# Evolving Unsupervised Deep Neural Networks for Learning Meaningful Representations

Yanan Sun, *Member, IEEE*, Gary G. Yen, *Fellow, IEEE*, and Zhang Yi, *Fellow, IEEE*

*Abstract*—**Deep Learning (DL) aims at learning the *meaningful representations*. A meaningful representation gives rise to significant performance improvement of associated Machine Learning (ML) tasks by replacing the raw data as the input. However, optimal architecture design and model parameter estimation in DL algorithms are widely considered to be intractable. Evolutionary algorithms are much preferable for complex and nonconvex problems due to its inherent characteristics of gradient-free and insensitivity to the local optimal. In this paper, we propose a computationally economical algorithm for evolving *unsupervised deep neural networks* to efficiently learn *meaningful representations*, which is very suitable in the current Big Data era where sufficient labeled data for training is often expensive to acquire. In the proposed algorithm, finding an appropriate architecture and the initialized parameter values for an ML task at hand is modeled by one computational efficient gene encoding approach, which is employed to effectively model the task with a large number of parameters. In addition, a local search strategy is incorporated to facilitate the exploitation search for further improving the performance. Furthermore, a small proportion labeled data is utilized during evolution search to guarantee the learned representations to be meaningful. The performance of the proposed algorithm has been thoroughly investigated over classification tasks. Specifically, error classification rate on MNIST with $1.15\%$ is reached by the proposed algorithm consistently, which is considered a very promising result against state-of-the-art unsupervised DL algorithms.**

*Index Terms*—**Deep learning, neural networks, representation learning, evolutionary algorithm, evolving neural networks.**

## I. INTRODUCTION

DEEP Learning (DL) algorithm, which is materialized by Deep Neural Networks (DNNs) for learning meaningful representations [1], is a very hot research area during recent years [2]–[4]. Meaningful representation refers to the outcome of the raw input data that goes through multiple nonlinear transformations in the DNNs, and the outcome could remarkably enhance the performance of the subsequent machine learning tasks. The hyper-parameter settings and

Yanan Sun is with the College of Computer Science, Sichuan University, Chengdu 610065 CHINA and with the School of Engineering and Computer Science, Victoria University of Wellington, Wellington 6140 NEW ZEALAND. E-mail:yanan.sun@ecs.vuw.ac.nz.

Gary G. Yen is with the School of Electrical and Computer Engineering, Oklahoma State University, Stillwater, OK 74075 USA. E-mail:gyen@okstate.edu. (*Corresponding author*)

Zhang Yi is with the College of Computer Science, Sichuan University, Chengdu 610065 CHINA. E-mail:zhangyi@scu.edu.cn.

parameter values in DNNs are substantially interrelated to the performance of DL algorithms. Specifically, hyper-parameters (such as the size of weights, types of nonlinear activation functions, a priori term types, and coefficient values) refer to the parameters that are needed to be assigned prior to training the models, and parameter values refer to the element values of the weights and are determined during the training phase. Due to the deficiencies of the current optimization techniques for searching for optimal hyper-parameter settings and parameter values, the power of DL algorithms cannot be shown fully. To this end, an effective and efficient approach concerning the hyper-parameter settings and parameter values has been proposed in this paper.

**Meaningful Representations** Typically, arbitrary DNNs can generate/learn Deep Representations (DRs). However, DRs are not necessarily meaningful, i.e., it is not true that all DRs contributed to the promising performance when they replace the raw data to be fed to machine learning algorithms (e.g., classification). In fact, DRs are the outcomes which have gone through nonlinear transformations from input data more than once [5], and are inspired by the mammalian hierarchical visual pathway [6]. Mathematically, the representations of the input data $X \in \mathbb{R}^m$ are formulated by (1)

$$
\begin{cases}
R_1 = & f_1(W_1 X) \\
R_2 = & f_2(W_2 R_1) \\
& \cdots \\
R_n = & f_n(W_n R_{n-1}) \\
R = & R_n
\end{cases}
\tag{1}
$$

where $f_1, \cdots, f_n$ denote a set of element-wise nonlinear activation functions, $W_1, \cdots, W_n$ refer to a series of connection weights and $R_1, R_2, \cdots, R_n$ are the learned representations (output) at the depth/layer $1, 2, \cdots,$ and $n$, among which $\mathbf{R} = \{R_i | 2 \leq i \leq n\}$ refers to the DRs. In addition, Fig. 1 shows the flowchart of deep representation learning and its role in machine learning tasks in a general case.
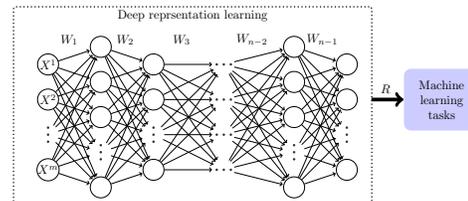


Fig. 1. An example to illustrate a general flowchart of deep representation learning and its relationship to machine learning tasks.

Obviously, multiple different DRs can be learned by varying $n$ in (1), while we only pay attention to the ones that give

the highest performance of the associated machine learning tasks. Based on literature reviews [7]–[9], these DRs are often called *meaningful representations*. Assuming $R_j$ are the meaningful representations, it is obvious that the hyper-parameter settings (e.g., the number of layers, $j$, and the chosen activation function types of $f_1, \cdots, f_j$) and parameter values (e.g., the values of each element in $\{W_1, \cdots, W_j\}$) would highly reflect the learned $R_j$ to be meaningful or not. To this end, the Back-Propagation algorithm (BP) [10] which relies on the gradient information is the widely employed algorithm in training parameter values. However, its performance is highly affected by the initialized setting due to its local search characteristics that could be easily trapped into local minima [11]. Although multiple implementations based on BP, such as Stochastic Gradient Descent (SGD), AdaGrad [12], RMSProp [13], and AdaDelta [14], have been presented to expectedly reduce the adverse impact of easily trapping into local minima, extra hyper-parameters (such as the initialization values of momentums and the balance factors) are introduced and also needed to be carefully tuned in advance. Furthermore, multiple algorithms [15], [16] have been proposed for optimizing the hyper-parameters, but they often require domain knowledge and are problem-dependent. To this end, the grid search method keeping its dominant position in selecting reasonable hyper-parameters was proposed [17]. However, the grid search method is an exhaustive approach, and would frequently miss the best hyper-parameter combinations when the hyper-parameters are continuous numbers.

**Deep Neural Networks** According to literature [18], [19], DL algorithms mainly include Convolutional NNs (CNNs), Deep Belief Networks (DBNs), and stacked Auto-Encoders (AEs). Specifically, CNNs are supervised algorithms for DL, and their numerous variants have been developed for various real-world applications [20]–[25]. Although these CNN algorithms have shown promising performance in some tasks, sufficient labeled training data, which is a must for successfully training them, are not easy to acquire. For example in the ImageNet benchmark [27], there are $10^9$ pictures that can be easily downloaded from the Google and Yahoo websites. It was reported that $48,940$ workers from $167$ countries are employed to label these photos. Therefore, the unsupervised NN approaches whose training processes rely solely on unlabeled data become preferable in this situation. DBNs [28] and stacked AEs [29], [30] are the mainly unsupervised DL algorithms [18], [19] for learning meaningful representations. Because of the unknown in training data targets during their training phase, learned representations from them are not necessarily to be *meaningful*. Therefore, *a priori* knowledge is needed to be incorporated into their training phase. For example, DBNs and stacked AEs trained with the sparsity constraint a priori with benefits of sparse coding [31] have been proposed in [32] and [33]. Furthermore, denoising AEs [34] have been proposed by artificially adding noise priori to input data for improving the ability to learn meaningful representations. In addition, Rifar *et al.* [35] have presented contractive AEs by introducing the term, which is the derivation of representations with respect to input data, for reducing the sensitivity a priori of representations.

**Evolutionary Algorithms for NNs** Evolutionary algorithms (EAs) are one class of population-based meta-heuristic optimization paradigms, and are motivated by the metaphors of biological evolution. During the period of evolution, individuals interact with each other and the beneficial traits are passed down to facilitate population adapting to the environment. Due to the nature of *gradient-free* and *insensitivity to local optima*, EAs are preferred in various problem domains [36]. Therefore, they have been extensively employed in optimizing NNs, which refers to the discipline of neuroevolution, such as for the connection weight optimization [37]–[39], the architecture setting [40]–[42] (more examples can be found in [36]). Generally, these algorithms employ direct or indirect methods to encode the optimized problems for the evolution. To be specific, each parameter in the connection weights is encoded by the binary numbers [37] or a real number [43] in the direct methods, which are effective for the small-scale problems. However, when they are used to encode the problems with a large number of parameters in connection weights, such as for processing the high-dimensional data, these methods become impractical due to the excessive length of the genotype explicitly representing each parameter no matter if coded in binary or real. To this purpose, Stanley and Miikkulainen have proposed the indirect-based Neural Evolution Augmenting Topologies (NEAT) method [44] for encoding connection weights and architectures with varying lengths of chromosomes. Because NEAT employs one unit to denote combinational information of one connection in the evolved NN, it still cannot effectively solve deep NNs where a large number of parameters exist. To this end, an improved version of NEAT (i.e., HyperNEAT) was proposed in [45] in which connection weights were evolved by composing different points in a fixed coordinate system with a series of predefined nonlinear functions. Although the indirect methods can reduce the length of the genotype representation, they limit the generalization of the neural networks and the feasible architecture space [36]. In 2015, Gong *et al.* [46] proposed a bi-objective evolutionary algorithm by using Differential Evolution [47] to concurrently consider the reconstruction error and sparsity of the AE, and chose the optimal sparsity from the knee area of the Pareto front. Recently, Liu *et al.* [48] presented a neural network connection pruning method by a multi-objective evolutionary algorithm to simultaneously consider the representation ability and the sparse measurement. Google [49] proposed their work on evolving CNNs for image classifications with a direct manner over 250 high performance servers for more than 400 hours. In this regard, the evolutionary approaches would surely be capable of evolving deep NNs, although the computational resource is not necessarily available to all interested researchers.

**Contributions** Based on the above investigations upon prospects of unsupervised deep NNs for learning meaningful representations and the EAs in evolving deep NNs, an effective and efficient approach named Evolving Unsupervised Deep Neural Networks (EUDNN) for learning meaningful representation through evolving unsupervised deep NNs, exactly evolving the building blocks of unsupervised deep NNs, has been proposed in this paper. In summary, the contributions of

this paper are documented as follows:

1) A computationally efficient gene encoding scheme of evolutionary approaches has been suggested, which is capable of evolving deep neural networks with a large number of parameters for addressing high-dimensional data with limited computational resources. With this design, the proposed algorithm can be smoothly implemented in academic environments with limited computational resources.

2) A fitness evaluation strategy has been employed to drive the unsupervised models towards usefulness in advance, which can drive the learned representations to be meaningful without any carefully designed a priori knowledge.

3) Deep neural networks with a large number of parameters involve a large-scale global optimization problem. As a result, the sole evolutionary scheme cannot generate the best results. To this end, the utilization of a local search strategy is proposed to be incorporated into the proposed algorithm to guarantee the desired performance.

**Organization** The remaining of this paper is organized as follows. First, related works and motivations of the proposed EUDNN are illustrated in Section II. Next, the details and discussions of the proposed algorithm are presented in Section III. To evaluate the performance, a series of experiments are performed by the proposed algorithm against selected peer competitors and the results measured by the chosen performance metric are analyzed in Section IV. Finally, conclusions and future work are drawn in Section V.

