# Beyond Gradient: Subspace Rotation Algorithm

Ci Lin, Tet Hin Yeap, Iluju Kiringa

**Abstract**

This study introduces the Subspace Rotation Algorithm (SRA), an innovative gradient-free method designed to discover the global optimal weight matrix. The SRA consists of two fundamental algorithms: the Left Subspace Rotation Algorithm (LSRA) and the Right Subspace Rotation Algorithm (RSRA). The combination of LSRA and RSRA, in two formats, LSRA-RSRA and RSRA-LSRA, can harness the advantages of both individual algorithms, resulting in enhanced performance.

Our observations reveal that shallow and wide Multilayer Perceptrons (MLP), trained using RSRA-LSRA, achieve higher training accuracy compared to the Backpropagation (BP) algorithm. Moreover, when combining SRA with the BP algorithm, a remarkable impact on training MLP models is observed. Our experiments demonstrate that the BP algorithm may become trapped in local optima, while RSRA-LSRA, though capable of escaping local optima, may not fully realize its potential when the number of hidden nodes is limited. The synergy of RSRA-LSRA and the BP algorithm allows for the utilization of the advantages from both approaches, achieving optimal MLP performance with fewer hidden nodes.

**Keywords:** Subspace Rotation Algorithm, Backpropagation, Global Optimization, Multilayer Perceptrons

## 1. Introduction

Deep learning has brought about a revolution in various fields, including computer vision [1], natural language processing [2], robotics [3], healthcare [4], and finance [5]. Training deep neural network models, such as the Multiple Layer Perception (MLP) [6], Convolutional Neural Network (CNN) [7], and Recurrent Neural Network (RNN) [8], involves solving optimization problems that aim to minimize objective functions based on the connection weights of the networks [9]. Unfortunately, these optimization problems come with inherent limitations, including non-linearity, non-convexity, the presence of numerous local minima, and broad regions linked to narrow ones.

The motivation behind this research is to address the challenges associated with gradient-based training algorithms. Many state-of-the-art methods rely on the error backpropagation (BP) algorithm, which utilizes gradients to update parameters. However, this approach encounters several issues, including the problems of vanishing and exploding gradients, difficulties in handling non-differentiable nonlinearities, and managing parallel weight updates across layers [10]. As a result, researchers have explored alternative training algorithms, including the Hilbert-Schmidt Independence Criterion (HSIC) bottleneck-based approach [11], the online alternating minimization approach [12], simultaneous perturbation stochastic approximation (SPSA) [13], and using the energy function as the objective function [14].

However, these newly developed methods also face challenges when it comes to finding global optimization solutions for non-convex problems. Drawing inspiration from the success of signal processing techniques [15, 16] in tasks such as dimensionality reduction [17], noise reduction [18], and feature extraction [19] and integrating the concepts of signal subspace [20], noise subspace [21], and the subspace rotation approach [22], this paper presents a gradient-free approach for finding global optimization solutions in neural networks. This approach allows for the differentiation and separation of the underlying data structure from unwanted variations or noise.

In this work, our contributions can be summarized as follows:

---

*corresponding_author@example.ca

(1) We introduce the Subspace Rotation Algorithm (SRA), a gradient-free approach designed for training neural network models. SRA addresses common issues associated with BP algorithms, such as gradient vanishing, explosion, or slow convergence, while seeking global optimal solutions.

(2) We conduct an in-depth exploration of the characteristics of the Left Subspace Rotation Algorithm (LSRA) and the Right Subspace Rotation Algorithm (RSRA). Our simulation results reveal that LSRA tends to exhibit greater stability when training MLP as the classifier. Additionally, by combining these two individual algorithms, LSRA-RSRA and RSRA-LSRA, we harness the advantages of both. Notably, our experimental findings demonstrate that, in general, RSRA-LSRA outperforms LSRA-RSRA.

(3) We achieve state-of-the-art performance for MLP models on the MNIST and Fashion-MNIST datasets by combining the RSRA-LSRA with the BP algorithm.

