in linear contexts). For example, the only sparse models referenced in the highly cited 071 review by Ribeiro et al. (2016) are linear. The most commonly used sparse methods (LASSO, EN, 072 and their variants) are linear regressors.<sup>2</sup> To address this limitation, later works consider sparse 073 nonlinear models. McConaghy (2011); Brewick et al. (2017); Sun & Braatz (2020) defined sets 074 of features consisting of polynomials (all works), interactions (all works), and/or non-polynomials 075 (McConaghy, 2011; Sun & Braatz, 2020). ALVEN (Sun & Braatz, 2020) uses an F-test for each feature (including the expanded set of features) to determine whether to keep a feature in the final EN 076 model, a filter approach. However, this F-test has very poor feature selectivity, as nearly all features 077 are selected when traditional values of  $\alpha$  (0.001  $\leq \alpha \leq 0.05$ ) are used. Furthermore, the ordering of the features with respect to their p-values does not follow their relevance, as many irrelevant features 079 are among those with the lowest p-values, and relevant features can be among those with the highest 080 p-values (including  $p \gg 0.05$ ). Other relevant methods include Fan & Li (2001), which introduced 081 the smoothly clipped absolute deviation (SCAD); Ravikumar et al. (2009), which introduced Sparse 082 Additive Models and for which Liu et al. (2022) used with  $L_0$  and  $L_2$  regularization; Zou & Zhang 083 (2009), which introduced the adaptive elastic net; Zhang (2010), which introduced the minimax 084 concave penalty (MCP/MC+); van de Geer et al. (2011), which introduced the thresholded LASSO; 085 Yamada et al. (2018), which used two forms of modified, nonlinear LASSO algorithms and achieves 086 high sparsity with the datasets tested; Bertsimas & Parys (2020), which introduced a novel cutting plane algorithm; and Xu et al. (2021), which introduced Bayesian-symbolic regression approaches. 087

880 More recently, the  $L_1$  regularization has been applied to neural networks for nonlinear feature se-089 lection. In its simplest form, group LASSO is applied to zero all the outputs of some neurons (spar-090 sifying the network and eliminating features when zeroing input-layer neurons) (Dinh & Ho, 2020; 091 Scardapane et al., 2017; Wang et al., 2021). LassoNet, a slightly modified version of this algorithm, 092 applies the  $L_1$  penalty only to the input layer and includes a skip-connection between that layer and the output layer (Lemhadri et al., 2021). More complex applications of this method include the 093 multi-modal neural networks of Zhao et al. (2015), the concrete autoencoders of Balin et al. (2019), and the teacher-student network of Mirzaei et al. (2020). The first two are notable for being at least partially unsupervised methods, suggesting they can select the most relevant features for a given 096 dataset no matter the task. Some works have also used approaches other than the  $L_1$  norm for neural network-based feature selection, such as the  $L_0$  norm (Yamada et al., 2020). While these deep learn-098 ing models are powerful tools, two considerable limitations are present: first, they do not provide any information on how the selected features are contributing to the final prediction, significantly 100 limiting interpretability. A posteriori methods to extract this information have been found unreliable 101 (Rudin, 2019), even if useful. Second, these complex model architectures may take "shortcuts" to 102 make apparently accurate predictions (Rosenzweig et al., 2021; Lapuschkin et al., 2019). However,

 <sup>&</sup>lt;sup>104</sup> <sup>1</sup>This last point is somewhat controversial in the literature, see Hebiri & Lederer (2013); Dalalyan et al. (2017), for example.

 <sup>&</sup>lt;sup>2</sup>Zhang et al. (2016) has claimed that, if thousands of samples are available, LASSO can consistently select features even in nonlinear settings, although the coefficients may be distorted. However, this consistency is disputed by Dinh & Ho (2020).

these "shortcuts" are not really relevant to the task, preventing proper generalization and human interpretation.

To create nonlinear, interpretable machine learning models with high predictive and descriptive 111 power, we propose the LASSO-Clip-EN (LCEN) algorithm. This algorithm generates an expanded 112 set of nonlinear features (such as in ALVEN) and performs feature selection and model fitting. This 113 feature set expansion, together with the Clip step, provide LCEN with the ability to do nonlin-114 ear predictions. The algorithm is tested on artificial and empirical data, successfully rediscovering 115 physical laws from data belonging to multiple different areas of knowledge with errors < 2% on the 116 coefficients, a value within the empirical noise of the datasets. On datasets from processes whose 117 underlying physical laws are not yet known, LCEN attains lower root mean square errors (RM-118 SEs) than many sparse and dense methods, leads to sparser models than all but one method tested, and simultaneously runs faster than most alternative methods. While formal theoretical proofs are 119 beyond the scope of this work due to space limitations, important previous works that proved desir-120 able theoretical properties of the thresholded LASSO (a LASSO-Clip model) include Zhou (2009), 121 Meinshausen & Yu (2009), Zhou (2010), and van de Geer et al. (2011). These works provide a 122 theoretical scaffold to partially justify the high performance of LCEN. 123

124 125

126

## 2 Methods

127 The LCEN algorithm (Algorithm 1) has five hyperparameters: *alpha*, which determines the regu-128 larization strength (as in the LASSO, EN, and similar algorithms); *l1\_ratio*, which determines how 129 much of the regularization of the EN step depends on the 1-norm as opposed to the 2-norm (as in the EN algorithm); degree, which determines the maximum degree for the basis expansion of the 130 data (Table A2); lag, which determines the maximum number of previous time steps from which X131 and y features are included (relevant only for dynamic models); and *cutoff*, which determines the 132 minimum value a scaled parameter needs to have to not be eliminated during the clip steps. Details 133 on the cross-validated hyperparameter values for all models are in Section A3. Three other hyperpa-134 rameters are relevant to the expansion of features (Algorithm 2) but do not interact with the LCEN 135 algorithm directly. The *trans\_type* hyperparameter controls what kinds of features are appended to 136 the data. It can be set to 'all' to include all transforms (see Table A2 for an example of what features 137 are included), 'poly' to include only polynomial and interaction terms, and 'simple\_interaction' to 138 include only interaction terms. The *interaction* hyperparameter is a boolean that controls whether 139 interaction terms are included in the feature expansion process. Finally, the *transform\_y* hyperpa-140 rameter, which is only relevant when lag > 0, is a boolean that determines whether the y features from previous time steps will also be transformed based on *trans\_type* or only the raw values of y 141 from previous time steps will be included. 142

143 The LCEN algorithm (Algorithm 1) begins with the LASSO step, which sets the *l1\_ratio* to 1. Five-144 fold cross-validation on the training set is employed among all combinations of *alpha*, *degree*, and 145 lag values. The training dataset is split randomly for each fold (as per sklearn's KFold function). 146 First, additional features are temporarily appended to the data based on the *degree*, *lag*, *trans\_type*, interaction, and transform\_y hyperparameters. We developed a custom algorithm (Algorithm 2) to 147 perform this feature expansion. Due to this dependency on the *degree* and *lag* hyperparameters, 148 this feature expansion occurs within the cross-validation process, creating a temporary augmented 149 dataset that is scaled to have mean = 0 and standard deviation = 1 and then input to the LASSO 150 method. For each hyperparameter combination and fold, a validation mean squared error (MSE) is 151 recorded. The values of *degree* and *lag* corresponding to the LASSO model with the lowest valida-152 tion MSE (averaged across all five folds) are recorded, and a LASSO model using this combination 153 of hyperparameters is fit using the training data to obtain scaled parameters (estimated coefficients). 154

The next step in the LCEN algorithm is the clip step, in which the features whose scaled LASSO parameters have absolute values smaller than the *cutoff* hyperparameter are recorded so they can be removed from the expanded dataset, and their coefficients are forced to 0.

The EN step involves cross-validation on the training set among all combinations of *alpha* and *l1\_ratio*, using the values of *degree* and *lag* obtained in the LASSO step. Once again, the training dataset is split randomly for each fold (as per sklearn's KFold function) and the features are expanded and scaled, then the features recorded in the Clip step are removed. For each hyperparameter combination and fold, a validation MSE is recorded. The values of *alpha* and *l1\_ratio* corresponding

to the EN model with the lowest validation MSE (averaged across all five folds) are recorded, and
 an EN model using this combination of hyperparameters is fit using the training data to obtain new
 scaled parameters.

A second clip step is done on these EN scaled parameters, zeroing the coefficients of the features whose scaled EN parameters have absolute values smaller than the *cutoff* hyperparameter. Lastly, some post-processing steps are done. The scaled coefficients are unscaled by multiplying the scaled coefficients by the standard deviation of the y training data and dividing by the standard deviation of each corresponding X feature. Then, a dot product of the train or test data and the unscaled coefficients is taken to obtain the final predictions. This procedure returns the trained EN model after the second clip step, which is interpretable and nonlinear, and the predictions made with the unscaled coefficients on the training and testing data.

173

174 Algorithm 1 LASSO-Clip-EN (LCEN) 175 Input: X and y data; lists of hyperparameters alpha, l1\_ratio, degree, lag; hyperparameters cutoff, 176 *trans\_type*, *interaction*, *transform\_y* 177 # LASSO step 178 Temporarily set *l1* ratio = 1. 179 for each hyperparameter combination in  $(alpha \times degree \times lag)$  do 180 Generate additional features based on the trans\_type, interaction, transform\_y, the current de-181 gree, and the current lag hyperparameters [Algorithm 2]. Temporarily append the new features to the X data for cross-validation. 182 Scale the data such that each feature's mean = 0 and its standard deviation = 1. 183 Perform five-fold cross-validation with random data splits and the LASSO method. 184 For each fold, record the validation MSE for this hyperparameter combination. 185 end for Obtain the combination of hyperparameters with the lowest average validation MSE from the 187 above cross-validation. Record the best *degree* and *lag* hyperparameters. 188 Fit a LASSO model on the scaled training data with these hyperparameters to obtain parameters. 189 # Clip step 190 **Record** all features whose scaled parameters have absolute values < *cutoff* from the training and 191 test data for removal during the EN training. 192 # EN step Restore *l1\_ratio* to its original list of values. 193 for each hyperparameter combination in (*alpha*  $\times$  *l1\_ratio*) do 194 Generate additional features based on the *trans\_type*, *interaction*, *transform\_y*, the optimal 195 *degree*, and the optimal *lag* hyperparameters [Algorithm 2]. 196 Temporarily append the new features to the X data for cross-validation. 197 Remove the features recorded in the Clip step. Scale the data such that each feature's mean = 0 and its standard deviation = 1. 199 Perform five-fold cross-validation with random data splits and the EN method. 200 For each fold, record the validation MSE for this hyperparameter combination. 201 end for 202 Obtain the combination of hyperparameters with the lowest average validation MSE from the 203 above cross-validation. Record the best *alpha* and *l1\_ratio* hyperparameters. 204 Fit an EN model on the scaled training data with these hyperparameters to obtain parameters. # Clip step II 205 Remove all features whose scaled parameters have absolute values < *cutoff*. 206 # Post-processing 207 Unscale the coefficients of the selected parameters based on the standard deviations of the data. 208 Obtain train/test predictions by performing a dot product of the unscaled coefficients with the 209 expanded train/test data containing only the selected features. 210 return the trained EN model and the predictions. 211 212 213 The rationale behind this sequence of steps – LASSO, then Clip, then EN, then a second Clip –

