# Deep Neural Network Hyperparameter Optimization with Orthogonal Array Tuning

Xiang Zhang, Xiaocong Chen, Lina Yao, Chang Ge, Manqing Dong

University of New South Wales, Australia xiang.zhang3@student.unsw.edu.au

Abstract. Deep learning algorithms have achieved excellent performance lately in a wide range of fields (e.g., computer version). However, a severe challenge faced by deep learning is the high dependency on hyperparameters. The algorithm results may fluctuate dramatically under the different configuration of hyper-parameters. Addressing the above issue, this paper presents an efficient Orthogonal Array Tuning Method (OATM) for deep learning hyper-parameter tuning. We describe the OATM approach in five detailed steps and elaborate on it using two widely used deep neural network structures (Recurrent Neural Networks and Convolutional Neural Networks). The proposed method is compared to the state-of-the-art hyper-parameter tuning methods including manually (e.g., grid search and random search) and automatically (e.g., Bayesian Optimization) ones. The experiment results state that OATM can significantly save the tuning time compared to the state-of-the-art methods while preserving the satisfying performance.

Keywords: orthogonal array, hyper-parameter, deep learning

# 1 Introduction

Deep learning has been recently attracting much attention in both academia and industry, due to its excellent performance on various research areas such as computer vision, speech recognition, natural language processing, and brain-computer interface [15]. Nevertheless, deep learning faces an important challenge that the performance of the algorithm highly depends on the selection of hyper-parameters. Compared with traditional machine learning algorithms, deep learning requires hyper-parameter tuning more urgently because deep neural networks: 1) have more hyper-parameters to be tunned; 2) have higher dependency on the configuration of hyper-parameters. [14] reports the deep learning classification accuracy dramatically fluctuates from 32.2% to 92.6% due to the different selection of hyper-parameters. Therefore, an effective and efficient hyper-parameter tuning method is necessary.

However, most of the existing hyper-parameter tuning methods have some drawbacks. In particular, grid search traverses all the possible combinations of different hyper-parameters, which is a time-consuming and ad-hoc process [2]. Random Search, which is developed based on grid research, set up a grid of

hyper-parameter values and selects random combinations to train the algorithm [2]. Random search method oversteps some disadvantages of grid search such as time-consuming but meanwhile brings a major disadvantage which cannot converge to the global optimum [1]. The randomly selected hyper-parameter combinations cannot guarantee a steady and competitive result. Apart from the manually tuning methods, automated tuning methods being more popular in recent years [10]. Bayesian Optimization, a most widely-used automated hyper-parameter tunning approach, attempts to find the global optimum in a minimum number of steps. Nevertheless, the results of Bayesian optimization are sensitive to parameters of the surrogate model and the performance is highly depending on the quality of the learning model [3].

To address the aforementioned issue, we propose the Orthogonal Array Tuning Method (OATM) which can achieve a trade-off of the less tuning time and competitive performance. In detail, the OATM manner is proposed based on Taguchi Approach [12]. The OATM is a highly fractional orthogonal design method that is based on a design matrix and allows the user to consider a selected subset of combinations of multiple factors at multiple levels. Additionally, the OATM is balanced to ensure that all possible values of all hyper-parameters are considered equally. Moreover, OATM has been commonly used as an experimental design method in a wide variety of domains like mechanical engineering [9] and electrical engineering [8]. To our best knowledge, our work is the first batch of work adopting orthogonal array into parameter tuning in deep learning.

The proposed OATM adopts the orthogonal array to extract the most representative and balanced combinations from the whole set of possible combinations. The proposed OATM will be explained in detail in the context of two popular deep learning structures (Section 5). In addition, the OATM is evaluated over three datasets, which demonstrate the universality and adaptability. We notice that source codes performing grid search, random search, and especially Bayesian Optimization on deep learning are hard to online acquire. Thus, we provide the reusable source codes and datasets for reproduction<sup>1</sup>.