The rest paper is structured as follows. Section 1 initiates our discussion by presenting an introduction to deep learning models and their associated BP algorithms. Additionally, we outline the motivation behind this research and highlight the specific contributions of our work. Section 2 delves into the core of our research, introducing two distinct SRA: the LSRA and the RSRA. Moreover, we provide a practical example to demonstrate the effectiveness of these algorithms. In Section 3, we meticulously design and conduct a series of four experiments aimed at evaluating the efficacy and efficiency of the SRA. This comprehensive assessment encompasses LSRA, RSRA, LSRA-RSRA, and RSRA-LSRA approaches, along with the BP algorithm. Our experiments involve the training of a variety of MLP models using diverse datasets, including MNIST, Fashion-MNIST, CIFAR10, and CIFAR100. Section 4, offers our concluding remarks. We summarize our key observations, noting that LSRA exhibits greater stability than RSRA in most scenarios. We also highlight the capacity of RSRA-LSRA and LSRA-RSRA to harness the strengths of individual algorithms, with RSRA-LSRA generally outperforming LSRA-RSRA. Most notably, we conclude that combining the RSRA-LSRA with the BP algorithm yields exceptional performance in training MLP models, surpassing the individual SRA and BP algorithm.

## 2. Subspace Rotation Algorithm: Discover the Optimal Solution without Gradient

### 2.1. Left Subspace Rotation Algorithm and Right Subspace Rotation Algorithm

The weight matrix of a neural network is connected to two layers of nodes, namely, the input layer of nodes and the output layer of nodes. If the input layer of nodes is used to evaluate and update the weight matrix, it is known as the LSRA. On the other hand, if it is the output layer of nodes that is used to evaluate and update the weight matrix, it is known as the RSRA. However, apart from this difference, these two algorithms are similar to each other. Hence, both these algorithms are discussed together in the following sections.

Algorithm 1 and 2 describe the LSRA and RSRA for MLP, respectively. These algorithms take as input the samples X, labels Y, and the number of hidden layers L in the MLP model, and output the weight matrices in the form of a list, known as Weight_List.

The algorithm initializes the weight list with orthogonal weight matrices and updates each weight matrix using the SRA (LSRA or RSRA) sequentially. In each update, it computes the output of the MLP up to a specific layer using the **ForwardPass** (**LeftForwardPass** or **RightForwardPass**) procedure and computes the reversed output of the MLP from the last layer up to the specific layer using the **BackwardPass** (**LeftBackwardPass** or **RightBackwardPass**) procedure. The algorithm executes the pseudoinverse operation on forward output and multiplies the obtained result with the backward output to compute a distance matrix. This matrix is decomposed into U, D, and V matrices using the Singular

---

**Algorithm 1** Left Subspace Rotation Algorithm For MLP

---

**Input**: Samples X, Labels Y, Number of Layers L
**Output**: weight matrix list Weight_List
**Initialization**: initialize orthogonal weight matrix and added it into the Weight_List one by one
**for** index = 0 **to** L **do**
  f_output = LeftForwardPass(X, index)
  b_output = LeftBackwardPass(Y, index, L)
  dist_matrix = PINV(f_output) × b_output
  U, D, V = SVD(dist_matrix)
  Weight_List[index] = U × V × Weight_LIST[index]
**end for**
**return** Weight_List

**LeftForwardPass(X, index)**
**Input**: Samples X, Layer Index: index
**Output**: Output of the Specific Layer
**if** index = 0 **then**
  f_output = X
**end if**
**for** ix = 0 **to** index **do**
  f_output = tanh( X × WeightList[ix])
  X = f_output
**end for**
**return** f_output

**LeftBackwardPass(X, index, L)**
**Input**: Labels Y, Layer Index: index, Length of Weight List: L
**Output**: Reversed Output of the Specific Layer
**for** ix = L **to** L - index **do**
  b_output = Y × WeightList[ix].T
  Y = atanh(b_output)
**end for**
**return** b_output

---

Value Decomposition Algorithm, where U is an orthogonal matrix that represents the left singular vectors of the distance matrix, D is a diagonal matrix that contains the singular values of the distance matrix, and V is the conjugate transpose of an orthogonal matrix that represents the right singular vectors of the distance matrix. Then, the U and V matrices are used to update the weight matrix for the specific layer.