The rationale behind this sequence of steps – LASSO, then Clip, then EN, then a second Clip –
 seeks to balance the algorithm's high accuracy and sparsity with a low runtime. As shown by the
 ablation experiments (Section A4), other combinations/variants do not achieve the same performance as LCEN. Specifically, the LASSO-Clip, EN-Clip, and LASSO-EN combinations tend to

- 216 have lower accuracy than LCEN. The EN-Clip combination is also much slower and less sparse 217 than LCEN, and the LASSO-EN combination is slower and slightly less sparse. The LASSO-Clip-218 LASSO combination is less accurate than LCEN, although it is slightly faster and more sparse. The 219 EN-Clip-EN combination achieves similar accuracy, but is much slower and less sparse than LCEN. 220 It is possible to add a debiasing step at the end by using ordinary least-squares (OLS) to estimate the coefficients of the features selected by LCEN (Belloni & Chernozhukov, 2013). This LCEN→OLS 221 variant model estimates coefficients more accurately for one of the ablation experiments, but per-222 forms worse in terms of test-set MSE on empirical datasets. 223
- The combinations that start with EN are slower than LCEN because EN has a greater number of hyperparameter combinations to be tested, and these combinations are tested with a higher number of features (as the full expanded feature set has not been subject to any selection via the LASSO and Clip steps). These combinations are also less sparse because the  $L_1$  and  $L_2$  norms compete during EN regularization, and a combination that prioritizes the  $L_2$  norm may have a lower cross-validation MSE. Beginning with the LASSO increases the algorithm's sparsity and speed at no accuracy cost.
- The use of hard-thresholding (Clip) steps improves LCEN's accuracy and sparsity. The increase in sparsity also lowers the algorithm's runtime. The second Clip step is less impactful than the first and does not affect the algorithm's runtime, but it still improves LCEN's feature selection capabilities.
  We highlight that the Clip steps operate on the scaled parameters; thus, any issues that may arise due to the (relative) magnitude of the coefficients are not significant.
- LCEN scales like EN. For an  $N \times P$  dataset under *F*-fold cross-validation with *A* potential  $\alpha$  values, *L* potential L<sub>1</sub> ratio values, and *C* potential *cutoff* values, LCEN scales as  $O(NP^2FALC)$ . The *degree* hyperparameter increases the number of features *P* in a supralinear way. The *lag* hyperparameter increases the number of features *P* in a linear way. Conversely, higher values for the *cutoff* hyperparameter decrease the number of features *P*.
- Multiple datasets (summarized in Table A1) are used to test the performance of the LCEN algorithm.
  These datasets can be divided into three categories: artificial data, empirical data from processes with known physical laws, and empirical data from processes with no known physical laws. Further
  description of these datasets and how the artificial datasets are generated is available in Section A2;
  the empirical datasets are also described in Sections A6.1 and 3.2. All models tested in this work
  had their hyperparameters selected by 5-fold cross-validation. The separation between training and
  testing sets varied depending on the dataset and is detailed in Section A2.
- 247 248

251

- 3 RESULTS
- - 3.1 ARTIFICIAL DATA HIGHLIGHT LCEN'S ROBUSTNESS TO NOISE, MULTICOLLINEARITY, AND HYPERPARAMETER VARIANCE
- 253 The first datasets used to validate the LCEN algorithm are multiple linear datasets ("Artificial lin-254 ear"). These datasets feature all combinations of  $\{100, 500, 1000\}$  samples  $\times \{100, 500, 1000\}$  true 255 features  $\times$  {0%} noise level  $\times$  {25%, 50%, 75%, 100%} additional false features. The noise level is 256 defined as mean(added noise/noiseless y)×100%. This added noise is Gaussian noise with  $\mu = 0$ 257 and a suitable  $\sigma$  to reach the desired noise level. For each combination, 3 repeats with different 258 random seeds were created. This experiment contains multiple challenging conditions, including 259 cases where the number of features was much larger than the number of samples ( $P \gg N$ ), cases 260 with a significant proportion of false features, and cases with a high noise (low correlation between the input X and the output y). 261
- 262 The methods LASSO, EN, fastSparseGAMs (FS-GAMs) (Liu et al., 2022), SCAD, MCP, symbolic 263 regression (SymReg) (implemented by Stephens et al. (2022)), and LCEN were tested on this feature 264 selection task. Overall, LCEN consistently outperformed all other methods in this task, as measured 265 by their Matthews Correlation Coefficients (MCC) and  $F_1$  scores (Figs. A1 to A12) Because most 266 tests contained more positive-class elements (true features) than negative-class elements (false fea-267 tures), it should be noted that the  $F_1$  score can be biased upwards, particularly if a model selects all features. EN was typically completely unselective, classifying all features as true except in a 268 few scenarios with N > P. SymReg performed marginally better, but still had a very low overall 269 performance. LASSO, FS-GAMs, SCAD, and MCP performed better than EN and SymReg, having

similar performances among themselves; notable exceptions were cases in which LASSO was completely unselective, and most scenarios with N > P, which allowed SCAD and MCP to perform perfect classification. LCEN performed perfect classification in the scenarios with N > P even more frequently than SCAD and MCP, and surpassed the other methods in the other scenarios in terms of absolute MCC by 19.8% on average when N = P and by 8.2% on average when N < P.

275 The first step of the LCEN algorithm uses LASSO, which has been claimed to underperform with 276 multicollinear data (Heinze et al., 2018; Zou & Hastie, 2005). Therefore, tests using multicollinear 277 data are done next. The goal is to verify whether LCEN can successfully determine the presence of 278 two different but correlated variables, as LASSO prefers to select only one variable in this scenario 279 (Zou & Hastie, 2005). Noise  $\epsilon_1$ , at different levels (as defined above), was added to the  $X_0$  variable 280 to create a correlated variable  $X_1$ . A second noise  $\epsilon_2$ , also at different levels, was added to the final y data. When  $\epsilon_1 = 0$ , both variables are equal and separation is not possible. However, at other  $\epsilon_1$ 281 values, the LCEN algorithm is very successful at identifying that two relevant variables exist and 282 assigning correct coefficients to them (Figure A13). Specifically, when the noise level  $\epsilon_1$  associated 283 with the X data (which indicates how different the X variables are, as highlighted by the variance 284 inflation factors [VIFs] in Fig. A13) is greater than the noise level  $\epsilon_2$  associated with the y data, 285 LCEN can separate both variables with coefficient errors  $\leq 5\%$ . When both noise levels are similar, 286 LCEN can separate both variables with coefficient errors between 5% and 10%. The X data used in 287 this experiment has very high multicollinearity (as shown by the VIFs); real data will typically have lower VIFs and thus be easier to separate using LCEN. 289

Next, a more complex equation is used to further validate LCEN. The "Relativistic energy" data 290 contain mass and velocity values used to calculate  $E^2 = c^4 m^2 + c^2 m^2 v^2$ . As before, datasets 291 with increasing noise levels are created. The *degree* hyperparameter is allowed to vary between 1 292 and 6 in this experiment. These *degree* values lead to expanded datasets with {8, 22, 42, 68, 100, 293 138} features respectively. Even though there are only two true features, there are a significant 294 number of false features, many of which have similar functional forms to the true features. LCEN 295 selected only relevant features for all noise levels tested ( $\leq 20\%$ ), and the coefficients were typically 296 equal to the ground truth (Table 1). The only major divergence happened at a noise level of 20%, as the coefficient for  $m^2$  had a 25% relative error. This error led to our hypothesizing that it is 297 challenging to distinguish among the features involving m (such as  $m, m^2$ , and  $m^3$ ) due to the low 298 range of the data. Thus, another dataset with the same properties but a larger range of values for 299 *m* is created. LCEN performed better on this dataset, again selecting only relevant features for all 300 noise levels tested (< 30%) and having much lower errors in the estimated coefficients (Table 1). 301 These experiments further highlight the robustness of LCEN and show how the range of the data 302 can affect predictions. To clarify our design choices and the relevance of each individual part of the 303 LCEN algorithm, ablation tests are performed with this dataset in Section A4 of the Appendix. 304

**Table 1:** Coefficient values and corresponding relative error to the ground truth for the "Relativistic energy" dataset at different noise levels. The first coefficient is for  $m^2$  and should be  $c^4 = 8.078 \times 10^{33} \text{ m}^4/\text{s}^4$ ; the second coefficient is for  $m^2v^2$  and should be  $c^2 = 8.988 \times 10^{16} \text{ m}^2/\text{s}^2$ . The left table is for the dataset with  $1 \le m < 10$ , and the right table is for the dataset with  $1 \le m < 100$ .

310

211	Noise	Coefficients	$\mathbf{Frror}(\mathcal{O}_{0})$	Noise	Coefficients	Error (%)
511	INDISC			0%	$8.078 \times 10^{33}$ , $8.987 \times 10^{16}$	0.001. 0.006
312	0%	$8.077 \times 10^{33}, 8.987 \times 10^{10}$	0.013, 0.009	5.01	<u>8 078 1033 8 086 1016</u>	0.005, 0.022
313	5%	8.081 $\times 10^{33}$ 8.969 $\times 10^{16}$	0.043 0.206	5%	$8.0/8 \times 10^{33}$ , $8.980 \times 10^{13}$	0.005, 0.022
	100	0.001/10 , 0.909/10	0.013, 0.200	10%	$8.078 \times 10^{33}, 8.984 \times 10^{16}$	0.009, 0.038
314	10%	$8.085 \times 10^{55}, 8.951 \times 10^{10}$	0.097, 0.410	15%	$8.070 \times 10^{33}$ 8.083 $\times 10^{16}$	0.013 0.054
315	15%	$8.089 \times 10^{33}$ , $8.935 \times 10^{16}$	0.139. 0.580	1570	$0.079 \times 10^{-3}, 0.903 \times 10^{-3}$	0.013, 0.034
	2007	6 0 27 × 1033 8 012 × 1016	25 20 0 844	20%	$8.079 \times 10^{33}, 8.981 \times 10^{10}$	0.017, 0.070
316	20%	$0.027 \times 10^{-3}$ , $8.912 \times 10^{-3}$	23.39, 0.844	30%	$8.080 \times 10^{33}$ 8.978 $\times 10^{16}$	0.025 0.103
047				5570	0.000/10 ,0.970/10	0.025, 0.105

317

Finally, LCEN is compared with the feature selection algorithm in ALVEN (Sun & Braatz, 2020), which uses the same basis function expansion, but uses f-tests for feature selection. The "4thdegree, univariate polynomial" dataset is created as per Sun & Braatz (2021), such that  $y = X + 0.5X^2 + 0.1X^3 + 0.05X^4 + \epsilon$ , 30 X points are available for training, and 1,000 X points are available for testing. These conditions simulate the scarcity of data potentially present in real datasets while ensuring test errors can be predicted with high confidence. Sun & Braatz (2021) created four types of ALVEN models for this prediction. On this same dataset and using the same four

types of models, LCEN attained median errors that are typically over 60% smaller than for ALVEN (Fig. 1). Discussion of these results are included in the Appendix (Section A5), including discussions on how LCEN consistently selected the correct *degree* hyperparameter via cross-validation despite the low number of training samples and high noise (Fig. A15).



