

SA-CNN: Application to text categorization issues using simulated annealing-based convolutional neural network optimization

Zihao Guo*

École Centrale de Nantes, France
zihao.guo@eleves.ec-nantes.fr

Yueying Cao

École Centrale de Nantes, France
yueying.cao@eleves.ec-nantes.fr

ABSTRACT

Convolutional neural networks (CNNs) are a representative class of deep learning algorithms including convolutional computation that perform translation invariant classification of input data based on their hierarchical architecture. However, classical convolutional neural network learning methods use the steepest descent algorithm for training, and the learning performance is greatly influenced by the initial weight settings of the convolutional and fully connected layers, requiring re-tuning to achieve better performance under different model structures and data. Combining the strengths of the simulated annealing algorithm in global search, we propose applying it to the hyperparameter search process in order to increase the effectiveness of convolutional neural networks (CNNs). In this paper, we introduce SA-CNN neural networks for text classification tasks based on Text-CNN neural networks and implement the simulated annealing algorithm for hyperparameter search. Experiments demonstrate that we can achieve greater classification accuracy than earlier models with manual tuning, and the improvement in time and space for exploration relative to human tuning is substantial.

CCS CONCEPTS

• **Computing methodologies** → Machine learning; Machine learning approaches; Neural networks.

KEYWORDS

Simulated Annealing Algorithm, Text Classification, Deep Learning, Self-optimization

ACM Reference Format:

Zihao Guo and Yueying Cao. 2022. SA-CNN: Application to text categorization issues using simulated annealing-based convolutional neural network optimization. In *2022 6th International Conference on Electronic Information Technology and Computer Engineering (EITCE 2022)*, October 21–23, 2022, Xiamen, China. ACM, New York, NY, USA, 8 pages. <https://doi.org/10.1145/3573428.3573788>

*Corresponding author

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

EITCE 2022, October 21–23, 2022, Xiamen, China

© 2022 Association for Computing Machinery.

ACM ISBN 978-1-4503-9714-8/22/10...\$15.00

<https://doi.org/10.1145/3573428.3573788>

1 INTRODUCTION

In recent years, significant breakthroughs have been achieved in the field of convolutional neural networks for text classification, and Yoon Kim [1] proposed a straightforward single-layer CNN architecture that can outperform traditional algorithms in a variety of uses. Rie Johnson and Tong Zhang [2] applied CNNs to high-dimensional text data and learned to embed small text regions for classification. Tong He and Weilin Huang [3] proposed a convolutional neural network that extracts regions and features related to text from image components. This type of model uses vectors to characterize each sentence in the text, which are then merged into a matrix and utilized as input for constructing a CNN network model.

Numerous experiments have shown, however, that the performance of neural networks is highly dependent on their architecture [4] [5] [6]. Due to the discrete nature of these parameters, accurate optimization algorithms cannot be used to resolve the architecture optimization problem [7]. Manually tuning the parameters of a model to optimize its performance for different tasks is not only inefficient, but also likely to miss the optimal parameters, resulting in a network architecture that does not achieve maximum performance, which is not advantageous in comparison to traditional classification algorithms [8]. In addition, the widely utilized grid search is an enumerative search, i.e., it tries every possibility by performing a cyclic traversal of all candidate parameter choices, which is marked by its high time consumption and limited globalization. Therefore, it is practical to use an algorithm to automatically and fairly rapidly determine the optimal architecture of a neural network.

It has been shown that tuning neural network hyperparameters with metaheuristic algorithms not only simplifies the network [9] [10], but also enhances its classification performance. In this paper, we use a simulated annealing algorithm to optimize the neural network architecture, and we model the neural network hyperparameter optimization problem as a dual-criteria optimization problem of classification accuracy and computational complexity. The resulting network achieves improved classification performance in the text classification task.

2 BACKGROUND AND RELATED WORK

2.1 The current utilization text classification in neural networks

Text classification was initially done by using knowledge engineering to build an expert system and perform classification, which is a laborious task with limited accuracy. After that, along with the development of statistical learning methods and machine learning

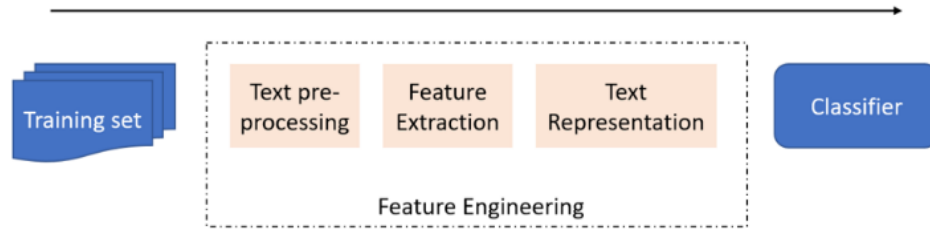


Figure 1: A diagrammatic representation of the process of shallow learning

disciplines, the classical approach of feature engineering plus shallow classification models developed gradually (Fig. 1). During this period, rule-based models: decision trees [11], probability-based models: Naive Bayes classification algorithms [12] [13], geometry-based models: SVM [14] and statistical models: KNN [15], etc. However, these models typically rely heavily on time-consuming feature engineering or large quantities of additional linguistic resources, and are ineffective at learning semantic information about words.

