
Automation for Interpretable Machine Learning Through a Comparison of Loss Functions to Regularisers

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 To increase the ubiquity of machine learning it needs to be automated. Automation
2 is cost-effective as it allows experts to spend less time tuning the approach, which
3 leads to shorter development times. However, while this automation produces
4 highly accurate architectures, they can be uninterpretable, acting as ‘black-boxes’
5 which produce low conventional errors but fail to model the underlying input-output
6 relationships—the ground truth. This paper explores the use of the Fit to Median
7 Error measure in machine learning regression automation, using evolutionary
8 computation in order to improve the approximation of the ground truth. When used
9 alongside conventional error measures it improves interpretability by regularising
10 learnt input-output relationships to the conditional median. It is compared to
11 traditional regularisers to illustrate that the use of the Fit to Median Error produces
12 regression neural networks which model more consistent input-output relationships.
13 The problem considered is ship power prediction using a fuel-saving air lubrication
14 system, which is highly stochastic in nature. The networks optimised for their
15 Fit to Median Error are shown to approximate the ground truth more consistently,
16 without sacrificing conventional Minkowski-r error values.

17 1 Development of Interpretable Machine Learning

18 Machine learning regression models are increasingly being used in industrial and engineering con-
19 texts for high-stakes decision making, automation and control. These methods often produce low
20 conventional error values, yet are known to produce physically inconsistent results which cannot
21 generalise off test set and so cannot be relied upon to model the ground truth of the system. The
22 models are designed to be accurate, not interpretable, and so a human cannot understand how changes
23 in the inputs change the prediction. In real-world applications, where the output can have a direct
24 effect on human life or the environment, model accuracy alone is not sufficient.

25 Trust is increased if a trained model approximates the true input-output relationships, performing
26 accurately within the bounds of the training data set and beyond it. It has been demonstrated that for
27 many applications minimising traditional error measures cannot guarantee an accurate approximation
28 of the ground truth (Willard et al. 2020). This is due to a poor inductive bias, the inherent prioritisation
29 of one solution over another (Battaglia et al. 2018), produced by conventional error measures which
30 are based on Minkowski-r metrics (Hanson & Burr 1987).

31 This trust can be increased by manually tuning to remove overfitting or to provide a solution that
32 makes more sense to the user. However, the expert knowledge and domain experience required to
33 properly tune a machine learning method manually are not always available in industry. Genetic
34 algorithms are therefore increasingly used to search a method’s hyperparameter space more efficiently

35 (Yang et al. 2021) (Kumar et al. 2021); which minimise conventional error measures on a test set,
36 often combined with lowering the complexity of the network. This automation exacerbates the
37 lack of interpretability, as models have a large flexibility, and prediction accuracy is prioritised,
38 a low conventional error is achieved without certainty that the method has modelled the correct
39 internal functions. Regularisation hyperparameters can be optimised alongside other neural network
40 parameters (Tani et al. 2021) (Luketina et al. 2016), which increases the search space and creates
41 more flexibility for methods to produce ‘accurate’ predictions and avoid overfitting.

42 Common regression regularisation methods are l1 and l2 regularisation and dropout. For l1 and l2
43 regularisation, large network weights are penalised in the loss function (Nowlan & Hinton 1992). The
44 absolute value of the weights is penalised in l1 regularisation and the squared value in l2, meaning
45 l1 encourages weights towards zero and l2 encourages weights to be small but non-zero. The l1, l2
46 and elastic net (l1+l2) regularisers improve a networks generality, increasing the applications where
47 the trained methods can be applied, by penalising complexity. Dropout, where a randomly selected
48 subset of weights are optimised at each epoch rather than the full set, improve the generality of the
49 trained models by preventing co-adaption of weight values (Srivastava et al. 2014). Dropout has been
50 shown to be equivalent to l2 regularisation after scaling by Fisher information (Wager et al. 2013),
51 suggesting that the two should not be used in unison. The neural network regularisation methods
52 discussed above aim to improve generality, reducing overfitting by simplifying the relationships
53 modelled by the networks.

54 Regularisers improve the modelling of the ground truth in scenarios adhering to the assumptions in the
55 proof in Bishop (1995), under which minimum Minkowski-r error values approximate the conditional
56 average of the dataset. This is because the inductive bias from the loss function guides the input-
57 output relationships towards the conditional average, while the regularisation stops overfitting by
58 simplifying the input-output relationships being modelled. However, these assumptions are restrictive
59 and it is noted that few regression applications adhere to them. For example, one assumption is that
60 the dataset is homoscedastic. In scenarios not adhering to these assumptions, network regularisation
61 simplifies the relationships being modelled but this does not necessarily improve the generality, or
62 model the ground truth.

63 The Fit to Median Error measure (Parkes et al. 2021) produces more interpretable regression, when
64 used in conjunction with conventional error measures. This is achieved by regularising the learnt
65 input-output relationships to the conditional median of the training dataset: the median output value,
66 conditioned on each isolated input variable in turn (Bishop 1995). For many regression applications
67 the conditional medians are a good approximation of the ground truth input-output relationships but
68 as yet it has not been explored as part of an automated approach.

69 A challenging regression problem is ship power prediction for a vessel using air lubrication to reduce
70 fuel consumption. It is chosen to be used in this study as it violates the assumptions in Bishop (1995),
71 where the noise in the output space is non-Gaussian and heteroscedastic. In this situation, correctly
72 modelling the ground truth and accurate prediction is required but there is limited understanding
73 of that ground truth (Parkes et al. 2018). The literature shows that shaft powering of a vessel can
74 be predicted with average accuracies of between 1.5-5% error with the use of a regression neural
75 network trained with high frequency data from the vessel (Pedersen & Larsen 2009), (Petersen et al.
76 2012), (Le et al. 2020), (Jeon et al. 2018), (Liang et al. 2019). All neural network applications to
77 ship power prediction in the literature use a combination of local searches and domain knowledge to
78 identify hyperparameter values. The addition of an air lubrication device increases the complexity of
79 the regression problem, as the system interacts with a number of interrelated input variables.