**Figure 1:** Test set median MSE for the "4th-degree, univariate polynomial" dataset. ALVEN results (left, reproduced from Sun & Braatz (2021) with permission) show that the error is monotonically increasing with noise and that the *degree* 4 "unbiased model" is the best at low noise levels, but is displaced by the *degree* 2 "biased model" at higher noise levels. On the other hand, LCEN results (right) show that the median errors converge at higher noises. Furthermore, the LCEN median errors are typically over 60% smaller than the ALVEN median errors, and the *degree* 4 "unbiased model" is always the best model no matter the noise. The "noise level" and "Noise variance  $\sigma^2$ " terms are equivalent in this figure. Fig. A14 contains interquartile ranges for the LCEN model's test MSEs.

# 3.2 LCEN SURPASSES MANY OTHER METHODS WHEN MAKING PREDICTIONS ON EMPIRICAL DATA

The applicability of an algorithm to real-world problems is judged only by its performance on real data, as data sparsity or real noise may affect the algorithm's capabilities. Tests done on empirical data generated by processes with known physical laws show that LCEN still displays exceptional feature selection capabilities, consistently selecting only the right features even when high hyperparameter variance – that is, scenarios where many potential combinations of hyperparameters exist (which increase the variance in a bias-variance context) – is present (Section A6.1 and Table 2).

**Table 2:** Summary of LCEN results for the empirical datasets from processes with known physical laws.

Dataset	Max. size	Only Correct Features	Coefficient Relative Error
CARMENES star data	$293 \times 350$	Yes	1.74%
Kepler's 3rd Law (1619)	$6 \times 18$	Yes	0.46%
Kepler's 3rd Law (Modern)	$8 \times 18$	Yes	0.07%

The final experiments to validate LCEN's performance involve comparisons to other algorithms on real datasets from processes with unknown physical laws. As there is no (computational) way to validate the feature selection by models trained on these datasets, the main focuses of this sec-tion are investigating prediction errors and sparsities of different models. The methods ordinary least squares (OLS), ridge regression (RR) (Tikhonov, 1963), partial least squares (PLS) (Wold, 1975a;b), LASSO, elastic net (EN) (Zou & Hastie, 2005), SCAD (Fan & Li, 2001), MCP (Zhang, 2010), sym-bolic regression (SymReg) (Stephens et al., 2022), random forest (RF) (Ho, 1995), gradient-boosted decision trees (GBDT) (Friedman, 2001), adaptive boosting (AdaB) (Freund & Schapire, 1997), support vector machine with radial-basis functions (SVM) (Boser et al., 1992), fastSparseGAMs (FS-GAMs) (Liu et al., 2022), multilayer perceptron (MLP), MLP with group LASSO (MLP-GL<sub>1</sub>) (Scardapane et al., 2017), and LassoNet (Lemhadri et al., 2021) were compared with LCEN. To clarify our design choices and the relevance of each individual part of the LCEN algorithm, ablation tests are performed with many of the datasets tested here in Section A4 of the Appendix.

The first dataset analyzed is the "Diesel Freezing Point" dataset (Hutzler & Westbrook, 2000), which is comprised of 395 diesel spectra measured at 401 wavelengths and used to predict the freezing point of these diesels. The dense, nonlinear methods SVM and MLP had the best prediction perfor-mance, with test RMSEs equal to 4.39 and 4.61 °C respectively (Table 3). They were followed by LCEN, RR, EN, LassoNet, LASSO, and MLP-GL<sub>1</sub>, which had test RMSEs between 4.83 and 4.92 °C. Other methods performed worse, with test RMSEs >  $5.0^{\circ}$ C. For comparison,  $5.0^{\circ}$ C is about 7.5% of the range of the test data, which contains diesels with freezing points between  $-59.5^{\circ}$ C and 6.6°C. The sparsest methods were LCEN, which selected only 36/401 features (9.0%) yet had a pre-diction error only 10.0% higher than that of the best dense method, and FS-GAMs, which selected only 2/401 features (0.5%) but had a prediction error 80.4% higher than the best dense method at its lowest cross-validation MSE. LCEN, FS-GAMs, and SVM were the only nonlinear methods that had a runtime faster than 10 seconds, a speed typically reserved for linear methods. EN and AdaB had runtimes  $\approx 20$  seconds, SCAD and MCP had runtimes  $\approx 30$  seconds, and all other models had runtimes on the order of hundreds or thousands of seconds. LCEN is the only method that combines a low test RMSE, interpretability, and a fast runtime. LASSO, the only other method that also has these properties, has a worse RMSE and sparsity. Finally, an end user could prioritize creating very sparse models, even at the expense of increasing these models' MSEs. To simulate such a scenario, the LCEN cutoff hyperparameter was increased from the value that minimizes the validation MSE to create sparser models. These models have much fewer features, yet their test set RMSEs typically increase by only small values. An LCEN model that selected 36 features had a test set RMSE that was 5.1% lower than that of an FS-GAM with the same number of features, whereas an LCEN model that selected 29 features had a test set RMSE that was 7.9% lower than that of a similar FS-GAM. This illustrates how LCEN can select the most critical features to make models with high sparsity and predictive power, and how these criteria can be prioritized by the end user. 

**Table 3:** Results of different models for the "Diesel Freezing Point" dataset. The number of features that minimizes the cross-validation MSE is 2 for FS-GAMs and 36 for LCEN.

Model	Test RMSE (°C)	Features	Runtime (s)
OLS	11.75	401	0.09
PLS	5.21	401	4.44
RR	4.83	401	4.87
EN	4.83	299	19.6
LASSO	4.90	39	2.06
SCAD	5.20	26	28.2
MCP	5.25	29	33.6
SymReg	5.72	6	696
RF	5.16	390	307
GBDT	5.40	354	2649
AdaB	5.60	300	22.0
SVM	4.39	401	5.33
MLP	4.61	401	121
MLP-GL <sub>1</sub>	4.92	80	569
LassoNet	4.83	401	474
	5.09	36	
	5.33	29	
FS-GAMs	5.26	13	3.75
	6.34	6	
	8.32	2	
	4.83	36	6.54
LCEN	4.91	29	6.27
	5.52	13	5.76
	7.40	6	5.53

The next dataset used is the "Abalone" dataset (Nash et al., 1995). LCEN surpassed all linear models and tied with the best nonlinear models in this task (Table A7 and Section A6.2). However, the other best nonlinear models all lack the interpretability of LCEN. By increasing the *cutoff* hyper-parameter, sparser LCEN models may be generated, which attain similar performances to the base

LCEN model. Moreover, LCEN models at the same or higher sparsities had lower test RMSEs than
 FS-GAMs, further highlighting the abilities of the LCEN algorithm.

LCEN is then tested on the "Concrete Compressive Strength" dataset (Yeh, 1998), which contains 435 the composition and age of 1,030 different types of concrete and their compressive strengths. The 436 relationship between these properties is nonlinear, and previous modeling attempts include algebraic 437 expressions and artificial neural networks (specifically, MLPs) (Yeh, 1998; 2006). These MLPs were 438 superior to the algebraic models, whereas the algebraic models provide interpretability on how the 439 properties of the concrete affect its compressive strength. LCEN is also considerably better than 440 the previously published algebraic models (Table 4), and its performance is competitive with that 441 of previously published ANNs without sacrificing interpretability. Furthermore, note that no type 442 of validation is mentioned in Yeh (1998), so the test and validation sets may be the same, making the MLP figures overoptimistic. LCEN has a validation RMSE of 4.66 MPa on this dataset, which 443 is slightly better than the RMSE of the MLP of Yeh (1998). Some other models surpass LCEN in 444 terms of test RMSE, but LCEN has the lowest validation RMSE out of all methods tested. 445

Table 4: Results of different models for the "Concrete Compressive Strength" dataset. All machinelearning models selected all 8 features except for FS-GAMs, which selected 4, and SymReg, which selected 6. No form of validation is mentioned in Yeh (1998), so the test and validation sets may be the same, making the ANN values overoptimistic.

Model	Test RMSE (MPa)
Algebraic expression (Yeh, 1998)	7.79
MLP (Yeh, 1998)	4.76
Linear+interactions model (Yeh, 2006)	7.43
SCAD = MCP	10.26
SymReg	10.51
RF	5.10
GBDT	7.29
AdaB	6.95
SVM	5.94
$MLP-GL_1$	5.47
LassoNet	5.53
FS-GAMs	11.39
LCEN	5.73

LCEN is also successful at predicting phenomena caused by human activity instead of physical laws. 465 In the modified "Boston housing" dataset (Harrison & Rubinfeld, 1978), LCEN attains a test RMSE 466 that is only 6% higher than that of a dense MLP (Table A8 and Section A6.2). Once again, LCEN 467 attains higher performance than many other methods in this dataset while also being completely 468 interpretable. Finally, the "GEFCom 2014" dataset was used to highlight the ability of LCEN to 469 predict in a complex and dynamic task (Hong et al., 2016). Two versions of the "GEFCom 2014" 470 dataset have been published: one that contains only energy consumption levels and another that 471 contains the same energy consumption data and also temperature data from multiple weather sta-472 tions. This work uses the former. "GEFCom 2014" is part of an energy forecasting competition 473 won by a LASSO-like model (Hong et al., 2016). More recently, deep learning has been applied to 474 this problem (Wilms et al., 2018; Gasparin et al., 2022), and deep learning models have achieved 475 strong 24-hour predictive performance (Gasparin et al., 2022). Despite the strong performance of multiple, complex ANN architectures, LCEN models obtain a 13.1% lower test RMSE on this fore-476 casting task than the state-of-the-art Seq2Seq model from Gasparin et al. (2022) (Table 5). Unlike 477 the ANNs, LCEN requires only a CPU for training and forecasting, and provides interpretable coef-478 ficients. LCEN can also be used for longer forecasts without significant increases in the prediction 479 error, further highlighting the robustness of the algorithm. 480

481

## 4 DISCUSSION

482 483

This work introduces the LASSO-Clip-EN (LCEN) algorithm for the creation of nonlinear, inter pretable machine learning models (Algorithm 1). LCEN is first validated using artificial data (Section 3.1), which provide an initial assessment of the algorithm's performance under multiple, in-

Table 5: Results of different models for the "GEFCom 2014" dataset. The deep learning models (TCN to Seq2Seq) and their results come directly from Gasparin et al. (2022).

489	Metric (mean)	TCN	RNN	LSTM	GRU	Seq2Seq			LCEN		
490	Hours Forecast	24	24	24	24	24	24	48	72	120	168
491	Test RMSE (MW)	17.2	18.0	19.5	19.0	17.1	14.9	18.9	21.0	23.4	24.7
492	Relative Error (%)	9.8	10.2	11.1	10.8	9.7	8.5	10.7	11.9	13.2	13.9

dependently controllable conditions. LCEN was then tested with data from processes with known physical laws (Section A6.1) and without known physical laws (Sections 3.2 and A6.2).

Overall, these experiments have demonstrated the applicability of LCEN to a multitude of scientific and nonscientific problems, even those with significant nonlinearities and complexity. On the real data from processes with known physical laws, LCEN successfully selected only the correct features with very low coefficient errors for all datasets used in this work, effectively rediscovering physical laws solely from data (Table 2). LCEN models were robust to defects in the real data, including noise, multicollinearity, or sample scarcity. LCEN models were typically as accurate as or more accurate than many alternative methods, yet were also much sparser. LCEN models are also trivial to interpret and display exactly how each input is contributing to the final output. This combina-tion of accuracy and interpretability is essential for the deployment of machine-learning models in performance-critical scenarios, from aviation to medicine. Moreover, the additional interpretability can assist in data or model refinement efforts and can make the models robust to changes in data or adversarial input. LCEN is free, open-source, and easy to use, allowing even non-specialists in machine learning to benefit from and use it. The main limitations of LCEN are that it is not a uni-versal function approximator, as it can model only the functions present in the expansion of dataset features, that the feature expansion algorithm is better suited to numerical data over image or text data, and that it sometimes is not as accurate as a dense deep learning method. If a GPU and enough time are available for model training, users in scenarios that focus on accuracy above anything else or with non-numerical data types may prefer to use a deep learning method. 