In recent years, research on deep learning that incorporates feature engineering into the process of model fitting has significantly enhanced the performance of text classification tasks. Kim [1] explored the use of convolutional neural networks with multiple windows for text classification, a method that has been widely adopted in industry due to its high computational speed and parallelizability. Yang et al [16] proposed HAN, a hierarchical attention mechanism network that mitigates the gradient disappearance problem caused by RNNs in neural networks. Johnson and Zhang [17] proposed a word-level deep CNN model that improves network performance by increasing network depth without significantly increasing computational burden.

Additionally, to convolutional neural networks and recurrent neural networks, numerous researchers have proposed more intricate models in recent years. The capsule networks-based text classification model proposed by Zhao et al [18], outperformed conventional neural networks. Google proposed the BERT model [19], which overcomes the problem that static word vectors cannot solve the problem of a word having multiple meanings. However, the parameters of each of the aforementioned deep learning models have a substantial effect on network performance and must be optimized for optimal network performance.

2.2 Current research status on the simulated annealing method and associated neural network optimization

The Simulated Annealing Technique (Fig. 2) is a stochastic optimization algorithm based on the Monte Carlo iterative solution approach that was first designed for combinatorial optimization and then adapted for general optimization. Its fundamental concept is based on the significance sampling approach published by Metropolis in 1953, but Kirkpatrick et al. [20] did not properly implement it into the field of combinatorial optimization until 1983.

The algorithm for simulated annealing is taken from the process of metal annealing [21], which may be summarized roughly

as follows. The simulated annealing technique begins with a high initial temperature, then as the temperature parameter decreases, it searches for the optimal solution among all conceivable possibilities. The SA method has a probability of accepting a solution that is worse than the initial one, i.e., the locally optimal solution can probabilistically jump out and finally converge to the global optimum. This likelihood of accepting a suboptimal answer decrease until the SA approaches the global optimal solution.

The standard optimization process of the simulated annealing algorithm can be described as follows.

Algorithm 1 Simulate Anneal Arithmetic

```

Input: Initial feasible solution:  $x_0$  ;
Initial temperature:  $T_0$  ;
Termination temperature :  $T_f$  ;
Iteration number:  $k$ 
 $x \leftarrow x_k, T \leftarrow T_0, k \leftarrow 0$ 
while  $k \leq k_{max}$  and  $T \geq T_f$  do
 $x_k \leftarrow NEIGHBOR(s)$ 
 $\delta f = f(x_k) - f(x)$ 
if  $\Delta f < 0$  or  $RANDOM(0, 1) \leq \Delta P(\Delta f, T)$  then
 $x \leftarrow x_k$ 
end if
 $T \leftarrow COOLING(T, k, k_{max})$ 
 $k \leftarrow k + 1$ 
end while

```

SA has been utilized extensively in VLSI [22], production scheduling [23], machine learning [24], signal processing [25], and other domains as a general-purpose stochastic search method. Boltzmann machine [26], the first neural network capable of learning internal expressions and solving complicated combinatorial optimization problems, utilizes the SA principle for optimization precisely, therefore the optimization potential of SA is evident.

RasdiRere et al. [27] employed simulated annealing to automatically construct neural networks and alter hyperparameters, and experimental findings revealed that the method could increase the performance of the original CNN, demonstrating the efficacy of this optimization technique. Mousavi et al. [28] updated a solar radiation prediction model by integrating artificial neural networks and simulated annealing with temperature cycling to increase ANN calibration performance. Choubin et al. [29] utilized the simulated

