

Improving Pre-Trained Weights Through Meta-Heuristics Fine-Tuning

Gustavo H. de Rosa, Mateus Roder, João Paulo Papa
Department of Computing
São Paulo State University
Bauru - SP, Brazil
{gustavo.rosa, mateus.roder, joao.papa}@unesp.br

Claudio F. G. dos Santos
Department of Computing
São Carlos Federal University
São Carlos - SP, Brazil
claudio.santos@ufscar.br

Abstract—Machine Learning algorithms have been extensively researched throughout the last decade, leading to unprecedented advances in a broad range of applications, such as image classification and reconstruction, object recognition, and text categorization. Nonetheless, most Machine Learning algorithms are trained via derivative-based optimizers, such as the Stochastic Gradient Descent, leading to possible local optimum entrapments and inhibiting them from achieving proper performances. A bio-inspired alternative to traditional optimization techniques, denoted as meta-heuristic, has received significant attention due to its simplicity and ability to avoid local optimums imprisonment. In this work, we propose to use meta-heuristic techniques to fine-tune pre-trained weights, exploring additional regions of the search space, and improving their effectiveness. The experimental evaluation comprises two classification tasks (image and text) and is assessed under four literature datasets. Experimental results show nature-inspired algorithms’ capacity in exploring the neighborhood of pre-trained weights, achieving superior results than their counterpart pre-trained architectures. Additionally, a thorough analysis of distinct architectures, such as Multi-Layer Perceptron and Recurrent Neural Networks, attempts to visualize and provide more precise insights into the most critical weights to be fine-tuned in the learning process.

Index Terms—Machine Learning, Meta-Heuristic Optimization, Weights, Fine-Tuning

I. INTRODUCTION

Intelligence-based systems brought better insights into decision-making tasks and withdrew part of the humans’ burden in recurring tasks, where most of these advances have arisen from research fostered by Artificial Intelligence (AI) [1] and Machine Learning (ML) [2]. They have been incorporated in a wide range of applications, such as autonomous driving [3], text classification [4], image and object recognition [5], and medical analysis [6], among others.

The increasing demand for more complex tasks and the ability to solve unprecedented problems strengthened an ML sub-area, denoted as Deep Learning (DL) [7]. DL algorithms are known for employing deep neural networks and millions of parameters to model the intrinsic nature of the human brain [8], i.e., learn how humans can process information

through their visual system and how they can communicate between themselves. Nevertheless, such learning is conditioned to the training data and the model’s parameters and often does not reproduce the real-world environment, leading to undesired behavior, known as underfitting/overfitting [9].

Even though proper training usually accompanies underfitting/overfitting, it is common to perceive a poor performance when the model is collated with unseen data. This discrepancy lies in the fact that the model “memorized” the training data instead of learning its patterns, thus not reproducing the desired outputs when applied to slightly-different data (test data). The best approach to overcome this problem would be to employ combinations of all possible parameters and verify whether they are suitable or not when applied to the testing data. Nevertheless, such an approach is unfeasible regarding DL architectures due to their vast number of parameters and exponential complexity [10].

On the other hand, a more feasible approach stands for optimization procedures, where parameters are optimized according to an objective function instead of being joined in all possible combinations. A recent technique, denoted as meta-heuristic, has attracted considerable attention in the last years, mainly due to its simple heuristics and ability to optimize non-differentiable functions. For instance, Rosa et al. [11] used the Harmony Search algorithm for fine-tuning Convolutional Neural Networks (CNN) hyperparameters, achieving improved results over the benchmark architectures. At the same time, Rodrigues et al. [12] explored single- and multi-objective meta-heuristic optimization in Machine Learning problems, such as feature extraction and selection, hyperparameter tuning, and unsupervised learning. Furthermore, Wang et al. [13] presented a fast-ranking version of the Particle Swarm Optimization algorithm to fine-tune CNN hyperparameters and remove the fitness function evaluation cost.

Nevertheless, most of the literature works focus on only optimizing the model’s hyperparameters (learning rate, number of units, momentum, weight decay, dropout) [14]–[17] instead of optimizing its parameters (layers’ weights and biases) [18]. Usually, parameters are optimized during the learning procedure through gradient-based approaches, such as the Stochastic Gradient Descent, yet they might benefit from the meta-heuristic techniques’ exploration and exploitation capabilities.

The authors are grateful to São Paulo Research Foundation (FAPESP) grants #2013/07375-0, #2014/12236-1, #2019/07665-4, #2019/02205-5, and #2020/12101-0, and to the Brazilian National Council for Research and Development (CNPq) #307066/2017-7 and #427968/2018-6.

This work proposes an additional fine-tuning after the model’s training, aiming to explore unknown search space regions that gradient-based optimizers could not find. Such an approach is conducted by exploring weights under pre-defined bounds and evaluating them according to an objective function (accuracy over the validation set). Therefore, the main contributions of this work are three-fold: (i) to introduce meta-heuristic optimization directly to the model’s parameters, (ii) to provide insightful analysis of whether gradient-based optimizers achieved feasible regions or not, and (iii) to fill the lack of research regarding meta-heuristic optimization applied to Machine Learning algorithms.