#### 2 Related Work

Currently, there are several widely used tuning methods such as grid search optimization, random search optimization, and Bayesian optimization. Grid search and random search require all possible values for each parameter whereas Bayesian optimization needs the range for each parameter. [4] proposed automated machine learning method based on the efficiency of Bayesian optimization and [11] applied multi-task Gaussian processes to Bayesian optimization to enhance the performance of Bayesian Optimization. However, these methods fail in deep learning architectures for which have larger amount of hyper-parameters and the performance highly rely on the configuration.

Apart from the aforementioned methods, the orthogonal array based hyperparameter tuning already used in a range of research areas such as mechanical

<sup>&</sup>lt;sup>1</sup> The link will be available after the paper is accepted

engineering and electrical engineering. J.A Ghani et al. [9] applied orthogonal array based approach to optimize the cutting parameters in the end milling. S.S. Mahapatra et al. [8] optimized wire electrical discharge machining (WEDM) process parameters by orthogonal array method.

**Summary.** The traditional methods are not suited for deep learning algorithms while the effectiveness of OATM has been demonstrated in many research topics. Intuitively, we adopt OATM for deep learning hyper-parameter tuning. To our best knowledge, our work is the first batch of studies in this area.

# 3 Orthogonal Array Tuning

In this section, we first provide the background knowledge of orthogonal array, namely, the definition, the compose principles, and the terminology. Then, we report the working procedure of OATM.

## 3.1 Background of Orthogonal Array

An Orthogonal Array is a table/array whose entries come from a fixed finite set of elements (typically, 1, 2, ..., n), arranged in a specific way that for every selection of two different columns of the table, all ordered 2-tuples of the elements appear for the same number of times. For example, Table 1 shows an Orthogonal Array whose entries come from a fixed finite set 1, 2, 3. In the Orthogonal Array, the column is called *factor* and each element in the finite set (or the element in each column) is called *level*.

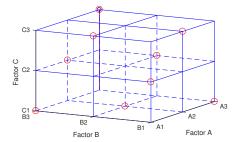
The Orthogonal Array holds two basic composition principles:

- First, in the same column (factor), different levels have the same appearing times. For example, in the first column of Table 1, each level (level 1, level 2, and level 3) appears for 3 times. Similarly, in the second and third columns, each level appears for 3 times.
- Second, in two randomly selected columns (factors), different level combinations are complete and balanced. The number of Orthogonal Array rows is determined by this principle. For example, in the first and second columns of Table 1, each column has 3 levels and there are totally 9 different ordered combinations: (1,2), (1,3), (1,3), (2,1), (2,2), (2,3), (3,1), (3,2), and (3,3). All the combinations are complete (every combination appears) and balanced (every combination appears once).

The essence of Orthogonal Array is a representative subset of the exhausting full set of the elements. We denote the exhausting combination of all the factors and all the levels (3 factors and 3 levels for the above example) as S. Apparently,  $Card(S) = 3^3 = 27$ . As shown in Table 1, the Orthogonal Array only has 9 rows. Let's say Card(O) = 9, where O denotes the set of the combinations in the OATM. Easy to know, O is a representative subset of S, or  $O \subseteq S$ . Intuitively, we can draw both sets out in a cube. In Figure 1,  $A_1, A_2, A_3$  represent 3 levels of factor A, while factors B, C are with the same tokens (factors are supposed to

#### 4 Xiang Zhang et.al

be statistically independent with each other). The total 27 nodes on the surface of the cube denote S while the 9 circled nodes represent the 9 combinations in O. It's easy to observe in Figure 1 that the combinations (circled node) sampled by OATM are uniformly distributed: each edge (totally 27 edges) of the cube has one circled node and each face (totally 9 faces) has three circled nodes.



Dow No	Factor No. Factor 1 Factor 2 Factor 3					
NOW INO.	Factor 1	Factor 2	Factor 3			
1	1	1	1			
2	1	2	2			
3	1	3	3			
4	2	1	2			
5	2	2	3			
6	2	3	1			
7	3	1	3			
8	3	2	1			
9	3	3	2			

Fig. 1: Orthogonal Array cube. The red Table 1: Orthogonal Array with circles are selected combinations by Orthogonal Array.