Additionally, the **ForwardPass** (**LeftForwardPass** or **RightForwardPass**) procedure takes as input the samples X and the index of the layer and computes the output of the MLP up to the indexed layer from the bottom to the top. The **BackwardPass** (**LeftBackwardPass** or **RightBackwardPass**) procedure takes as input the labels Y, the layer index, and the length of the weight list L, and computes the reversed output of the MLP up to the indexed layer from the top to the bottom. It is essential to note that the **tanh** and **atanh** are inverse functions. Moreover, the domain of **atanh** is (-1, 1), and during the calculation of backward values, it is possible that the backward value cannot be

---

**Algorithm 2** Right Subspace Rotation Algorithm For MLP

---

**Input**: Samples X, Labels Y, Number of Layers L
**Output**: weight matrix list Weight_List
**Initialization**: initialize orthogonal weight matrix and added it into the Weight_List one by one
**for** index = 0 **to** L **do**
   f_output = RightForwardPass(X, index)
   b_output = RightBackwardPass(Y, index, L)
   dist_matrix = PINV(f_output) × b_output
   U, D, V = SVD(dist_matrix)
   Weight_List[index] = Weight_LIST[index] × U × V
**end for**
**return** Weight_List

 

**RightForwardPass(X, index)**
**Input**: Samples X, Layer Index: index
**Output**: Output of the Specific Layer
**for** ix = 0 **to** index+1 **do**
   f_output = X × WeightList[ix]
   X = tanh(f_output)
**end for**
**return** f_output

 

**RightBackwardPass(X, index, L)**
**Input**: Labels Y, Layer Index: index, Length of Weight List: L
**Output**: Reversed Output of the Specific Layer
**if** index = L -1 **then**
   b_output = Y
**end if**
**for** ix = L-1 **to** L - index **do**
   b_output = atanh(Y × WeightList[ix].T)
   Y = b_output
**end for**
**return** b_output

---

calculated because the input is outside the domain of (-1, 1), which correspond to $-\infty$ or $+\infty$. In such situation, we can clip the output value of **atanh** in the range (-100, 100).

### 2.2. **Mathematical Analysis for the Subspace Rotation Algorithm**

Assuming there are $n$ samples, each represented as $x_i \in \mathbb{R}^p$, and all $n$ samples can be combined into a matrix $X \in \mathbb{R}^{p \times n}$. Let initial weight matrix be $W \in \mathbb{R}^{p \times r}$. If the Frobenius norm is used to calculate the distance between the original matrix and the transformed matrix, the issue can be reformulated as an optimization problem:

$$\min_{W \in O(p,r)} \|X - WW^T X\|_F \tag{2.1}$$

Solution: Without loss of generality, assuming $p \geq r$, $B = WW^T X$. According to the rank properties of matrices, we can obtain $rank(B) = r$. Equation 2.1 can be rewritten in the following form:

$$\min_{rank(B) \leq r} ||X - B||_F \tag{2.2}$$

According to the Eckart–Young–Mirsky theorem, let $X = UDV^T \in \mathbb{R}^{p \times n}, \quad n \geq p$. If we aim to find a matrix B with rank $r$, where $r < p$, we can decompose U, D, and V in the following way:

$$U = [U_1, U_2], D = \begin{bmatrix} D_1 & 0 \\ 0 & D_2 \end{bmatrix}, V = [V_1, V_2]$$

where $U_1$ is $m \times r$, $D_1$ is $r \times r$, and $V_1$ is $n \times r$. Then the rank-r matrix B, obtained from the truncated singular value decomposition is:

$$B = U_1 D_1 V_1^T \tag{2.3}$$

Now, let us return to the original target function, where we can make an observation:

$$B = WW^T X = WW^T UDV^T = U_1 D_1 V_1^T \tag{2.4}$$

In this case, the above equation will be successful only when $W = U_1$. Thus, $U_1$ is the analytic solution we are looking for.

## 3. Experiment and Analysis

### 3.1. Sample Preparation

Before proceeding to the simulation part, it is necessary to describe the approach of preparing the samples. For the samples from the MNIST dataset, the Fashion-MNIST dataset, the CIFAR10 dataset, and CIFAR100 dataset, the RGB color, which is ranged from 0 to 255, will be scaled to the range of -1 to 1.

### 3.2. Experiment Design

Four experiments have been purposefully designed to assess the efficacy and effectiveness of the following algorithms: LSRA, RSRA, LSRA-RSRA, RSRA-LSRA, and BP (with different optimizers).

(1) **Robustness of SRA**: The first experiment involves designing three MLP models with different architectures, which are trained using LSRA, RSRA, LSRA-RSRA, and RSRA-LSRA on the MNIST dataset. The objective of this experiment is to identify the most robust and optimal SRA among LSRA, RSRA, LSRA-RSRA, and RSRA-LSRA.