There are at least two clear future directions for this work. The first involves using the LCEN algorithm in classification tasks, as many important problems in science and engineering involve classification. A comprehensive analysis of the performance of LCEN in classification tasks will follow this paper. The second involves applying the LCEN algorithm to automatically generate physical equations for hybrid model architectures (such as physics-constrained or physics-guided ML), which have high potential for scientific applications (Peng et al., 2021; Willard et al., 2022).

# 540 REFERENCES

552

553

554

555

563

564

565

566 567

568

569

581

582

583

- H. Akaike. A new look at the statistical model identification. *IEEE Transactions on Automatic Control*, 19(6):716–723, 1974. URL https://doi.org/10.1109/TAC.1974.1100705.
- Muhammed Fatih Balın, Abubakar Abid, and James Zou. Concrete autoencoders: Differentiable feature selection and reconstruction. In *Proceedings of the 36th International Conference on Machine Learning*, volume 97, pp. 444–453, 2019. URL https://proceedings.mlr. press/v97/balin19a.html.
- Alexandre Belloni and Victor Chernozhukov. Least squares after model selection in highdimensional sparse models. *Bernoulli*, 19(2):521–547, 2013. URL https://doi.org/10. 3150/11-BEJ410.
  - Dimitris Bertsimas and Bart Van Parys. Sparse high-dimensional regression: Exact scalable algorithms and phase transitions. *The Annals of Statistics*, 48(1):300 323, 2020. URL https://doi.org/10.1214/18-AOS1804.
- Bernhard E. Boser, Isabelle M. Guyon, and Vladimir N. Vapnik. A training algorithm for optimal
   margin classifiers. In *Proceedings of the Fifth Annual Workshop on Computational Learning Theory*, pp. 144–152, 1992. URL https://doi.org/10.1145/130385.130401.
- Patrick T. Brewick, Sami F. Masri, Biagio Carboni, and Walter Lacarbonara. Enabling reduced-order data-driven nonlinear identification and modeling through naïve elastic net regularization. *International Journal of Non-Linear Mechanics*, 94:46–58, 2017. URL https://doi.org/10.1016/j.ijnonlinmec.2017.01.016.
  - D. Clark, Z. Schreter, and A. Adams. A quantitative comparison of dystal and backpropagation. In *Proceedings of the Seventh Australian Conference on Neural Networks*, pp. 132–137, 1996. URL https://www.tib.eu/de/suchen/id/BLCP%3ACN016972815.
  - Arnak S. Dalalyan, Mohamed Hebiri, and Johannes Lederer. On the prediction performance of the Lasso. *Bernoulli*, 23(1):552–581, 2017. URL https://doi.org/10.3150/15-BEJ756.
- Vu C Dinh and Lam S Ho. Consistent feature selection for analytic deep neural networks. In Advances in Neural Information Processing Systems, volume 33, pp. 2420–2431, 2020. URL https://proceedings.neurips.cc/paper\_files/paper/2020/file/1959eb9d5a0f7ebc58ebde81d5df400d-Paper.pdf.
- Jianqing Fan and Runze Li. Variable selection via nonconcave penalized likelihood and its oracle properties. *Journal of the American Statistical Association*, 96(456):1348–1360, 2001. URL https://doi.org/10.1198/016214501753382273.
- Yoav Freund and Robert E Schapire. A decision-theoretic generalization of on-line learning and an application to boosting. *Journal of Computer and System Sciences*, 55(1):119–139, 1997. URL https://doi.org/10.1006/jcss.1997.1504.
  - Jerome H. Friedman. Greedy function approximation: A gradient boosting machine. *The Annals of Statistics*, 29(5):1189–1232, 2001. URL https://doi.org/10.1214/aos/ 1013203451.
- Alberto Gasparin, Slobodan Lukovic, and Cesare Alippi. Deep learning for time series forecasting: The electric load case. CAAI Transactions on Intelligence Technology, 7(1):1–25, 2022. doi: https://doi.org/10.1049/cit2.12060. URL https://ietresearch.onlinelibrary. wiley.com/doi/abs/10.1049/cit2.12060.
- <sup>589</sup> D Jr Harrison and D L Rubinfeld. Hedonic housing prices and the demand for clean air. J. Environ. Econ. Manage., 5:81–102, 1978. URL doi.org/10.1016/0095-0696(78)90006-2.
- Mohamed Hebiri and Johannes Lederer. How correlations influence lasso prediction. *IEEE Transactions on Information Theory*, 59(3):1846–1854, 2013. URL https://doi.org/10.1109/TIT.2012.2227680.

594 Georg Heinze, Christine Wallisch, and Daniela Dunkler. Variable selection - A review and rec-595 ommendations for the practicing statistician. Biometrical Journal, 60(3):431-449, 2018. URL 596 https://doi.org/10.1002/bimj.201700067. 597 Tin Kam Ho. Random decision forests. In Proceedings of 3rd International Conference on Doc-598 ument Analysis and Recognition, volume 1, pp. 278–282, 1995. URL doi.org/10.1109/ ICDAR.1995.598994. 600 601 Sungsoo Ray Hong, Jessica Hullman, and Enrico Bertini. Human factors in model interpretabil-602 ity: Industry practices, challenges, and needs. Proceedings of the ACM on Human-Computer 603 Interaction, 4(68):1-26, May 2020. URL https://doi.org/10.1145/3392878. 604 Tao Hong, Pierre Pinson, Shu Fan, Hamidreza Zareipour, Alberto Troccoli, and Rob J. Hyndman. 605 Probabilistic energy forecasting: Global energy forecasting competition 2014 and beyond. Inter-606 national Journal of Forecasting, 32(3):896-913, 2016. doi: https://doi.org/10.1016/j.ijforecast. 607 2016.02.001. URL https://www.sciencedirect.com/science/article/pii/ 608 S0169207016000133. 609 610 Scott A. Hutzler and S. R. Westbrook. Estimating chemical and bulk properties of middle dis-611 tillate fuels from near-infrared spectra. Technical report, Defense Technical Information Center, U.S. Army TARDEC, Warren, Michigan, 2000. URL https://apps.dtic.mil/sti/ 612 citations/ADA394209. Report TFLRF No. 348. 613 614 J. Kepler, E. J. Aiton, A. M. Duncan, and J. V. Field. The Harmony of the World, pp. 418, 422. 615 American Philosophical Society, 1997. URL https://books.google.com/books?id= 616 rEkLAAAAIAAJ. 617 Sebastian Lapuschkin, Stephan Wäldchen, Alexander Binder, Grégoire Montavon, Wojciech Samek, 618 and Klaus-Robert Müller. Unmasking clever hans predictors and assessing what machines really 619 learn. Nature Communications, 10:1096, 3 2019. doi: 10.1038/s41467-019-08987-4. 620 621 Ismael Lemhadri, Feng Ruan, and Rob Tibshirani. LassoNet: Neural networks with feature spar-622 sity. In Proceedings of The 24th International Conference on Artificial Intelligence and Statis-623 tics, volume 130, pp. 10-18, 2021. URL https://proceedings.mlr.press/v130/ 624 lemhadri21a.html. 625 Jiachang Liu, Chudi Zhong, Margo Seltzer, and Cynthia Rudin. Fast sparse classification for gen-626 eralized linear and additive models. In Gustau Camps-Valls, Francisco J. R. Ruiz, and Isabel 627 Valera (eds.), Proceedings of The 25th International Conference on Artificial Intelligence and 628 Statistics, volume 151, pp. 9304–9333, 2022. URL https://proceedings.mlr.press/ 629 v151/liu22f.html. 630 631 Trent McConaghy. FFX: Fast, Scalable, Deterministic Symbolic Regression Technology, 632 pp. 235–260. Springer, New York, 2011. URL https://doi.org/10.1007/ 978-1-4614-1770-5\_13. 633 634 Nicolai Meinshausen and Bin Yu. Lasso-type recovery of sparse representations for high-635 dimensional data. The Annals of Statistics, 37(1):246-270, 2009. URL https://doi.org/ 636 10.1214/07-AOS582. 637 638 Ali Mirzaei, Vahid Pourahmadi, Mehran Soltani, and Hamid Sheikhzadeh. Deep feature selection using a teacher-student network. Neurocomputing, 383:396-408, 2020. doi: https://doi.org/ 639 10.1016/j.neucom.2019.12.017. URL https://www.sciencedirect.com/science/ 640 article/pii/S0925231219317199. 641 642 Warwick Nash, Tracy Sellers, Simon Talbot, Andrew Cawthorn, and Wes Ford. Abalone. UCI 643 Machine Learning Repository, 1995. URL https://doi.org/10.24432/C55C7W. 644 Menelaos Pavlou, Gareth Ambler, Shaun Seaman, Maria De Iorio, and Rumana Z Omar. Review 645 and evaluation of penalised regression methods for risk prediction in low-dimensional data with 646 few events. Statistics in Medicine, 35(7):1159-1177, 2016. URL https://doi.org/10. 647 1002/sim.6782.