80 This paper explores the automation of neural network training to a new problem, with a focus
81 on producing a network which accurately models the ground truth. It compares the ground truth
82 representation of a neural network when a genetic algorithm optimises the network’s hyperparameters
83 to reduce the Mean Fit to Median Error measure and compares it to standard regularization using l1,
84 l2 and dropout, and to a network optimised to minimise the Maximum Absolute Error. It is illustrated
85 that neural network regularisation methods (l1, l2 and dropout) can be replaced by the use of the
86 Mean Fit to Median performance measure as an objective in the genetic algorithm, reducing the
87 complexity of the search space and producing networks which more consistently model the ground
88 truth.

89 **2 Neural Networks Parameters**

90 Previous applications of neural networks to ship power prediction use between 1 and 3 hidden layers
 91 (Leifsson et al. 2008) (Parkes et al. 2019), and between 5 and 300 neurons in each hidden layer (Jeon
 92 et al. 2018). To provide a sufficiently large search space to allow verification, or otherwise, of these
 93 parameters a maximum of 4 hidden layers and 1000 neurons in each layer are used. The majority
 94 of the literature treats the problem as time-invariant and use feed-forward networks, so no recurrent
 95 parameters are optimised. As the optimiser or activation functions are rarely documented in the
 96 literature, the state-of-the-art optimisers and activation functions available in the Keras framework
 97 (Chollet et al. 2015) are used in the optimisation, Table 1.

Table 1: Selected Neural Network Hyperparameters

Hyperparameter	Value or set
Layers	[1,4]
Neurons in each layer	[1,1000]
Epochs	Increasing from 1-20 for increasing generations
Early stopping patience	5
Loss function	Mean Absolute Error
Performance measures	Mean Absolute Relative Error, Maximum Absolute Relative Error, Mean Fit to Median Error
Optimiser	SGD, Adam (Kingma & Ba 2014), Nadam (Dozat 2016), RMSprop (Hinton et al. 2012), Adagrad (Duchi et al. 2011), Adadelta (Zeiler 2012), Adamax (Kingma & Ba 2014)
Activation function	ReLU, sigmoid, softmax, softplus, softsign, tanh, selu, elu
l1 & l2 Rates	0, 0.01, 0.001, 0.0001, 0.00001
Dropout	[0,0.9)
Initialiser	Random Normal ($\mu = 0, \sigma = 0.1$)

98 The number of epochs and early stopping procedure are not optimised, as there was a need for
 99 predictable compute requirements and allowing the optimisation of these parameters leads to unpre-
 100 dictable run times. The number of epochs to train each network increases for increasing generation
 101 number in the genetic algorithm, from 1 epoch in the first 15 generations to 20 in the final 15. This
 102 was also implemented to reduce compute and it was validated that when more than 20 epochs were
 103 allowed, that the early stopping, with a patience of 5, stopped the training within 20 epochs for the
 104 majority of networks. The loss function is similarly not optimised, the Mean Absolute Error is used,
 105 as the conditional medians are closer to the ground truth input-output relationships in these datasets
 106 than the conditional means.

107 The performance measures, or the genetic algorithm’s fitness functions, are the Mean Absolute
 108 Relative Error, the Maximum Absolute Relative Error and the Mean Fit to Median Error. Different
 109 combinations of these, alongside the use of regularisation parameters in the search space are compared
 110 to illustrate the effect of different types of regularisation.

111 **3 cMLSGA Parameters**

Table 2: Selected cMLSGA Hyperparameters

Hyperparameter	Value or set
Algorithm at Individual Level	HEIA, IBEA
Crossover Type & Rate	SBX & DE, 1
Mutation Type & Rate	Polynomial, 0.08
Number of eliminated collectives	1
Generations between elimination	10
Population size	1000
Generations	300
Proportion elite	10%

112 In this study cMLSGA¹ is selected as it shows the top performance on a range of evolutionary
 113 benchmarking problems (Grudniewski & Sobey 2021) and practical problems (Grudniewski & Sobey
 114 2019). Genetic algorithms are increasing used to tune neural network hyperparameters including
 115 regularisation parameters for use on new problems (Jin et al. 2004). Many approaches have multiple
 116 genetic algorithm objectives, although these all minimise an error measure and a measure of network
 117 complexity (Wang et al. 2019) and (Smith & Jin 2014). The use of multiple different performance
 118 measures as objectives is yet to be explored in the literature.

Table 3: Genetic Algorithm Approaches

Approach	Objective(s)	Network Regularisation
GAi	Mean Absolute Error	l1, l2 and dropout
GAii	Mean Absolute Error Maximum Absolute Error	l1, l2 and dropout
GAiii	Mean Fit to Median Error Mean Absolute Error	None
GAiv	Mean Absolute Error Maximum Absolute Error	None

119 Four approaches are investigated in this study, summarised in Table 3, for approach (GAi) and
 120 (GAii) the genetic algorithm cMLSGA optimises all variables in Table 2, including the l1 and l2
 121 regularisation rate and the dropout rate of the networks. Although it is advised that l2 regularisation
 122 and dropout are not used in the same network the genetic algorithms are provided with zero options
 123 for all regularisation parameters, to identify if one is preferable in this scenario.

124 Approach (GAi) is a single objective genetic algorithm optimising the Mean Absolute Error which is
 125 compared to a multi-objective formulation where the (GAii) approach optimises both Mean Absolute
 126 Error and Maximum Absolute Error. For approaches (GAiii) and (GAiv) no network regularisation
 127 parameters are optimised: l1, l2 and dropout rates are all set permanently to zero. They avoid
 128 producing networks that have overfitted by the use of two performance metrics as multi-objectives,
 129 (GAiii) uses the Mean Fit to Median and Mean Absolute Errors to be minimised and (GAiv) uses the
 130 Maximum Absolute and Mean Absolute. All approaches use 40 CPUs with 2.0 GHz Intel Skylake
 131 processors and 192 GB of DDR4 memory, and take less than 3 days, this setup may not be feasible
 132 for widespread industrial application, although it is suggested it is within reach of some industries.