## II. RELATED WORKS AND MOTIVATIONS

We will detail the unsupervised DL models that motive our work in this paper, highlight their deficiencies in learning *meaningful* representations, and rationalize our motivations in Subsection II-A. With this same detailed manner, the evolutionary algorithms which demonstrate the potential for evolving deep NNs will be documented in Subsection II-B.

### A. Unsupervised Deep Learning Models

In this subsection, the unsupervised DL models are reviewed first (Subsection II-A1). Then, their building blocks are introduced (Subsection II-A2). Next, the mechanisms guaranteeing the learned representations to be meaningful are formulated and commented (Subsection II-A3). Finally, the motivations of the proposed algorithm in reducing the adverse impact of their deficiencies are elaborated (Subsection II-A4).

*1)* Unsupervised DL models cover DBNs [28] and variants of stacked AEs (i.e., stacked sparse AEs (SAEs) [32], [33], stacked denoising AEs (DAEs) [34], and stacked contract AEs (CAEs) [35]). Moreover, the building block of DBNs is a Restricted Boltzmann machine (RBM) [50], and that of stacked AEs is an AE. Furthermore, the parameter values in DBNs and stacked AEs are optimized by the greedy layer-wise training method, which is composed of two phases [51]: pre-training and fine-tuning. Conveniently, Fig. 2a depicts the pre-training phase, where a set of three-layer (the input layer, the hidden layer, and the output layer) NNs with varying numbers

of units are individually trained by minimizing reconstruction errors. In the fine-tuning phase which is illustrated by Fig. 2b, these hidden layers are first sequentially stacked together with the parameter values trained in the pre-training phase, then a classification layer (i.e., the classifier) is added to the tail to perform the fine-tuning by optimizing the corresponding loss function determined by the particular task at hand.
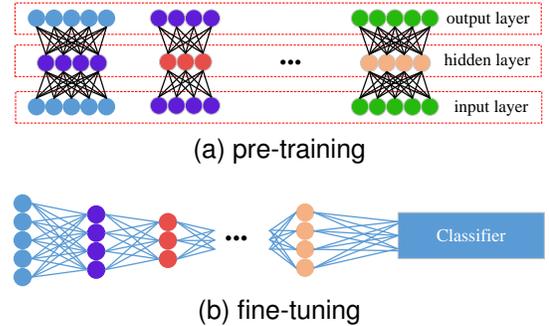


Fig. 2. The training process of unsupervised deep neural networks.

*2)* Unsupervised DL algorithms are considerably preferred mainly due to their requirements upon fewer labeled data especially in the current Big Data era[1]. However, a major issue of training these models is how to guarantee the learned representations to be meaningful. Specifically in the pre-training phase for training one NN unit (see Fig. 3 as an example), let $X \in R^n$ denote the input data, $W \in R^{n \times k}$ denote the connection weight matrix from the input layer to the hidden layer, while $W' \in R^{k \times n}$ denote the connection weight matrix from the hidden layer to the output layer. The NN unit is trying to minimize the reconstruction error $L$ between the input data $X$ and the output data $X'$ by (2)[2]
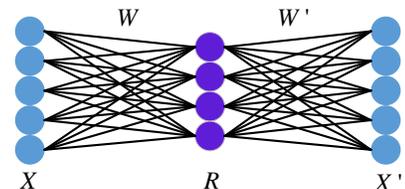


Fig. 3. An example of unsupervised deep neural network unit model.

$$\begin{cases} R & = & f(WX) \\ X' & = & f(W'R) \\ L & = & l(X, X') \end{cases} \quad (2)$$

In (2), $R$ denotes the learned representations (i.e., the output of the hidden layer), $f$ denotes the activation function, and $l$ denotes the function to measure the differences between $X$ and $X'$.

---

[1]Even data is abundant in the Big Data era, most raw data collected is unlabeled for a classification task, e.g., the ImageNet classification benchmark that has been discussed in Section I.

[2]Bias terms, which are another kind of connection weights widely existing in NNs, are incorporated into $W$ and $W'$ here for simplicity.

*3)* It is obvious that the learned representations $R$ are not necessarily meaningful only by minimizing $L$ due to no information of the associated classification task existing in this phase and arbitrary $R$ will lead to a minimal $L$, while $R$ is meaningful only when they could improve the performance of the associated classification task. To this end, literature have presented unsupervised DL algorithms with different a priori knowledge [31], [33]–[35] which is denoted as $\Theta$, and then the reconstruction error is transformed to $L = l(X, X') + \lambda\Theta$ where $\lambda$ denotes a balance factor to determine the weight of the associated a priori term. Although a prior knowledge would help the learned representations to be meaningful, major issues remain:

- The prior knowledge is designed with different assumptions, which do not necessarily satisfy the current situations.
- The prior knowledge is presented specifically for general tasks, while it is hopeful that the performance would be improved on particular tasks.
- It is difficult to choose the most suitable a priori term for the current task.
- The balance factor $\lambda$ is a hyper-parameter whose value is not easily to be assigned [35].

*4)* Considering this problem, the method that has been developed in our previous work [9] is employed in this proposed algorithm. To be specific, a small proportion of labeled data is employed during the fitness evaluation of EAs, and the learned representations are directly quantified based on the classification task that is employed in the fine-tuning phase. With the environmental selection in EAs, individuals that have the positive effect on the classification task survive into the current generation and are expected to generate offspring with better performance in the next generation, which ultimately leads to the learned representations to be meaningful. Because the employed labeled data can be injected from the fine-tuning phase, and the classification task is the same as that in the fine-tuning phase, this strategy for learning meaningful representations would not introduce extra cost.

### B. Evolutionary Algorithms for Evolving Neural Networks

Although multiple related literature for evolving NNs have been mentioned in Section I, only the works in [44], [45] (i.e., the NEAT and the HyperNEAT) will be concerned here because our proposed algorithm aims at evolving *deep* NNs[3]. In the following, the details of NEAT, as well as HyperNEAT and their deficiencies in evolving deep NNs are documented in Subsections II-B1 and II-B2, respectively. Combined with the challenge of EAs in evolving deep NNs, i.e., the upper bound encoding problem, the motivations of the proposed EUDNN are presented in Subsection II-B3. In addition, another challenge, i.e., EAs cannot fully solve the optimization problems with a large number of parameters, and the corresponding motivations are given in Subsection II-B4.

---

[3]The works in [36]–[42] were proposed two decades ago and cannot be applied for deep NNs, the work in [48] concerned only the weight pruning, and the work in [49] employed a direct way for evolving and did not have a general meaning.

*1)* The NEAT [44] has been proposed with an indirect method for adaptively increasing the complexity of the evolved NNs. Specifically, two types of genes, i.e., the node genes and the connection genes, exist in the NEAT. The node genes, which are used to represent all the units of the evolved NN, are encoded with the type of the unit (i.e., the input unit, the hidden unit, or the output unit) and one identification number. The connection genes that are employed to denote the connection information between the node genes, and one node gene is encoded with five elements (the numbers of the input and output units, the value of the connection, one bit indicating whether the connection is activated or not, and one innovation number which records the index of the connection gene with an increased manner). During the evolution process, the individuals are first initialized only with the input and output units of the network, and the random connections between these units. Then, individuals are recombined and mutated. To be specific, there are two types of mutations including the connection mutations and the node mutations. When the connection mutations occur, one connection gene will be added to the list of the connection genes to denote that a pair of node genes is connected. While for the node mutations, one hidden node is generated, then the corresponding connection gene is created to split one existed connection into two parts. Although the NEAT is flexible to evolve NNs, a deterministic number of the output is required, which is impractical in the DL. Furthermore, due to each connection and unit in NEAT are explicitly encoded, it is not suitable for evolving deep NNs that often have a large number connections and units. For remedying this deficiency regarding the incapacity of evolving deep NNs, the connective compositional pattern producing networks (CPNN) [45], [52] has been presented and led to the HyperNEAT.

*2)* The HyperNEAT has been proposed by combining the NEAT with the CPNN encoding scheme. Particularly, the CPNN employs one low-dimensional coordinate system to generate connections for the NEAT by a list of predefined nonlinear functions. To be specific, any point in the coordinate system is picked up, and then fed into a series of compositional functions from the list to complete the transformation from the genotype to the phenotype. Because any number of points can be selected from the low-dimensional coordinate system, numerous connections would be represented with a low computational cost. In this regard, the HyperNEAT has the most potential for evolving a deep NN, while the size of the output still needs to be set in advance, which faces the same problem to NEAT in practice. Furthermore, all the values of the connections in the HyperNEAT are generated by the genetic operators during the evolution, which cannot guarantee the best performance in evolving a deep NN due to the nature of the large-scale global problem. In addition, the recurrent connections or the connections between the same layers are involved in this algorithm, which is also not suitable for learning compact meaningful representations.

*3)* As we have discussed in Section I, the performance of DL algorithms is highly affected by the hyper-parameter settings and the parameter values. In the pre-training phases, one of the key hyper-parameters is the size of hidden layers.

One problem would be naturally raised when EA approaches are employed to search for the sizes, that is how we can ensure the upper bound of the hidden layer sizes given a fixed-length gene encoding strategy. Although the indirect encoding scheme can alleviate this situation somewhat, it limits the generalization of the evolved NNs and the feasible architecture space [37]. On the other hand, if we employ a larger number as the upper bound, it is difficult to determine how large it is reasonable because too large a number would consume more computational resources, otherwise deteriorate the model performance. Excitingly, Yang *et al.* [53] have mathematically pointed out that the meaningful representations of the input data lie at its original space. Supposed that the input data is with $n$ dimension, the size of the associated hidden layer should be no more than $n$. Furthermore, we know that $n$ orthogonal $n$-dimensional basis vectors are sufficient to span a $n$-dimensional space based on Theorem 1. Consequently, we only need to compute one basis $r_1$ of $n$-dimensional space, and the other $(n-1)$ $n$-dimensional basis vectors can be explicitly computed by (3) to find the null space[4]. To this end, we can efficiently model the problem with $n^2$ parameters by employing a genetic algorithm to explicitly encode about $n$ parameters, which is a computational efficient gene encoding approach.

**Theorem 1.** *A set of orthogonal vectors $b_i \in R^n$ ($i = 1, \cdots, n$) is sufficient to span the space $S \in R^n$.*

$$\text{null space}(r_1) = \{x \in R^n | r_1 x = 0\} \tag{3}$$

*4)* Here, we would point out another challenge to inspire our motivation for evolving deep NNs by employing GAs. In our proposed algorithm, the computationally efficient gene encoding strategy mentioned above is employed to model unsupervised deep NNs where a large number of parameters exist. Although the length of the encoded parameters has been reduced appreciably in this regard, the number of the parameters in the original problems remains constant no matter what encoding method is employed. In fact, the effects of one gene in the employed encoding strategy is equivalent to that of multiple parameters in the original problems. For example, for an NN which has $100,000$ parameters, only $1,000$ genes are employed by the computationally efficient gene encoding strategy proposed herein. As a result, one gene represents $100$ parameters in average, and if one gene is changed with the crossover and mutation operators, it will involve the changes of $100$ parameters. Moreover, it is well known that performances of EAs are guaranteed by their exploration search (given by mutation operators) and exploitation search (given by crossover operators) which introduce the global search and local search abilities, respectively. Because a slight change of one gene in the proposed algorithm will lead to the changes of many parameters which affect the global behavior, it can be viewed as that EAs lack of the local search from the problem to be solved. In addition, the

data which are processed by DL algorithms is common with high dimension, which leads to a large number of decision variables in the encoded chromosomes of EAs, although our employed encoding strategy has saved much space compared to existing approaches. Extensive experiments have quantified that EAs are difficult to reach the best performance upon the problems with high input dimensions. To address this issue, we incorporate a local search strategy into the proposed algorithm for assuring the desirable performance.