9 rows, 3 factors and each factor has 3 levels

#### 3.2 Orthogonal Array Tuning Method

In this section, we propose the Orthogonal Array Tuning Method inspired by the basic principles of orthogonal array. Although deep learning algorithms can achieve good performance in many research areas, tuning the hyper-parameters (e.g., the number of layers, the number of nodes in each layer and the learning rate) is time-consuming and dependent on user's expertise.

In OATM, the hyper-parameters are regarded as factors and different values of each hyper-parameter are regarded as levels. The procedure is listed as follows.

- Step 1: Build the F-L (factor-level) table. Determine the number of to-betuned factors and the number of levels for each factor. The levels should be determined by experience and literature. We further suppose each factor has the same number of levels $^2$ .
- Step 2: Construct Orthogonal Array Tuning table. The constructed table should obey the basic composition principles. Here<sup>3</sup> shows some commonly used tables. An alternative way is to use the software. The Orthogonal Array Tuning table can be generated by software such as Weibull++4 and SPSS<sup>5</sup>,

 $<sup>^{2}</sup>$  For the sake of simplicity, we consider all the factors with the same number of levels. More advanced knowledge can be found in [12] for more complex situations.

 $<sup>^3</sup>$  https://www.york.ac.uk/depts/maths/tables/taguchi\_table.htm

 $<sup>^4~{\</sup>tt http://www.reliasoft.com/Weibull/index.htm}$ 

<sup>&</sup>lt;sup>5</sup> https://www.ibm.com/analytics/au/en/technology/spss/

more details in this link<sup>6</sup>. The Orthogonal Array Tuning table is marked as  $L_M(h^k)$  which has k factors, h levels, and totally M rows.

- Step 3: Run the experiments with the hyper-parameters determined by the Orthogonal Array Tuning table.
- Step 4: Range analysis. This is the key step of OATM. Based on the experiment results in the previous step, range analysis method is employed to analyze the results and figure out the optimal levels and importance of each factor. The importance of a factor is defined by its influence on the results of the experiments. Note that range analysis optimizes each factor and combines the optimal levels together, which means that the optimized hyper-parameter combination is not restricted to the existing Orthogonal Array table.
- Step 5: Run the experiment with the optimal hyper-parameters setting.

The OATM is enabled to optimize the hyper-parameters by utilizing a very small set of highly representative hyper-parameter combinations. The high efficiency can be demonstrated by a simple sample in Figure 1. The OATM only takes 9 combinations (red cycles) which means the hyper-parameters can be optimized by running the experiment for 9 times. In contrast, the grid search requires trying all the 27 combinations (27 nodes in the cube). Therefore, through OATM, we can save about 67% (0.67 = 1 - 9/27) work in the tuning procedure.

# 4 Experimental Setting

To evaluate the proposed OATM, we design extensive experiments to tune the hyper-parameters of two most widely used deep learning structures, i.e., the Recurrent Neural Networks (RNNs) and the Convolutional Neural Networks (CNNs). Both of the two deep learning structures are employed on three real-world applications: 1) a human intention recognition task based on the Electroencephalography (EEG) signals; 2) activity recognition based on wearable sensors like Inertial Measurement Unit (IMU); 3) activity recognition based on pervasive sensors like Radio Frequency IDentification (RFID).

# 4.1 Data Setting

The proposed OATM is evaluated over three different tasks on three benchmark datasets. Each dataset is divided into a training set (80%) and a testing set (20%).

**EEG-based Intention Recognition.** We select the widely used EEG dataset from PhysioNet eegmmidb database<sup>7</sup> which contains 5 different categories. In this paper, we choose a subset of eegmmidb which contains 28,000 EEG samples. Every sample is a vector with 64 elements corresponding to 64 channels.