133 4 Data

134 The data used in this study are from a large vessel equipped with the Silverstream® Air Lubrication
 135 System. The air lubrication system works through use of fluid sheering to create an air microbubble
 136 carpet directly captured within the boundary layer on the ship hull bottom. The bubble carpet reduces
 137 the frictional resistance thereby increasing the speed and reducing the shaft power. Compressors
 138 provide a constant supply of air to the hull bottom to maintain a uniform bubble carpet operated at
 139 the optimal compressor power that maximises the energy balance. The study is performed on both
 140 system on and system off datasets, however for brevity only results for system off are presented as
 141 they show similar performance. This prediction is required for a baseline determination of how the
 142 system is working, but the relationships between the power, weather, ocean and operating conditions
 143 are complex and difficult to model.

144 The variables considered in this study are the shaft power, speed through water, relative wind speed
 145 and direction, draught and trim, with shaft power the target variable. These are selected based on
 146 a detailed study into variable selection for shaft power prediction (Parkes et al. 2019). The speed
 147 through water is selected over the speed over ground, for use as an input variable, as it is more
 148 hydrodynamically relevant and its accuracy is validated by comparison to the speed over ground.
 149 The dataset is cleaned by removing rows with missing or non-physical values and all datapoints
 150 below 0.05 normalised shaft power are removed. The dataset is split into two using the air lubrication
 151 system status: system on and system off, where system on is defined as air lubrication system power
 152 greater than zero. The system on dataset contains 352,690 datapoints and system off contains 237,962.
 153 The data is split into training, testing and validation sets of 70%, 15% and 15% respectively. Each

¹The code for cMLSGA is available at https://www.bitbucket.org/*****.

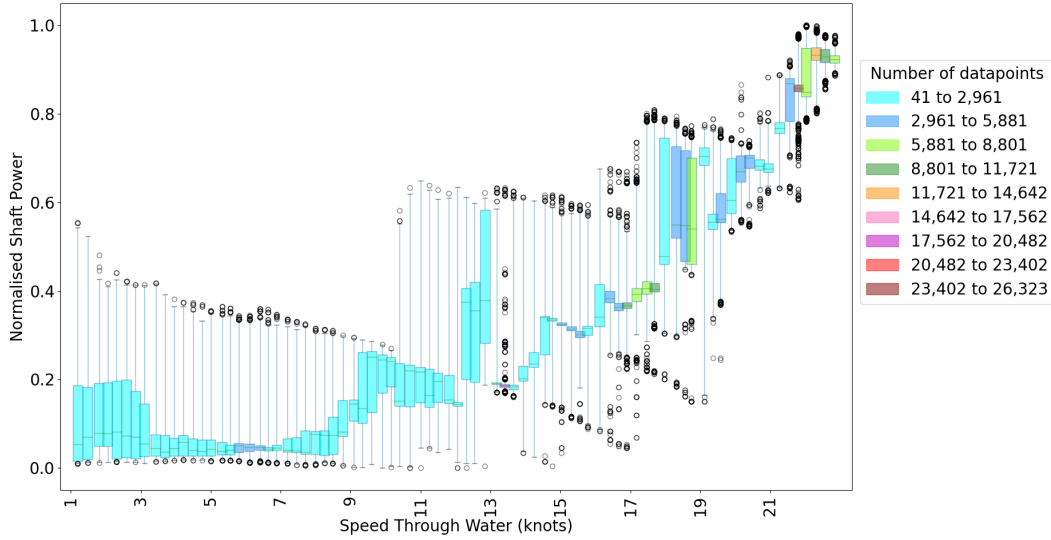


Figure 1: The distribution of the observed shaft powers for half knot bins of speed through the water for dataset where the system is off. In the box and whisker plots the boxes contain 50% of the distribution and the whiskers extend to the datum which is at 1.5 times the interquartile range.

154 network in the genetic algorithm trains on a randomly sampled 35,000 datapoints from the training
 155 set and uses randomly sampled sets of size 7,500 from validation and testing sets for validation during
 156 training, and testing to produce the fitness of the network for the genetic algorithm. The errors stated
 157 in the paper are from networks on the Pareto fronts of each approach, which are validated on the full
 158 testing set.

159 The datasets contain large regions of sparse data in all input variable domains, this is exemplified
 160 by the ship speed domain where each half-knot interval below 16 knots contains less than 0.8% of
 161 the data, which accounts for more than half the speed domain, Figure 1. In addition, the boxplot
 162 ranges and outliers show high heteroscedicity with idiosyncratic noise caused by situations where the
 163 angle of the propeller blades is varied to achieve the required speed. This highlights the complexity
 164 in developing models of the powering of this vessel, as the dataset also contains the effects from other
 165 latent variables, such as piloting behaviour and route taken.

166 5 Optimisation including regularisation parameters: (GAi) and (GAii)

167 Previous studies predicting ship powering using neural networks report that l_1 , l_2 and elastic net
 168 increase both test set and off-test set errors and that optimal values for both l_1 and l_2 are zero.
 169 Therefore the genetic algorithm setup is biased towards low and zero values of regularisation rates by
 170 using a set of exponentially decreasing values and an explicit zero option.