The remainder of this paper is organized as follows. Section II presents a theoretical background regarding the employed ML architectures, e.g., Multi-Layer Perceptrons and Recurrent Neural Networks. Section III introduces a brief explanation about meta-heuristic optimization, as well as the Genetic Algorithm and Particle Swarm Optimization. Section IV presents the mathematical formulation of the proposed approach, its complexity analysis, the employed datasets, and the experimental setup. Finally, Section V discusses the experimental results while Section VI states the conclusions and future works.

II. MACHINE LEARNING

This section introduces brief concepts regarding the Multi-Layer Perceptron and the Long Short-Term Memory.

A. Multi-Layer Perceptron

Multi-Layer Perceptron has arisen from the traditional Perceptrons and represents a type of feed-forward artificial neural network. Instead of having a single intermediate layer and a linear function as the Perceptron has, the MLP architecture comprises multiple Perceptrons arranged in hidden layers and followed by non-linear activations, which allows it to distinguish non-linearly separable data. Figure 1 illustrates the standard architecture of an MLP.

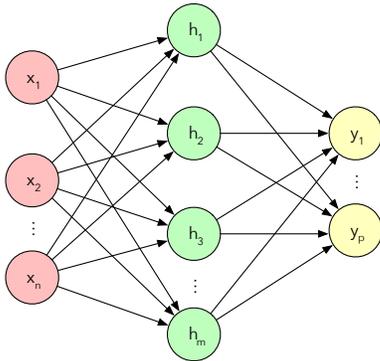


Fig. 1. Multi-Layer Perceptron standard architecture.

Commonly, MLP-based networks are employed in supervised learning tasks and trained through the Backpropagation [19] algorithm, which calculates the output error and corrects the model’s weights according to the derivative of

the activation function and part of the error. Modern MLP architectures use more sophisticated activation functions, such as the Rectified Linear Unit (ReLU) [20], instead of only relying on traditional ones, e.g., sigmoid and hyperbolic tangents. Additionally, they have been used as the foundation of several Deep Learning architectures, such as VGG [21] and Inception [22] ones.

B. Long Short-Term Memory

Hochreiter et al. [23] proposed the Long Short-Term Memory networks, which are particular types of Recurrent Neural Networks designed to learn information through long periods. Their main difference when compared to traditional RNNs lies in their hidden layer, which employs a cell (unit) with four gate mechanism that interacts between themselves. Figure 2 illustrates such a cell architecture.

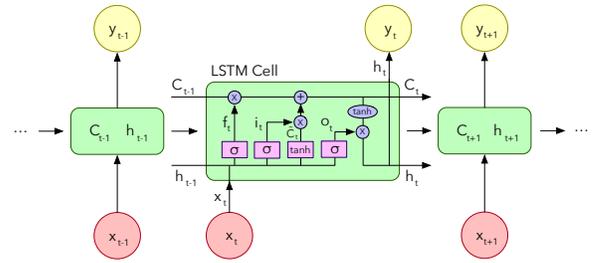


Fig. 2. Long Short-Term Memory cell architecture.

An LSTM cell is represented by the following variables: cell state $C_t \in \mathbb{R}^k$, input neurons $x_t \in \mathbb{R}^n$, hidden layer $h_t \in \mathbb{R}^k$ and output neurons $y_t \in \mathbb{R}^p$. The cell state works as a conveyor, running through the belt and suffering linear combinations. Additionally, LSTM’s cell may add or remove self-contained information through its gate mechanisms.

The gates allow or deny the flow of information and usually constitute non-linear neural layers, such as sigmoid activation and point-wise operations. The output of a sigmoid function creates a real number between 0 and 1, which describes the amount of information the gate should propagate. Note that 0 stands for no-information, while 1 stands for the full-information. Commonly, an LSTM cell is regulated by three gates: input, forget, and output.

III. META-HEURISTIC OPTIMIZATION

Optimization problems consist of maximizing or minimizing mathematical functions through potential values, while the optimization procedure aims to find those values given a pre-defined domain. Traditional optimization methods [24], such as the combinatorial and iterative methods, such as the Grid-Search, the Newton method, the Quasi-Newton, the Gradient Descent, Interpolation methods, use the evaluation of gradients and Hessians, thus, being only practical when applied to differentiable functions. Additionally, they elevate the computational burden due to calculating first- and second-order derivatives.

Alternatively, an approach known as meta-heuristic has been applied to solve optimization tasks. Meta-heuristic techniques [25] consists of high-level algorithms designed to create and select heuristics capable of producing feasible solutions to the optimization problem. Hence, meta-heuristic optimization is a procedure that connects notions of *exploration*, used to conduct extensive searches throughout the space, and *exploitation*, used to improve potential solutions based on their neighborhoods.

A. Genetic Algorithm

Genetic Algorithm is a traditional evolutionary-based algorithm inspired by the process of natural selection. It commonly relies on biological operators, such as selection, mutation, and crossover, and generates feasible solutions for optimization tasks. Each individual is represented by an n -dimensional position array \mathbf{x} , where each dimension stands for a decision variable, and a fitness value associated with this particular position, i.e., $f(\mathbf{x})$.