(2) **Pros and Cons of RSRA-LSRA and BP**: The second experiment features a two-layer MLP architecture (with 10000 hidden nodes), trained using RSRA-LSRA and BP algorithms on four different datasets: MNIST, Fashion-MNIST, CIFAR10, and CIFAR100. For the BP algorithm, we utilize the SGD optimizer with a learning rate of $5e-3$ and a momentum of 0.2, implementing early stopping to prevent overfitting. This experiment aims to uncover the strengths and weaknesses of RSRA-LSRA and BP algorithms in training shallow and wide MLP models.

(3) **Model Complexity and Overfitting**: The third experiment designs four two-layer MLP models with varying numbers of hidden nodes, which are trained using RSRA-LSRA and BP algorithms on the CIFAR100 dataset. For the BP algorithm, we employ the SGD optimizer with a learning rate of $1e-2$ and a momentum of 0.9. Early stopping is not implemented, and the training epoch is fixed at 500. The purpose of this experiment is to explore the relationship between model complexity and overfitting for RSRA-LSRA and BP algorithms.

(4) **Combining SRA and BP Algorithm**: The fourth experiment introduces a two-layer MLP model (with 500 hidden nodes), trained using RSRA-LSRA, BP, and COM (combining RSRA-LSRA and BP algorithms) on two datasets: MNIST and Fashion-MNIST. For the BP algorithm, we employ the Adam optimizer with a learning rate of $1e - 3$, without early stopping. The aim of this experiment is to compare the performance of RSRA-LSRA, BP, and COM algorithms.

### 3.2.1. Explore LSRA, RSRA and the combination of LRSA and RSRA

*Table 1.* The Accuracy of MLP Trained by LSRA, RSRA, LSRA-RSRA, and RSRA-LSRA

| | LSRA | | RSRA | | LSRA-RSRA | | RSRA-LSRA | |
|---|---|---|---|---|---|---|---|---|
| | Train Acc | Test Acc | Train Acc | Test Acc | Train Acc | Test Acc | Train Acc | Test Acc |
| MLP[1] | 98.01 | 96.92 | 44.87 | 45.58 | 98.27 | 97.18 | 98.32 | 97.16 |
| MLP[2] | 86.91 | 86.78 | 43.20 | 43.11 | 83.55 | 84.05 | 87.81 | 87.66 |
| MLP[3] | 94.47 | 94.33 | 94.56 | 94.20 | 93.97 | 93.57 | 94.64 | 94.16 |

[1] The architecture of MLP is 784-5000-10, with a total of $3.97 \times 10^6$ parameters.
[2] The architecture of MLP is 784-600-10-600-10, with a total of $4.88 \times 10^5$ parameters.
[3] The architecture of MLP is 784-1200-1400-1200-1000-800-600-400-200-10, with a total of $7.10 \times 10^6$ parameters.

SRA is a newly developed algorithm that has not been fully explored. It essentially comprises two variants: the LSRA and the RSRA. These two algorithms can be used individually or in combination. Before comparing them with the BP algorithm, it is necessary to discuss how to achieve optimal performance with SRA.

In this section, our aim is to evaluate the effectiveness of LSRA, RSRA, and the combination of LSRA and RSRA. To do this, we intentionally designed three MLP models and trained them on the MNIST dataset using LSRA, RSRA, LSRA-RSRA, and RSRA-LSRA, respectively.

The first MLP model (784-5000-10) is a shallow model with a wider hidden layer. It is used to test whether SRA can perform better with a shallow model. The second MLP model (784-600-10-600-10) is slightly deeper, with a bottleneck in the middle hidden layer, to determine whether the SRA can handle non-standard MLP architectures effectively. The third MLP model (784-1200-1400-1200-1000-800-600-400-200-10) is a deep neural network with smoothly changing hidden layers. It aims to assess the performance of SRA on deep neural networks.

The experimental results are quite inspiring. For MLP with smoothly changing hidden layers, where the number of hidden nodes gradually changes, both LSRA and RSRA exhibit almost the same performance, as shown in the third row of Table 1. However, for MLP with rapidly changing hidden layers, LSRA outperforms RSRA, as demonstrated in the first and second rows of Table 1.