171 The single objective (GAi) fails to identify that zero regularisation rates produce the lowest errors,
 172 favouring networks with the highest possible rate of l_2 (0.01), Figure 2b. (GAi) produces networks
 173 with the highest Mean Absolute Relative Errors of all the approaches, $(5.19 \pm 0.00)\%$ from Figure
 174 4a. In contrast, (GAii) favours lower l_1 and l_2 rates of 0 or 0.00001, Figures 2c and 2d, which results
 175 in networks with the lowest Mean Absolute Relative Errors of all four approaches, on average, with a
 176 value of $(2.87 \pm 0.45)\%$, shown in Figure 4a. This is around 0.5% higher than the lowest documented
 177 error for ship power prediction.

178 It is posited that the high error for the single objective problem is directly related to the use of a
 179 large l_2 regularisation rate, as noted in previous studies for ship power prediction. It is possible
 180 that the use of a multi-objective search algorithm for a single-objective problem means that the
 181 optimal hyperparameters can't be found, resulting in large errors. The implementation also requires
 182 restrictions in the number of epochs used for training in the initial generations, it is possible this
 183 biases (GAi) towards certain size networks, where higher l_2 rates are preferable. This hypothesis is

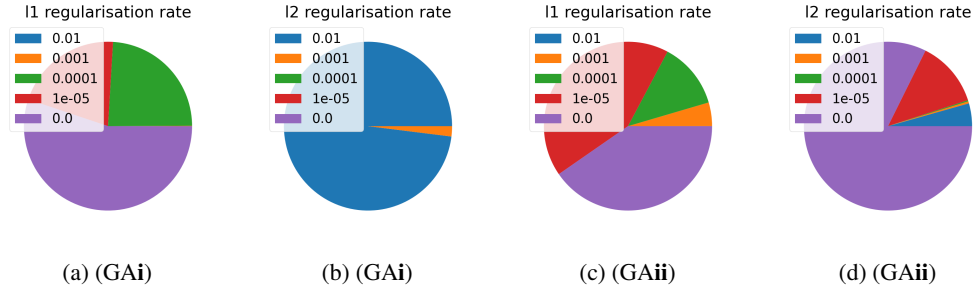


Figure 2: Distribution of regularisation rates for networks in the last 15 generations of (GAi) cMLSGA with multi-objectives of minimising Maximum and Mean Absolute Error for (a) I1 and (b) I2 and (GAii) cMLSGA with the single objective of minimising Mean Absolute Error for (c) I1 and (d) I2

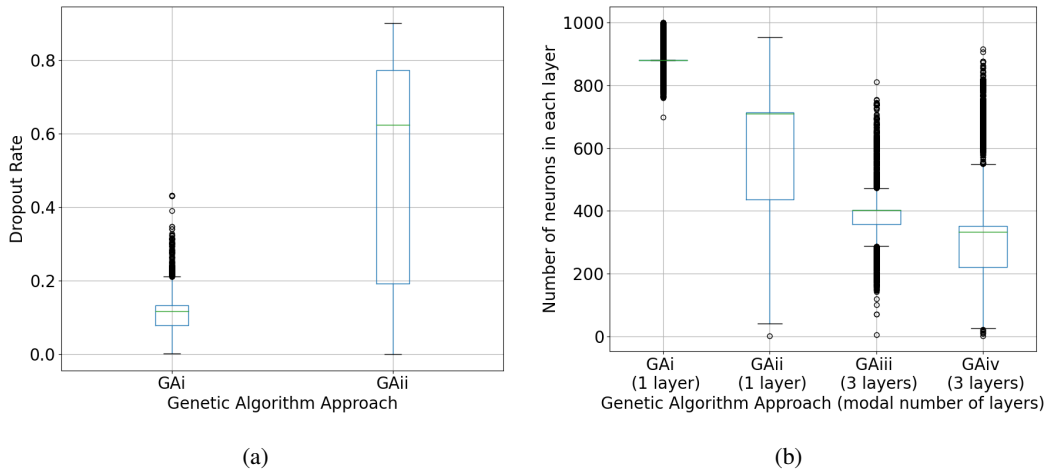


Figure 3: (a) Dropout rate for networks in the last 15 generations of cMLSGA with (GAi) the single objective of minimising Mean Absolute Error and (GAii) multi-objectives of minimising Maximum and Mean Absolute Error and (b) the number of neurons in each layer for networks in the last 15 generations of cMLSGA with (GAi), (GAii), (GAiii) and (GAiv).

184 supported by the fact that 74.4% of networks in the first 15 generations of (GAi) have 1 hidden layer,
 185 and that over 99.8% of the networks in the final 15 generations have 1 hidden layer, with 880 ± 17
 186 neurons in this layer, Figure 3b. This is significantly more neurons than those in the hidden layer of
 187 networks in the final 15 generations of (GAii) which range from 3-952 with a median value of 709,
 188 Figure 3b. The added objective of minimising Maximum Absolute Error in (GAii) may cause these
 189 slightly smaller networks to be more attractive as they are in a sense regularised by their size, as they
 190 have reduced modelling flexibility therefore are less likely to overfit and produce high Maximum
 191 Absolute Errors.

192 Another explanation for the difference in I2 rates chosen by (GAi) and (GAii) is the equivalence
 193 of I2 and dropout. Since I2 and dropout are equivalent up to a Fisher transformation, their use in
 194 conjunction is not recommended. The evidence for this is that (GAi) favours the highest I2 rate and
 195 has a median dropout rate in the final 15 generations of 0.116, whereas (GAii) favours the zero I2 rate
 196 and has a median dropout rate of 0.624, Figure 3a. This illustrates that the genetic algorithms will
 197 chose either I2 or dropout to minimise the Mean Absolute Relative Error. The I1 rates also support
 198 this hypothesis, as chosen rates for I1 regularisation in the final 15 generations are more comparable
 199 for (GAi) and (GAii).