0.40	
648	Grace C. Y. Peng, Mark Alber, Adrian Buganza Tepole, William R. Cannon, Suvranu De, Savador
649	Dura-Bernal Krishna Garikinati George Kamiadakis William W Lytton Paris Perdikaris Linda
650	Data Definar, Histina Guindputt, Goola medaling machina lagrang: What can ya lagra?
651	reizold, and Ellen Kunt. Multiscale findering inees machine learning. What can we learning
1001	Archives of Computational Methods in Engineering, 28:1017–1037, 2021. URL doi.org/10.
652	1007/s11831-020-09405-5.
653	
65/	Pradeep Ravikumar, John Lafferty, Han Liu, and Larry Wasserman. Sparse Additive Models. Jour-
0.0	nal of the Royal Statistical Society Series B: Statistical Methodology, 71(5):1009–1030, 2009.
655	IIRL https://doi.org/10.1111/j.1467-9868.2009.00718 x
656	
657	Marco Tulio Ribeiro Sameer Singh and Carlos Guestrin Model-agnostic interpretability of ma-
CE0	ching lagring. In ICMI Workshop on Human Interpretability in Maching Lagring pp 01 05
000	conte talling. In <i>TextL workshop on Linual accession water the tearning</i> , pp. 91–95,
659	2010. UKL https://arxiv.org/abs/1606.05386.
660	Luis Decomposity Loophing Schooling Hawkey Michael Meels and Menny Abile Detab
661	Juna Rosenzweig, Joachim Sicking, Sebastian Houben, Michael Mock, and Maram Akila. Pach
001	shortcuts: Interpretable proxy models efficiently find black-box vulnerabilities. In <i>IEEE/CVF</i>
002	Conference on Computer Vision and Pattern Recognition Workshops, pp. 56–65, 2021. doi: 10.
663	1109/CVPRW53098.2021.00015.
664	
665	Cynthia Rudin. Stop explaining black box machine learning models for high stakes decisions and
005	use interpretable models instead Nature Machine Intelligence 1:206-215 2019 JIBI https://
666	(Act and 10 1020 (2000) 10 10 10 10 10 10 10 10 10 10 10 10 10
667	//doi.org/10.1038/542256-019-0048-x.
668	Endil Contact and William W. Company Linear investigation of hand limited and stime science service
0000	Facili Santosa and William W. Symes. Linear Inversion of Band-Immited reflection seismograms.
669	SIAM Journal on Scientific and Statistical Computing, 7(4):1307–1330, 1986. URL https:
670	//doi.org/10.1137/0907087.
671	
670	Simone Scardapane, Danilo Comminiello, Amir Hussain, and Aurelio Uncini. Group sparse reg-
012	ularization for deep neural networks. <i>Neurocomputing</i> , 241:81–89, June 2017. URL http:
673	//dx_doi_org/10_1016/i_neucom_2017_02_029
674	// dx.doi.org/10.1010/ j.nedeom.2017.02.025
675	Gideon Schwarz Estimating the dimension of a model The Annals of Statistics 6(2):461–464
676	1079 IDI https://doi.org/10.1014/ooc/1176244126
070	1978. UKL https://doi.org/10.1214/a05/11/0344150.
677	A Schweitzer V M Passagger C Cifuentes V J S Béier M Cortés Contraras I A Caballero
678	A. Schweizer, V. M. Fassegger, C. Chuennes, V. J. S. Bejar, M. Contes-Contrelas, J. A. Cabaneto,
679	C. del Burgo, S. Czesia, M. Kurster, D. Montes, M. R. Zapatero Osorio, I. Ribas, A. Reiners,
0.0	A. Quirrenbach, P. J. Amado, J. Aceituno, G. Anglada-Escudé, F. F. Bauer, S. Dreizler, S. V. Jef-
680	fers, E. W. Guenther, T. Henning, A. Kaminski, M. Lafarga, E. Marfil, J. C. Morales, J. H. M. M.
681	Schmitt, W. Seifert, E. Solano, H. M. Tabernero, and M. Zechmeister. The CARMENES search
682	for exoplanets around M dwarfs. Different roads to radii and masses of the target stars. As-
683	then Astrophysical Science Marcian Difference in the second science of the second science $(25, 26, 20, 20, 10, 10, 10, 10, 10, 10, 10, 10, 10, 1$
000	101. Astrophys. 023. Add, May 2019. OKL https://doi.org/10.1051/0004-0301/
684	201834965.
685	Calif. Showedli, Ta canalain and and dist? Statistical Science, 25(2):280, 210, 2010, UDI, but here
686	Gant Snmuen. To explain of to predict? Stansneal Science, 25(5):289–510, 2010. URL https:
607	//doi.org/10.1214/10-STS330.
007	
688	Gary Smith. Step away from stepwise. Journal of Big Data, 5:32, 2018. URL https://doi.
689	org/10.1186/s40537-018-0143-6.
690	
601	Trevor Stephens et al. Genetic programming in Python with a scikit-learn inspired API: gplearn,
091	2022. URL https://github.com/trevorstephens/gplearn.
692	
693	Weike Sun and Richard D. Braatz. ALVEN: Algebraic learning via elastic net for static and dynamic
694	nonlinear model identification Computers & Chemical Engineering 143:107103 2020 URL
001	https://doi.org/10.1016/j.compchomong.2020.107103
690	neeps.//doi.org/10.1010/j.compenemeng.2020.10/103.
696	Weike Sun and Richard D Braatz Smart process analytics for predictive modeling Comput
697	and R. Chamiagi Engineering 144.107124 2021 LIDI https://doi.org/10.1010//
698	ers & Chemical Engineering, 144:10/134, 2021. UKL https://doi.org/10.1016/j.
600	compcnemeng.2020.10/134.
023	
700	Shaonan Iian, Yan Yu, and Hui Guo. Variable selection and corporate bankruptcy forecasts.
701	<i>Journal of Banking &amp; Finance</i> , 52:89–100, 2015. URL https://doi.org/10.1016/j.jbankfin.2014.12.003.

702 Robert Tibshirani. Regression shrinkage and selection via the lasso. Journal of the Royal Statistical 703 Society: Series B (Methodological), 58(1):267-288, 1996. URL https://doi.org/10. 704 1111/j.2517-6161.1996.tb02080.x. 705 A. N. Tikhonov. Solution of incorrectly formulated problems and the regularization method. Dok-706 lady Akademii Nauk SSSR, 4:1035–1038, 1963. 707 708 Sara van de Geer, Peter Bühlmann, and Shuheng Zhou. The adaptive and the thresholded Lasso 709 for potentially misspecified models (and a lower bound for the Lasso). Electronic Journal of 710 Statistics, 5:688 – 749, 2011. URL https://doi.org/10.1214/11-EJS624. 711 712 Jian Wang, Huaqing Zhang, Junze Wang, Yifei Pu, and Nikhil R. Pal. Feature selection using a neural network with group Lasso regularization and controlled redundancy. IEEE Transactions 713 on Neural Networks and Learning Systems, 32:1110-1123, 2021. URL doi.org/10.1109/ 714 TNNLS.2020.2980383. 715 716 S. Waugh. Extending and Benchmarking Cascade-Correlation: Extensions to the Cascade-717 Correlation Architecture and Benchmarking of Feed-forward Supervised Artificial Neural Net-718 works. PhD thesis, University of Tasmania, 1995. URL https://api.semanticscholar. 719 org/CorpusID:53803349. 720 721 Mark J. Whittingham, Philip A. Stephens, Richard B. Bradbury, and Robert P. Freckleton. Why do we still use stepwise modelling in ecology and behaviour? Journal of Animal Ecology, 75(5): 722 1182-1189,2006.URL https://doi.org/10.1111/j.1365-2656.2006.01141.x. 723 724 Jared Willard, Xiaowei Jia, Shaoming Xu, Michael Steinbach, and Vipin Kumar. Integrating scien-725 tific knowledge with machine learning for engineering and environmental systems. ACM Comput. 726 Surv., 55(4):66, Nov 2022. URL https://doi.org/10.1145/3514228. 727 728 Henning Wilms, Marco Cupelli, and Antonello Monti. Combining auto-regression with exogenous variables in sequence-to-sequence recurrent neural networks for short-term load forecasting. In 729 IEEE 16th International Conference on Industrial Informatics, pp. 673–679, 2018. URL doi. 730 org/10.1109/INDIN.2018.8471953. 731 732 Herman Wold. 11 - Path models with latent variables: The NIPALS approach. In H. M. Blalock, 733 A. Aganbegian, F. M. Borodkin, Raymond Boudon, and Vittorio Capecchi (eds.), Quantitative 734 Sociology, pp. 307–357. Academic Press, New York, 1975a. URL https://doi.org/10. 735 1016/B978-0-12-103950-9.50017-4. 736 Herman Wold. Soft modelling by latent variables: The non-linear iterative partial least squares 737 (NIPALS) approach. Journal of Applied Probability, 12(S1):117-142, 1975b. URL doi.org/ 738 10.1017/S0021900200047604. 739 740 Wolfram Alpha LLC. Wolfram Alpha, 2022. 741 742 Kai Xu, Akash Srivastava, Dan Gutfreund, Felix Sosa, Tomer Ullman, Josh Tenenbaum, 743 and Charles Sutton. A Bayesian-symbolic approach to reasoning and learning in intuitive 744 physics. In Advances in Neural Information Processing Systems, volume 34, pp. 2478–2490, 2021. URL https://proceedings.neurips.cc/paper\_files/paper/2021/ 745 file/147540e129e096fa91700e9db6588354-Paper.pdf. 746 747 Makoto Yamada, Jiliang Tang, Jose Lugo-Martinez, Ermin Hodzic, Raunak Shrestha, Avishek 748 Saha, Hua Ouyang, Dawei Yin, Hiroshi Mamitsuka, Cenk Sahinalp, Predrag Radivojac, Filippo 749 Menczer, and Yi Chang. Ultra high-dimensional nonlinear feature selection for big biological 750 data. IEEE Transactions on Knowledge and Data Engineering, 30(7):1352–1365, 2018. URL 751 doi.org/10.1109/TKDE.2018.2789451. 752 Yutaro Yamada, Ofir Lindenbaum, Sahand Negahban, and Yuval Kluger. Feature selection using 753 stochastic gates. In Proceedings of the 37th International Conference on Machine Learning, 754 volume 119, pp. 10648-10659, 2020. URL https://proceedings.mlr.press/v119/ 755 yamada20a.html.

756 757 758	IC. Yeh. Modeling of strength of high-performance concrete using artificial neural networks. <i>Cement and Concrete Research</i> , 28(12):1797–1808, 1998. URL https://doi.org/10.1016/S0008-8846(98)00165-3.
759 760 761 762	I-Cheng Yeh. Analysis of strength of concrete using design of experiments and neural networks. <i>Journal of Materials in Civil Engineering</i> , 18(4):597–604, 2006. URL https://doi.org/ 10.1061/(ASCE)0899–1561(2006)18:4(597).
763 764 765	I-Cheng Yeh. Concrete Compressive Strength. UCI Machine Learning Repository, 2007. URL https://doi.org/10.24432/C5PK67.
766 767	Cun-Hui Zhang. Nearly unbiased variable selection under minimax concave penalty. <i>The Annals of Statistics</i> , 38(2):894–942, 2010. URL https://doi.org/10.1214/09-AOS729.
768 769 770 771	Yue Zhang, Weihong Guo, and Soumya Ray. On the consistency of feature selection with lasso for non-linear targets. In <i>Proceedings of The 33rd International Conference on Machine Learn-</i> <i>ing</i> , volume 48, pp. 183–191, 2016. URL https://proceedings.mlr.press/v48/ zhangal6.html.
772 773 774 775	<ul> <li>Lei Zhao, Qinghua Hu, and Wenwu Wang. Heterogeneous feature selection with multi-modal deep neural networks and sparse group LASSO. <i>IEEE Transactions on Multimedia</i>, 17(11):1936–1948, 2015. URL doi.org/10.1109/TMM.2015.2477058.</li> </ul>
776 777 778 779	Shuheng Zhou. Thresholding procedures for high dimensional variable selection and sta- tistical estimation. In Advances in Neural Information Processing Systems, volume 22, 2009. URL https://proceedings.neurips.cc/paper_files/paper/2009/ file/92fb0c6d1758261f10d052e6e2c1123c-Paper.pdf.
780 781 782	Shuheng Zhou. Thresholded lasso for high dimensional variable selection and statistical estimation, 2010. URL https://arxiv.org/abs/1002.1583.
783 784 785	Hui Zou and Trevor Hastie. Regularization and Variable Selection Via the Elastic Net. <i>Journal</i> of the Royal Statistical Society Series B: Statistical Methodology, 67(2):301–320, 2005. URL https://doi.org/10.1111/j.1467-9868.2005.00503.x.
786 787 788 790 791 792 793 794 795 796 797 798 799 800 801 802 803 804 803 804 805 806 807 808 809	Hui Zou and Hao Helen Zhang. On the adaptive elastic-net with a diverging number of parameters. <i>The Annals of Statistics</i> , 37(4):1733 – 1751, 2009. URL https://doi.org/10.1214/08-AOS625.