<sup>6</sup> https://www.youtube.com/watch?v=C7PIcOX1WQg

<sup>7</sup> https://www.physionet.org/pn4/EEGmmidb/

**IMU-based Activity Recognition.** This dataset is collected by 9 participants [5], which contains 1200000 samples. 8 ADLs are selected as a subset of our paper. The activity is measured by 3 IMUs and each IMU collects sensor signal with 14 dimensions including two 3-axis accelerometers, one 3-axis gyroscopes, one 3-axis magnetometers, and one thermometer.

RFID-based Activity Recognition. We collect the signals from passive RFID tags [13] and have 3100 samples in total. 21 activities, including 18 ADLs (Activity of Daily Living) and 3 abnormal falls, are performed by 6 subjects. Each sample has 12 dimensions corresponding to 12 RFID tags. RSSI measures the power present in a received radio signal, which is a convenient environmental measurement in ubiquitous computing.

## 4.2 Deep Learning Structures

In this section, we briefly describe RNN and CNN structures and then introduce the key hyper-parameters that will be tuned in the experiments.

RNN Structure The RNN [7], one of the most widely-used deep neural networks, is generally employed to explore the feature dependencies over time dimension through an internal state of the network. Unlike feed-forward neural networks, RNNs can use their internal memory to process arbitrary sequences of inputs and exhibit dynamic temporal behavior. Such characteristic ensures RNNs achieve excellent performance in time-series tasks such as speech recognition and natural language processing.

The RNN structure used in this paper is shown in Figure 2. In the hidden layer, to implement the recurrent function, two LSTM (Long Short-Term Memory) layer is concentrated. LSTM is a simple cell structure which can be used to build a recurrent neural network. Different from other fully connected layers, LSTM layer is composed of cells (shown as rectangles) instead of neural nodes (shown as circles).

In this RNN structure, based on the deep learning hyper-parameters tuning experience, the learning rate, the regularization, and the number of nodes in each hidden layer are key factors affecting the algorithm performance. The loss is calculated by cross-entropy function, and the regularization method is  $\ell_2$  norm with the coefficient  $\lambda$ , The loss is finally optimized by the AdamOptimizer algorithm. In summary, we choose four factors as to-be-tuned hyper-parameters: the learning rate lr, the regularization coefficient  $\lambda$ , the number of hidden layers  $n_l$ , and the number of nodes<sup>8</sup> in each hidden layer  $n_n$ .

CNN Structure The CNN is another popular deep neural network, which shows strong ability to capture the latent spatial relevance of the input data

<sup>&</sup>lt;sup>8</sup> Assume all the hidden layers have the same fixed number of nodes.

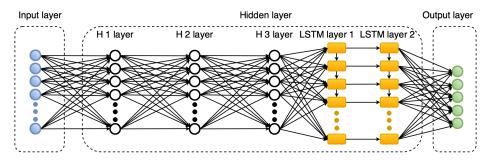


Fig. 2: The schematic diagram of RNN structure. 'H' denotes Hidden, where, for example, the  $H\ 1\ layer$  denotes the first hidden layer.

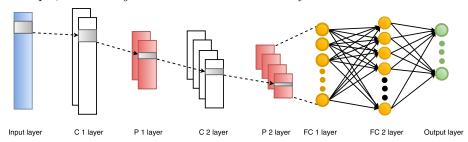


Fig. 3: The schematic diagram of CNN structure. C, P, and FC denote convolutional layer, pooling layer, and fully connected layer, respectively.

and has been demonstrated in a wide range of research topics such as computer vision [6]. The CNN structure contains three categories of components: the convolutional layer, the pooling layer, and the fully connected layer. Each component may appear one or multiple times in a CNN.

As shown in Figure 3, the schematic diagram of CNN is stacked in the following order: the input layer, the first convolutional layer, the first pooling layer, the second convolutional layer, the second pooling layer, the first fully connected layer, the second fully connected layer, and the output layer. The loss function, regularization method, and optimizer are the same as those in the RNN structure. Based on hyper-parameters tuning experience on CNN, we choose four most crucial factors to be tuned by OATM: the learning rate lr', the filter size f', the number of convolutional and pooling layers  $n'_{l}$ , and the number of nodes  $n'_{n}$  in the second fully connected layer.