The Genetic Algorithm’s main objective is to evolve a population of m individuals in an iterative way, where the so-called biological operators are applied over the population to create a more fit population, e.g., lower fitness values in minimization problems. During each iteration/generation, the population is evaluated, and a set of $p_s \times m$ individuals are stochastically selected from it, where p_s stands for the selection proportion selection and m the number of individuals.

Furthermore, the selected individuals are divided into pairs to form the “parents” and bred into offsprings according to a crossover probability, denoted as p_c . The offsprings are new individuals who share characteristics inherited from their parents, i.e., they have randomly selected positions from their mother and father. Afterward, the offsprings are mutated according to a mutation probability p_m , which occasionally adds a noise value to one of the offsprings’ positions.

Finally, the population is re-evaluated, and the iterative process continues until a stop criterion is satisfied, such as an epsilon or a maximum number of generations. Combining the biological operators’ explorability and exploitability allows the population to convergence to more appropriate values, thus producing feasible solutions to optimization tasks.

B. Particle Swarm Optimization

Particle Swarm Optimization is a nature-inspired algorithm that designs each agent as a bird that belongs to a swarm and searches for optimal food sources. Each agent is represented by a (\mathbf{x}, \mathbf{v}) tuple, where \mathbf{x} stands for its position and \mathbf{v} for its velocity. The initial position \mathbf{x} is described as an n -dimensional randomly vector, while the velocity \mathbf{v} is an n -dimensional vector of zeros, where each dimension stands for the decision variable. Additionally, the objective corresponds to searching the most likely decision variables, which maximizes or minimizes a target function.

Let \mathbf{v}_i^t be the velocity of an agent i at iteration t , belonging to a swarm of size M , such that $i \in \{1, 2, \dots, M\}$. One can update its velocity according to Equation 1, as follows:

$$\mathbf{v}_i^{t+1} = w\mathbf{v}_i^t + c_1r_1(\mathbf{x}_i^* - \mathbf{x}_i^t) + c_2r_2(\mathbf{g} - \mathbf{x}_i^t), \quad (1)$$

where \mathbf{x}_i^* stands for the best position obtained by agent i , and \mathbf{g} denotes the current best solution. Additionally, w , c_1 , and c_2 stand for the inertia weight, the social parameter, and the cognitive ratio, respectively. Finally, r_1 and r_2 are uniformly distributed random numbers in the range $[0, 1]$.

Furthermore, let \mathbf{x}_i^t be the position of an agent i at iteration t . One can update its position according to Equation 2, as follows:

$$\mathbf{x}_i^{t+1} = \mathbf{x}_i^t + \mathbf{v}_i^{t+1}. \quad (2)$$

IV. METHODOLOGY

This section presents a brief discussion regarding the proposed approach, its complexity analysis, the employed datasets, and the experimental setup.

A. Proposed Approach

The proposed approach aims to pre-train an architecture through its standard pipeline, e.g., stochastic gradient optimization across a training set, followed by a fine-tuning using meta-heuristic optimization across a validation set (post-trained). The idea is to use meta-heuristic techniques to explore the search space better and intensify a promising solution found by the traditional optimization algorithm.

Let θ be the weights of a pre-trained neural network, where $\theta \in \mathbb{R}^n$ and n stands for the number of weights. Additionally, let Δ be a defined value which stands for the search bounds around θ . The initial solutions are randomly sampled from the $[\theta - \Delta, \theta + \Delta]$ interval¹ and feed to the meta-heuristic technique, which will explore the search space and find the most suitable solutions given a fitness function, i.e., accuracy over the validation set. At the end of the optimization procedure, both post-trained and pre-trained networks are evaluated over testing sets and compared. Figure 3 illustrates an overview of the proposed approach pipeline.

B. Complexity Analysis

Let $O(\iota)$ and $O(\zeta)$ be the complexity of training and validating a network for each epoch, respectively. We can observe that in the proposed approach, we opted to pre-train the whole network using T_n iterations, with additional validations for every epoch. Hence, the whole pre-training procedure complexity is depicted by Equation 3, as follows:

$$O(\iota) \cdot T_n + O(\zeta) \cdot T_n = O(\iota + \zeta) \cdot T_n. \quad (3)$$

The proposed approach intends to provide an additional optimization step after the network’s pre-training, where agents will encode the pre-trained weights and biases as their positions, search for better solutions throughout the space and evaluate the fitness function (validation). Let T_o be the number of optimization iterations, m the number of agents,

¹Note that although the search interval (Δ) is equal for each variable in θ , they may have a different value.