#### 810 **APPENDIX – FEATURE EXPANSION ALGORITHM** A1 811 812 Algorithm 2 Feature expansion for LCEN 813 814 **Input:** X and y data; hyperparameters *degree*, *lag*, *trans\_type*, *interaction*, *transform\_y* 815 if lag > 0 then 816 Append X data from the previous *lag* time steps (samples) to each time step. **if** *transform y* == True **then** 817 Append y data from the previous *lag* time steps (samples) to each time step. 818 end if 819 Discard the first *lag* time steps. 820 end if 821 Using sklearn's PolynomialFeatures function, generate polynomial (and interaction if the hyper-822 parameter *interaction* is True) transforms of the X data for the given *degree*. 823 **if** *trans\_type* == 'all' **then** 824 Generate logarithm transforms. 825 Generate square root transforms for the features with all values > 0. Generate inverse $[1/X_k]$ transforms. 827 if degree > 2 then Generate transforms with noninteger degrees $[(X_k)^{N+1/2}$ for integers N such that |N| <828 *degree*] for the features whose every sample is > 0. 829 Generate the log-inverse transforms $[(\ln X_k)^N/(X_k)^M]$ for natural numbers N and M 830 such that N + M < degree]. 831 Generate the log-sqrt-inverse transforms $[(\ln X_k)^N/(X_k)^{M-1/2}]$ for natural numbers N 832 and M such that N + M < degree - 1 for the features with all values > 0. 833 end if 834 end if 835 return the transformed X features 836

## A2 APPENDIX – DESCRIPTION OF DATASETS USED IN THIS WORK

837 838

839

840 Three types of data are used in this work: artificial data ["Artificial Linear", "Multicollinear data", 841 "Relativistic energy", and "4th-degree, univariate polynomial"], empirical data from processes with known physical laws ["CARMENES star data" and "Kepler's 3rd Law"], and empirical data from 842 processes with no known physical laws ["Diesel Freezing Point", "Abalone", "Concrete Compres-843 sive Strength", "Boston housing", and "GEFCom 2014"]. The artificial data are generated by us as 844 described in the next paragraph. These artificial data are used for an initial assessment of the LCEN 845 algorithm and to investigate how properties of the data, such as noise or data range, affect its feature 846 selection capabilities. Empirical data from processes with known physical laws are described in 847 Section A6.1 and used to verify whether the LCEN algorithm can rediscover known physical laws 848 using data with real properties. Empirical data from processes with no known physical laws are 849 described in Sections 3.2 and A6.2, and used to compare the performance of the LCEN algorithm 850 against other linear and nonlinear models, including deep learning models.

851 The "Artificial Linear" datasets were created by drawing numbers from a uniform distribution be-852 tween -10 and 10 in intervals of 0.1 for the samples X and coefficients k and summing to generate 853 the outputs  $y = \sum_{i=1}^{n\_samples} k_i X_i$ . These datasets feature all combinations of {100, 500, 1000} sam-854 ples  $\times$  {100, 500, 1000} true features  $\times$  {0%} noise level  $\times$  {25%, 50%, 75%, 100%} additional 855 false features. The noise level is defined as mean(added noise/noiseless y) $\times 100\%$ . The "Multi-856 collinear data" dataset was created by drawing numbers from a uniform distribution between 1 and 10 to create one variable  $X_0$ , which was used together with a small amount of noise to create a 858 correlated variable  $X_1 = X_0 + \epsilon_1$ ; finally, they were summed such that  $y = 2X_0 + 2X_1 + \epsilon_2$ . The 859 "Relativistic energy" dataset was created by drawing numbers from a uniform distribution between 1 and 10 or 1 and 100 for masses, and  $5 \times 10^7$  and  $2.5 \times 10^8$  for velocities, which represent the energy 860 of a body as  $E^2 = c^4 m^2 + c^2 m^2 v^2$ . With these velocity numbers, relativistic effects are respon-861 sible for 20.4% of the total squared energy on average. The "4th-degree, univariate polynomial" 862 dataset was created by drawing numbers from a normal distribution with mean 0 and variance 5 and 863 transforming them into the polynomial  $y = X + 0.5X^2 + 0.1X^3 + 0.05X^4 + \epsilon$ .

All models tested in this work had their hyperparameters selected by 5-fold cross-validation. The separation between training and testing sets varied depending on the dataset. None of the artificial datasets or datasets containing empirical data from processes with known physical laws have a sep-arate test set, as they are used to investigate the capability of the LCEN algorithm to select correct features (which occurs based on the training set). For the "Diesel freezing point" dataset, 30% of the dataset was randomly separated to form the test set. For the "Abalone" dataset, the last 1,044 entries (25%) were used as the test set as per Waugh (1995); Clark et al. (1996). For the "Concrete Compressive Strength" dataset, 25% of the dataset was randomly separated to form the test set as per Yeh (1998). For the "Boston housing" dataset, 20% of the dataset was randomly separated to form the test set. For the "GEFCom 2014" dataset, the data from task 1 were used as the training set and all data from tasks 2-15 were used as the test set. 

Table A1: Datasets used in this work and their sources. The artificial datasets are used in Section
3.1; the real datasets from processes with known physical laws are used in Section A6.1; and the
real datasets from processes with unknown physical laws are used in Section 3.2.

Dataset Name	Source
Artificial Linear	Artificial data generated by us
Multicollinear data	Artificial data generated by us
Relativistic energy	Artificial data generated by us
4th-degree, univariate polynomial	Artificial data generated by us
CARMENES star data	Schweitzer et al. (2019) [link to dataset]
Kepler's 3rd Law	Kepler et al. (1997) (Original from 1619) Wolfram Alpha LLC (2022) (Modern)
Diesel Freezing Point	Hutzler & Westbrook (2000) [link to dataset]
Abalone	Nash et al. (1995)
Concrete Compressive Strength	Yeh (1998) [dataset: Yeh (2007)]
Boston housing (modified by us)	Harrison & Rubinfeld (1978) [link to dataset]
GEFCom 2014	Hong et al. (2016) [link to dataset]

A3 APPENDIX – LIST OF HYPERPARAMETERS USED IN THIS WORK

All possible permutations of the hyperparameters below were cross-validated.

- 1. For the LASSO and Ridge regression models:  $\alpha = 0$  and 20 log-spaced values between -4.3 and 0 (as per np.logspace (-4.3, 0, 20)).
- 2. For the elastic net (EN) models:  $\alpha$  as above and L1 ratios equal to [0, 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 0.95, 0.97, 0.99].
- 3. For the SCAD models:  $\alpha$  as above and the *a* parameter (also written as  $\gamma$ ) equal to 3.7, the default value. According to Fan & Li (2001), SCAD is invariant to changes in *a*.
- 4. For the MCP models:  $\alpha$  as above and  $\gamma$  equal to [1, 1.5, 2, 2.5, 3, 3.5, 4].
- 5. For the symbolic regression (SymReg) models: most hyperparameters were set to their default values as per Stephens et al. (2022), except for the following. *population\_size* was increased to 2,000, *p\_crossover* equal to [0.7, 0.8, 0.9, 0.95], *p\_subtree\_mutation* equal to [0.01, 0.025, 0.05, 0.1, 0.15], *p\_hoist\_mutation* equal to [0.01, 0.025, 0.05, 0.1], and *p\_point\_mutation* equal to [0.01, 0.025, 0.05, 0.1, 0.25, 0.05, 0.1, 0.15] were tested. Because the sum of these probabilities must be  $\leq 1$ , some combinations are not feasible.
- 6. For the partial least squares (PLS) models: a number of components equal to all integers between 1 and a limit were used. This limit is either the number of features or 80% of the number of samples, whichever is smaller.
- 9137. For the LCEN models:  $\alpha$  and L1 ratios as above. *degree* values equal to [1, 2, 3] were typ-914ically used, except when otherwise indicated (such as in the "Relativistic energy" dataset).915lag = 0 was used, except for the "GEFCom 2014" dataset, which used lag = 168. *cutoff*916values between  $1 \times 10^{-3}$  and  $5.5 \times 10^{-1}$  were used; higher values were used only when in-917tentionally creating models with fewer selected features. A *cutoff* = 0 is used in the ablation tests for the LASSO-EN model (Section A4).

918	8	For the random forest (RF) and gradient-boosted decision tree (GRDT) models: [10, 25
919	0.	50, 100, 200, 300] trees, maximum tree depth equal to $[2, 3, 5, 10, 15, 20, 40]$ , mini-
920		mum fraction of samples per leaf equal to [0.01, 0.02, 0.05, 0.1], and minimum fraction of
921		samples per tree equal to [0.1, 0.25, 0.333, 0.5, 0.667, 0.75, 1.0]. For the GBDT models,
922		learning rates equal to [0.01, 0.05, 0.1, 0.2] were also used.
923	9.	For the AdaBoost (AdaB) models: [10, 25, 50, 100, 200, 300] trees/estimators and learning
924		rates equal to [0.01, 0.05, 0.1, 0.2] were used.
925	10	For the support vector machine (SVM) models: C values equal to [0.01, 0.1, 1, 10, 50]
926	10.	100], epsilon values equal to [0.01, 0.025, 0.05, 0.075, 0.1, 0.15, 0.2, 0.3], and gamma
927		values equal to $[1/50, 1/10, 1/5, 1/2, 1, 2, 5, 10, 50]$ divided by the number of features in a
928		dataset were used.
929	11.	For the fast sparse GAMs (FS-GAMs): the $L_0L_2$ penalty, num gamma = 20, gamma min
930		$= 5 \times 10^{-5}$ , gamma max = 1, and a max support size equal to the larger of 20% of the
931		features or 8 were used.
932	12	For the multilayer perceptron (MLP) MLP with group LASSO (MLP-GL) and LassoNet
933	12.	models: the hidden layer sizes varied for each dataset. Representing an MLP with one
934		hidden layer as [X] and an MLP with two as [X, Y], hidden layer sizes of {[800], [400],
935		[200], [100], [800, 800], [800, 400], [400, 400], [400, 200], [200, 200], [200, 100], [100,
936		100]} were used with the "Diesel Freezing Point" dataset, {[18], [9], [4], [18, 18], [18, 9],
937		[9, 9], [9, 4], [9, 2], [4, 4]} were used with the "Abalone" dataset, {[16], [8], [4], [48, 16],
938		[48, 8], [40, 24], [40, 16], [40, 8], [32, 16], [32, 8], [24, 16], [24, 8], [16, 16], [16, 8], [8, 8],
939		$[8, 4]$ were used with the "Concrete Compressive Strength" dataset, and $\{[26], [13], [13], [6], [13],$
940		[78, 26], [65, 39], [65, 26], [65, 13], [52, 39], [52, 26], [52, 13], [39, 39], [39, 26], [39, 13], [52, 26], [52, 12], [12,
941		[20, 20], [20, 15], [15, 15], [15, 0]} were used with the Boston housing dataset. Learning
942		size of 32 100 enochs and a cosine scheduler with a minimum learning rate equal to 1/16
943		size of 52, 100 epoens, and a cosine scheduler with a minimum rearming rate equal to 1710
		of the original learning rate with 10 epochs of warm-up were also used. For the MLP-GL <sub>1</sub>
944		of the original learning rate with 10 epochs of warm-up were also used. For the MLP-GL <sub>1</sub> and LassoNet, regularization parameters equal to $[1 \times 10^{-4}, 1 \times 10^{-3}, 1 \times 10^{-2}]$ were used.
944 945		of the original learning rate with 10 epochs of warm-up were also used. For the MLP-GL <sub>1</sub> and LassoNet, regularization parameters equal to $[1 \times 10^{-4}, 1 \times 10^{-3}, 1 \times 10^{-2}]$ were used.