In summary, the difficulties of deep unsupervised NNs for learning meaningful representations and EAs for evolving deep NNs have been clarified first, and then addressed by our motivations in this section. In the next section, the technical details will be implemented based on these motivations.

## III. PROPOSED ALGORITHM

In this section, the details of the proposed EUDNN are presented. To be specific, the framework which is composed of two distinct stages is depicted at first (Subsection III-A). Next the specifics of each stage are elaborated, respectively (Subsections III-B and III-C). Furthermore, the over-fitting problem preventing mechanism of EUDNN and the significant differences against its peer competitor are discussed (Subsection III-D).

### A. Framework of EUDNN

In this subsection, the framework of the proposed EUDNN is presented. For convenience of the development, it is assuming that the learned representations are for a classification task in which the *meaningful* representations can improve its performance in term of a higher Correct Classification Rate (CCR) (the CCR upon the training data is collected during the training/optimization phase, and that upon the test data during the test/experimental phase). Moreover, given a set of data $D$ in this classification task, a portion of $D$ which is denoted by $D_{train} = \{(x_1, y_1), \cdots, (x_k, y_k)\}$ is considered as the training data in which $x_i$ denotes the input data and $y_i$ is the corresponding label, while the remaining data is regarded as the test data $D_{test}$ for checking whether the learned representations are meaningful. Furthermore, the flowchart of the proposed EUDNN is illustrated in Fig. 4, which clearly shows the two stages of the design: 1) finding the optimal architectures in deep NNs, the desirable initialization of connection weight, and the activation functions (pre-training), and 2) fine-tuning all of the parameter values in connection weights from the desirable initialization.

To this end, one genetic approach with an efficient strategy introduced in Subsection II-B is employed to encode the potential architectures and the associated large numbers of parameters in connection weights by a set of individuals, and then the EA is utilized to evolve and select the individual who has the best performance based on the fitness measures. For warranting the learned representations being meaningful, the method introduced in Subsection II-A is employed, i.e., a small part of data $D_f$ from $D_{train}$ is randomly selected, and the representations of $D_f$ are learned based on the models encoded by the individuals, then they are fed with

---

[4]Theoretically, multiple solutions could be found in computing the bases of the null space. In practice, we only accept the orthonormal basis for the corresponding null space obtained from the singular value decomposition.
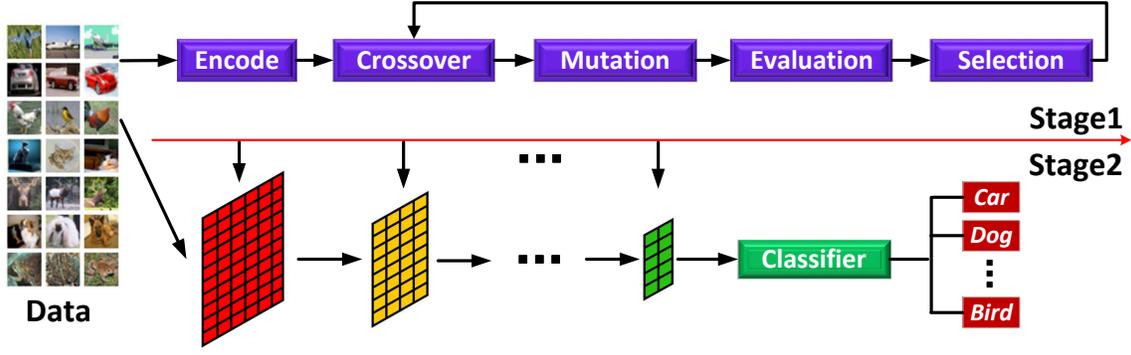
Fig. 4. The flowchart of the proposed algorithm that is composed of two distinct stages. Especially, the first stage is for finding optimal architectures as well as desirable initializations of the connection weight parameter values. The second stage is to fine-tune them for a potentially better performance.

---

**Algorithm 1:** Framework of the Proposed EUDNN

**Input:** Training data $D_{train}$; maximum number $p$ of layers; classifier $C(\cdot)$; test data $D_{test}$.

**Output:** Predicted labels of $D_{test}$.

1 $i \leftarrow 0$;
2 **while** $i < p$ **do**
3    $i \leftarrow i + 1$;
4    $W_j, f_j(\cdot) \leftarrow$ Obtain the optimal connection weight and the corresponding activation function via evolving;
5 **end**
6 Fine-tune all the connection weights $W_1, \cdots, W_p$;
7 $Y_{test} = C(f_p(W_p \times \cdots f_2(W_2 \times f_1(W_1 \times D_{test}))))$;
8 **Return** $Y_{test}$.

---

**Algorithm 2:** Obtain the Optimal Connection Weight and Activation Function

**Input:** Input data; size of population $m$; probability of crossover $\rho$; probability of mutation $\mu$.

**Output:** Optimal connection weight $W$; activation function $f(\cdot)$.

1 Initialize the population $P$ with the size $m$;
2 **while** *stopping criteria are not satisfied* **do**
3    Evaluate the fitness of individuals in $P$;
4    $Q \leftarrow$ Generate new offspring with the probability $\rho$ from two parents selected with binary tournament selection;
5    $Q \leftarrow$ Mutate all the individuals in $Q$ with the probability $\mu$;
6    $S \leftarrow$ Select the individual with the best fitness from $P \cup Q$;
7    $P \leftarrow S \cup$ Select $(m-1)$ individuals from $(P \cup Q) \setminus S$ with binary tournament selection;
8 **end**
9 Evaluate the fitness of the individuals in $P$;
10 $ind_{best} \leftarrow$ Select the individual with the best fitness from $P$;
11 **Return** $W$ and $f(\cdot)$ represented by $ind_{best}$.

---

the associated classification task to select the ones which give the higher CCR for evolution. Based on the investigations in Subsection II-B, a fine-tuning approach additionally, which introduces the exploitation local search, is utilized in the second stage to archive the best performance ever found, which complements with the exploration global search in the first stage. In summary, these two stages collectively ensure the learned representations to be meaningful through unsupervised deep NNs.

In addition, the framework of the proposed EUDNN is presented in Algorithm 1. Specifically, lines 2-5 describe the first stage, while line 6 defines the second stage. Finally, the predicted labels of the test data are calculated and returned in lines 7 and 8. Next, the details of these two stages are documented, respectively.

*B. Obtaining Optimal Connection Weights and Activation Functions via Evolving*

The process of obtaining all the optimal connection weights and their corresponding activation functions contains a series of repeated subprocesses. In this subsection, we first in Algorithm 2 propose how to obtain one optimal connection weight and its activation function. Then, the entire process is described.

To be specific in Algorithm 2, $m$ individuals that encode the information of potential optimal connection weights and their

corresponding activation functions are initialized first (line 1). Then, the evolution takes effect (lines 2-8) until the stopping conditions, such as exceeding the maximum generations, are met. During each generation, the fitness of all the individuals are evaluated first (line 3). Next, new offspring are generated with the probability $\rho$, and their parents are selected from $P$ with the binary tournament selection (line 4). Then, all the offspring in $Q$ are mutated with the probability $\mu$ (line 5). Furthermore, lines 6-7 describe the environmental selection in which the best individual is preserved first for the elitism, then $m - 1$ individuals are selected from the remaining solutions in $P \cup Q$ with binary tournament selection. Specifically, two individuals are randomly selected from $(P \cup Q) \setminus S$ first. Then the one with better CCR is chosen, and the other is put back. With the same process, this operation is repeated $m - 1$ times.

When the evolution terminates, the best solution is selected from the current population for transforming the optimal
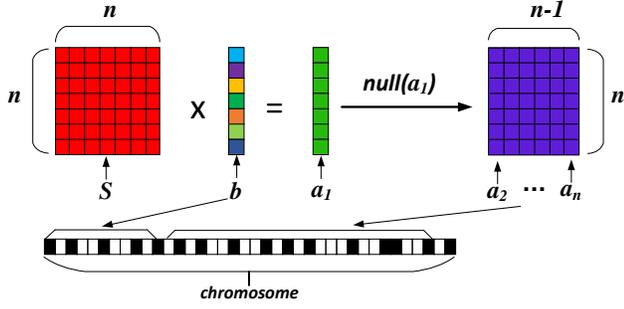
Fig. 5. A flowchart describes the process of encoding the potential connection weight and activation function. First, a set of basis vectors $S$ is given in the original space with $n$-dimension. Then, a set of coefficients $b$ is generated to represent the vector $a_1$ by linear combining the basis vectors. Then, the orthogonal complements $\{a_2, \cdots, a_n\}$ of $a_1$ are computed. Finally, all the information of computing $a_1$, indicating whether the basis from $\{a_2, \cdots, a_n\}$ is selected, and the activation functions are encoded into the chromosomes that are used to evolve to obtain the optimal connection weight and activation function.

connection weight and the activation function (lines 9-10). Next, the details of the employed gene encoding strategy will be discussed, although its fundamental principles have been documented in Subsection II-B. It has been pointed out in [53] that the potential connection weight for obtaining the meaningful representations likely lies in a subspace of the original space. As a consequence, the search for the optimal connection weight can be constrained in the space of input data. Specifically, it is assuming that the input data is $n$-dimensional. First, a set of basis $S = [s_1, \cdots, s_n]$ which can span a $n$-dimensional space is given, e.g., any $n$ linear independent $n$-dimensional vectors. Then the vector $a_1$ is linearly combined by the bases in $S$ with the coefficients $b = [b_1, \cdots, b_n]$ that are randomly specified. Next, the orthogonal complements $\{a_2, \cdots, a_n\}$ of $a_1$ are computed by (3). It is obvious that $\{a_1, a_2, \cdots, a_n\}$ are capable of spanning the space of input data. Finally, a part of these bases, which span a subspace of the original space, are selected for constructing the optimal connection weight by a binary encoded string indicating whether the corresponding basis is available. Furthermore, the corresponding activation function is also encoded into the chromosome. Specifically, a list of selected activation functions with different nonlinear capacities is given, then their indexes in this list are chosen to indicate which one is selected. Moreover, Fig. 5 is provided to intuitively illustrate our intention on efficiently encoding the connection weight and activation function. When the optimal connection weight $W_i$ and its corresponding activation function $f_i$ are found for the $i$-th layer with Algorithm 2, then that for the $(i + 1)$-th layer can be optimized with the same algorithm by setting the input data as $f_i(W_i \times R_i)$ where $R_i$ denote the representations at the $i$-th layer. In the employed gene encoding approach, each coefficient of $b$ is represented with nine bits in which the leftmost bit denotes the positive or negative of the coefficient. Then, one bit is used to indicate whether the basis $a_j$ $(j \in [2, \cdots, n])$ is selected for the connection weight. Finally, two bits are utilized to represent the activation function. In addition to

the well-adopted sigmoid and hyperbolic tangent functions, rectifier function [54], which is reported recently to have a superior performance in some applications, is also considered as one candidate. As a consequence, one chromosome needs $10n + 1$ bits for the $n$-dimensional input data. If the real number encoding method is employed here, a multiple of eight memory space would be taken, which is the major reason that the proposed EUDNN employs the binary encoding method being a contribution to the so claimed computational efficient gene encoding strategy.