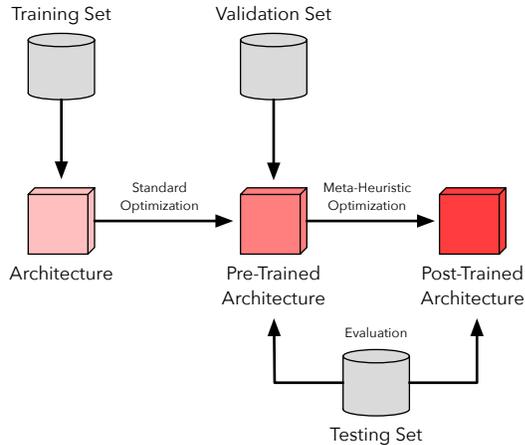


Fig. 3. Proposed approach pipeline.

and the whole optimization procedure complexity described by Equation 4, as follows:

$$O(\zeta) \cdot T_o \cdot m. \quad (4)$$

Therefore, summing both Equations 3 and 4 together, it is possible to achieve the proposed approach complexity, as described by Equation 5:

$$O(\iota + \zeta) \cdot T_n + O(\zeta) \cdot T_o \cdot m. \quad (5)$$

Note that if $O(\zeta) \rightarrow 0$, both approaches will have the same complexity. Unfortunately, this pattern will not often happen, mainly due to the complexity of larger models, such as CNNs and RNNs. However, the advantage is that $O(\zeta)$ is significantly smaller than $O(\iota)$, thus adding only a small burden over the standard pre-training.

C. Datasets

We considered four datasets in the experimental section, being two image- and two text-based:

- CIFAR-10 [26]: is a subset image database from the “80 million tiny images” dataset. Composed of 60,000 32x32 colour images divided in 10 classes, with 6,000 images per class. It is divided into five training batches and one test batch, each containing 10,000 images;
- CIFAR-100 [26]: is almost like the CIFAR-10 dataset, yet it provides a more challenging problem. Composed of 60,000 32x32 colour images divided in 100 classes, with 600 images per class. It is also split into five training batches and one test batch, each containing 10,000 images;
- IMDB Large Movie Reviews [27]: is a dataset for the task of sentiment analysis (binary classification) and is composed of 50,000 raw-text movie reviews, equally split into training and testing sets;
- Stanford Sentiment Treebank (SST) [28]: is composed of three-level sentiment analysis (positive, negative, and

neutral) of syntactically plausible phrase in thousands of sentences from Rotten Tomatoes movie reviews (more than 200,000).

D. Experimental Setup

The proposed approach is evaluated amongst two distinct tasks: image classification (CIFAR-10 and CIFAR-100) and sentiment analysis (SST and IMDB). Regarding the former task, we employed a standard Multi-Layer Perceptron (input, hidden, and output layers), while the latter task uses a Long Short-Term Memory (embedding, hidden, and output layers). Table I describes the hyperparameters used in the image- and text-based datasets, respectively.

TABLE I
HYPERPARAMETERS CONFIGURATION USED IN IMAGE CLASSIFICATION (MLP) AND SENTIMENT ANALYSIS (LSTM) TASKS.

Hyperparameter	CIFAR-10	CIFAR-100	IMDB	SST
n_i (input units)	3,072	3,072	–	–
n_e (embedding units)	–	–	128	64
n (hidden units)	1,024	2,048	512	128
n_o (output units)	10	100	2	3
e (epochs)	50	50	10	20
bs (batch size)	100	100	128	8
η (learning rate)	0.0001	0.0001	0.01	0.0001

Regarding the meta-heuristic techniques, we opted to use an evolutionary-based algorithm denoted as Genetic Algorithm and a swarm-based one, known as Particle Swarm Optimization. Both algorithms are available in the Opyoptimizer [16] package and the paper’s source code at GitHub². For each meta-heuristic, we employed three search space configurations, as follows:

- α : 10 agents optimized over 5 iterations;
- β : 50 agents optimized over 25 iterations;
- γ : 100 agents optimized over 50 iterations.

Additionally, Table II describes the meta-heuristics parameters configuration.

TABLE II
META-HEURISTICS PARAMETERS CONFIGURATION.

Meta-heuristic	Parameters
GA	$p_s = 0.75 \mid p_c = 0.5 \mid p_m = 0.25$
PSO	$w = 0.7 \mid c_1 = 1.7 \mid c_2 = 1.7$

Finally, to provide a more robust analysis, we conduct 10 runnings with different seeds (different splits) for each (metaheuristic, dataset) pair, followed by a statistical analysis according to the Wilcoxon signed-rank test [29] with a 0.05 significance.

²The source code is available at https://github.com/gugarosa/mh_fine_tuning.

V. EXPERIMENTS AND DISCUSSION

This section presents the experimental results concerning the employed datasets and tasks. Furthermore, we present additional discussions regarding the effect on weights’ optimization and how the search space’s bounds influence the fine-tuning.

A. Overall Discussion

Table III describes the experimental results over the image classification task, which comprehends both CIFAR-10 (top) and CIFAR-100 (bottom) datasets. The most striking point to elucidate is that almost every meta-heuristic model could slightly improve the baseline accuracy, i.e., every meta-heuristic on the CIFAR-10 dataset, as well as α -GA-MLP and β -PSO-MLP on the CIFAR-100 dataset. Additionally, considering both datasets, α -GA-MLP obtained the best accuracy and recall metrics amongst the evaluated models (underlined cells), yet every evaluated model was statistically similar according to the Wilcoxon signed-rank test (bolded cells).