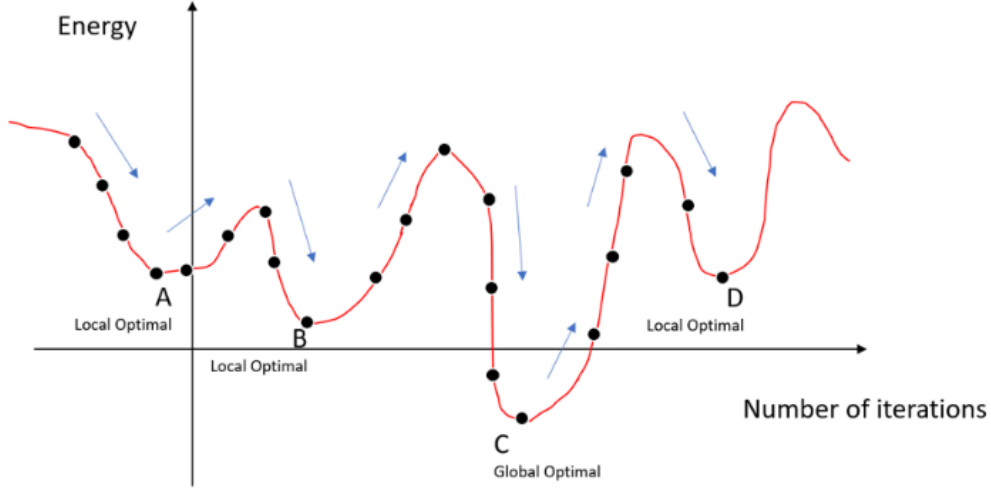


Figure 2: Simulated annealing method picture schematic.

annealing (SA) feature selection approach to find the influential factors of PM modeling based on current air detection machine learning models for the spatial risk assessment of PM10 in Barcelona.

According to the study, however, there is no research on the combination of simulated annealing and neural networks for text categorization tasks. This research conducts tests on neural networks employing simulated annealing to enable automated hyperparameter search, based on the fact that neural networks now generate superior outcomes in text categorization.

3 METHODS

3.1 Convolutional neural networks for text processing tasks (Text-CNN)

Convolutional neural networks (CNN) (Fig. 3) originated in the field of computer vision; however, with the recent deformation of the CNN input layer, this neural network structure has been steadily transferred to the field of natural language processing, where it is often referred to as Text CNN. The schematic is shown below.

A text statement consists of n words, therefore the text is separated according to words and each word is mapped to the Word Embedding layer as a k -dimensional vector. At this point, for this input model, the text may be regarded as a $n \times k$ single-channel picture. During the processing of the convolution layer, the convolution is used to extract the relationships between tuples containing different numbers of words, i.e., to generate new features and obtain different feature maps, when the width of the convolution kernel is equal to the dimension k of the word vector.

The feature map in the form of an n -dimensional vector is then sampled using max-pooling to determine the maximum value, and the data is pooled and utilized as input to the fully connected layer. The SoftMax function [Eq. 1] is then employed to convert these probabilities into discrete 0 or 1 class labels in order to solve this

classification challenge.

$$y = \text{softmax}(Wh + b1) \quad (1)$$

As demonstrated in Eq. 2), a cross-entropy loss function is frequently utilized for the classification job in model training.

$$\text{Loss} = - \sum_{i=1}^n y_i \times \log(y'_i) \quad (2)$$

where y_i is the label value corresponding to the real probability distribution of the i^{th} sample, y' is the prevalence measure corresponding to the projected probability distribution of the i^{th} sample, and n is the number of samples in the training set.

The hyperparameter optimization process of the neural network by simulated annealing method is shown in (Fig. 4).

3.2 Hyperparametric optimization method based on multi-objective simulated annealing method

In this study, we use the MOSA algorithm proposed by Gülüc and Kus [30] to optimize the hyperparameters of the convolutional neural network in order to find the most suitable parameters quickly and achieve a higher accuracy rate. Similar to the single-objective SA algorithm, we extend the SA algorithm, which only considers the error rate of the network implementation, to consider the two objectives of the number of FLOPs required by the network and the error rate of the network, respectively, and define the stopping criterion of the simulated annealing method as the number of iterations.

3.2.1 Main flow of MOSA algorithm. The MOSA algorithm primarily uses Smith's Pareto dominance rule [31] due to complications such as the need for two target values to be on the same scale, followed by the application of decision rules to aggregate the probabilities, and the need to maintain different temperatures due to the