200 **6 Optimisation using multiple performance measures: (GAiii) and (GAiv)**

201 For approaches (GAiii) and (GAiv) all neural network regularisation parameters are set to zero.
 202 The regularisation is performed by minimising different network performance measures, the Mean
 203 Absolute and Mean Fit to Median for (GAiii), and the Mean Absolute and Maximum Absolute for
 204 (GAiv). The trade-off between the two objectives produces regularised neural networks, without
 205 explicitly changing the architecture or loss function. The Mean Fit to Median is chosen as it indicates
 206 how close the relationships modelled by a network are to the conditional averages of the dataset, in
 207 many regression examples this is akin to the ground truth input-output relationships (Parkes et al.
 208 2021). The Maximum Absolute is chosen as for many industrial applications of machine learning the
 209 maximum prediction error is more pertinent than the mean error. The Mean Absolute Error is used
 210 instead of the Mean Squared Error in both approaches, as the conditional medians are closer to the
 211 ground truth input-output relationships in these datasets than the conditional means.

212 Differently shaped networks are favoured by (GAiii) and (GAiv), compared to (GAi) and (GAii),
 213 focusing on networks with 3 hidden layers and on average less than 400 neurons in each layer, Figure
 214 3b. These networks have 51 times the number of connections than the networks chosen in (GAi) and
 215 (GAii). Apart from (GAi), (GAiii) has the most consistently sized networks in the final 15 generations,
 216 with an interquartile range of 46 neurons, compared to (GAiv) which have an interquartile range of
 217 131 neurons. It is suggested that as the Mean Fit to Median Error biases networks towards specific
 218 input-output relationships, there is a smaller range of potential network architectures which habitually
 219 model these relationships. Whereas networks which minimise the Maximum Absolute Error are less
 220 restricted and can model a wider range of input and output relationships.

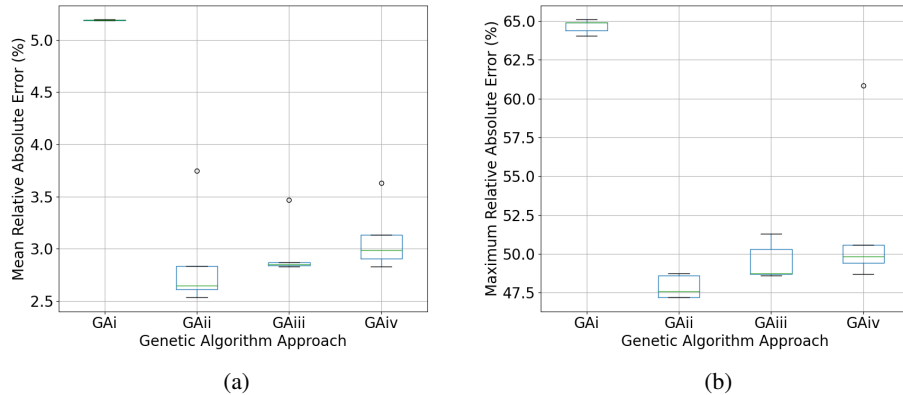


Figure 4: Mean Relative Absolute Error (a) and Maximum Absolute Error (b) from cMLSGA with (GAi) the single objective of minimising Mean Absolute Error and (GAii) multi-objectives of minimising Maximum and Mean Absolute Error, both optimising the parameters for 11, 12 regularisation and dropout in the networks, and (GAiii) and (GAiv) which do not use network regularisation but minimise Mean Fit to Median and Maximum Absolute Error respectively, alongside Mean Absolute Error

221 The Mean Absolute Relative Errors from networks in the Pareto fronts are $(2.97 \pm 0.25)\%$ for (GAiii)
 222 and $(3.10 \pm 0.28)\%$ for (GAiv). It is expected that (GAiv) would produce higher Mean Absolute
 223 Relative Errors as discussed above, minimising the Maximum Absolute Error should bias predictions
 224 towards the midpoint of the conditional output distributions, whereas minimising the Mean Absolute
 225 Error should bias predictions towards the median of these distributions. As it is established that noise
 226 in the output distribution is non-Gaussian, Figure 1, these values will not align so some sacrifice
 227 in Mean Absolute Error is expected from (GAiv). Both (GAiii) and (GAiv) produce comparable
 228 Maximum Absolute Errors, of $(49.5 \pm 1.1)\%$ and $(51.9 \pm 4.5)\%$. It is suggested that this is because,
 229 although the conditional median output value and conditional midpoint output value do not align
 230 for the majority of the input domain, they are sufficiently close to produce comparable Maximum
 231 Absolute Errors.

232 Across all four approaches, the genetic algorithm producing networks with the highest Mean Absolute
 233 Error is the approach which does not provide extra weighting to sparse areas of data. The approaches
 234 minimising Maximum Absolute Error are implicitly biased away from networks which predict the
 235 majority of the testing datapoints correctly, but predict one datapoint poorly, favouring networks
 236 which predict all testing datapoints to a moderate degree of error. Approach (GAiii) more explicitly
 237 weights prediction in sparse areas of data by favouring networks which model the conditional median
 238 of the dataset across all input domains, irrespective of the quantity of data across each input domain.
 239 The regression problem of ship power prediction is chosen in part because of its irregular data
 240 distribution; more than 9% of the dataset lies in less than a 0.5 knot interval of ship speed, Figure 1.
 241 This provides an explanation for the high testing errors from (GAi), where only the Mean Absolute
 242 Error is minimised, there is little incentive for the genetic algorithm to produce networks which
 243 generalise across the full range of the input domain well.

244 7 Comparison of the interpretability

245 To assess the interpretability of the networks selected by the four different approaches the learnt
 246 relationship between an input, the ship speed, and the output, shaft power, for the networks in
 247 the Pareto front of each approach are visualised, Figure 5. These are extracted with the following
 248 procedure: set all but one input variable to be constant at the mode; cycle the remaining variable from
 249 its minimum to its maximum recorded values with 150 points evenly spaced along the domain and
 250 run the new dataset through the trained network.