Furthermore, when MLP are trained using LSRA-RSRA (alternating training with LSRA and RSRA until all layers are updated) or RSRA-LSRA (alternating training with RSRA and LSRA until all layers are updated), they generally exhibit mixed performance, as shown in Table 1. Essentially, RSRA-LSRA and LSRA-RSRA combine the strengths of LSRA and RSRA. Therefore, their effectiveness is similar to LSRA, which proves to be a more stable algorithm than RSRA in almost all cases. Notably, RSRA-LSRA subtly outperforms LSRA-RSRA, as indicated in Table 1.

### 3.2.2. Subspace Rotation Algorithm Versus Backpropagation Algorithm

As discussed in Section 3.2.1, the SRA performed surprisingly well on shallow and wide MLP models, which runs counter to conventional wisdom in the field of deep learning. Now,
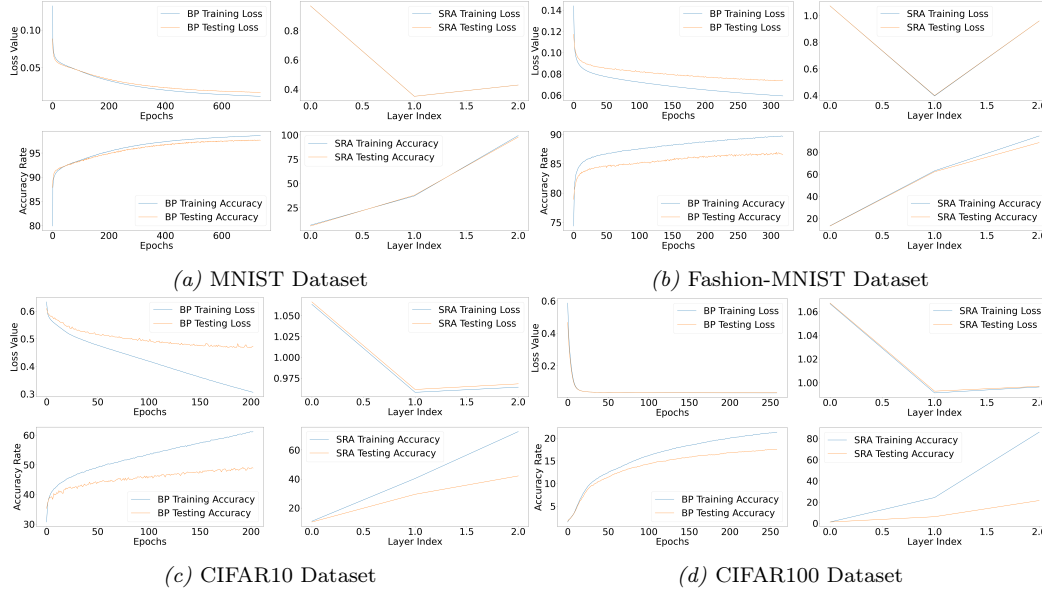
*Figure 1.* Loss Value and Accuracy of the Model on Different Datasets

*Table 2.* The Performance of MLP Trained on MNIST, Fashion-MNIST, CIFAR10, and CIFAR100 Dataset by RSRA-LSRA and BP

| | RSRA-LSRA | | | | BP | | | |
|---|---|---|---|---|---|---|---|---|
| | Train Acc | Train Loss | Test Acc | Test Loss | Train Acc | Train Loss | Test Acc | Test Loss |
| MNIST [1] | 99.345 | 0.4291 | 97.82 | 0.4291 | 98.59 | 0.0119 | 97.64 | 0.0172 |
| Fashion MNIST [1] | 94.44 | 0.9592 | 88.47 | 0.9592 | 89.68 | 0.0596 | 86.86 | 0.0736 |
| CIFAR10 [2] | 72.68 | 0.9641 | 42.11 | 0.9681 | 67.55 | 0.1705 | 49.53 | 0.2365 |
| CIFAR100 [2] | 86.00 | 0.9962 | 21.48 | 0.9968 | 23.02 | 0.0319 | 18.26 | 0.0329 |

[1] The architecture of MLP is 784-10000-10, with a total of $7.94 \times 10^6$ parameters. The activation function is tanh function.