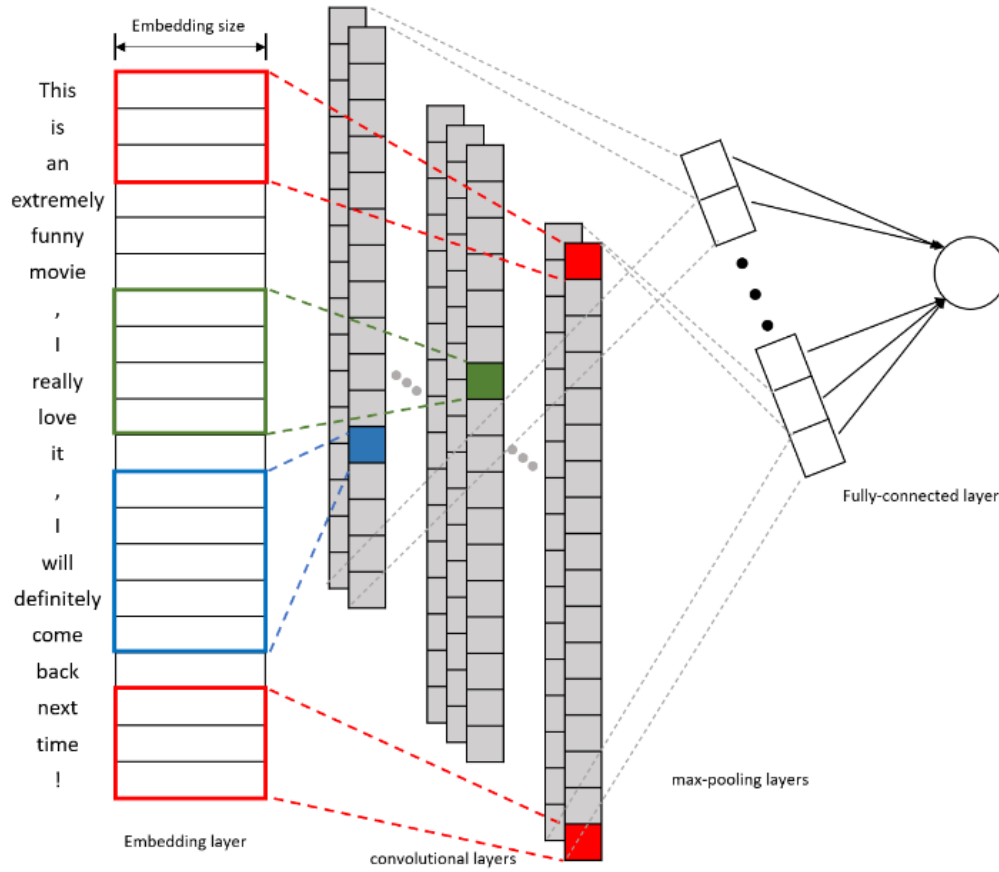


Figure 3: Schematic diagram of Text-CNN.

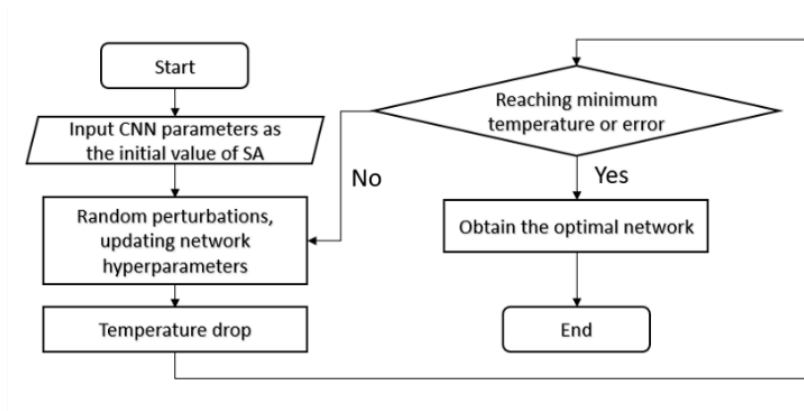


Figure 4: Flowchart of neural network hyperparameter tuning using simulated annealing.

different probabilities evaluated for each target. All non-dominated solutions encountered during the search are stored in an external archive, A , when the first iteration begins at the initial temperature. As new solutions are accepted, A is updated (by inserting the new solution X' and removing all solutions dominated by it) and a superior solution is eventually formed as the Pareto frontier. As depicted

in the following flowchart, whenever a new solution X' is obtained, X and A are updated based on the dominance relationship between the current solution X , the new solution X' , and the solution in the external archive A . This process of re-visiting previously visited archive solutions is known as the return-to-base strategy. In contrast to the single-target SA algorithm, the ΔF calculation used to