Furthermore, the linear Support Vector Machine (SVM) [55] is employed for evaluating the quality of individuals due to its promising computational efficiency and its linear nature for better discriminating power whether the learned representations are meaningful or not. Next, we will give the details of the fitness evaluation by using SVM based on the design principle described in Subsection II-A4. For convenience of the development, let $D_{train} = \{X_{train}, Y_{train}\}$ denote the training set where $X_{train}$ are the data and $Y_{train}$ are the corresponding labels, and the selected individual for fitness evaluation is denoted by $ind_i$. Firstly, a small fraction of data denoted by $D_{eval} = \{X_{eval}, Y_{eval}\}$ is randomly selected from $D_{train}$. Secondly, the corresponding model is transformed from the encoded individual $ind_i$. Thirdly, the representations (denoted by $F_{eval}$) of $X_{eval}$ are calculated based on the formulas in (1). Fourthly, $\{F_{eval}, Y_{eval}\}$ are fed to SVM and the CCR on $X_{eval}$ is estimated. Finally, the CCR is used as the fitness of $ind_i$.

### C. Fine-tuning Connection Weights

To further improve the performance, an exploitation mechanism implemented by local search strategy is incorporated into the second stage to fine-tune parameter values in connection weights. In this stage, the architecture is fixed with the evolved activation functions and the initialization values of the connection weights, and then a local search method is used to tune the connection weights further. Fig. 6 shows an example of this process. Specifically, when all the connection weights and activation functions have been optimized in the first stage, all the hidden layers are connected to a list based on their orders in the first stage by adding one input layer at the top of this list. Then, the connection weights in this list are initialized with the values confirmed in the first stage. Finally, a classifier is added to the tail of this list to perform the fine-tuning process. Note here that the BP algorithm is employed for the fine-tuning. Actually, any local search algorithm can be used in the second stage. The reasons for employing BP are largely due to two aspects: 1) the gradient information in the loss function is always analytical and the BP that is based on the gradient is naturally employed in most designs; 2) multiple libraries of BP have been implemented for accelerating the computation with the Graphics Processing Units (GPUs) and the computational cost can be reduced remarkably, especially in the situations of processing high-dimensional data. Furthermore, when the rectifier activation function that is not differentiable at the point 0 is selected, the value 0 is assigned according to the convention of the community [56].
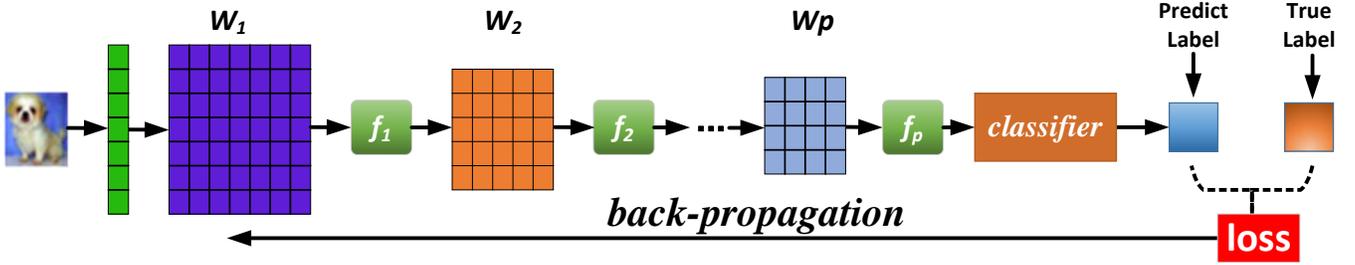
Fig. 6. The flowchart of the second stage in the proposed EUDNN. Especially, the predicted label is computed with the connection weights and activation functions for the input data. Then the loss of the classifier is formulated between the predicted label and the true label. Next, the error is back propagated and the parameter values of the connection weights are updated.

*D. Discussions*

In this subsection, we mainly discuss the over-fitting problem preventing mechanism utilized by the proposed EUDNN, and the significant differences of the proposed EUDNN against the Direct Evolutionary Feature Extraction algorithm (DEFE) [57] that employs a similar gene encoding strategy to EUDNN.

The over-fitting problem implies the poor generalization ability of models, i.e., the trained model reaches a better CCR upon training data at the cost of a worsen CCR upon test data. Because the goal in training a classification model is for obtaining a higher CCR upon test data, the over-fitting problem should be prevented by some mechanisms. Commonly, given a number of models which are all capable of solving a particular classification task, the model with a smaller Vapnik Chervonenkis (VC) dimension[5] [7] usually has a better generalization ability, which does not lead to an over-fitting problem. Because the number of parameters is positive to the value of a VC dimension, and deep NN architectures are generally with the numerous number of parameters, the over-fitting problem easily occurs in these models.
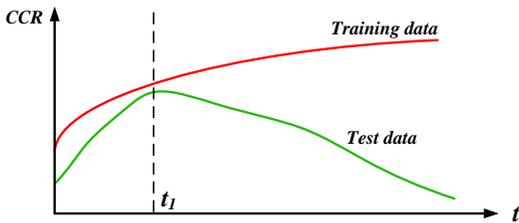


Fig. 7. Correct classification rates of training data and test data as training process continues.

More specifically, Fig. 7 illustrates a typical instance in CCR on training data (red curve) and CCR of checking on test data (green curve) as the training process continues. Especially, CCR on both data are continuously growing until the time $t_1$, and CCR on the training data continues to increase while CCR on the test data begins to drop when the training time

is greater than $t_1$, which obviously indicates the presence of an over-fitting problem. As we have claimed that the best performance of the proposed EUDNN cannot be guaranteed during the training in the first stage, and the second stage is introduced to expectedly help the proposed EUDNN arrive at the best performance. To this end, it is concluded that the over-fitting problem will not occur in the first stage of the proposed EUDNN (because the first stage terminates prior to the time $t_1$, while the over-fitting problem might occur after the time $t_1$), but may occur in the second stage. Consequently, some rules need to be utilized to prevent this problem only in the second stage. Here, the "early stop" approach is utilized for this purpose, i.e., a group of data $D_{validate}$ is uniformly selected from $D_{train}$ as the validate data to replace the checking upon test data in Fig. 7, when we first observe the CCR of validate data begins to decrease while the CCR of training is still increasing (i.e., the particular time $t_1$ is found), the fine-tuning in the second stage is terminated and the optimal model that gives the best performance is obtained. Next, the second concern, i.e., the differences between the proposed EUDNN and the DEFE, will be discussed.

It has been observed that 1) DEFE learns only linear representations and 2) shallow representations of input data. These two observations cause that DEFE cannot learn the meaningful representations [28]. Next, the details of these conclusions are discussed. To be specific, the learned representations $R$ of DEFE can be formulated as $R = WX$ [57] where $W$ is the transformation matrix (i.e., the connection weight in deep NN models) and $X$ is the input data. It is evident that there is no nonlinear transformation upon $WX$. Consequently, only linear representations would be learned by DEFE, while in the proposed EUDNN, a list of nonlinear activation functions with different nonlinear transformation abilities is incorporated into the evolution for performing nonlinear representation learning. Furthermore, although multiple transformations like that in the proposed EUDNN can be implemented by DEFE to learn deep representations, deep linear transformations are equivalent to a one layer linear representation.

In summary, DEFE cannot be employed for learning meaningful representations due to its linear nature, while the success of deep NNs is mainly caused by the meaningful representations learned by deep nonlinear transformations, which have been explicitly implemented by the proposed EUDNN.

[5]Generally, the VC dimension can be viewed as an indicator measuring the complexity of multiple models which are capable of solving one particular task [58]. The smaller the VC dimension, the more simplicity is the corresponding model, and a more simplicity model is with better generalization [59]. Commonly, a large number and magnitude of elements in the transformation matrixes are positive to the VC dimension.

## IV. Experiments

In order to examine the performance of the proposed EUDNN, experiments based on a set of image classification benchmarks against selected peer competitors are performed. During the comparisons, the chosen performance metric considers the CCR on the test data. In the following, the employed benchmarks are outlined first. Then the chosen peer competitors are reviewed, and the justification for selecting them is explained further. This is followed by the descriptions of the performance metric chosen and the specifics of parameter settings employed by these compared algorithms. Finally, the quantitative as well as the qualitative experimental results are illustrated and comprehensively analyzed.

### A. Benchmark Test Datasets

Benchmarks used by compared algorithms are the handwritten digits benchmark test dataset MNIST [21], basic MNIST dataset (MNIST-basic) [60], a rotated version of MNIST (MNIST-rot) [60], MNIST with random noise background (MNIST-back-rand) [60], MNIST with random image background (MNIST-back-image) [60], MNIST-rot with random image background (MNIST-rot-back-image) [60], tall and wide rectangles dataset (Rectangles) [60], rectangles dataset with random image background (Rectangles-image) [60], convex sets recognition dataset (Convex) [60], and the gray version of Canadian Institute for Advanced Research object recognition dataset [61] (Cifar10-bw) over 10 classes, i.e., airplane, automobile, bird, cat, deer, dog, frog, horse, ship, and truck.



Fig. 8. A group of digit samples $(0 - 9)$ from the MNIST benchmark test dataset.

Briefly, these benchmark test datasets are categorized into three different classes based on the object types that they intend to recognize. The first one is about the hand-written digits and covers the MNIST, MNIST-basic, MNIST-rot, MNIST-back-rand, MNIST-back-image, and MNIST-rot-back-image benchmarks. Examples from the MNIST benchmark are depicted in Fig. 8 for reference. The second one is to classify the geometries and the rectangles, such as the Rectangles, Rectangles-image, and the Convex benchmarks. The last one is to identify the natural objects in Cifar10-bw. Different variants in MNIST and rectangles datasets present the algorithms dissimilar difficulties from the aspects of perturbations, the small number of training dataset, and the large testing dataset size. Furthermore, the dimensions, number of classes, and the sizes of training set and test set of the chosen benchmark datasets are shown in Table I.

### B. Performance Metric

Technically speaking, it is difficult to directly evaluate whether the learned representations are meaningful or not because they are intermediate outcomes. A general practice

TABLE I
THE CONFIGURATIONS OF THE CHOSEN BENCHMARK DATASETS.

| Benchmark | Dimension | # of class | Size of training set | Size of test set |
|---|---|---|---|---|
| MNIST | $28 \times 28$ | 10 | 50,000 | 10,000 |
| MNIST-basic | $28 \times 28$ | 10 | 12,000 | 50,000 |
| MNIST-rot | $28 \times 28$ | 10 | 12,000 | 50,000 |
| MNIST-back-rand | $28 \times 28$ | 10 | 12,000 | 50,000 |
| MNIST-back-image | $28 \times 28$ | 10 | 12,000 | 50,000 |
| MNIST-rot-back-image | $28 \times 28$ | 10 | 12,000 | 50,000 |
| Rectangles | $28 \times 28$ | 2 | 1,200 | 50,000 |
| Rectangles-image | $28 \times 28$ | 2 | 12,000 | 50,000 |
| Convex | $28 \times 28$ | 2 | 8,000 | 50,000 |
| Cifar10-bw | $32 \times 32$ | 10 | 50,000 | 10,000 |

for this is to feed these learned representations to a particular classification task, and then to investigate the CCR by a classifier. Commonly, a higher CCR implies that the learned representations are more meaningful. Because the benchmarks employed in these experiments are multi-class classification tasks, the softmax regression classifier [62] is employed here to measure the corresponding CCR according to the convention adopted in the community.

It is assumed that a set of training data and their corresponding labels with $k$ distinct integer values are denoted as $\{x_1, \cdots, x_m\}$, and $\{y_1, \cdots, y_m\}$, respectively, where $x_i \in \mathcal{R}^n$ and $y_i \in \{1, \cdots, k\}$. To be specific, the label of the sample $x_i$ $(i \in \{1, \cdots, m\})$ is predicted by (4) with the softmax regression,

$$\arg\max_{j} \ p_j(x_i) = \frac{exp(\theta_j^T x_i)}{\sum_{l=1}^{k} exp(\theta_l^T x_i)} \qquad (4)$$

where $\Theta = [\theta_1, \cdots, \theta_k]^T$ are obtained by minimizing

$$J(\Theta) = -\frac{1}{m} \left[ \sum_{i=1}^{m} \sum_{j=1}^{k} f(y_i, j) log \frac{exp(\theta_j^T x_i)}{\sum_{l=1}^{k} exp(\theta_l^T x_i)} \right]$$

in which $f(y_i, j) = 1$ if $y_i = j$, otherwise $f(y_i, j) = 0$.