# 5 Results and Analysis

In this section, we present the hyper-parameter tuning results by OATM and compare to the state-of-the-art methods over a very comprehensive scenario which contains two deep learning structures working on three datasets. For simplification, we set the same hyper-parameter ranges for the three datasets. All

<sup>&</sup>lt;sup>9</sup> We consider each convolutional layer and the following pooling layer as whole.

the codes are open-sourced, please check the code for the training details which are not presented here due to page limitation,

#### 5.1 Overall Comparison

In this section, we compare the proposed OATM with the most competitive state-of-the-art hyper-parameter tuning approaches including two manually methods (grid search and random search) and an automated one (Bayesian Optimization). It's easy to compute that there are  $81=3^4$  exhausted combinations in grid search since we have four factors with three levels of the hyper-parameters. Thus, grid search requires 81 runnings to get the optimal hyper-parameters. On the other hand, our method requires only 9 runnings described in the corresponding orthogonal array table (detailed in Section 5.2). Due to the numbers of runnings in random search and Bayesian Optimization are manually set, they are set as 9 runnings which is same with our method in order to keep fair comparison. The baselines are introduced here:

- Grid search simply goes through all the possible combinations according to the values provided which is exhaustive [2].
- Random search randomly picks combinations from all possible ones. It may not find a decent combination but is widely adopted in industry for the high-efficiency [1].
- Bayesian optimization uses a Gaussian process to minimize the loss function in order to maximize performance [10].

The hyper-parameter levels are selected based on empirical values. For grid search, random search, and our OATM, the empirical values are discrete as listed in Table 3 (take eegmmidb as an example). For Bayesian Optimization, the hyper-parameter ranges from the maximum and minimum of each factor. For instance, the lr ranges from [0.005, 0.015]. All the experiments are implemented in NVIDIA Titian X (Pascal) GPU and each reported value is the average of five runnings under the same setting.

The comparison results are shown in Table 2. It can be observed that:

- under the same running numbers (9 runnings), our method outperforms the random search and Bayesian Optimization over all the datasets and deep learning architectures;
- our method performs slightly lower than grid search but still competitive, however, take EEG dataset with RNN as an example, our approach saves 88% tuning time which is indicated from that the OATM only requires 9 runnings and costs 821.9s while grid search requires 81 runnings and 6853.6s;
- the optimal factors selected by our method approximate to the global optimal factors selected by grid search.

#### 5.2 Case Study in RNN and CNN

In this section, we take EEG classification as an example to present the detailed procedure of OATM in RNN and CNN architecture. The overall paradigm can be divided into five steps.

Table 2: Comparison with the state-of-the-art methods over three datasets and two deep learning architectures. The F1  $\sim$  F4 represent four tuning factors. Acc, Prec and F-1 denote accuracy, precision and F-1 score, respectively.