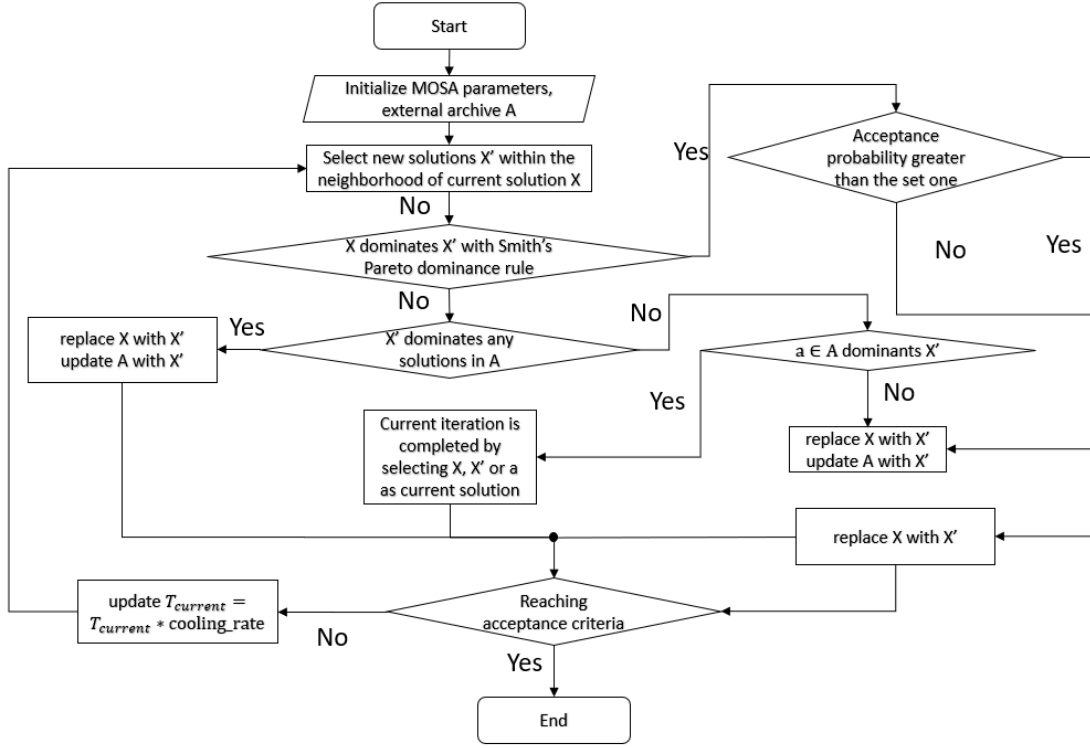


Figure 5: Schematic diagram of MOSA algorithm flow.

determine the probability of acceptance is different in this method. For calculation purposes, a single temperature will be maintained regardless of the number of targets.

3.2.2 Setting and calculation of temperature and probability for SA method. First, the initial and final temperatures. Theoretically, the higher the initial temperature setting, the better, but since the time required for convergence increases correspondingly with the temperature, it is necessary to make a compromise between convergence time and convergence accuracy. To define T_{init} , we apply a real-time initial temperature selection strategy. In this strategy, rather than using Eq. 3) to calculate the initial probability value for a given initial temperature (where ΔF is the amount of deterioration in the objective function and T_{cur} is the current temperature).

$$P_{acc} = \min\{1, \exp(\Delta F / T_{cur})\} \quad (3)$$

We use Eq. 4) to calculate the initial temperature with an initial probability value (where F_{ave} is the average amount of deterioration penalty calculated during short-term “aging”). F_{final} is also defined by this method of real-time temperature adjustment)

$$T_{init} = \Delta F_{ave} / \ln(P_{acc}) \quad (4)$$

The application flow chart is shown in (Fig. 5).

3.2.3 Acceptance criteria in searching parameters. In this experiment, the number of iterations is used to define the stopping criterion for the simulated annealing method. The rationale is as follows: if poor early performance is observed on the validation set, the current training process is terminated and the next search round is

initiated. This approach has the benefit of introducing noise while decreasing the total running time of the HPO method. In each iteration of the simulated annealing method, the newly generated configuration is trained on the training set of the original training set and evaluated on the validation set (i.e., the test split of the original training set). We apply the Xavier weight initial value setting term, a learning rate of 0.001, the Rmsprop optimizer, and a batch size of 32.

3.2.4 Optimization of the Simulated Annealing method. Starting with the initial solution I and initial temperature T_{init} , the iterative process “generate a new solution to calculate the objective function difference accept or reject” is repeated for the current solution and temperature. T_{cur} is gradually decayed, and if the optimal error rate achieved on the validation set does not improve after three consecutive calendar hours, the training procedure is terminated, indicating that the current solution is the approximate optimal solution, i.e., the optimal state of the network has been reached.

4 EXPERIMENTS AND RESULTS

4.1 Introduction to the experimental data set

This experiment utilized two short text categorization datasets, the MR dataset and the TREC dataset. MR is a dataset for the sentiment analysis of movie reviews, with each sample classified as positive sentiment or negative sentiment. The CR dataset consists of reviews of five electronic products, and these sentences have been manually tagged with the sentiment of the reviews. TREC is a question

Table 1: Statistics for the text classification dataset

Dataset	C	AvgLen	DSize	V	Vpre	Test
MR	2	20	10662	18765	16448	CV
CR	2	19	3775	5340	5046	CV
TREC	6	10	5952	9592	9125	50

* C: number of target categories, AvgLen: average sentence length, DSize: dataset size, |V|: number of words, |Vpre|: number in pretrained word vector, Test: test set size (CV means there was no standard train/test split and thus 10-fold CV was used)

classification dataset in which each data point represents a question description, and the job entails categorizing a question into six question kinds, including person, place, numerical information, abbreviation, description, and entity.