### C. Compared Algorithms

Because of the proposed EUDNN aiming at evolving *unsupervised deep neural networks* for learning *meaningful representations*, algorithms related to evolving deep NNs (NEAT [44], HyperNEAT [45]), unsupervised deep NNs (DBNs [51], and variants of stacked AEs [33]) that have been discussed in Section I should be all employed as peer competitors. However, the NEAT and the HyperNEAT cannot be used to learn meaningful representations due to the reasons that have been discussed in Section I and further analyzed in Section II. As a result, they are excluded from the selected compared algorithms. To this end, DBNs and variants of stacked AEs are employed for performing the comparison experiments. Because RBMs [50] and AEs [10], [29], [30] are the building blocks to train DBNs and stacked AEs, respectively, these two types of algorithms are considered as the peer competitors in our experiments to compare the performance of the learned representations against that of the proposed algorithm (i.e., we will evolve RBMs and AEs as the unsupervised deep NN models, which are named EUDNN/RBM

and EUDNN/AE, respectively, to perform the comparisons against considered peer competitors). Specifically, the variants of AEs, i.e., the Sparse AEs (SAEs) [31], the Denoising EAs (DAEs) [34], and the Contractive AEs (CAEs) [35], have been proposed with different regularization terms for learning meaningful representations in recent years and also have obtained comparable performance in multiple tasks. As a consequence, they are also included as the peer competitors in the experiments, in addition to the DBNs.

### D. Parameter Settings

For a fair comparison, multiple parameters in the second stage of the proposed EUDNN and the competing ones are the same. As a consequence, we will first give details of these generic parameter settings in this subsection. Then, the particular parameter settings are individually introduced. Because the best performance of the compared algorithms often strongly depends on the particular benchmark dataset and the corresponding parameter settings, in order to do a fair comparison, we first test these parameters from the range widely used in the community upon the corresponding training data, then the best performance upon test data of each compared algorithm is selected for comparisons.

*1) Learning Rate and Batch Size* The Stochastic Gradient Descent (SGD) algorithm is chosen as the algorithm to train the SAE, the DAE, the CAE, and the softmax regression, and its learning rates as well as the batch sizes vary in $\{0.0001, 0.001, 0.01, 0.1\}$ and $\{10, 100, 200\}$, respectively, according to the community convention.

*2) Number of Runs and Stop Criteria* All the compared algorithms are independently performed 30 runs. In addition, a performance monitor is injected into each epoch in training the softmax regression to record the best CCR over the test dataset as the best performance of the algorithm that feeds the HLlearned representations to the softmax regression.

*3) Unit Number and Depth* The number of the units for the SAE, the DAE, the CAE, and the RBM in each layer is set to be from 200 to 3,000 using a $log$ function with an interval 0.5 as recommended by [63], and the maximum depth is set to be 5 (this depth is excluded from the input layer, i.e., the maximum number of hidden layers).

*4) Statistical Significance* The results measured by the selected performance metric need to be statistically compared due to the heuristic natures of the first stage in the proposed EUDNN. In these experiments, the Mann-Whitney-Wilcoxon rank-sum test [64] with a $5\%$ significant level is employed for this purpose according to the community convention.

In addition, the sparsity of the SAE, the binary corrupted level of the DAE, and the coefficient of the contractive term in the CAE are set to be $10\%, 30\%, 50\%$ and $70\%$, respectively. Because of the nature of the RBM, the CD-$k$ algorithm [65] is selected as its training algorithm and $k$ is set to be 1 based on the suggestion in [63]. In order to speed up the proposed algorithm in the first stage, a proportion (i.e., $20\%$) of the training dataset is randomly selected in each generation for the fitness evaluation. In addition, the connection weights and the biases are respectively set to be

between $[-4 \times 6/\sqrt{n_{number}}, 4 \times 6/\sqrt{n_{number}}]$ with a uniform sampling and 0, respectively [66], if required, where $n_{number}$ denotes the total number of the units in two adjacent layers based on the experiences suggested in [66].

Because parameter settings in the second stage of the proposed EUDNN are the same as that of the peer competitors, parameter settings of the evolution related parameters in the first stage are declared next. Conveniently, one chromosome in this stage can be divided into three parts: main basis related coefficients (Part 1) which are used to represent the vector $a_1$ in Fig. 5, projected space related coefficients (Part 2) which are employed to indicate which bases are selected for the connection weight, and the coefficients (Part 3) which denote the type of activation functions. Because Parts 1 and 2 have strong effects on the quality of the connection weight, it is hopefully that crossover operation should be promoted in these two parts for improving the exploitation local search that provides much better performance based on the exploration global search. As a consequence, one point crossover operator is employed in Parts 1 and 2. In addition, three widely used nonlinear activation functions are considered in the proposed algorithm and one is to be selected for the corresponding connection weight. Therefore, it is hopefully that the information representing the activation function is not modified often since it is hard to determine which one is the best. Consequently, Parts 2 and 3 are considered as one part to participate in the crossover operation. It is noted here that, when the value in Part 3 is invalid, a random one is chosen to reset it. Noting that the polynomial mutation [67] is used here as the mutation operator (distribution index is set to be 20). In addition, the population size is set to be 50. As for the crossover probability and the mutation probability in the proposed algorithm, both of them are set to be the same as that of the community convention (i.e., 0.9 for crossover and 0.1 for mutation). A proportion of 10% is randomly selected from the training set for the fitness evaluation. Codes of the proposed EUDNN can be made available upon request through the first author.

### E. Experimental Results

Based on the motivation of our design, the proposed EUDNN 1) employs evolutionary algorithm and local search strategy to ensure the learned representations through deep NNs to be meaningful, 2) employs evolutionary approach in the first stage to help the deep NNs find the optimal architectures and the good initialized weights, which give a better starting position for the second stage, and 3) employs the local search strategy in the second stage to improve the intended performance much further. Consequently, a series of experiments are carefully crafted to evaluate the performance of the proposed design.

*1) Performance of the Proposed Algorithm* In order to quantify whether the representations learned by the proposed EUDNN are meaningful, a series of experiments are well-designed and comparisons are performed. Specifically, EUDNN/AE and EUDNN/RBM are two implementations of the proposed algorithm over the unsupervised neural network models (i.e., AEs and RBMs, respectively). Then they are used

to learn the representations together with the selected peer competitors employing the configurations introduced above. Next, the softmax regression metric is employed to measure whether the learned representations improve the associated classification tasks through CCR, which in turn indicates the learned representations being meaningful or not.

Particularly, the mean values and standard derivations of CCR resulted by these compared algorithms over 30 independent runs are listed in Table II in which the best results over the same benchmark are highlighted in boldface. In addition, the symbols "+," "-," and "=" denote whether the CCR of the proposed algorithm upon the corresponding benchmarks are statistically better than, worse than, and equal to that of the associated peer competitors, respectively, with the employed rank-sum test[6]. Furthermore, the summarizations, how many times over the considered benchmarks the proposed EUDNN are better than, worse than, and equal to the corresponding peer competitor, are listed in the last row of Table II.

In Table III, the first column shows the names of the chosen benchmark datasets, the second column provides the corresponding best CCRs obtained, while the third column presents the numbers of neurons of the deep models (excluding the the classifier layer) with which the best CCRs are reached on the corresponding benchmark dataset. As we have claimed in Subsection IV-D that the maximum number of building blocks investigated in this paper is set to be five. Therefore, the number of layers, which include the input layer and hidden layers, shown in Table III for each benchmark dataset does not exceed six. For the first row in Table III as an example, it indicates that the best CCR of 98.85% on the MNIST benchmark dataset is achieved with only four building blocks where the input layer is with 784 neurons, and hidden layers are with 400, 202, 106, and 88 neurons, respectively.

It is clearly shown in Table II[7] that the proposed EUDNN/AE obtains the best mean values upon the MNIST-rot, the MNIST-rot-back-image, the Convex, and the Cifar10-bw benchmarks, and the best rank-sum results upon the MNIST-rot, the Convex, and the Cifar10-bw benchmarks. Moreover, the proposed EUDNN/RBM wins both the best mean values and the rank-sum results upon the MNIST, and the MNIST-back-image benchmarks. Although the best result of the proposed EUDNN (obtained by the EUDNN/AE) over the MNIST-basic benchmark is a little worse than that of the SAE, which is the winner of the best mean value and rank-sum results, EUDNN/AE outperforms all the other peer competitors. Furthermore, the SAE obtains the best mean values upon the MNIST-basic and the MNIST-back-rand benchmarks, but the best result of the proposed algorithm (obtained by the EUDNN/AE) is statistically equal to that of the SAE upon the MNIST-back-rand benchmark, and also outperforms other competing algorithms. Upon the Rectangles-image benchmarks, the best result of the proposed algorithm

(obtained by the EUDNN/RBM) is worse than that of the CAE and the SAE, while the EUDNN/RBM and CAE have the same results statistically. In addition, the best results of the proposed algorithm upon the MNIST-rot-back-image (obtained by the EUDNN/AE) and the Rectangles (obtained by the EUDNN/RBM) benchmarks are all statistically equivalent to that of the DBN, while the best mean values upon these two benchmarks are obtained by the EUDNN/AE and the EUDNN/RBM, respectively. Note here that the MNIST is a widely used classification benchmark for quantifying the performance of deep learning models, and the best results are frequently obtained by supervised models, which require sufficient labeled training data during their training phases. To our best knowledge, the CCR with 98.85% obtained by the proposed algorithm (EUDNN/RBM), which is an unsupervised approach is a very promising result among unsupervised deep learning models. In summary, the proposed algorithm totally wins 34 times over the 40 comparisons against the selected peer competitors, which reveals the superior performance of the proposed algorithm in learning *meaningful representations* with *unsupervised neural network models*.

*2) Performance Analysis Regarding the First Stage* Since we have claimed that the first stage of the proposed algorithm helps the unsupervised NN-based models learn optimal architectures and better-initialized parameter values, component-wise experiments over the optimal architectures and the initialized parameter values should be performed to investigate their respective effects to justify our designs. However, the initialized parameter values are dependent on the architectures. This leads to the specific experiment by varying only the architecture configurations on investigating how the learned architectures solely affect the performance is difficult to design. Hence, the performance regarding the initialized parameter values is mainly investigated here.

To this end, we first record the architecture configurations (see Table III) with which the proposed algorithm presents the promising performance in best mean values of EUDNN/AE and EUDNN/RBM upon each benchmark from Table II. Then experiments are re-performed by peer competitors with the recorded architecture configurations and randomly initialized parameter values. Finally, the learned representations are fed to the considered performance metric to measure whether these representations are meaningful. Specifically, experimental results are depicted in Fig. 9 in which the vertical axis denote the CCR while A-J in the horizontal axis represent the benchmarks MNIST, MNIST-basic, MNIST-rot, MNIST-back-rand, MNIST-back-image, MNIST-rot-back-image, Rectangles, Rectangles-image, Convex, and Cifar10-bw, respectively.