Table IV describes the experimental results over the sentiment classification task, which comprehends both IMDB (top) and SST (bottom) datasets. The meta-heuristics performance was marginally inferior to the baseline architecture (best results marked by underlined cells) considering all metrics and models. Such behavior is possibly explained by the fact that meta-heuristics optimized the last fully-connected layer, which is not the most descriptive layer in a recurrent network. Even though they did not achieve outstanding results, every model has been statistically similar to the baseline experiment according to the Wilcoxon signed-rank test (bolded cells) and reinforces the meta-heuristics capacity in attempting to search for more reasonable minimum points.

Nonetheless, we think that meta-heuristics’ performance was hindered by their high dimensional space (thousands of features), which often bottlenecks the exploration and exploitation phases as a low number of agents (roughly 10^2 smaller order of magnitude) are not capable of traversing local optima and hence not converging to the most proper locations. On the other hand, the experimental results showed that meta-heuristics could sway the pre-trained weights in search of better values, which fosters feasible improvements and the necessity of additional experimentations.

B. How Weights Influence the Fine-Tuning?

Figure 4 illustrates the weights distribution concerning LSTM and γ -GA-LSTM models over the IMDB dataset, while Figure 5 depicts the weights concerning LSTM and γ -PSO-LSTM models over the SST dataset. In both figures, it is hard to observe differences between plots (a) and (b), although plot (b) has slightly lighter colors when compared to plot (a). Such behavior indicates that the meta-heuristic algorithms found local optima when searching for weights values inferior to the pre-trained ones.

Furthermore, the absence of a critical difference between both plots indicates that the meta-heuristics could not adequately explore the search space, only attaining minimal

superior or inferior values. As previously mentioned, this may happen due to an optimization in a not so “important” layer, hence marginally improving or decreasing the optimized networks’ performance.

C. Bounding the Search Space

Finally, we opted to conduct the last experiment to verify the influence of distinct search bounds over the proposed methodology. Table V describes the accuracy results between distinct search bounds ($\Delta = 0.0001$ and $\Delta = 0.001$) over CIFAR-10 (top) and CIFAR-100 (bottom) datasets. Even though a larger Δ provided the best results (underlined cells), every model has been statistically similar to each other according to the Wilcoxon signed-rank test (bolded cells). Such performance strengthens the hypothesis that optimizing high dimensional search spaces requires more complex structures, such as increased agents and iterations, to accomplish an adequate convergence.

Moreover, Table VI describes the accuracy results between distinct search bounds ($\Delta = 0.0001$ and $\Delta = 0.001$) over IMDB (top) and SST (bottom) datasets. In such an experiment, it is possible to observe the same behavior depicted by Table V, where every meta-heuristic has been statistically similar to each other. Notwithstanding, one can also perceive that γ -based models with $\Delta = 0.001$ achieved the best results concerning both datasets, supporting that LSTMs may benefit from broader than narrower searches.

VI. CONCLUSION

This work presented an early draft of how to fine-tune neural networks’ performance through meta-heuristic techniques. Essentially, after training architectures through standard gradient descent algorithms, meta-heuristic algorithms attempt to explore the trained search space and find more suitable positions.

The experimental results showed that Multi-Layer Perceptron applied to image classification tasks could benefit from weights fine-tuning. While the meta-heuristics (GA and PSO) performance were statistically similar to the standard trained architecture, they could attain a higher mean accuracy and a lower standard deviation over the CIFAR-10 dataset (52.54 ± 0.52 against 52.40 ± 0.62) while obtaining a higher mean accuracy and higher standard deviation over the CIFAR-100 dataset (24.96 ± 0.33 against 24.93 ± 0.30). Such results strengthen the ability of meta-heuristics to fine-tune and better adjust the learned weights.

Regarding the sentiment analysis task conducted by the LSTMs, one can perceive that none meta-heuristic could obtain better metrics than the standard architecture, yet all results were statistically similar according to the Wilcoxon signed-rank test. Such behavior might be explained due to only fine-tuning the last fully-connected layer (layer before the Softmax activation) instead of fine-tuning the recurrent layer and ignoring the biases fine-tuning, which may help meta-heuristics in achieving more competitive results.

Finally, glimpsing through future works, we aim to extend the current framework by using new sets of meta-heuristics

TABLE III
EXPERIMENTAL RESULTS (%) OVER CIFAR-10 (TOP) AND CIFAR-100 (BOTTOM) DATASETS.