**Table A2:** Additional features included for each value of the *degree* hyperparameter for a dataset with three features labeled  $X_0$ ,  $X_1$ , and  $X_2$  when the *lag* hyperparameter is set to 0, the *trans\_type* hyperparameter is set to 'all', and the *interaction* hyperparameter is set to True. If the *trans\_type* hyperparameter were set to 'poly', only the features of the form  $(X_k)^n$  and the interaction terms (if *interaction* were still set to True) would be present. A *degree* of n (any natural number) also includes all features from degrees 1 to n - 1.

953	Degree	Sample new features included [for all k]	Features after expansion
954	1	intercept, $X_k$ , ln $X_k$ , $(X_k)^{1/2}$ , $1/X_k$	13
955	2	$(X_k)^2$ , 2-way interactions, $(\ln X_k)^2$ , $(X_k)^{3/2}$ , $\frac{1}{(X_k)^2}$ , $\frac{\ln X_k}{X_k}$	37
956	3	$(X_k)^3$ , 3-way interactions, $(\ln X_k)^3$ , $(X_k)^{5/2}$ , $\frac{1}{(X_k)^3}$ , $\frac{(\ln X_k)^2}{(X_k)^2}$ , $\frac{\ln X_k}{(X_k)^2}$	75
957	4	$[\dots]$	129
958	5	[]	201
959			
960			
961			
962			
963			

- 964 965
- 966
- 967
- 968
- 969
- 970 971

### 972 A4 APPENDIX – ABLATION TESTS

To better clarify the design choices of the LCEN algorithm and highlight the relevance of each individual part of the algorithm, ablation tests are performed. Three ablated algorithms - LASSO-Clip (LC), EN-Clip (ENC), and LASSO-EN (LEN) – are compared with the original LCEN algorithm. Three variant algorithms, LASSO-Clip-LASSO (LCL), EN-Clip-EN (ENCEN), and regular LCEN followed by OLS for debiasing (LCEN $\rightarrow$ OLS), are also compared. The "Relativistic en-ergy", "Diesel Freezing Point", "Abalone", and "Concrete Compressive Strength" datasets are used in the ablation tests. Tests with the "Relativistic energy" dataset show that models with a Clip step had some degree of success with selecting only relevant features (Table A3). However, the ablated algorithms (LC, ENC, and LEN) had much higher prediction errors for the coefficients of the relevant features, even though LC and ENC were able to select only the relevant features. The variant algorithms (LCL, ENCEN, and LCEN $\rightarrow$ OLS) had performances closer to that of LCEN, but LCL was slightly worse in terms of error. LCEN $\rightarrow$ OLS was able to estimate the coefficients with a lower error than LCEN, but this improvement in coefficient estimation performance comes with an increase in test-set MSEs in other tasks. The models that begin with EN (ENC and EN-CEN) are approximately one order of magnitude slower than LCEN on all datasets (Tables A3–A6), whereas LC was approximately 50% faster than LCEN. However, LCEN consistently had the low-est validation RMSE in all datasets, and had the lowest test-set RMSE in all but one dataset. As highlighted by Table A4, LCEN built the sparsest and most accurate models out of all ablated and variant algorithms trained with the "Diesel Freezing Point" dataset. Overall, these ablation exper-iments highlight how LCEN is the optimal algorithm to maximize accuracy and selectivity while maintaining a low runtime. 

**Table A3:** Relative error to the ground truth for the "Relativistic energy" dataset with  $1 \le m < 100$  at different noise levels for ablated and variant LCEN algorithms. The first coefficient is for  $m^2$  and the second coefficient is for  $m^2v^2$ . Compare with the right-side table of Table 1 and a runtime of 4.79 seconds for LCEN.

	Noise Level	LC Error (%)	ENC Error (%)	LEN Error (%)
	0%	36.62, 18.08	41.52, 20.91	43.75, 22.32
	5%	37.68, 18.85	41.30, 21.16	43.93, 23.45
	10%	18.92, 1.647	44.31, 23.70	1.137, 0.468
	15%	39.71, 20.61	44.31, 23.70	46.65, 26.40
	20%	39.65, 21.13	39.99, 21.95	45.87, 27.23
	30%	22.35, 2.603	22.93, 2.660	7.649, 1.036
	Runtime (s)	3.70	37.1	5.40
Γ	Noise Level	LCL Error (%)	ENCEN Error (%	) LCEN→OLS
	Noise Level 0%	LCL Error (%) 0.007, 0.012	ENCEN Error (% 0.001, 0.006	) LCEN $\rightarrow$ OLS 0, 0
	Noise Level 0% 5%	LCL Error (%) 0.007, 0.012 0.011, 0.029	ENCEN Error (% 0.001, 0.006 0.005, 0.022	) LCEN→OLS 0, 0 0.004, 0.016
	Noise Level 0% 5% 10%	LCL Error (%) 0.007, 0.012 0.011, 0.029 0.015, 0.045	ENCEN Error (% 0.001, 0.006 0.005, 0.022 0.009, 0.038	) LCEN $\rightarrow$ OLS 0, 0 0.004, 0.016 0.008, 0.032
	Noise Level 0% 5% 10% 15%	LCL Error (%) 0.007, 0.012 0.011, 0.029 0.015, 0.045 0.019, 0.061	ENCEN Error (% 0.001, 0.006 0.005, 0.022 0.009, 0.038 0.013, 0.054	) LCEN $\rightarrow$ OLS 0, 0 0.004, 0.016 0.008, 0.032 0.012, 0.048
	Noise Level 0% 5% 10% 15% 20%	LCL Error (%) 0.007, 0.012 0.011, 0.029 0.015, 0.045 0.019, 0.061 0.023, 0.077	ENCEN Error (% 0.001, 0.006 0.005, 0.022 0.009, 0.038 0.013, 0.054 0.017, 0.070	$\begin{array}{c c} & & \\ \hline & & \\ \hline & & \\ 0,0 \\ & & \\ 0.004, 0.016 \\ & & \\ 0.008, 0.032 \\ & & \\ 0.012, 0.048 \\ & & \\ 0.016, 0.064 \end{array}$
	Noise Level 0% 5% 10% 15% 20% 30%	LCL Error (%) 0.007, 0.012 0.011, 0.029 0.015, 0.045 0.019, 0.061 0.023, 0.077 0.031, 0.109	ENCEN Error (% 0.001, 0.006 0.005, 0.022 0.009, 0.038 0.013, 0.054 0.017, 0.070 0.025, 0.103	$\begin{array}{c c} & & \\ \hline & & \\ \hline & & \\ 0, 0 \\ & & \\ 0.004, 0.016 \\ & & \\ 0.008, 0.032 \\ & & \\ 0.012, 0.048 \\ & & \\ 0.016, 0.064 \\ & & \\ 0.024, 0.096 \end{array}$

Algorithm	Test RMSE (°C)	Features	Runtime (s)
IC	4.84	37	4.39
LC	5.15	29	4.31
FNC	4.80	263	20.6
LINC	5.00	257	20.7
LEN	4.87	39	6.85
ICI	4.93	33	4.74
LUL	5.01	28	4.62
ENCEN	4.83	191	31.1
ENCEN	4.90	173	30.1
	5.02	36	6.54
LUEN-JULS	5.07	29	6.27

Table A4: Results of different ablated and variant LCEN algorithms for the "Diesel Freezing Point" dataset. Compare with Table 3.

**Table A5:** Results of different ablated and variant LCEN algorithms for the "Abalone" dataset.Compare with Table A7 and a runtime of 19.2 seconds for LCEN.

Algorithm	Test RMSE (rings)	Features	Runtime (s)
LC	2.1	8	11.8
ENC	2.1	8	297
LEN	2.1	8	26.3
LCL	2.0	8	12.9
ENCEN	2.1	8	308

Table A6: Results of different ablated and variant LCEN algorithms for the "Concrete Compressive Strength" dataset. All models selected all 8 features, but a varying number of transforms of these features. Compare with Table 4 and a runtime of 39.7 seconds for LCEN.

Algorithm	Test RMSE (MPa)	Runtime (s)
LC	5.53	24.6
ENC	16.5	800
LEN	5.50	44.9
LCL	5.92	26.4
ENCEN	6.12	863



1119Figure A1: Plots of the Matthews Correlation Coefficients (MCCs) [top] and precision-recall (P-R)1120curves [bottom] for models tested under the "Artificial Linear" dataset. Scenario conditions are as1121written in each subfigure's title. The black dotted lines in the P-R curves are  $F_1$  score isolines. The1122 $F_1$  score corresponding to each isoline is shown on the right.



Figure A2: Plots of the Matthews Correlation Coefficients (MCCs) [top] and precision-recall (P-R) curves [bottom] for models tested under the "Artificial Linear" dataset. Scenario conditions are as written in each subfigure's title. The black dotted lines in the P-R curves are  $F_1$  score isolines. The  $F_1$  score corresponding to each isoline is shown on the right.



Figure A3: Plots of the Matthews Correlation Coefficients (MCCs) [top] and precision-recall (P-R) curves [bottom] for models tested under the "Artificial Linear" dataset. Scenario conditions are as written in each subfigure's title. The black dotted lines in the P-R curves are F1 score isolines. The F<sub>1</sub> score corresponding to each isoline is shown on the right.



Figure A4: Plots of the Matthews Correlation Coefficients (MCCs) [top] and precision-recall (P-R) curves [bottom] for models tested under the "Artificial Linear" dataset. Scenario conditions are as written in each subfigure's title. The black dotted lines in the P-R curves are F1 score isolines. The F<sub>1</sub> score corresponding to each isoline is shown on the right.



Figure A5: Plots of the Matthews Correlation Coefficients (MCCs) [top] and precision-recall (P-R) curves [bottom] for models tested under the "Artificial Linear" dataset. Scenario conditions are as written in each subfigure's title. The black dotted lines in the P-R curves are F<sub>1</sub> score isolines. The F<sub>1</sub> score corresponding to each isoline is shown on the right.



Figure A6: Plots of the Matthews Correlation Coefficients (MCCs) [top] and precision-recall (P-R) curves [bottom] for models tested under the "Artificial Linear" dataset. Scenario conditions are as written in each subfigure's title. The black dotted lines in the P-R curves are F1 score isolines. The F<sub>1</sub> score corresponding to each isoline is shown on the right.



Figure A7: Plots of the Matthews Correlation Coefficients (MCCs) [top] and precision-recall (P-R) curves [bottom] for models tested under the "Artificial Linear" dataset. Scenario conditions are as written in each subfigure's title. The black dotted lines in the P-R curves are F1 score isolines. The F<sub>1</sub> score corresponding to each isoline is shown on the right.



Figure A8: Plots of the Matthews Correlation Coefficients (MCCs) [top] and precision-recall (P-R) curves [bottom] for models tested under the "Artificial Linear" dataset. Scenario conditions are as written in each subfigure's title. The black dotted lines in the P-R curves are  $F_1$  score isolines. The F<sub>1</sub> score corresponding to each isoline is shown on the right.