[2] The architecture of MLP is 3072-10000-10, with a total of $3.082 \times 10^7$ parameters. The activation function is tanh function.

we aim to compare RSRA-LSRA with BP algorithm for training shallow and wide MLP models. These models have an architecture of $X - 10000 - Y$, where X represents the input nodes of the MLP, and Y represents the output nodes of the MLP.

All four datasets, including MNIST, Fashion-MNIST, CIFAR10, and CIFAR100, are used to train the MLP models. From a model accuracy perspective, for MNIST, Fashion-MNIST, and CIFAR100 datasets, both the training and testing accuracy of the RSRA-LSRA are better than that of the BP algorithm. However, for the CIFAR10 dataset, the RSRA-LSRA achieves higher training accuracy but fails to outperform the BP algorithm in terms of testing accuracy. The explanation for this behavior is related to the signal subspace and noise subspace. In relatively clean datasets like MNIST and Fashion-MNIST, the RSRA-LSRA can effectively extract the major signal subspace from the training data, resulting in good performance on the testing dataset. In other words, the signal subspace from the training data has a significant overlap with the signal subspace from the testing data. However, for CIFAR10 and CIFAR100 datasets, which contain a considerable amount of noise, the subspace extracted by the RSRA-LSRA from the training data also includes a substantial amount of noise subspace. Consequently, the MLP models do not perform well

on the testing dataset, and for CIFAR10 and CIFAR100, there is a notable drop in accuracy from training to testing.

Although CIFAR10 and CIFAR100 share some similarities, they each have distinct characteristics. As shown in Table 2, in the case of CIFAR10, with 10 output nodes, the SRA struggles to accommodate the noise subspace, leading to better training accuracy (72.68%) than the BP algorithm (67.55%), but lower testing accuracy (42.11%) compared to the testing accuracy (49.53%) achieved by BP algorithm. However, in the case of CIFAR100, which has 100 output nodes, the subspace extracted by the RSRA-LSRA is more extensive and can better accommodate the noise subspace. The RSRA-LSRA achieves significantly better training accuracy (86.00%) than the BP algorithm (23.03%). While there is a substantial drop in accuracy from training to testing, the RSRA-LSRA still outperforms the BP algorithm, achieving a testing accuracy of 21.48%, which is higher than 18.25% obtained from the BP algorithm. However, due to the early stopping technique we applied in BP algorithm, we can not conclude that BP algorithm could not perform better than RSRA-LSRA on CIFAR100 dataset.

In terms of loss value, for the BP algorithm, the loss value is closely related to accuracy. However, for the RSRA-LSRA, the loss value achieved by the mean square error (MSE) algorithm does not accurately reflect accuracy. In general, the loss value (MSE value) for the RSRA-LSRA is larger than the loss value for the BP algorithm. For the BP algorithm, if there is a significant difference between training and testing accuracy, this is also reflected in the loss value. For the RSRA-LSRA, regardless of the difference between training and testing accuracy, the training and testing loss values remain similar, as shown in Figure 1. This suggests that the models trained by the RSRA-LSRA can be fine-tuned with the BP algorithm and improve the accuracy further.

### 3.2.3. Model Complexity and Overfitting



*Figure 2.* MLP Models with Different Complexity Trained on CIFAR100 Dataset

For the clean datasets, such as MNIST and Fashion-MNIST, the growing complexity of the model does not seem to have a significant impact on overfitting. However, when it comes to noisy dataset, such as CIFAR10 and CIFAR100 datasets, the presence of noise in the samples amplifies the influence of model complexity on testing performance.

As depicted in the first chart (top-left) of Figure 2 and in Table 3, which represents MLP models with varying complexities trained by RSRA-LSRA. It is observed that as

*Table 3.* The Performance of MLP Models with Different Complexity Trained on CIFAR100 Dataset by RSRA-LSRA and BP

|  | RSRA-LSRA | | | | BP | | | |
|---|---|---|---|---|---|---|---|---|
|  | Train Acc | Train Loss | Test Acc | Test Loss | Train Acc | Train Loss | Test Acc | Test Loss |
| MLP [1] | 47.07 | 0.9892 | 11.57 | 0.9909 | 39.602 | 0.0301 | 20.45 | 0.4029 |
| MLP [2] | 69.65 | 0.9838 | 20.31 | 0.9871 | 54.04 | 0.0247 | 23.02 | 0.4002 |
| MLP [3] | 86.00 | 0.9962 | 21.48 | 0.9968 | 53.45 | 0.0260 | 22.57 | 0.4008 |
| MLP [4] | 96.36 | 0.9975 | 19.75 | 0.9978 | 51.11 | 0.0277 | 21.39 | 0.4071 |

[1] The architecture of MLP is 3072-1000-100, with a total of $3.172 \times 10^6$ parameters. The activation function is tanh function.