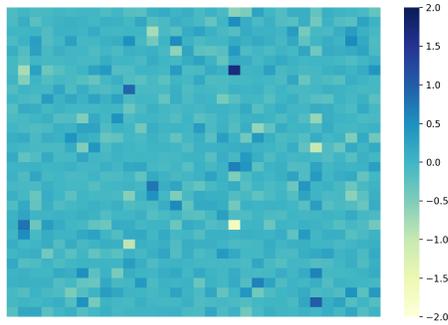
Model	Accuracy	Precision	Recall	F1-Score
MLP	52.40 ± 0.62	52.73 ± 0.57	52.40 ± 0.62	52.45 ± 0.60
α -GA-MLP	<u>52.54 ± 0.52</u>	<u>52.76 ± 0.54</u>	<u>52.54 ± 0.52</u>	<u>52.54 ± 0.52</u>
α -PSO-MLP	52.51 ± 0.56	52.72 ± 0.57	52.51 ± 0.56	52.50 ± 0.56
β -GA-MLP	52.49 ± 0.57	52.70 ± 0.57	52.49 ± 0.57	52.48 ± 0.58
β -PSO-MLP	52.53 ± 0.56	52.73 ± 0.57	52.53 ± 0.56	52.52 ± 0.56
γ -GA-MLP	52.52 ± 0.59	52.73 ± 0.62	52.52 ± 0.59	52.51 ± 0.60
γ -PSO-MLP	52.52 ± 0.59	52.73 ± 0.61	52.52 ± 0.59	52.51 ± 0.60
MLP	24.93 ± 0.30	25.59 ± 0.35	24.93 ± 0.30	24.74 ± 0.29
α -GA-MLP	<u>24.96 ± 0.33</u>	25.64 ± 0.40	<u>24.96 ± 0.33</u>	24.77 ± 0.32
α -PSO-MLP	24.93 ± 0.31	25.59 ± 0.33	24.93 ± 0.31	24.73 ± 0.30
β -GA-MLP	24.91 ± 0.33	25.63 ± 0.40	24.91 ± 0.33	24.75 ± 0.33
β -PSO-MLP	<u>24.96 ± 0.34</u>	25.61 ± 0.37	<u>24.96 ± 0.34</u>	24.76 ± 0.32
γ -GA-MLP	24.92 ± 0.31	<u>25.65 ± 0.38</u>	24.92 ± 0.31	<u>24.78 ± 0.32</u>
γ -PSO-MLP	24.92 ± 0.31	25.52 ± 0.39	24.92 ± 0.31	24.70 ± 0.31

TABLE IV
EXPERIMENTAL RESULTS (%) OVER IMDB (TOP) AND SST (BOTTOM) DATASETS.

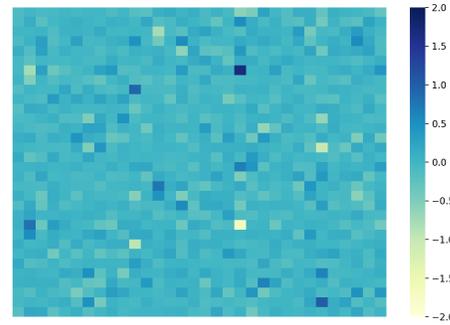
Model	Accuracy	Precision	Recall	F1-Score
LSTM	<u>49.21 ± 2.05</u>	<u>49.17 ± 1.69</u>	<u>49.21 ± 1.23</u>	<u>48.54 ± 1.98</u>
α -GA-LSTM	48.98 ± 1.76	48.91 ± 1.75	48.98 ± 1.62	48.14 ± 1.81
α -PSO-LSTM	48.97 ± 1.54	48.93 ± 1.60	48.92 ± 1.81	48.10 ± 1.90
β -GA-LSTM	48.98 ± 1.93	48.91 ± 1.90	48.98 ± 1.97	48.14 ± 1.77
β -PSO-LSTM	48.98 ± 1.88	48.90 ± 1.78	48.97 ± 1.99	48.16 ± 1.73
γ -GA-LSTM	49.03 ± 1.88	48.96 ± 1.77	49.03 ± 1.90	48.18 ± 1.81
γ -PSO-LSTM	48.99 ± 1.84	48.95 ± 1.73	48.97 ± 1.71	48.13 ± 1.77
LSTM	<u>55.31 ± 2.64</u>	<u>50.24 ± 1.48</u>	<u>48.89 ± 1.71</u>	<u>48.41 ± 1.92</u>
α -GA-LSTM	54.90 ± 2.85	49.95 ± 1.36	48.44 ± 1.89	47.95 ± 2.08
α -PSO-LSTM	54.93 ± 2.86	49.98 ± 1.38	48.47 ± 1.93	47.99 ± 2.13
β -GA-LSTM	54.91 ± 2.87	49.96 ± 1.38	48.45 ± 1.91	47.96 ± 2.10
β -PSO-LSTM	54.92 ± 2.89	49.97 ± 1.36	48.45 ± 1.93	47.98 ± 2.13
γ -GA-LSTM	54.90 ± 2.87	49.96 ± 1.39	48.44 ± 1.93	47.96 ± 2.12
γ -PSO-LSTM	54.95 ± 2.88	50.00 ± 1.38	48.48 ± 1.94	48.00 ± 2.14

algorithms, including fine-tuned biases, and extending the number of fine-tuned layers. Additionally, we aim at extending the current work to other neural network architectures, such as Convolutional Neural Networks and Transformers. We believe that the key to proper fine-tuning lies in selecting the most

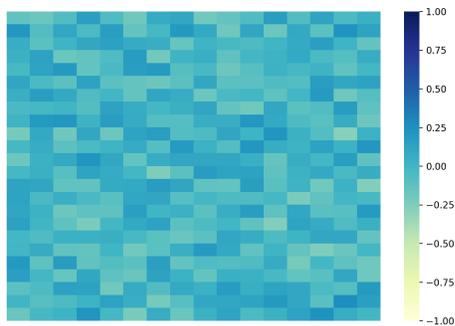
informative layer (a layer that is responsible for extracting the most important information) and a not so overwhelming search space (fewest feature spaces).