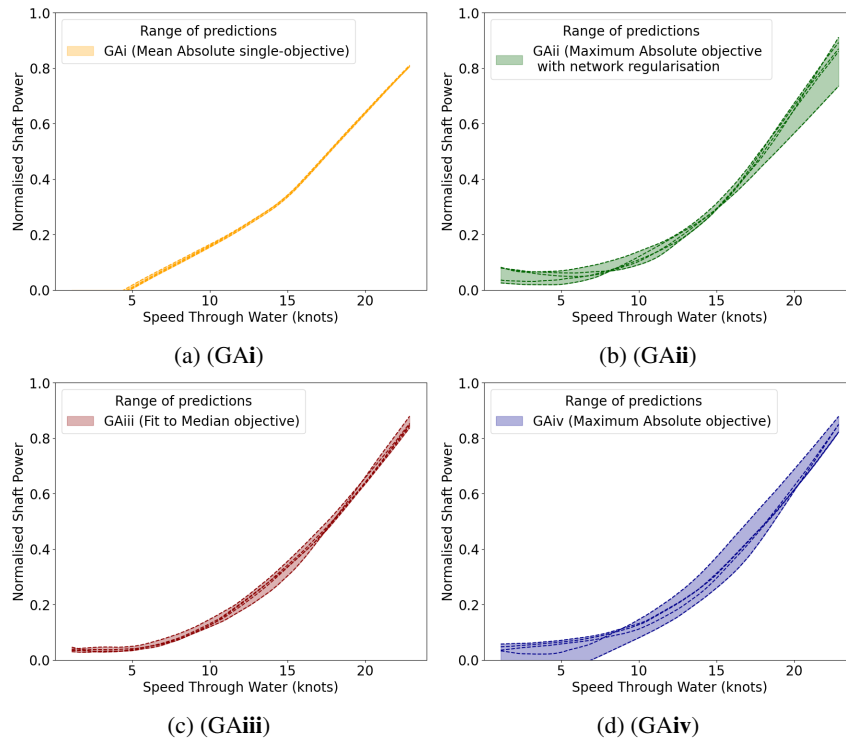


Figure 5: The learnt speed-power curves from 5 networks on the Pareto fronts of (GAii), (GAiii), (GAiv) and the 5 networks producing lowest Mean Absolute Relative Error from (GAi).

251 The approach which produces the most consistent speed-power relationships is (GAi), with an average
 252 variation of $1.8\%^2$, Figure 5a. However, the relationship modelled by the 5 networks with the lowest
 253 Mean Absolute Relative Error in (GAi) all approximate a piece-wise linear relationship which clearly
 254 underfits the dataset in Figure 1. The expected trend between ship speed through the water and shaft

²Average variation in just the speed-power curves are discussed in this section, but it is verified that all input-power curves follow the same trends with variation around 0.5% across input variables.

255 power is a cubic polynomial, therefore as well as producing the highest Mean Absolute Relative
256 Errors, networks chosen by (GAi) model the ground truth input-output relationships the worst out of
257 the four approaches. Both (GAii) and (GAiv) produce 5 fairly consistent speed-power curves, with
258 average variations of 5.9% and 10% respectively, Figures 5b and 5d. Both approaches approximate
259 smooth polynomial curves, although the degrees of the polynomials might differ, as multiple curves
260 intersect at various points along the speed axis. The spread of learnt relationships is greater at the
261 highest and lowest speeds for (GAii), with a decrease in spread for speeds of around 15 knots, where
262 many of the curves intersect. The curves from (GAiv) show equal spread across the speed domain.

263 The approach with both accurate and consistent learnt speed-power curves is (GAiii), with limited
264 intersections of curves and an average spread of 3.0%. It is suggested that the reason using the
265 Mean Fit to Median Error as an objective in a multi-objective genetic algorithm produces more
266 interpretable results, or more consistent learnt relationships, is because instead of encouraging the
267 networks to model more simple relationships. It encourages the networks to model the conditional
268 median functions of the dataset, supported by the increase in network connections. Whereas the other
269 approaches leave room for networks to fail to model the conditional averages, especially in irregularly
270 distributed and non-normally distributed datasets. The Mean Absolute Error values from networks
271 selected by (GAiii) are on average 0.1% higher than those from (GAii), and the Maximum Absolute
272 Error values are 1.6% higher.

273 A limitation of the approach is that the Fit to Median Error measure will likely perform best at
274 improving interpretability on datasets which violate the assumptions in Bishop (1995); the ship
275 powering example is chosen to illustrate this as it provides a clearly heteroscedastic dataset. For
276 applications where noise profiles are Gaussian, and there is no effect from latent or interrelated input
277 variables, the Fit to Median Error will not improve interpretability, but will perform the same as
278 conventional Minkowski-r metrics, either Mean Squared or Mean Absolute Error depending on the
279 convexity of input-output relationships.

280 Interestingly, the approach producing the lowest Mean Absolute and Maximum Absolute Errors
281 does not model the ground truth the most accurately. This creates a potential for negative societal
282 impacts, as the standard performance metrics for regression neural networks do not provide a full
283 picture of performance or expected behaviour. Interpretability of trained methods is essential for
284 safe application of machine learning in the real world, especially when automated methods are
285 used to replace experienced professionals. (GAii) demonstrates the same accuracy of approach as
286 those with standard network regularisation, but with a better fit to the ground truth. This approach
287 bypasses the need to use, and therefore to optimise the parameters of the regularisation methods.
288 If evolutionary computation is already being used to optimise network parameters, then compute
289 is saved by removing the network regularisation parameters l1, l2 and dropout. (GAiii) completed
290 300 generations in 46hours whereas (GAii) required 12 hours more computation to complete 300
291 generations.