It is shown in Fig. 9 that most of the peer competitors employing the chosen architecture configurations listed in Table III obtain worse CCR upon the considered benchmarks compared to the proposed algorithm. Specifically, the proposed algorithm shows these best CCR upon MNIST, MNIST-rot, MNIST-back-image, MNIST-rot-back-image, Convex, and Cifar10-bw benchmarks, which is consistent with the findings listed in Table II. In addition, the proposed algorithm wins the best CCR upon MNIST-basic and MNIST-back-rand

---

[6]To do this statistically test, we first select the better CCR generated by EUDNN/AE and EUDNN/RBM with the same benchmark, then the selected results are used to do the rank-sum test.

[7]In this paper, the statistical results biases the results generated by the statistical significance toolkit, i.e., the Mann-Whitney-Wilcoxon rank-sum test [67] with a 5% significant level.

| Benchmark | EUDNN | | DAE | CAE | SAE | DBN |
|---|---|---|---|---|---|---|
| | AE | RBM | | | | |
| MNIST | 0.9878(0.00751) | **0.9885(0.00255)** | 0.9820(0.00506)(+) | 0.9843(0.00699)(+) | 0.9832(0.00891)(+) | 0.9771(0.00959)(+) |
| MNIST-basic | 0.9674(0.00616) | 0.9633(0.00473) | 0.9580(0.00352)(+) | 0.9635(0.00831)(+) | **0.9776(0.00585)(-)** | 0.9658(0.00550)(+) |
| MNIST-rot | **0.7952(0.00917)** | 0.7549(0.00286) | 0.7274(0.00757)(+) | 0.7706(0.00754)(+) | 0.7852(0.00380)(+) | 0.7639(0.00568)(+) |
| MNIST-back-rand | 0.8843(0.00076) | 0.8386(0.00054) | 0.7725(0.00531)(+) | 0.5741(0.00779)(+) | **0.8851(0.00934)(=)** | 0.8221(0.00130)(+) |
| MNIST-back-image | 0.4325(0.00569) | **0.4830(0.00469)** | 0.4022(0.00012)(+) | 0.4010(0.00337)(+) | 0.4638(0.00162)(+) | 0.4587(0.00794)(+) |
| MNIST-rot-back-image | **0.8925(0.00906)** | 0.8879(0.00815) | 0.8691(0.00127)(+) | 0.6574(0.00913)(+) | 0.8733(0.00632)(+) | 0.8830(0.00098)(=) |
| Rectangles | 0.9627(0.00311) | **0.9681(0.00829)** | 0.9232(0.00166)(+) | 0.6275(0.00602)(+) | 0.9408(0.00263)(+) | 0.9622(0.00154)(=) |
| Rectangles-image | 0.7521(0.00689) | 0.7716(0.00048) | 0.7598(0.00451)(+) | **0.7810(0.00784)(=)** | 0.7725(0.00002)(-) | 0.7628(0.00913)(+) |
| Convex | **0.8113(0.00052)** | 0.8085(0.00826) | 0.7930(0.00538)(+) | 0.8016(0.00996)(+) | 0.8053(0.00878)(+) | 0.7895(0.00443)(+) |
| Cifar10-bw | **0.4798(0.00107)** | 0.4331(0.00962) | 0.4309(0.00005)(+) | 0.4860(0.00775)(+) | 0.4423(0.00817)(+) | 0.4598(0.00869)(+) |
| +/-/= | | | 10/0/0 | 9/0/1 | 7/2/1 | 8/0/2 |

benchmarks as well, with these architecture configurations. In addition to the proposed algorithm in which the initialized parameter values are set by the proposed evolutionary approach, all the results illustrated in Fig. 9 are obtained by the compared algorithms with the same architecture configurations and commonly used parameter initializing methods for the second stage. As we all know that the performance of local search strategies is strongly rely on their starting position, therefore, it is reasonable to conclude that the evolutionary scheme employed by the first stage of the proposed algorithm has substantially helped the learned representations to be meaningful.

TABLE III

THE BEST CORRECT CLASSIFICATION RATE (CCR) OF THE PROPOSED ALGORITHM UPON MNIST, MNIST-BASIC, MNIST-ROT, MNIST-BACK-RAND, MNIST-BACK-IMAGE, MNIST-ROT-BACK-IMAGE, RECTANGLES, RECTANGLES-IMAGE, CONVEX, CIFAR10-BW BENCHMARKS AND THE CORRESPONDING ARCHITECTURE CONFIGURATIONS.

| Benchmark | Best CCR | Architecture configurations |
|---|---|---|
| MNIST | 0.9885 | 784, 400, 202, 106, 88 |
| MNIST-basic | 0.9674 | 784, 400, 211, 120 |
| MNIST-rot | 0.7952 | 784, 400, 233, 133, 100, 81 |
| MNIST-back-rand | 0.8843 | 784, 397, 202, 123 |
| MNIST-back-image | 0.4830 | 784, 386, 191, 1088, 100 |
| MNIST-rot-back-image | 0.8925 | 784, 378, 205, 106 |
| Rectangles | 0.9681 | 784, 397, 205, 113, 100, 75 |
| Rectangles-image | 0.7716 | 784, 402, 214, 122, 89 |
| Convex | 0.8113 | 784, 394, 200, 110, 55, 49 |
| Cifar10-bw | 0.4798 | 1024, 502, 253, 141, 130 |

*3) Performance Analysis Regarding the Second Stage* In this experiment, we mainly investigate whether the local search strategy employed in the second stage promotes the integral performance of the proposed algorithm compared to only the evolutionary methods used in the first stage. For this purpose, we first pick up the promising CCR obtained by the proposed algorithm from Table II in which the results of the proposed algorithm are collectively achieved by the evolutionary method employed in the first stage and the local search strategy employed in the second stage. Then we select the corresponding results performed without the local search
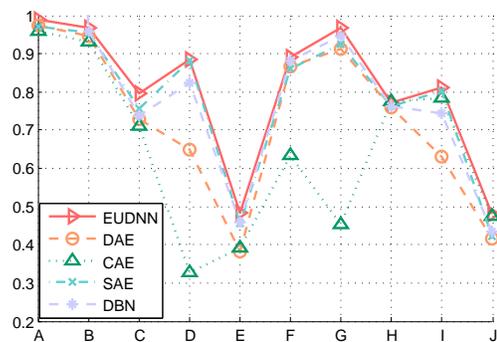


Fig. 9. The performance of the proposed algorithm against DAE, CAE, SAE, and DBN with the configurations on which the proposed algorithm obtains the best correct classification rates over benchmarks measured by softmax regression. Especially, A-J denote the benchmarks MNIST, MNIST-basic, MNIST-rot, MNIST-back-rand, MNIST-back-image, MNIST-rot-back-image, Rectangles, Rectangles-image, Convex, and Cifar10-bw, respectively.

strategy (i.e., the results obtained by the proposed algorithm during the first stage). Finally, these results are illustrated in Fig. 10 for quantitative comparisons. Specifically in Fig. 10 the vertical axis denotes the CCR, while A-J in the horizontal axis represent the benchmarks MNIST, MNIST-basic, MNIST-rot, MNIST-back-rand, MNIST-back-image, MNIST-rot-back-image, Rectangles, Rectangles-image, Convex, and Cifar10-bw, respectively, and the bars in blue denote the results obtained by the proposed algorithm without the second stage, while the bars in red refer to that with the second stage.

It is clearly shown in Fig. 10 that the performance has been improved with the second stage of the proposed EUDNN over all the considered benchmarks compared to the algorithm that only the first stage is employed. Especially, the CCR have been significantly improved by about 20% upon the MNIST-rot, MNIST-back-rand, MNIST-back-image, MNIST-rot-back-image, and Cifar10-bw benchmarks and 12.83% on the MNIST benchmark. In summary, it is concluded from these experimental results that the local search strategy utilized in the second stage helps the performance of the proposed algorithm to be improved much further, which promotes the
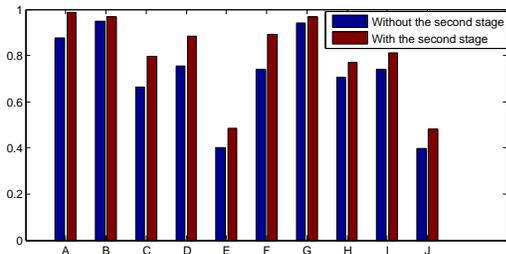
Fig. 10. Correct classification rate (CCR) comparisons of the proposed algorithm without (denoted by blue bars) and with (denoted by red bars) the second stage upon the MNIST, MNIST-basic, MNIST-rot, MNIST-back-rand, MNIST-back-image, MNIST-rot-back-image, Rectangles, Rectangles-image, Convex, and Cifar10-bw benchmarks, which are denoted by A-J, respectively.

learned representations to be meaningful and satisfies our motivation of this design.

*F. Visualizations of Learned Representations*

In Subsection IV-E, a series of quantitative experiments has been given to highlight the performance of the proposed algorithm in learning meaningful representations with unsupervised deep NN-based models. Here, a qualitative experiment is provided for comprehensively understanding what the representations are learned from the proposed algorithm via visualizations, which is a common approach employed by related works [7]–[9], [34], [35] to intuitively investigate the learning representations. For this purpose, the activation maximization approach [68] is utilized to visualize the learned representations of the proposed algorithm over MNIST dataset and a number of 100 randomly selected visualizations of the patches are illustrated [8] in Fig. 11. Furthermore, the SGD is employed during the optimization of the activation maximization with $10,000$ iterations and a fixed learning rate of 0.1. To be specific, Fig. 11a shows the learned representations on depth 1 in which the visualization is commonly describable [68]. It is clear in Fig. 11a that some strokes are learned in most patches and a part of the representations is similar to that of the RBM [68], which can be viewed as the effectiveness of the proposed algorithm, because these similar representations over MNIST dataset have been reported in multiple kinds of literature [8], [9]. The visualizations of the representations on depths 2 and 3 are depicted in Figs. 11b and 11c, respectively. However, these representations are difficult to understand intuitively and be interpretable due to the high-level hierarchical nature [68]. But it still can be concluded that the proposed algorithm has learned the meaningful representations by comparing them to the experiments simulated in [68] that learned representations herein resemble those of the DAE to some extent. Noting that multiple learned features shown in Fig. 11a seem to be random. The reason is that not all the neurons in the corresponding hidden layer have learned the meaningful features. Specifically, the visualization of features is from the 100 neurons randomly selected from the

---

[8]Because visualizations of representations learned from the depth larger than one are difficult to interpret, and that from the depth larger than three have no reference for comparisons, only representations with depths 1, 2, and 3 are visualized here.

313,600 (this number can be calculated from Table III), and it is not necessary that all the 313,600 neurons have learned the meaningful features. In summary, these visualizations give a qualitative observation to highlight that the meaningful representations have been effectively learned by the proposed algorithm.