Data Model	Models	s Methods	Optimal Factors			ors	Metrics					
	wiodels		F1	$\mathbf{F2}$	F3	F4	#-Runnings	Time (s)	Acc	Prec	Recall	F-1
		Grid	0.005	0.004	6	64	81	6853.6	0.9251	0.9324	0.9139	0.9231
	RNN	Random	0.01	0.008	6	32	9	766.8	0.7941	0.8003	0.7941	0.7947
		во	0.0135	0.0049	5	32	9	703.4	0.718	0.7246	0.6474	0.6838
EEG		Ours	0.005	0.004	6	64	9	821.9	0.925	0.9335	0.9223	0.9279
EEG		Grid	0.005	4	3	192	81	31891.5	0.828	0.8137	0.8256	0.8269
	CNN	Random	0.003	2	1	128	9	662.8	0.7268	0.7277	0.7269	0.7266
	CIVIN	во	0.001	4	3	139	9	721.9	0.7244	0.7302	0.7244	0.7263
		Ours	0.003	4	1	128	9	680.4	0.797	0.7969	0.8112	0.8003
		Grid	0.005	0.004	6	96	81	3027.2	0.9936	0.9909	0.9976	0.9971
	RNN	Random	0.015	0.004	4	32	9	1008.5	0.9139	0.9209	0.9412	0.9156
	IUININ	во	0.0132	0.0041	4	48	9	1078.8	0.9872	0.9877	0.9851	0.9863
$\mathbf{IMU}$		Ours	0.005	0.004	6	64	9	1138.2	0.9913	0.9924	0.9905	0.9919
INIO	CNN	Grid	0.003	2	1	128	81	41804.9	0.9732	0.9708	0.9708	0.9707
		Random	0.003	2	2	128	9	7089.2	0.9692	0.9691	0.9692	0.9691
	CIVIV	во	0.0012	2	2	192	9	6559.7	0.9696	0.9702	0.9701	0.9701
		Ours	0.003	2	2	128	9	6809.8	0.9702	0.9699	0.9703	0.9702
	RNN	Grid	0.005	0.008	6	96	81	2846.1	0.9342	0.9388	0.9201	0.9252
RFID—		Random	0.005	0.012	4	32	9	642.3	0.8891	0.9138	0.8826	0.8895
	IUININ	во	0.0142	0.0093	6	79	9	452.2	0.9071	0.8556	0.8486	0.8436
		Ours	0.005	0.008	6	64	9	497.1	0.9134	0.9138	0.9029	0.9162
		Grid	0.005	4	2	192		7890.8	0.9316	0.9513	0.9316	0.9375
	CNN	Random	0.005	2	1	128	9	1210.3	0.8683	0.9113	0.8684	0.8779
		во	0.005	5	3	64	9	872.9	0.9168	0.9058	0.9194	0.9086
		Ours	0.005	4	3	192	9	980.3	0.9235	0.9316	0.9188	0.9326

Step 1: Build the F-L table According to the description in Section 4.2, OATM will work on four different hyper-parameters (factors): the learning rate lr, the l-2 norm coefficient  $\lambda$ , the number of hidden layers  $n_l$ , and the number of nodes  $n_n$ . The number of levels h is set to be 3 which could be much larger in real-world applications. Based on the related work and tuning experience [14], the empirical values are shown in Table 3.

**Step 2: OATM table** Then, choosing a suitable Orthogonal Array table which contains 4 factors and 3 levels for our experiments in this link<sup>10</sup> wich contains 9 combinations. The OATM table should satisfy two basic principles: i) in each column, different levels have the same appear times; ii) in any two randomly-selected columns, nine differently-ordered element combinations are completed and balanced.

**Step 3: Run the experiments** Following the OATM table, run the 9 experiments and record the classification accuracy. In our case, each experiment runs 5 times with the corresponding average accuracy recorded. Each experiment is trained for 1,000 iterations to guarantee the convergence.

**Step 4: Range analysis** This is the key step of Orthogonal Array Tuning. The overall range analysis procedure and results are shown in Table 4. The first 9 rows are measured and recorded in Step 3.  $R_{leveli}$  denotes the sum of accuracy

https://www.york.ac.uk/depts/maths/tables/taguchi\_table.htm

Table 3: Factor-Level table of RNN and CNN.

		Factor 1 (lr)	Factor 2 $(\lambda)$	Factor 3 $(n_l)$	Factor 4 $(n_n)$
RNN	Level 1	0.005	0.004	4	32
	Level 2	0.01	0.008	5	64
	Level 3	0.015	0.012	6	96
		Factor 1 $(lr')$	Factor 2 $(f')$	Factor 3 $(n'_l)$	Factor 4 $(n'_n)$
CNN	Level 1	0.001	[1,2]	1	64
	Level 2	0.003	[1,4]	2	128
	Level 3	0.005	[1.6]	3	192