292 8 Conclusion

293 Interpretable and accurate methods are required for widespread application of machine learning
294 to real-world regression problems. To automate the training of neural networks so that the result
295 is interpretable, three different genetic algorithm approaches are compared: one to minimise the
296 Maximum Absolute Error of the networks, which includes standard regularization using l1, l2 and
297 dropout; and two which do not use any network regularisation, one minimising the Mean Fit to
298 Median and one to minimise the Maximum Absolute Error. The results show that all three approaches
299 give similar Mean Absolute Errors from networks on their Pareto fronts, from 2.9% for the approach
300 with regularisation to 3.1% for the approach minimising Maximum Absolute Error. However, the
301 Mean Fit to the Median approach shows a considerably better interpretability, or fit to the ground
302 truth, with a spread in predicted input-output curves of 3% compared to a spread of 6% for the
303 approach using regularisation and 10% when minimising the Maximum Absolute Error.

304 References

305 Battaglia, P. W., Hamrick, J. B., Bapst, V., Sanchez-Gonzalez, A., Zambaldi, V., Malinowski, M.,
306 Tacchetti, A., Raposo, D., Santoro, A., Faulkner, R., Gulcehre, C., Song, F., Ballard, A., Gilmer,
307 J., Dahl, G., Vaswani, A., Allen, K., Nash, C., Langston, V., Dyer, C., Heess, N., Wierstra, D.,

- 308 Kohli, P., Botvinick, M., Vinyals, O., Li, Y. & Pascanu, R. (2018), ‘Relational inductive biases,
309 deep learning, and graph networks’.
- 310 Bishop, C. (1995), *Neural Networks for Pattern Recognition*, Oxford University Press, chapter 6,
311 pp. 194–225.
- 312 Chollet, F. et al. (2015), ‘Keras’, <https://keras.io>.
- 313 Dozat, T. (2016), ‘Incorporating nesterov momentum into adam’, *International Conference on*
314 *Learning Representations 2016* .
- 315 Duchi, J., Hazan, E. & Singer, Y. (2011), ‘Adaptive subgradient methods for online learning and
316 stochastic optimization’, *Journal of Machine Learning Research* **12**(61), 2121–2159.
- 317 Grudniewski, P. A. & Sobey, A. J. (2019), Do general genetic algorithms provide benefits when
318 solving real problems?, in ‘2019 IEEE Congress on Evolutionary Computation (CEC)’, IEEE,
319 pp. 1822–1829.
- 320 Grudniewski, P. A. & Sobey, A. J. (2021), ‘cMLSGA: a co-evolutionary multi-level selection genetic
321 algorithm for multi-objective optimization’.
- 322 Hanson, S. & Burr, D. (1987), Minkowski-r back-propagation: learning in connectionist models with
323 non-euclidian error signals., in ‘Neural Information Processing Systems (NIPS 1987)’.
- Hinton, G., Srivastava, N. & Swersky, K. (2012), ‘Lecture 6a:overview of mini-batch gradient
descent’.
URL: http://www.cs.toronto.edu/~tijmen/csc321/slides/lecture_slides lec6.pdf
- 324 Jeon, M., Noh, Y., Shin, Y., Lim, O., Lee, I. & Cho, D. (2018), ‘Prediction of ship fuel consumption
325 by using an artificial neural network’, *Journal of Mechanical Science and Technology* **32**(12), 5785–
326 5796.
- 327 Jin, Y., Okabe, T. & Sendhoff, B. (2004), Neural network regularization and ensembling using
328 multi-objective evolutionary algorithms, in ‘Proceedings of the 2004 Congress on Evolutionary
329 Computation (IEEE Cat. No.04TH8753)’, Vol. 1.
- 330 Kingma, D. P. & Ba, J. (2014), ‘Adam: A method for stochastic optimization’, *arXiv preprint*
331 *arXiv:1412.6980* .
- 332 Kumar, P., Batra, S. & Raman, B. (2021), ‘Deep neural network hyper-parameter tuning through
333 twofold genetic approach’, *Soft Computing* pp. 1–25.
- 334 Le, L., Lee, G., Park, K. & Kim, H. (2020), ‘Neural network-based fuel consumption estimation for
335 container ships in korea’, *Maritime Policy & Management* pp. 1–18.
- 336 Leifsson, L., Sævarsdóttir, H., Sigurðsson, S. & Vésteinsson, A. (2008), ‘Grey-box modeling of an
337 ocean vessel for operational optimization’, *Simulation Modelling Practice and Theory* **16**(8), 923–
338 932.
- 339 Liang, Q., Tsvete, H. A. & Brinks, H. W. (2019), Prediction of vessel propulsion power using machine
340 learning on ais data, ship performance measurements and weather data, in ‘Journal of Physics:
341 Conference Series’, Vol. 1357, p. 012038.
- 342 Luketina, J., Berglund, M., Greff, K. & Raiko, T. (2016), ‘Scalable gradient-based tuning of continu-
343 ous regularization hyperparameters’.
- 344 Nowlan, S. J. & Hinton, G. E. (1992), ‘Simplifying neural networks by soft weight sharing’, *Neural*
345 *Computation* .
- 346 Parkes, A. I., Savasta, T. D., Sobey, A. J. & Hudson, D. A. (2019), Efficient vessel power prediction
347 in operational conditions using machine learning., in ‘Practical Design of Ships and Other Floating
348 Structures(PRADS), September 2019, Yokohama, Japan’.
- 349 Parkes, A. I., Sobey, A. J. & Hudson, D. A. (2018), ‘Physics-based shaft power prediction for large
350 merchant ships using neural networks’, *Ocean Engineering* **166**, 92–104.