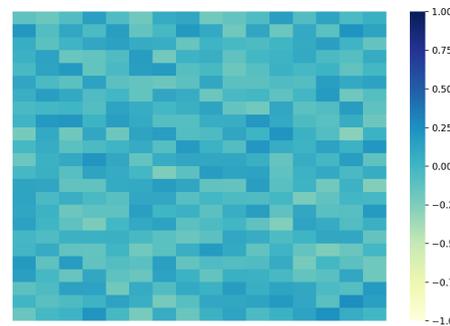
(a)



(b)

Fig. 4. IMDB weights matrix comparison between: (a) LSTM and (b) γ -GA-LSTM.

(a)



(b)

Fig. 5. SST weights matrix comparison between: (a) LSTM and (b) γ -PSO-LSTM.

REFERENCES

- [1] Y. Lu, "Artificial intelligence: a survey on evolution, models, applications and future trends," *Journal of Management Analytics*, vol. 6, no. 1, pp. 1–29, 2019.
- [2] R. Boutaba, M. A. Salahuddin, N. Limam, S. Ayoubi, N. Shahriar, F. Estrada-Solano, and O. M. Caicedo, "A comprehensive survey on machine learning for networking: evolution, applications and research opportunities," *Journal of Internet Services and Applications*, vol. 9, no. 1, p. 16, 2018.
- [3] E. Yurtsever, J. Lambert, A. Carballo, and K. Takeda, "A survey of autonomous driving: Common practices and emerging technologies," *IEEE Access*, vol. 8, pp. 58 443–58 469, 2020.
- [4] K. Kowsari, K. Jafari Meimandi, M. Heidarysafa, S. Mendu, L. Barnes, and D. Brown, "Text classification algorithms: A survey," *Information*, vol. 10, no. 4, p. 150, 2019.
- [5] P. Druzhkov and V. Kustikova, "A survey of deep learning methods and software tools for image classification and object detection," *Pattern Recognition and Image Analysis*, vol. 26, no. 1, pp. 9–15, 2016.
- [6] G. Litjens, T. Kooi, B. E. Bejnordi, A. A. A. Setio, F. Ciompi, M. Ghahfoorian, J. A. Van Der Laak, B. Van Ginneken, and C. I. Sánchez, "A survey on deep learning in medical image analysis," *Medical image analysis*, vol. 42, pp. 60–88, 2017.
- [7] S. Pouyanfar, S. Sadiq, Y. Yan, H. Tian, Y. Tao, M. P. Reyes, M.-L. Shyu, S.-C. Chen, and S. Iyengar, "A survey on deep learning: Algorithms, techniques, and applications," *ACM Computing Surveys (CSUR)*, vol. 51, no. 5, pp. 1–36, 2018.
- [8] N. Kriegeskorte and T. Golan, "Neural network models and deep learning," *Current Biology*, vol. 29, no. 7, pp. R231–R236, 2019.
- [9] M. Belkin, D. J. Hsu, and P. Mitra, "Overfitting or perfect fitting? risk bounds for classification and regression rules that interpolate," in *Advances in neural information processing systems*, 2018, pp. 2300–2311.
- [10] H. Y. Xiong, Y. Barash, and B. J. Frey, "Bayesian prediction of tissue-regulated splicing using rna sequence and cellular context," *Bioinformatics*, vol. 27, no. 18, pp. 2554–2562, Sep. 2011.
- [11] G. Rosa, J. Papa, A. Marana, W. Scheirer, and D. Cox, "Fine-tuning convolutional neural networks using harmony search," in *Iberoamerican Congress on Pattern Recognition*. Springer, 2015, pp. 683–690.
- [12] D. Rodrigues, J. P. Papa, and H. Adeli, "Meta-heuristic multi- and many-objective optimization techniques for solution of machine learning problems," *Expert Systems*, vol. 34, no. 6, p. e12255, 2017.
- [13] Y. Wang, H. Zhang, and G. Zhang, "cpsy-cnn: An efficient pso-based algorithm for fine-tuning hyper-parameters of convolutional neural networks," *Swarm and Evolutionary Computation*, vol. 49, pp. 114 – 123, 2019.
- [14] J. Bergstra, B. Komer, C. Eliasmith, D. Yamins, and D. D. Cox, "Hyperopt: a python library for model selection and hyperparameter optimization," *Computational Science & Discovery*, vol. 8, no. 1, p. 014008, 2015.
- [15] C. Yao, D. Cai, J. Bu, and G. Chen, "Pre-training the deep generative models with adaptive hyperparameter optimization," *Neurocomputing*, vol. 247, pp. 144–155, 2017.
- [16] G. H. de Rosa and J. P. Papa, "Soft-tempering deep belief networks parameters through genetic programming," *Journal of Artificial Intelligence and Systems*, vol. 1, no. 1, pp. 43–59, 2019.
- [17] J.-H. Han, D.-J. Choi, S.-U. Park, and S.-K. Hong, "Hyperparameter optimization using a genetic algorithm considering verification time in

TABLE V
ACCURACY RESULTS (%) BETWEEN DISTINCT SEARCH BOUNDS (Δ) OVER
CIFAR-10 (TOP) AND CIFAR-100 (BOTTOM) DATASETS.