[2] The architecture of MLP is 3072-5000-100, with a total of $1.586 \times 10^7$ parameters. The activation function is tanh function.

[3] The architecture of MLP is 3072-10000-100, with a total of $3.172 \times 10^7$ parameters. The activation function is tanh function.

[4] The architecture of MLP is 3072-20000-100, with a total of $6.344 \times 10^7$ parameters. The activation function is tanh function.

the complexity increases, indicated by the number of hidden nodes, the training accuracy of models trained by RSRA-LSRA also rises, ranging from 47.07% to 96.36%. In theory, a similar trend should be expected for models trained by the BP algorithm. However, considering the actual learning rate, within 500 epochs, the MLP with 10000 hidden nodes performed the best among all models trained by the BP algorithm, achieving a training accuracy of 54.04% and a testing accuracy of 23.02%. Nonetheless, when considering only the training accuracy, MLP models trained by RSRA-LSRA can outperform MLP models models trained by BP, as shown in the first and second charts in Figure 2.

However, when we account for testing accuracy, the performance of MLP models trained by BP tends to surpass MLP models trained by RSRA-LSRA. When combining both training and testing accuracy, it becomes evident that RSRA-LSRA demonstrates strong learning capabilities compared to the BP algorithm. Thus, the RSRA-LSRA may be more sensitive to the general noise present in the sample space, while it is expected to be robust to outliers within the samples.

### 3.2.4. **Combining SRA and BP Algorithms**



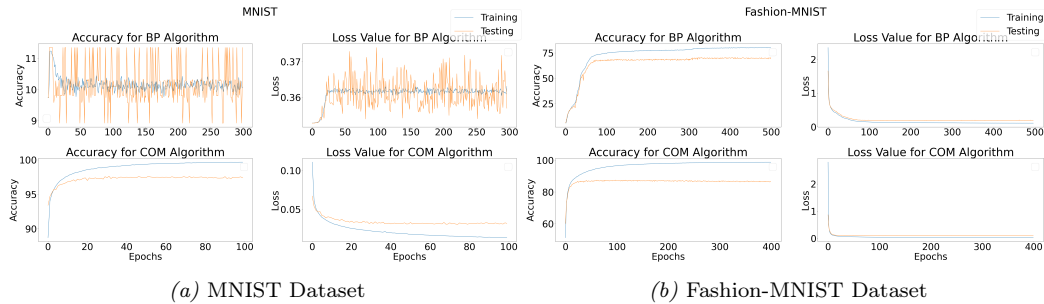*(a)* MNIST Dataset      *(b)* Fashion-MNIST Dataset

*Figure 3.* Loss Value and Accuracy of the Model on MNIST and Fashion-MNIST Dataset (Please note, for BP algorithm, the loss values of the first 3 epochs are omitted since it is too large)

As analyzed in the previous section, the SRA approach can achieve the global optimal subspace for the training dataset. Consequently, when the MLP model is sufficiently complex, it can achieve impressive performance in terms of training accuracy. However, this can also lead to overfitting. In contrast, the BP algorithm focuses on finding local optimal

*Table 4.* The Performance of MLP Trained by Different Algorithms[1]

|  | RSRA-LSRA | | | | BP | | | | COM[2] | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
|  | Train A | Train L[3] | Test A | Test L | Train A | Train L | Test A | Test L | Train A | Train L | Test A | Test L |
| MNIST | 72.76 | 55.26 | 73.25 | 56.19 | 10.07 | 0.3631 | 11.35 | 0.3569 | 99.67 | 0.0134 | 97.53 | 0.0319 |
| F-MNIST | 79.31 | 87.98 | 77.86 | 87.74 | 80.17 | 0.1105 | 69.33 | 0.1913 | 98.75 | 0.0221 | 86.52 | 0.1029 |

[1] The architecture of MLP is 784-500-10, with a total of $3.97 \times 10^5$ parameters. The activation function is tanh function.

[2] The MLP model is adjusted by the RSRA-LSRA for three times, and then trained by BP algorithm with Adam optimizer.