- 351 Parkes, A. I., Sobey, A. J. & Hudson, D. A. (2021), ‘Towards error measures which influence a
352 learners inductive bias to the ground truth’.
- 353 Pedersen, B. P. & Larsen, J. (2009), Prediction of full-scale propulsion power using artificial neural
354 networks, in ‘Proceedings of the 8th international conference on computer and IT applications in
355 the maritime industries (COMPIT’09), Budapest, Hungary May’, pp. 10–12.
- 356 Petersen, J. P., Jacobsen, D. J. & Winther, O. (2012), ‘Statistical modelling for ship propulsion
357 efficiency’, *Journal of marine science and technology* **17**(1), 30–39.
- 358 Smith, C. & Jin, Y. (2014), ‘Evolutionary multi-objective generation of recurrent neural network
359 ensembles for time series prediction’, *Neurocomputing* **143**, 302–311.
- 360 Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I. & Salakhutdinov, R. (2014), ‘Dropout:
361 a simple way to prevent neural networks from overfitting’, *The Journal of Machine Learning
362 Research* **15**(1), 1929–1958.
- 363 Tani, L., Rand, D., Veelken, C. & Kadastik, M. (2021), ‘Evolutionary algorithms for hyperparameter
364 optimization in machine learning for application in high energy physics’, *The European Physical
365 Journal C* **81**(2), 1–9.
- 366 Wager, S., Wang, S. & Liang, P. (2013), ‘Dropout training as adaptive regularization’, *arXiv preprint
367 arXiv:1307.1493*.
- 368 Wang, B., Sun, Y., Xue, B. & Zhang, M. (2019), Evolving deep neural networks by multi-objective
369 particle swarm optimization for image classification, in ‘Proceedings of the Genetic and Evolution-
370 ary Computation Conference’, GECCO ’19, p. 490–498.
- 371 Willard, J., Jia, X., Xu, S., Steinbach, M. & Kumar, V. (2020), ‘Integrating physics-based modeling
372 with machine learning: a survey’.
- 373 Yang, S., Tian, Y., He, C., Zhang, X., Tan, K. C. & Jin, Y. (2021), ‘A gradient-guided evolutionary
374 approach to training deep neural networks’, *IEEE Transactions on Neural Networks and Learning
375 Systems* pp. 1–15.
- 376 Zeiler, M. D. (2012), ‘Adadelta: an adaptive learning rate method’.

377 Checklist

378 The checklist follows the references. Please read the checklist guidelines carefully for information on
379 how to answer these questions. For each question, change the default **[TODO]** to **[Yes]**, **[No]**, or
380 **[N/A]**. You are strongly encouraged to include a **justification to your answer**, either by referencing
381 the appropriate section of your paper or providing a brief inline description. For example:

- 382 • Did you include the license to the code and datasets? **[Yes]** See Section ??.
- 383 • Did you include the license to the code and datasets? **[No]** The code and the data are
384 proprietary.
- 385 • Did you include the license to the code and datasets? **[N/A]**

386 Please do not modify the questions and only use the provided macros for your answers. Note that the
387 Checklist section does not count towards the page limit. In your paper, please delete this instructions
388 block and only keep the Checklist section heading above along with the questions/answers below.

389 1. For all authors...

- 390 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
391 contributions and scope? **[Yes]** The scope of the paper (regression problem automation)
392 is made clear in the abstract. The claims of improved interpretability using the Fit to
393 Median Error measure are reflected in Figure 5, as well as the body of text.
- 394 (b) Did you describe the limitations of your work? **[Yes]** The limitations of the Fit to
395 Median are discussed in Section 7, and the limitations relating to compute are discussed
396 in Section 3.

- 397 (c) Did you discuss any potential negative societal impacts of your work? [Yes] See
398 Section 6
- 399 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
400 them? [Yes]
- 401 2. If you are including theoretical results...
- 402 (a) Did you state the full set of assumptions of all theoretical results? [No] Citation of
403 (Bishop 1995) is included instead of explicit statement of all assumptions, although the
404 most pertinent assumptions are discussed in Section 1
- 405 (b) Did you include complete proofs of all theoretical results? [N/A] No theorems are
406 posited, the paper explores a practical comparison of common methods.
- 407 3. If you ran experiments...
- 408 (a) Did you include the code, data, and instructions needed to reproduce the main ex-
409 perimental results (either in the supplemental material or as a URL)? [Yes] Instruc-
410 tions needed to replicate the results on a different dataset are provided in the main
411 text; Sections 3 and 2 specify the experimental setup. The code is based on the
412 cMLSGA algorithm, available at https://www.bitbucket.org/***** (redacted
413 for anonymity) which is provided in Section 3. The data is proprietary so is not pro-
414 vided.
- 415 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
416 were chosen)? [Yes] See Sections 3 and 2
- 417 (c) Did you report error bars (e.g., with respect to the random seed after running ex-
418 periments multiple times)? [Yes] All results are reported to plus/minus the standard
419 deviation, and all figures show multiple runs with respect to either a random seed or an
420 evolutionary method.
- 421 (d) Did you include the total amount of compute and the type of resources used (e.g., type
422 of GPUs, internal cluster, or cloud provider)? [Yes] See Section 3
- 423 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 424 (a) If your work uses existing assets, did you cite the creators? [Yes] The creators of the
425 genetic algorithm cMLSGA are cited in Section 3 and the data is acknowledged to
426 belong to Silverstream Technologies Ltd in Section 4
- 427 (b) Did you mention the license of the assets? [Yes] The code is licensed under the GNU
428 General Public License. The data is not licensed.
- 429 (c) Did you include any new assets either in the supplemental material or as a URL? [No]
- 430 (d) Did you discuss whether and how consent was obtained from people whose data
431 you're using/curating? [No] The work is performed in collaboration with Silverstream
432 Techonologies Ltd.
- 433 (e) Did you discuss whether the data you are using/curating contains personally identifiable
434 information or offensive content? [No] The shaft power variable is normalised to
435 remove any possibility of identifying the vessel in question.
- 436 5. If you used crowdsourcing or conducted research with human subjects...
- 437 (a) Did you include the full text of instructions given to participants and screenshots, if
438 applicable? [N/A]
- 439 (b) Did you describe any potential participant risks, with links to Institutional Review
440 Board (IRB) approvals, if applicable? [N/A]
- 441 (c) Did you include the estimated hourly wage paid to participants and the total amount
442 spent on participant compensation? [N/A]