## V. CONCLUSION

In order to warrant the representations learned by *unsupervised* deep neural networks to be meaningful, the existing approaches for learning them need optimal combinations of hyper-parameters, appropriate parameter values, and sufficient labeled data as the training data. These approaches generally employ the exhaustive grid search method to directly optimize hyper-parameters due to their unavailable gradient information, which give an unaffordable computational complexity that increases with an order of magnitude as the number of hyper-parameter grows. Furthermore, the gradient-based training algorithms in these existing algorithms are easy to be trapped into the local minima, which cannot guarantee them the best performance. In addition, in the current era of Big Data, the volume of labeled data is limited and obtaining sufficient data with labels is expensive, if not impossible. To address these concerning issues, we have proposed an evolving unsupervised deep neural networks method which heuristically searches for the best hyper-parameter settings and the global minima to learn the meaningful representations without sufficient labeled data. To be specific, two stages are composed in the proposed algorithm. In the first stage, all the information regarding hyper-parameter and parameter settings are encoded into the individual chromosome and the best one is selected when they go through a series of crossover, mutation, and selection operations. Furthermore, the activation functions that provide the nonlinear ability to the learning algorithm are also incorporated into the individual chromosome to go through the evolutions of obtaining the promising performance. In addition, the orthogonal complementary techniques are employed in the proposed algorithm to reduce the computational complexity for effectively learning the deep representations. Specifically, only a limited number of labeled data is needed in the proposed algorithm to direct the search to learn representations with meaningfulness. For further improving the performance, the second stage is introduced with a local search strategy to complement with the ability of the exploitation search for training the proposed algorithm with the architecture and the activation function optimized in the first stage. These two stages collectively promote the proposed algorithm effectively learning the meaningful representations with unsupervised deep neural network-based models. To evaluate the meaningfulness of the learned representations, a series of experiments are given against peer competitors over multiple benchmarks related classification tasks. The results measured by the softmax regression show the considerable competitiveness of the proposed algorithm in learning meaningful representations. In near future, we will place more focus on the efficient encoding methods as well as the way measuring the quality of the representation during the evolution of a larger scale and higher
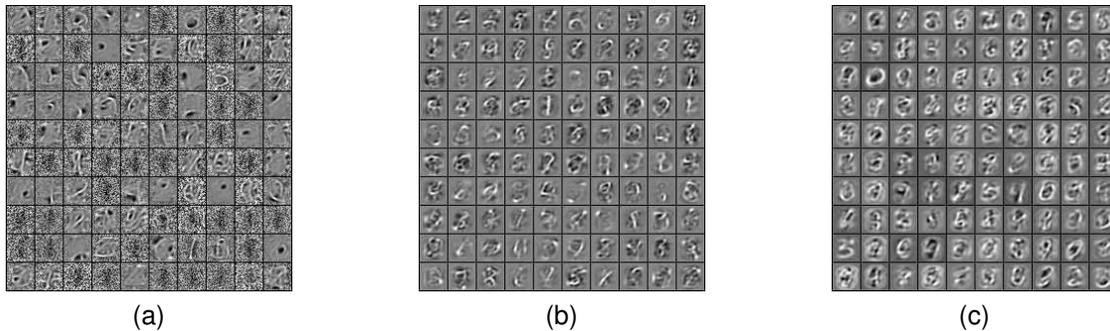
Fig. 11. Visualizations of the proposed algorithm over MNIST dataset with depths 1 (Fig. 11a), 2 (Fig. 11b), and 3 (Fig. 11c) by activation maximization method.

dimensional data. In addition, we would also investigate how to effectively evolve deep supervised neural networks, such as CNNs.

## REFERENCES

[1] Y. Bengio, A. Courville, and P. Vincent, "Representation learning: A review and new perspectives," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 35, no. 8, pp. 1798–1828, 2013.

[2] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Advances in Neural Information Processing Systems*, 2012, pp. 1097–1105.

[3] C. Farabet, C. Couprie, L. Najman, and Y. LeCun, "Learning hierarchical features for scene labeling," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 35, no. 8, pp. 1915–1929, 2013.

[4] J. J. Tompson, A. Jain, Y. LeCun, and C. Bregler, "Joint training of a convolutional network and a graphical model for human pose estimation," in *Advances in Neural Information Processing Systems*, 2014, pp. 1799–1807.

[5] O. Delalleau and Y. Bengio, "Shallow vs. deep sum-product networks," in *Advances in Neural Information Processing Systems*, 2011, pp. 666–674.

[6] D. H. Hubel and T. N. Wiesel, "Receptive fields, binocular interaction and functional architecture in the cat's visual cortex," *The Journal of physiology*, vol. 160, no. 1, pp. 106–154, 1962.

[7] Y. Bengio, "Learning deep architectures for ai," *Foundations and Trends in Machine Learning*, vol. 2, no. 1, pp. 1–127, 2009.

[8] Y. Sun, H. Mao, Q. Guo, and Z. Yi, "Learning a good representation with unsymmetrical auto-encoder," *Neural Computing and Applications*, pp. 1–7, 2015.

[9] Y. Sun, H. Mao, Y. Sang, and Z. Yi, "Explicit guiding auto-encoders for learning meaningful representation," *Neural Computing and Applications*, pp. 1–8, 2015.

[10] D. E. Rumelhart, G. E. Hinton, and R. J. Williams, "Learning representations by back-propagating errors," *Cognitive Modeling*, vol. 5, no. 3, p. 1, 1988.

[11] R. S. Sutton, "Two problems with backpropagation and other steepest-descent learning procedures for networks," in *Proc. 8th annual conf. cognitive science society*, 1986, pp. 823–831.

[12] J. Duchi, E. Hazan, and Y. Singer, "Adaptive subgradient methods for online learning and stochastic optimization," *Journal of Machine Learning Research*, vol. 12, no. Jul, pp. 2121–2159, 2011.

[13] T. Tieleman and G. Hinton, "Rmsprop," *COURSERA: Lecture*, 2012.

[14] M. D. Zeiler, "Adadelta: an adaptive learning rate method," *arXiv preprint arXiv:1212.5701*, 2012.

[15] J. S. Bergstra, R. Bardenet, Y. Bengio, and B. Kégl, "Algorithms for hyper-parameter optimization," in *Advances in Neural Information Processing Systems*, 2011, pp. 2546–2554.

[16] J. Bergstra, D. Yamins, and D. D. Cox, "Making a science of model search: Hyperparameter optimization in hundreds of dimensions for vision architectures." *ICML (1)*, vol. 28, pp. 115–123, 2013.

[17] P. Lerman, "Fitting segmented regression models by grid search," *Applied Statistics*, pp. 77–84, 1980.

[18] Y. Bengio, Y. LeCun *et al.*, "Scaling learning algorithms towards ai," *Large-scale kernel machines*, vol. 34, no. 5, pp. 1–41, 2007.

[19] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," *Nature*, vol. 521, no. 7553, pp. 436–444, 2015.

[20] Y. LeCun, B. Boser, J. S. Denker, D. Henderson, R. E. Howard, W. Hubbard, and L. D. Jackel, "Backpropagation applied to handwritten zip code recognition," *Neural Computation*, vol. 1, no. 4, pp. 541–551, 1989.

[21] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.

[22] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, "Going deeper with convolutions," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 1–9.

[23] M. D. Zeiler and R. Fergus, "Visualizing and understanding convolutional networks," in *European Conference on Computer Vision*. Springer, 2014, pp. 818–833.

[24] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.

[25] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," *arXiv preprint arXiv:1512.03385*, 2015.

[26] C. W. Tsai, C. F. Lai, H. C. Chao, and A. V. Vasilakos, "Big data analytics: a survey," *Journal of Big Data*, vol. 2, no. 1, p. 21, 2015.

[27] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *Computer Vision and Pattern Recognition, 2009. CVPR 2009. IEEE Conference on*. IEEE, 2009, pp. 248–255.

[28] G. E. Hinton and R. R. Salakhutdinov, "Reducing the dimensionality of data with neural networks," *Science*, vol. 313, no. 5786, pp. 504–507, 2006.

[29] H. Bourlard and Y. Kamp, "Auto-association by multilayer perceptrons and singular value decomposition," *Biological Cybernetics*, vol. 59, no. 4-5, pp. 291–294, 1988.

[30] G. E. Hinton and R. S. Zemel, "Autoencoders, minimum description length, and helmholtz free energy," *Advances in Neural Information Processing Systems*, pp. 3–3, 1994.

[31] B. A. Olshausen and D. J. Field, "Sparse coding with an overcomplete basis set: A strategy employed by v1?" *Vision Research*, vol. 37, no. 23, pp. 3311–3325, 1997.

[32] H. Lee, C. Ekanadham, and A. Y. Ng, "Sparse deep belief net model for visual area v2," in *Advances in Neural Information Processing Systems*, 2008, pp. 873–880.

[33] Y. Bengio, P. Lamblin, D. Popovici, H. Larochelle *et al.*, "Greedy layer-wise training of deep networks," *Advances in Neural Information Processing Systems*, vol. 19, p. 153, 2007.

[34] P. Vincent, H. Larochelle, Y. Bengio, and P.-A. Manzagol, "Extracting and composing robust features with denoising autoencoders," in *Proceedings of the 25th International Conference on Machine Learning*. ACM, 2008, pp. 1096–1103.

[35] S. Rifai, P. Vincent, X. Muller, X. Glorot, and Y. Bengio, "Contractive auto-encoders: Explicit invariance during feature extraction," in *Proceedings of the 28th International Conference on Machine Learning*, 2011, pp. 833–840.

[36] X. Yao, "Evolving artificial neural networks," *Proceedings of the IEEE*, vol. 87, no. 9, pp. 1423–1447, 1999.

[37] D. Whitley, T. Starkweather, and C. Bogart, "Genetic algorithms and

neural networks: Optimizing connections and connectivity," *Parallel Computing*, vol. 14, no. 3, pp. 347–361, 1990.

[38] L. D. Whitley, "The genitor algorithm and selection pressure: Why rank-based allocation of reproductive trials is best." in *ICGA*, vol. 89, 1989, pp. 116–123.

[39] D. J. Montana and L. Davis, "Training feedforward neural networks using genetic algorithms." in *International Joint Conference on Artificial Intelligence*, vol. 89, 1989, pp. 762–767.

[40] S. F. Christian and C. Lebiere, "The cascade-correlation learning architecture," in *Advances in Neural Information Processing Systems 2*. Citeseer, 1990.

[41] M. Frean, "The upstart algorithm: A method for constructing and training feedforward neural networks," *Neural Computation*, vol. 2, no. 2, pp. 198–209, 1990.

[42] J. Sietsma and R. J. Dow, "Creating artificial neural networks that generalize," *Neural Networks*, vol. 4, no. 1, pp. 67–79, 1991.

[43] R. Zi-wu and S. Ye, "Improvement of real-valued genetic algorithm and performance study [j]," *Acta Electronica Sinica*, vol. 2, p. 017, 2007.

[44] K. O. Stanley and R. Miikkulainen, "Evolving neural networks through augmenting topologies," *Evolutionary Computation*, vol. 10, no. 2, pp. 99–127, 2002.

[45] K. O. Stanley, "Compositional pattern producing networks: A novel abstraction of development," *Genetic Programming and Evolvable Machines*, vol. 8, no. 2, pp. 131–162, 2007.

[46] M. Gong, J. Liu, H. Li, Q. Cai, and L. Su, "A multiobjective sparse feature learning model for deep neural networks," *IEEE transactions on neural networks and learning systems*, vol. 26, no. 12, pp. 3263–3277, 2015.

[47] R. Storn and K. Price, "Differential evolution a simple and efficient heuristic for global optimization over continuous spaces," *Journal of Global Optimization*, vol. 11, no. 4, pp. 341–359, 1997.

[48] J. Liu, M. Gong, Q. Miao, X. Wang, and H. Li, "Structure learning for deep neural networks based on multiobjective optimization," *IEEE Transactions on Neural Networks and Learning Systems*, 2017.

[49] E. Real, S. Moore, A. Selle, S. Saxena, Y. L. Suematsu, Q. Le, and A. Kurakin, "Large-scale evolution of image classifiers," *arXiv preprint arXiv:1703.01041*, 2017.

[50] P. Smolensky, "Information processing in dynamical systems: Foundations of harmony theory," DTIC Document, Tech. Rep., 1986.

[51] G. E. Hinton, S. Osindero, and Y.-W. Teh, "A fast learning algorithm for deep belief nets," *Neural Computation*, vol. 18, no. 7, pp. 1527–1554, 2006.

[52] K. O. Stanley, "Exploiting regularity without development," in *Proceedings of the AAAI Fall Symposium on Developmental Systems*. AAAI Press Menlo Park, CA, 2006, p. 37.