Model	$\Delta = 0.0001$	$\Delta = 0.001$
α -GA-MLP	52.53 \pm 0.57	<u>52.54 \pm 0.52</u>
α -PSO-MLP	52.53 \pm 0.58	52.51 \pm 0.56
β -GA-MLP	52.53 \pm 0.58	52.49 \pm 0.57
β -PSO-MLP	52.52 \pm 0.56	52.53 \pm 0.56
γ -GA-MLP	<u>52.54 \pm 0.57</u>	52.52 \pm 0.59
γ -PSO-MLP	52.51 \pm 0.58	52.52 \pm 0.59
α -GA-MLP	24.94 \pm 0.31	<u>24.96 \pm 0.33</u>
α -PSO-MLP	24.94 \pm 0.32	24.93 \pm 0.31
β -GA-MLP	24.93 \pm 0.31	24.91 \pm 0.33
β -PSO-MLP	24.92 \pm 0.31	<u>24.96 \pm 0.34</u>
γ -GA-MLP	24.93 \pm 0.31	24.92 \pm 0.31
γ -PSO-MLP	24.94 \pm 0.32	24.92 \pm 0.31

TABLE VI
ACCURACY RESULTS (%) BETWEEN DISTINCT SEARCH BOUNDS (Δ) OVER
IMDB (TOP) AND SST (BOTTOM) DATASETS.

Model	$\Delta = 0.0001$	$\Delta = 0.001$
α -GA-LSTM	48.98 \pm 1.69	48.98 \pm 1.76
α -PSO-LSTM	48.98 \pm 1.72	48.97 \pm 1.54
β -GA-LSTM	48.98 \pm 1.85	48.98 \pm 1.93
β -PSO-LSTM	48.98 \pm 1.73	48.98 \pm 1.88
γ -GA-LSTM	48.98 \pm 1.77	<u>49.03 \pm 1.88</u>
γ -PSO-LSTM	48.97 \pm 1.69	48.99 \pm 1.84
α -GA-LSTM	54.91 \pm 2.88	54.90 \pm 2.85
α -PSO-LSTM	54.91 \pm 2.88	54.93 \pm 2.86
β -GA-LSTM	54.91 \pm 2.88	54.91 \pm 2.87
β -PSO-LSTM	54.91 \pm 2.88	54.92 \pm 2.89
γ -GA-LSTM	54.91 \pm 2.88	54.90 \pm 2.87
γ -PSO-LSTM	54.91 \pm 2.88	<u>54.95 \pm 2.88</u>

a convolutional neural network,” *Journal of Electrical Engineering & Technology*, vol. 15, no. 2, pp. 721–726, 2020.

- [18] N. M. Nawari, A. Khan, M. Rehman, R. Naseem, and J. Uddin, “Studying the effect of optimizing weights in neural networks with meta-heuristic techniques,” in *Proceedings of the International Conference on Data Engineering 2015 (DaEng-2015)*. Springer, 2019, pp. 323–330.
- [19] D. E. Rumelhart, G. E. Hinton, and R. J. Williams, “Learning representations by back-propagating errors,” *nature*, vol. 323, no. 6088, pp. 533–536, 1986.
- [20] V. Nair and G. E. Hinton, “Rectified linear units improve restricted boltzmann machines,” in *ICML*, 2010.
- [21] K. Simonyan and A. Zisserman, “Very deep convolutional networks for

- large-scale image recognition,” *arXiv preprint arXiv:1409.1556*, 2014.
- [22] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna, “Rethinking the inception architecture for computer vision,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016.
- [23] S. Hochreiter and J. Schmidhuber, “Long short-term memory,” *Neural Comput.*, vol. 9, no. 8, pp. 1735–1780, Nov. 1997. [Online]. Available: <http://dx.doi.org/10.1162/neco.1997.9.8.1735>
- [24] D. P. Bertsekas, *Nonlinear programming*. Athena Scientific, 1999.
- [25] X.-S. Yang, “Review of meta-heuristics and generalised evolutionary walk algorithm,” *International Journal of Bio-Inspired Computation*, vol. 3, no. 2, pp. 77–84, 2011.
- [26] A. Krizhevsky, “Learning multiple layers of features from tiny images,” Citeseer, Tech. Rep., 2009.
- [27] A. Maas, R. E. Daly, P. T. Pham, D. Huang, A. Y. Ng, and C. Potts, “Learning word vectors for sentiment analysis,” in *Proceedings of the 49th annual meeting of the association for computational linguistics: Human language technologies*, 2011, pp. 142–150.
- [28] R. Socher, A. Perelygin, J. Wu, J. Chuang, C. D. Manning, A. Y. Ng, and C. Potts, “Recursive deep models for semantic compositionality over a sentiment treebank,” in *Proceedings of the 2013 conference on empirical methods in natural language processing*, 2013, pp. 1631–1642.
- [29] F. Wilcoxon, “Individual comparisons by ranking methods,” *Biometrics Bulletin*, vol. 1, no. 6, pp. 80–83, 1945.