Table 4: Range analysis of RNN

Row No.	Factor 1 (lr)	Factor 2 $(\lambda)$	Factor 3 (n <sub>l</sub> )	Factor 4 $(n_n)$	Acc		
1	0.005	0.004	4	32	0.875		
2	0.005	0.008	5	64	0.8		
3	0.005	0.012	6	96	0.521		
4	0.01	0.004	5	96	0.888		
5	0.01	0.008	6	32	0.797		
6	0.01	0.012	4	64	0.451		
7	0.015	0.004	6	64	0.897		
8	0.015	0.008	4	96	0.335		
9	0.015	0.012	5	32	0.471		
$R_{level1}$	2.196	2.66	1.661	2.143			
$R_{level2}$	2.136	1.932	2.159	2.148			
$R_{level3}$	1.703	1.443	2.215	1.744			
$A_{level1}$	0.732	0.887	0.554	0.714			
$A_{level2}$	0.712	0.644	0.720	0.716			
$A_{level3}$	0.568	0.481	0.738	0.581			
Lowest Acc	0.568	0.481	0.554	0.581			
Highest Acc	0.732	0.887	0.738	0.716			
Range	0.164	0.406	0.184	0.135			
Importance	$lambda > n_l > lr > n_n$						
Best Level	Level 1	Level 1	Level 3	Level 2			
Optimal Value	0.005	0.004	6	64	0.925		

under level i. For example,  $R_{level1}$  in factor 1 is the sum of the accuracy in the first 3 rows (0.196 = 0.875 + 0.8 + 0.521), where factor 1 is on level 1.  $A_{leveli}$  denotes the average accuracy of level i, calculated by  $A_{leveli} = R_{leveli}/h$ . In the above example, we calculate  $A_{level1}$  as 0.732 = 2.196/3. Lowest and highest accuracy values, measuring the maximum and minimum of  $A_{leveli}$  respectively, are used to calculate the range of  $A_{leveli}$ . The importance denotes how important the factor is, which is ranked by the range value. Best level is the selected optimal level based on the  $Highest\ Acc$  while  $Optimal\ Value$  represents the corresponding value of the best level.

Step 5: Run the optimal setting Since the best level is given by the range analysis in the previous step, we run the experiment with the optimal hyperparameters (lr = 0.004,  $\lambda = 0.005$ ,  $n_l = 6$ , and  $n_n = 64$ ) and finally we got the optimal accuracy as 0.925. We can observe that:

Row No.	Factor 1 $(lr')$	Factor 2 $(f')$	Factor 3 $(n'_l)$	Factor 4 $(n'_n)$	) Acc
1	0.001	[1,2]	1	64	0.707
<b>2</b>	0.001	[1,4]	2	128	0.771
3	0.001	[1,6]	3	192	0.775
4	0.003	[1,2]	2	192	0.779
5	0.003	[1,4]	3	64	0.752
6	0.003	[1,6]	1	128	0.797
7	0.005	[1,2]	3	128	0.784
8	0.005	[1,4]	1	192	0.782
9	0.005	[1,6]	2	64	0.756
$R_{level1}$	2.253	2.27	2.993	2.215	
$R_{level2}$	2.328	2.305	2.306	2.352	
$R_{level3}$	2.322	2.328	2.311	2.336	
$A_{level1}$	0.751	0.757	0.998	0.738	
$A_{level2}$	0.776	0.768	0.769	0.784	
$A_{level3}$	0.774	0.776	0.770	0.779	
Lowest Acc	0.751	0.757	0.769	0.738	
Highest Acc	0.776	0.776	0.998	0.784	
Range	0.025	0.019	0.229	0.046	
Importance		$n_l' > n_n'$	> lr' > f'		
Best Level	Level 2	Level 3	Level 1	Level 2	
Optimal Value	0.003	[1,6]	1	128	0.797

Table 5: Range analysis of CNN

- The optimal accuracy 0.925 is higher than the maximum of the accuracy (0.897) in the OATM experiments, which demonstrates that the OATM is enabled to approximate the global optimal instead of the local optimal.
- The importance of each factor is ranked through the range analysis:  $lambda > n_l > lr > n_n$ , which can guide the researcher to grasp the dominating variable in the RNN structure and be helpful in the future development.