The table 1 displays the statistical information of the three data sets.

Several samples from each of the two datasets are provided below.

- It's a square, sentimental drama that satisfies, as comfort food often can. [MR Dataset, tags: positive].
- The sort of movie that gives tastelessness a bad rap. [MR Dataset, tags: negative].
- This camera is so easy to use! [CR Dataset, tags: positive].
- The sound level is also not as high as I would have expected. [CR Dataset, tags: negative].
- What is Australia's national flower? [TREC Dataset, tags: place]
- Who was the first man to fly across the Pacific Ocean? [TREC Dataset, tags: person]

4.2 Introduction to the comparison model

Comparing the model in the article to the experimental model in Kim's [1] study for experimentation.

- CNN-rand: All word vectors are initialized at random before being utilized as optimization parameters during training.
- CNN-static: All word vectors are directly acquired with the Word2Vec tool and are fixed.
- CNN-multichannel: A mix of CNN-static and CNN-non-static, i.e., two types of inputs.
- DCNN [32]: Dynamic Convolutional Neural Network with k-max pooling.
- MV-RNN [33]: Matrix-Vector Recursive Neural Network with parse trees.

4.3 SA-CNN parameter setting

4.3.1 Parameter setting of MOSA. In this research, we employed the identical parameter settings for the simulated annealing approach as Gülcü and Kuş [30], set the starting initial probability value to 0.5, and derived $T_{init} \approx 0.577, T_{final} = 0.12$ in a similar fashion. The link between the number of exterior and internal iterations estimated for different cooling rate values is depicted in the table 2.

The number of outer iterations defines the number of search networks, whereas the number of inner iterations determines the number of single network training iterations. As the cooling rate,

we chose 0.95, where the number of external cycles is greater than the number of internal cycles, to ensure that as many network structures as possible are searched for and to avoid repeated training of a single network structure as much as possible, thereby avoiding becoming trapped in a local optimal solution of network selection.

4.3.2 Search range of neural network hyperparameters. In this study, we utilize a 300-dimensional word2vec trained by Mikolov [34] as the initial representation and continually update the word vector during the training process, similar to Kim's [1] approach. In this paper, the empirical range of hyperparameters to be tuned in the network is provided so that the simulated annealing technique can search for a new solution from the existing one. Expanding the range of searchable hyperparameters may result in improved experimental outcomes, provided computational resources permit.

1. Conv:

– kernelCount: [32, 64, 96, 100, 128, 160, 256]

– dropoutRate: [0.1, 0.2, 0.3, 0.4, 0.5]

2. fullyConnected:

– unitCount: [16, 32, 64, 128, 256, 512]

– dropoutRate: [0.1, 0.2, 0.3, 0.4, 0.5]

3. learningProcess:

– activation: [relu, leaky relu, elu, tanh, linear]

– learningRate: [0.0001, 0.001, 0.01, 0.0002, 0.0005, 0.0008, 0.002, 0.004, 0.005, 0.008]

– batchSize: [64, 128, 256]

4. seedNumber: 40

5. ratioInit: 0.9

4.4 Experimental results and discussion

4.4.1 Comparison of model accuracy results. The following table shows the accuracy of different CNN models for text classification tasks on MR, CR and TREC datasets.

As shown in the table 3, on the MR and CR datasets, the model presented in this paper outperformed other neural network structures, leading the authors to conclude that, due to a lack of computational resources, the parameter search range was restricted and the optimal network was not found. Nevertheless, the performance of the convolutional neural network was utilized, and on the TREC dataset, the model SA-CNN achieved the highest accuracy rate. Under the assumption of using the same model structure, the experimental results demonstrate that using the simulated annealing algorithm to find the optimal hyperparameters not only reduces the tediousness of manual parameter tuning, but also yields better parameters than manual tuning if the search range is correct, thereby achieving a high test accuracy.