[53] J. Yang, A. F. Frangi, J.-y. Yang, D. Zhang, and Z. Jin, "Kpca plus lda: a complete kernel fisher discriminant framework for feature extraction and recognition," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 27, no. 2, pp. 230–244, 2005.

[54] X. Glorot, A. Bordes, and Y. Bengio, "Deep sparse rectifier neural networks." in *Aistats*, vol. 15, no. 106, 2011, p. 275.

[55] C. Cortes and V. Vapnik, "Support-vector networks," *Machine Learning*, vol. 20, no. 3, pp. 273–297, 1995.

[56] A. L. Maas, A. Y. Hannun, and A. Y. Ng, "Rectifier nonlinearities improve neural network acoustic models," in *Proc. ICML*, vol. 30, no. 1, 2013.

[57] Q. Zhao, D. Zhang, and H. Lu, "A direct evolutionary feature extraction algorithm for classifying high dimensional data," in *Proceedings of the National Conference on Artificial Intelligence*, vol. 21, no. 1. Menlo Park, CA; Cambridge, MA; London; AAAI Press; MIT Press; 1999, 2006, p. 561.

[58] V. N. Vapnik, "The nature of statistical learning theory," *IEEE Transactions on Neural Networks*, vol. 8, no. 6, pp. 1564–1564, 1997.

[59] V. Vapnik, *Statistical learning theory*. DBLP, 2010.

[60] H. Larochelle, D. Erhan, A. Courville, J. Bergstra, and Y. Bengio, "An empirical evaluation of deep architectures on problems with many factors of variation," in *Proceedings of the 24th international conference on Machine learning*. ACM, 2007, pp. 473–480.

[61] A. Krizhevsky and G. Hinton, "Learning multiple layers of features from tiny images," 2009.

[62] J. Engel, "Polytomous logistic regression," *Statistica Neerlandica*, vol. 42, no. 4, pp. 233–252, 1988.

[63] G. Hinton, "A practical guide to training restricted boltzmann machines," *Momentum*, vol. 9, no. 1, p. 926, 2010.

[64] R. G. Steel, D. JH Dickey *et al.*, *Principles and procedures of statistics a biometrical approach*. WCB/McGraw-Hill, 1997, no. 519.5 S813 1997.

[65] M. A. Carreira-Perpinan and G. E. Hinton, "On contrastive divergence learning." in *AISTATS*, vol. 10. Citeseer, 2005, pp. 33–40.

[66] X. Glorot and Y. Bengio, "Understanding the difficulty of training deep feedforward neural networks." in *Aistats*, vol. 9, 2010, pp. 249–256.

[67] K. Deb, *Multi-objective optimization using evolutionary algorithms*. John Wiley & Sons, 2001, vol. 16.

[68] D. Erhan, Y. Bengio, A. Courville, and P. Vincent, "Visualizing higher-layer features of a deep network," *University of Montreal*, vol. 1341, 2009.

**Yanan Sun** (S'15-M'18) received a Ph.D. degree in engineering from the Sichuan University, Chengdu, China, in 2017. From 2015.08-2017.02, he is a jointly Ph.D. student financed by the China Scholarship Council in the School of Electrical and Computer Engineering, Oklahoma State University (OSU), USA. He is currently a Postdoc Research Fellow in the School of Engineering and Computer Science, Victoria University of Wellington, Wellington, New Zealand. His research topics are many-objective optimization and deep learning.

**Gary G. Yen** (S'87-M'88-SM'97-F'09) received a Ph.D. degree in electrical and computer engineering from the University of Notre Dame in 1992. Currently he is a Regents Professor in the School of Electrical and Computer Engineering, Oklahoma State University (OSU). Before joined OSU in 1997, he was with the Structure Control Division, U.S. Air Force Research Laboratory in Albuquerque. His research interest includes intelligent control, computational intelligence, conditional health monitoring, signal processing and their industrial/defense applications.

Dr. Yen was an associate editor of the *IEEE Control Systems Magazine, IEEE Transactions on Control Systems Technology*, *Automatica, Mechantronics*, *IEEE Transactions on Systems, Man and Cybernetics, Parts A and B* and I*EEE Transactions on Neural Networks*. He is currently serving as an associate editor for the *IEEE Transactions on Evolutionary Computation* and the *IEEE Transactions on Cybernetics*. He served as the General Chair for the *2003 IEEE International Symposium on Intelligent Control* held in Houston, TX and *2006 IEEE World Congress on Computational Intelligence* held in Vancouver, Canada. Dr. Yen served as Vice President for the Technical Activities in 2005-2006 and then President in 2010-2011 of the IEEE Computational intelligence Society. He was the founding editor-in-chief of the *IEEE Computational Intelligence Magazine*, 2006-2009. In 2011, he received Andrew P Sage Best Transactions Paper award from *IEEE Systems, Man and Cybernetics Society* and in 2014, he received Meritorious Service award from *IEEE Computational Intelligence Society*.

**Zhang Yi** (F'16) received a Ph.D. degree in mathematics from the Institute of Mathematics, The Chinese Academy of Science, Beijing, China, in 1994. Currently, he is a Professor at the Machine Intelligence Laboratory, College of Computer Science, Sichuan University, Chengdu, China. He is the co-author of three books: *Convergence Analysis of Recurrent Neural Networks* (Kluwer Academic Publishers, 2004), *Neural Networks: Computational Models and Applications* (Springer, 2007), and *Subspace Learning of Neural Networks* (CRC Press, 2010). He was an Associate Editor of *IEEE Transactions on Neural Networks and Learning Systems* (2009 2012), and He is an Associate Editor of *IEEE Transactions on Cybernetics* (2014 ). His current research interests include Neural Networks and Big Data. He is a fellow of IEEE.

## SUPPLEMENTAL MATERIALS

### VI. AN EXAMPLE OF ENCODING STRATEGY

In this section, the main steps of the encoding scheme are given, and then an illustrating example based on these steps is provided. For the convenience of the development, assuming the input data is with the dimension of $n \times n$. The main steps of this encoding scheme are detailed below:

1) Randomly generate $n$ orthogonal vectors and each vector is $n$-dimensional, these vectors are denoted by $S = [s_1, \cdots, s_n]$;
2) Randomly generate $n$ real numbers that are denoted by $b = [b_1, \cdots, b_n]$;
3) Compute $a_1 = b_1 \times s_1 + \cdots + b_n \times s_n$;
4) Compute the bases $a_2, \cdots, a_n$ of the null space of $a_1$;
5) Initialize a chromosome with a length of $2n + 1$;
6) Copy the values of $b_1, \cdots, b_n$ into the first position to the $n$-th position of this chromosome (i.e., they are used to denote the value of $a_1$);
7) Randomly generate $n - 1$ numbers from $\{0, 1\}$, copy them to the $(n + 1)$-th to $(2n - 1)$-th position of this chromosome (they are used to represent whether the corresponding basis from $\{a_2, \cdots, a_n\}$ would be selected as the subspace or not);
8) Randomly generate a number from $\{1, 2, 3\}$, convert it to the binary format with the length of 2, and copy it to the $2n$-th to the $(2n + 1)$-th position of this chromosome (they are used to denote the index of the chosen activation function).

Supposed that $n$ is equal to 5, an example based on the description above is given as follow.

1) Randomly generate vectors $S = [s_1, s_2, s_3, s_4, s_5]$ where $S \in R^{5 \times 5}$;

$$S = \begin{bmatrix} -0.4861 & -0.6498 & 0.2718 & 0.1572 & 0.4927 \\ -0.4617 & -0.2830 & -0.1205 & -0.6073 & -0.5686 \\ -0.4438 & 0.3468 & 0.5339 & 0.4669 & -0.4240 \\ -0.4721 & 0.6142 & -0.0597 & -0.3831 & 0.4995 \\ -0.3614 & 0.0075 & -0.7893 & 0.4916 & -0.0681 \end{bmatrix}$$

2) Randomly generate 5 numbers stored into $b$;

$$b = \begin{bmatrix} 0.7303 & 0.4886 & 0.5785 & 0.2373 & 0.4588 \end{bmatrix}$$

3) Compute the linear combination $a_1$ of $S$ and $b$;

$$a_1 = \begin{bmatrix} 1.2642 & -1.4589 & -1.0880 & 1.2815 & -0.1746 \end{bmatrix}^T$$

4) Compute the bases ($a_2, a_3, a_4$, and $a_5$) of the null space of $a_1$;

$$a_2 = \begin{bmatrix} -0.8108 & 0.4590 & 0.0400 & 0.0336 & -0.3596 \end{bmatrix}^T$$

$$a_3 = \begin{bmatrix} 0.0600 & 0.0400 & 0.9970 & -0.0025 & 0.0266 \end{bmatrix}^T$$

$$a_4 = \begin{bmatrix} 0.0504 & 0.0336 & -0.0025 & 0.9979 & 0.0224 \end{bmatrix}^T$$

$$a_5 = \begin{bmatrix} -0.5388 & -0.3596 & 0.0266 & 0.0224 & 0.7611 \end{bmatrix}^T$$

5) Initialize a chromosome with a length of 11;

| Null | Null | Null | Null | Null | Null | Null | Null | Null | Null | Null |
|---|---|---|---|---|---|---|---|---|---|---|

6) Copy the elements in $b$ into the $1^{\text{rd}}$—$5^{\text{th}}$ positions.

| 0.7303 | 0.4886 | 0.5785 | 0.2373 | 0.4588 | Null | Null | Null | Null | Null | Null |
|---|---|---|---|---|---|---|---|---|---|---|

| 0.7303 | 0.4886 | 0.5785 | 0.2373 | 0.4588 | 0 | 1 | 1 | 0 | Null | Null |
|---|---|---|---|---|---|---|---|---|---|---|

7) Copy $\{0, 1, 1, 0\}$ that are 4 randomly generated numbers from $\{0, 1\}$ into the $6^{\text{th}}$—$9^{\text{th}}$ positions.
8) A randomly generated number 2 from $\{1, 2, 3\}$, convert 2 to its binary form 10, and copy 10 to the $10^{\text{th}}$—$11^{\text{th}}$ positions.

| 0.7303 | 0.4886 | 0.5785 | 0.2373 | 0.4588 | 0 | 1 | 1 | 0 | 1 | 0 |
|---|---|---|---|---|---|---|---|---|---|---|

### VII. STEPS OF CROSSOVER AND MUTATION

In the following, we will give the steps of the crossover operation on two parent solutions in the proposed algorithm.

1) Assuming that the parent solutions are denoted by $ind_1$ and $ind_2$;
2) Randomly generate a number from $[0, 1]$, and if the generated number is below the predefined crossover probability, perform steps 3)-5), otherwise go to step 6);
3) Calculate the length (denote by $l_1$) of the first two parts, and the length (denoted by $l_2$) of the third part of the individual (the information of these three parts can be seen in Section IV-D of the manuscript);
4) Randomly generate an integer number (denoted by $i_1$) from $[1, l_1]$, and another integer number (denoted by $i_2$) from $[1, l_2]$;
5) Exchange the first two parts of $ind_1$ and $ind_2$ on the position $i_1$ with the one point crossover operator, and exchange the third part of $ind_1$ and $ind_2$ on the position $i_2$ with the one point crossover operator;
6) Return $ind_1$ and $ind_2$.

Next, we will give the steps of the mutation operation on the individual $ind_1$.

1) Randomly generate a number from $[0, 1]$, and if this number is below the predefined mutation probability, performed steps 2), otherwise go to step 3);
2) For each position in $ind_1$, randomly generate a number from $[0, 1]$, if this number is less than $0.5$, perform the polynomial mutation on the current position, otherwise skip to next position.
3) Return $ind_1$.