The OATM paradigm of CNN is similar to RNN. Here, we only report the F-L table (Table 3) and the range analysis table (Table 5).

# 6 Discussion and Conclusion

In this paper, we present an efficient OATM approach for hyper-parameter tuning in the context of deep learning. The proposed OATM is evaluated over two popular deep learning structures(RNN and CNN) over three real-world datasets. The experiment results show that our approach outperforms state-of-the-art hyper-parameter tuning methods.

One disadvantage of OATM is that it requires the empirical values as prerequisites. The values of the F-L table should be chosen appropriately. However, this is the common drawback of all the tuning methods. For instance, the hyper-parameter ranges in Bayesian Optimization are also pre-defined based on empirical values.

## References

- Andradóttir, S.: A review of random search methods. In: Handbook of Simulation Optimization, pp. 277–292. Springer (2015)
- Bergstra, J.S., Bardenet, R., Bengio, Y., Kégl, B.: Algorithms for hyper-parameter optimization. In: Advances in Neural Information Processing Systems 24, pp. 2546– 2554 (2011)
- 3. Calandra, R., Gopalan, N., Seyfarth, A., Peters, J., Deisenroth, M.P.: Bayesian gait optimization for bipedal locomotion. In: Learning and Intelligent Optimization. pp. 274–290 (2014)
- 4. Feurer, M., Klein, A., Eggensperger, K., Springenberg, J., Blum, M., Hutter, F.: Efficient and robust automated machine learning. In: Advances in Neural Information Processing Systems 28 (2015)
- 5. Fida, B., Bibbo, D., Bernabucci, I., et al.: Real time event-based segmentation to classify locomotion activities through a single inertial sensor. In: Proceedings of the 5th EAI International Conference on Wireless Mobile Communication and Healthcare. pp. 104–107 (2015)
- Gu, J., Wang, Z., Kuen, J., Ma, L., Shahroudy, A., Shuai, B., Liu, T., Wang, X., Wang, G., Cai, J., et al.: Recent advances in convolutional neural networks. Pattern Recognition 77, 354–377 (2018)
- Li, K., Xu, H., Wang, Y., Povey, D., Khudanpur, S.: Recurrent neural network language model adaptation for conversational speech recognition. INTERSPEECH, Hyderabad pp. 1–5 (2018)
- 8. Mahapatra, S., Patnaik, A.: Optimization of wire electrical discharge machining (wedm) process parameters using taguchi method. The International Journal of Advanced Manufacturing Technology 34(9), 911–925 (2007)
- Nalbant, M., Gökkaya, H., Sur, G.: Application of taguchi method in the optimization of cutting parameters for surface roughness in turning. Materials & design 28(4), 1379–1385 (2007)
- Snoek, J., Larochelle, H., Adams, R.P.: Practical bayesian optimization of machine learning algorithms. In: Advances in Neural Information Processing Systems 25, pp. 2951–2959. Curran Associates, Inc. (2012)
- Swersky, K., Snoek, J., Adams, R.P.: Multi-task bayesian optimization. In: Advances in Neural Information Processing Systems 26, pp. 2004–2012 (2013)
- 12. Taguchi, G., Taguchi, G.: System of experimental design; engineering methods to optimize quality and minimize costs. Tech. rep. (1987)
- 13. Yao, L., Sheng, Q.Z., Li, X., Gu, T., Tan, M., Wang, X., Wang, S., Ruan, W.: Compressive representation for device-free activity recognition with passive rfid signal strength. IEEE Transactions on Mobile Computing 17(2), 293–306 (2017)
- Zhang, X., Yao, L., Huang, C., Sheng, Q.Z., Wang, X.: Intent recognition in smart living through deep recurrent neural networks. In: International Conference on Neural Information Processing (ICONIP). pp. 748–758. Springer (2017)
- 15. Zhang, X., Yao, L., Sheng, Q.Z., Kanhere, S.S., Gu, T., Zhang, D.: Converting your thoughts to texts: Enabling brain typing via deep feature learning of eeg signals (2018)