Table 2: Relationship Between The Number Of External And Internal Iterations Calculated For Different Values Of Cooling Rate

Iteration budget	Tinit	Tfinal	CR*	OI*	II*
250	0.577	0.12	0.99	156.2	1.6
250	0.577	0.12	0.95	30.6	8.1
250	0.577	0.12	0.9	14.9	26.7
250	0.577	0.12	0.85	9.6	25.8
250	0.577	0.12	0.8	7.0	35.5

*CR: Cooling rate, OI: Outer iterations, II: Inner iterations

*The table is cited from Gülüc and Kus’s experiment [30] on the effect of MOSA cooling rate on the number of iterations

Table 3: Accuracy of different models on the dataset

Model	MR	CR	TREC
CNN-rand	76.1	79.8	91.2
CNN-static	81.0	84.7	92.8
CNN-non-static	81.5	84.3	93.6
CNN-multichannel	81.1	85.0	92.2
D-CNN	-	-	93.0
MV-RNN	79.0	-	-
SA-CNN(This article)	80.7	83.2	93.8

Table 4: TOP3 hyperparameters on TREC dataset.

Hyperparameters	Top1	Top2	Top3
filter num of win 3	100	100	32
filter num of win 4	64	64	50
filter num of win 5	32	64	50
activation function	Tanh	Relu	Relu.85
Learning Rate	0.002	0.002	0.002
Dropout Rate	0.4	0.1	0.4
Batch Size	64	64	64

1) Discussion of experimental results: In order to comprehend the characteristics of the ideal hyperparameters discovered by the simulated annealing technique, the table 4 list the top 3 optimal hyperparameters sought by the algorithm in the TREC dataset.

Compared to manual tuning, the simulated annealing algorithm may search for hyperparameter combinations that one would not ordinarily consider; for instance, the number of CNN convolutional kernels for the Top1 model on the TREC dataset is 100, 64, and 32 for three different strings, which is a combination that one would not ordinarily consider. Therefore, by utilizing the simulated annealing process to optimize the neural network’s hyperparameters, it is theoretically possible to acquire hyperparameter combinations that have not been considered or disregarded based on prior experience, and so achieve improved performance.

5 CONCLUSION

In this article, we proposed a machine learning method combining simulated annealing and convolutional neural networks. The main

goal is to adjust the hyperparameters of neural networks using simulated annealing in order to prevent manual tuning parameters into local optima, thus failing to enhance neural network performance. The experimental results demonstrate that the method of implementing simulated annealing to tune the hyperparameters of a neural network is effective in overcoming the constraints of manual parameter tuning, is practical, and can be theoretically applied to additional natural language processing problems. Due to the limited resource space and the different cooling rate for defining the initialization, which may result in different time costs and architecture, the final solution may only be an approximate optimal solution, and this approximation may differ. Consequently, the simulated annealing method may be integrated with other algorithms or the multi-objective simulated annealing algorithm may be further optimized, thereby further enhancing the efficiency of the simulated annealing algorithm on neural network optimization.

REFERENCES

- [1] Yoon Kim. Convolutional neural networks for sentence classification. In Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 1746–1751, Doha, Qatar, October, 2014. Association for Computational Linguistics.
- [2] Rie Johnson and Tong Zhang. Effective use of word order for text categorization with convolutional neural networks. arXiv preprint arXiv:1412.1058, 2014.
- [3] Tong He, Weilin Huang, Yu Qiao, and Jian Yao. Text-attentional convolutional neural network for scene text detection. IEEE transactions on image processing, 25 (6): 2529–2541, 2016.
- [4] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In Proceedings of the IEEE conference on computer vision and pattern recognition, pages 1–9, 2015.
- [5] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In Proceedings of the IEEE conference on computer vision and pattern recognition, pages 770–778, 2016.