Figure A9: Plots of the Matthews Correlation Coefficients (MCCs) [top] and precision-recall (P-R) curves [bottom] for models tested under the "Artificial Linear" dataset. Scenario conditions are as written in each subfigure's title. The black dotted lines in the P-R curves are F1 score isolines. The F<sub>1</sub> score corresponding to each isoline is shown on the right.



Figure A10: Plots of the Matthews Correlation Coefficients (MCCs) [top] and precision-recall (P-R) curves [bottom] for models tested under the "Artificial Linear" dataset. Scenario conditions are as written in each subfigure's title. The black dotted lines in the P-R curves are  $F_1$  score isolines. The  $F_1$  score corresponding to each isoline is shown on the right.



Figure A11: Plots of the Matthews Correlation Coefficients (MCCs) [top] and precision-recall (P-R) curves [bottom] for models tested under the "Artificial Linear" dataset. Scenario conditions are as written in each subfigure's title. The black dotted lines in the P-R curves are F1 score isolines. The F<sub>1</sub> score corresponding to each isoline is shown on the right.



Figure A12: Plots of the Matthews Correlation Coefficients (MCCs) [top] and precision-recall (P-R) curves [bottom] for models tested under the "Artificial Linear" dataset. Scenario conditions are as written in each subfigure's title. The black dotted lines in the P-R curves are F<sub>1</sub> score isolines. The F<sub>1</sub> score corresponding to each isoline is shown on the right.



1746Figure A13: LCEN model output at different X-data noise levels  $\epsilon_1$  and y-data noise levels  $\epsilon_2$ .1747Bright red squares indicate both variables were selected and their coefficients had errors  $\leq 5\%$ . Light1748red squares indicate that both variables were selected and their coefficients had  $5\% < \text{errors} \leq 10\%$ .1749White squares indicate that both variables were selected and their coefficients had  $10\% < \text{errors} \leq 20\%$ . Light blue squares indicate that both variables were selected and their coefficients had  $10\% < \text{errors} \leq 20\%$ . Bright blue squares indicate that only one of the variables was selected.

1782 For the "4th-degree, univariate polynomial" dataset, Sun & Braatz (2020) created four models: one 1783 that always uses degree = 4 ("unbiased model"), one that always uses degree = 2 ("biased model"), 1784 one that selects a *degree* between 1 and 10 based on cross-validation ("cv"), and one that selects a 1785 degree equal to 2 or 4 based on cross-validation ("cv limited order"). Sun & Braatz (2021) noted 1786 that the *degree* 4 "unbiased model" was the best at low noise levels, but its error quickly increases, leading to the degree 2 "biased model" becoming the best for noise levels > 75 (Fig. 3 of Sun & 1787 Braatz (2021); reproduced with permission here as the left subfigure of Fig. 1). The model with 1788 degree equal to 2 or 4 "cv limited order" was typically very close in performance to the best model 1789 at all noise levels, whereas the model with a *degree* between 1 and 10 "cv" had lower performance. 1790 Sun & Braatz (2021) explain these observations with the bias-variance tradeoff: at low noise levels, 1791 models should follow the ground truth as closely as possible; thus, the *degree* 4 "unbiased model" 1792 was the best. However, at sufficiently high noise levels, it becomes impossible to obtain enough 1793 signal to compensate for the additional degrees of freedom (variance) in a 4th degree model; thus, 1794 the degree 2 "biased model" becomes the best. The degree between 1 and 10 "cv" model had lower 1795 performance due to its greater hyperparameter variance, and the *degree* equal to 2 or 4 "cv limited 1796 order" model struck a balance between the "unbiased model" and the "biased model".

1797 Similarly to the models generated using ALVEN, the LCEN model with a *degree* between 1 and 1798 10 "cv" had the lowest performance and the LCEN model with degree equal to 2 or 4 "cv limited 1799 order" had a performance between the *degree* 4 "unbiased model" and the *degree* 2 "biased model". However, the degree 4 "unbiased model" was always the best model, no matter the noise level used. 1801 We attribute this considerable reduction in median test MSEs and the superiority of the *degree* 4 "unbiased model" created by LCEN to the improved feature selection algorithm, which is able to 1803 better resist variance due to noise and a large number of hyperparameters. This is corroborated by how the model with a *degree* between 1 and 10 "cv" tended to select *degree* = 4 at lower noise levels and degree = 2 at higher noise levels (Figure A15), showing how LCEN can automatically follow 1805 the bias-variance tradeoff hypothesis.



1833 1834

**Figure A14:** 25% (squares) and 75% quartile (triangles) test set MSEs for the LCEN model trained for the "4th-degree, univariate polynomial" dataset. The trends tend to match those from Fig. 1.





1880

## A6 APPENDIX – ADDITIONAL RESULTS WITH EMPIRICAL DATA

### 1881 1882 A6.1 DATASETS FOR WHICH PHYSICAL LAWS ARE AVAILABLE

The first test of an empirical dataset from a process with a known physical law uses the "CARMENES star data" dataset from Schweitzer et al. (2019). This dataset contains information on temperature (T), radius (R), and luminosity (L) of 293 white dwarf stars. These features are linked together by the Stefan-Boltzmann equation,  $L = 4\pi R^2 \sigma T^4$ , where  $\sigma$  is a constant. Normalizing this equation to values from another star (typically, the Sun), conveniently sets the constant terms to 1. This normalization is applied to the "CARMENES star data" dataset. LCEN with *degrees* from 1 to 10 was applied to this normalized dataset. Despite the very large number of potential features (due to the high *degree* values used), LCEN correctly selected only the  $R^2T^4$  feature. The coefficient assigned to  $R^2T^4$  is 0.9826, which is well within the 2–3% error on these data (as reported by Schweitzer et al. (2019)). LCEN retained high performance for this real data in a high-hyperparameter variance scenario.

A potential limitation in real datasets is data scarcity. To evaluate the LCEN algorithm in a low-1894 data scenario, the "Kepler's 3rd Law" datasets are created. The first version uses the original data obtained by Kepler, first published in 1619 and republished in Kepler et al. (1997). From only 1896 6 (slightly inaccurate) measurements, Kepler was able to derive the eponymous Kepler's 3rd law, which states that the period T of a celestial body is related to the semi-major axis of its orbit a by 1898  $T = ka^{3/2}$ . The constant k depends on the masses of the central and orbiting bodies; however, 1899 as the mass of the central body is typically much larger, the mass of the orbiting body is ignored. 1900 In this and Kepler's works, T is measured in Earth days, so the constant k is  $\sim$  365.25 when using modern data and  $\sim$ 365.15 when using Kepler's original data. Despite the low number of data points, 1901 LCEN correctly selected only the  $a^{3/2}$  feature. Moreover, the coefficient assigned to that feature was 1902 366.82, an error of only 0.46% relative to Kepler's k = 365.15. 1903

1904 1905 1906 1906 1907 1907 1907 1908 1909 LCEN is then evaluated using a modern version of the same dataset, which contains 8 points (as Uranus and Neptune were discovered after Kepler's observations) whose data were measured with greater accuracy. On this modern "Kepler's 3rd Law" dataset, LCEN again selects only the  $a^{3/2}$ feature. The coefficient assigned to the  $a^{3/2}$  feature is 365.00, an error of only 0.07% relative to the modern value k = 365.25. LCEN did perfect feature selection in these data-scarce scenarios, with parameter estimates minimally affected by experimental noise.

- 1910
- 1911 A6.2 DATASETS FOR WHICH NO PHYSICAL LAW IS AVAILABLE 1912

1913 Abalone (Haliotis sp.) are sea snails whose age can be determined by cutting their shells, staining them, and counting the stained shell rings under a microscope. This process is laborious and 1914 error-prone. An alternative is to estimate the number of rings based on readily available physical 1915 characteristics, such as weight and size. As before, LCEN was compared with other dense and 1916 sparse machine learning models, and LCEN models with increased sparsity were also generated 1917 (Table A7). In this problem, OLS, PLS, RR, LASSO, and EN all converged to the OLS solution 1918 (that is, no regularization), selecting all 8 linear features and having an RMSE of 2.1 rings. On the 1919 other hand, LCEN automatically detected that 2nd degree features would be relevant. The LCEN 1920 algorithm model also selected all 8 features and had an RMSE of 2.0 rings. Nonlinear models had 1921 test RMSEs between 2.0 and 2.5 rings, but most all lack the interpretability of LCEN. By increasing 1922 the cutoff hyperparameter, sparser LCEN models may be generated. An LCEN model with only 3 features had an RMSE of 2.1 rings, and another with only 2 features had an RMSE of 2.2 rings. 1923 1924 This experiment further illustrates LCEN's robust feature selection, and how very sparse LCEN models retain significant performance. Furthermore, LCEN models with the same or lower number 1925 of selected features had a lower test set RMSE than FS-GAMs. 1926

1927 The "Boston housing" dataset contains the median value of owner-occupied houses and many inter-1928 nal and external measurements, such as the per-capita crime rate of the region, the average number of rooms, and the concentration of nitric oxides in the area (Harrison & Rubinfeld, 1978). We modified 1929 this dataset to detransform the B variable into its raw value; samples in which this detransformation 1930 led to multiple possible values were discarded. In terms of test RMSE, the linear models (OLS, 1931 PLS, RR, EN, LASSO) tended to perform equal to each other and quite poorly on this dataset. RF 1932 and SVM performed relatively well, but GBDT and AdaB had the two worst performances among 1933 the nonlinear models. A dense MLP was the best model in terms of test RMSE, and the MLP-GL<sub>1</sub> 1934 and LassoNet performed similarly but slightly worse. LCEN had a very high performance on this 1935 regression task, reaching a test RMSE only 6% higher than that of the dense MLP. LCEN also had 1936 the lowest validation RMSE, which was 54% lower than that of the dense MLP.

1937 1938

## 1939 A7 APPENDIX – COMPUTATIONAL RESOURCES USED

1940

All experiments were done in a personal computer equipped with a 13th Gen Intel<sup>®</sup> Core<sup>™</sup> i5-13600K CPU, 64 GB of DDR4 RAM, and an NVIDIA GeForce RTX 4090 GPU. Runtimes for the models trained on the "Diesel Freezing Point" dataset are provided in Table 3, and runtimes for LCEN and ablated algorithms are provided in the tables of Section A4.

Table A7: Results of different models for the "Abalone" dataset. The number of features that minimizes the cross-validation MSE is 6 for FS-GAMs and 8 for LCEN. These data and results are discussed in Section 3.2.

Model	Test RMSE (rings)	Features
OLS = PLS = RR = LASSO = EN	2.1	8
SCAD	2.1	8
MCP	2.1	8
SymReg	2.3	3
RF	2.1	8
GBDT	2.2	8
AdaB	2.3	8
SVM	2.0	8
MLP	2.0	8
$MLP-GL_1$	2.0	8
LassoNet	2.0	8
	2.1	8
FS-GAMs	2.2	6
	2.4	2
	2.0	8
LCEN	2.1	3
	2.2	2

Table A8: Results of different models for the "Boston housing" dataset.

Model	Test RMSE (Thousands USD)
OLS	6.38
PLS	6.39
RR = EN	6.39
LASSO	6.38
SCAD	6.38
MCP	6.38
SymReg	6.70
RF	5.02
GBDT	5.91
AdaB	5.67
SVM	5.15
MLP	4.51
$MLP-GL_1$	4.88
LassoNet	4.76
FS-GAMs	7.32
LCEN	4.78