- [6] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. *Advances in neural information processing systems*, 25, 2012.
- [7] Francisco Erivaldo Fernandes Junior and Gary G Yen. Particle swarm optimization of deep neural networks architectures for image classification. *Swarm and Evolutionary Computation*, 49:62–74, 2019.
- [8] Deng K J, Deng C M, Li C *et al*. Application of genetic algorithm in text sentiment classification. *Journal of Sichuan University (Natural Science Edition)*, 56(1):45–49, 2019.
- [9] Amr AbdelFatah Ahmed, Saad M Saad Darwish, and Mohamed M ElSherbiny. A novel automatic cnn architecture design approach based on genetic algorithm. In *International Conference on Advanced Intelligent Systems and Informatics*, pages 473–482. Springer, 2019.
- [10] Sehla Loussaief and Afef Abdelkrim. Convolutional neural network hyper-parameters optimization based on genetic algorithms. *International Journal of Advanced Computer Science and Applications*, 9(10), 2018.
- [11] Jiang Su and Harry Zhang. A fast decision tree learning algorithm. In *Aaai*, volume 6, pages 500–505, 2006.
- [12] Andrew McCallum, Kamal Nigam, *et al*. A comparison of event models for naive bayes text classification. In *AAAI-98 workshop on learning for text categorization*, volume 752, pages 41–48. Madison, WI, 1998.
- [13] Jingnian Chen, Houkuan Huang, Shengfeng Tian, and Youli Qu. Feature selection for text classification with naïve bayes. *Expert Systems with Applications*, 36(3):5432–5435, 2009.
- [14] Zi-Qiang Wang, Xia Sun, De-Xian Zhang, and Xin Li. An optimal svm-based text classification algorithm. In *2006 International Conference on Machine Learning and Cybernetics*, pages 1378–1381. IEEE, 2006.
- [15] Zhou Yong, Lishi Youwen, and Xia Shixiong. An improved knn text classification algorithm based on clustering. *Journal of computers*, 4(3):230–237, 2009.
- [16] Zichao Yang, Diyi Yang, Chris Dyer, Xiaodong He, Alex Smola, and Eduard Hovy. Hierarchical attention networks for document classification. In *Proceedings of the 2016 conference of the North American chapter of the association for computational linguistics: human language technologies*, pages 1480–1489, 2016.
- [17] Rie Johnson and Tong Zhang. Deep pyramid convolutional neural networks for text categorization. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 562–570, 2017.
- [18] Wei Zhao, Jianbo Ye, Min Yang, Zeyang Lei, Suofei Zhang, and Zhou Zhao. Investigating capsule networks with dynamic routing for text classification. *arXiv preprint arXiv:1804.00538*, 2018.
- [19] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- [20] Scott Kirkpatrick, C Daniel Gelatt Jr, and Mario P Vecchi. Optimization by simulated annealing. *science*, 220(4598):671–680, 1983.
- [21] Weixun Yong, Jian Zhou, Danial Jahed Armaghani, MM Tahir, Reza Tarinejad, Binh Thai Pham, and Van Van Huynh. A new hybrid simulated annealing-based genetic programming technique to predict the ultimate bearing capacity of piles. *Engineering with Computers*, 37(3):2111–2127, 2021.
- [22] DF Wong, Hon Wai Leong, and HW Liu. *Simulated annealing for VLSI design*, volume 42. Springer Science & Business Media, 2012.
- [23] Peter JM Van Laarhoven, Emile HL Aarts, and Jan Karel Lenstra. Job shop scheduling by simulated annealing. *Operations research*, 40(1):113–125, 1992.
- [24] Christophe Andrieu, Nando De Freitas, Arnaud Doucet, and Michael I Jordan. An introduction to mcmc for machine learning. *Machine learning*, 50(1):5–43, 2003.
- [25] Sheng Chen and Bing Lam Luk. Adaptive simulated annealing for optimization in signal processing applications. *Signal Processing*, 79(1):117–128, 1999.
- [26] Geoffrey E Hinton. Boltzmann machine. *Scholarpedia*, 2(5): 1668, 2007.
- [27] LM Rasdi Rere, Mohamad Ivan Fanany, and Aniati Murni Arymurthy. Simulated annealing algorithm for deep learning. *Procedia Computer Science*, 72:137–144, 2015.
- [28] Seyyed Mohammad Mousavi, Elham S Mostafavi, and Pengcheng Jiao. Next generation prediction model for daily solar radiation on horizontal surface using a hybrid neural network and simulated annealing method. *Energy conversion and management*, 153:671–682, 2017.
- [29] Bahram Choubin, Mahsa Abdolshahnejad, Ehsan Moradi, Xavier Querol, Amir Mosavi, Shahaboddin Shamshirband, and Pedram Ghamisi. Spatial hazard assessment of the pm10 using machine learning models in Barcelona, Spain. *Science of The Total Environment*, 701:134474, 2020.
- [30] Ayla Gu'lcu" and Zeki Kuş, Multi-objective simulated annealing for hyper-parameter optimization in convolutional neural networks. *PeerJ Computer Science*, 7:e338, 2021.
- [31] Kevin I Smith, Richard M Everson, Jonathan E Fieldsend, Chris Murphy, and Rashmi Misra. Dominance-based multiobjective simulated annealing. *IEEE Transactions on Evolutionary computation*, 12(3): 323–342, 2008.
- [32] Nal Kalchbrenner, Edward Grefenstette, and Phil Blunsom. A convolutional neural network for modeling sentences. *arXiv preprint arXiv:1404.2188*, 2014.
- [33] Richard Socher, Brody Huval, Christopher D Manning, and Andrew Y Ng. Semantic compositionality through recursive matrix-vector spaces. In *Proceedings of the 2012 joint conference on empirical methods in natural language processing and computational natural language learning*, pages 1201–1211, 2012.
- [34] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. Distributed representations of words and phrases and their compositionality. *Advances in neural information processing systems*, 26, 2013.