[3] Here, "A" means "Accuracy" and "L" means "Loss".

solutions. While it may not discover the global optimum, it can fine-tune the model and steadily decrease the loss value until it reaches a local minimum. The idea behind combining the SRA and the BP algorithm is to leverage the strengths of both approaches and find a global optimization solution for the model.

To illustrate the advantages of combining the BP and SRA, we design a standard MLP model (784-500-10) with weight matrix initialized uniformly in the range [0, 1]. In this scenario, if we were to apply the BP algorithm to train the model, it would become trapped in a local minimum. As shown in Table 4 and Figure 3, we train the MLP model using the Adam Optimizer with a learning rate set to $1e-3$.

For the MNIST dataset, the Adam optimizer cannot escape local minimization even after 300 epochs. The training accuracy fluctuates around 10.07%, and the testing accuracy remains at approximately 11.35%, essentially at the level of random guessing. Both the training and testing loss values hover around 0.36. For the Fashion-MNIST dataset, the situation is slightly better. After 500 epochs, the model achieves a training accuracy of 80.17% and a testing accuracy of 69.33% with loss values of 0.1105 and 0.1913, respectively.

As shown in Table 4, the RSRA-LSRA, while capable of escaping local minimization, cannot fully realize its potential and reach optimal performance. For the MNIST dataset, it achieves a training accuracy of 72.76% and a testing accuracy of 73.25%. For the Fashion-MNIST dataset, it obtains a training accuracy of 79.31% and a testing accuracy of 77.86%. Although the loss value is not directly related to accuracy for the RSRA-LSRA, it is notably higher than usual.

When we initialize the weight matrix using the same scheme, the MLP model, after being adjusted by RSRA-LSRA three times, is trained by the BP algorithm with the Adam optimizer. For the MNIST dataset, within 100 epochs, the training accuracy reaches 99.67%, and the testing accuracy is 97.53%. For the Fashion-MNIST dataset, within 400 epochs, the training accuracy reaches 98.75%, surpassing most advanced neural network architectures with highly complex structures, and the testing accuracy is 86.52%. For both the MNIST and Fashion-MNIST datasets, the combination of training algorithms achieves better performance than individual BP or RSRA-LSRA.

## 4. Conclusion and Future Work

This study introduces a gradient-free SRA to train the neural network. The algorithm leverages SVD to calculate the rotation matrix and update the weight matrix to achieve an optimal representation of sample patterns.

Regarding SRA, two fundamental types exist: LSRA and RSRA. LSRA and RSRA can be applied separately or in combination to train MLP models. Each approach has a distinct effect on the performance of MLP models. The experiments demonstrate that RSRA-LSRA can achieve optimal solutions and is more stable than other approaches when training a variety of MLP models, including shallow MLP, deep MLP, and deep MLP with

a bottleneck layer. Therefore, when compared with the BP algorithm, the RSRA-LSRA scheme is selected.

Two-layer MLP models with 10000 hidden nodes are designed to evaluate the effectiveness of the RSRA-LSRA and the BP algorithm. It is observed that for shallow and wide MLP models, the RSRA-LSRA can achieve higher training accuracy than the BP algorithm. For less noisy datasets, such as MNIST and Fashion-MNIST, the testing accuracy achieved by the RSRA-LSRA is better than the model trained by the BP algorithm. However, for noisy datasets, such as CIFAR10 and CIFAR100, since the existence of overfit, the testing accuracy of the model trained by RSRA-LSRA is worse than model trained by BP algorithm.

To study the relationship between model complexity and overfitting, the most noisy dataset, CIFAR100 is chosen. Four two-layer MLP models with hidden node numbers: 1000, 5000, 10000, and 20000, are designed and trained using the RSRA-LSRA and the BP algorithm with the SGD optimizer. With the RSRA-LSRA, the training accuracy grows quickly with increasing complexity, from 47.07% to 96.36%. However, the testing accuracy does not show the same trend due to the impact of overfitting. The best testing accuracy is for the MLP with 10000 hidden nodes, which achieves 21.48% testing accuracy. For the BP algorithm, the MLP model with 5,000 hidden nodes achieves optimal performance, with 54.04% training accuracy and 23.02% testing accuracy. In this case, both algorithms cause overfitting.

The combination of SRA and the BP algorithm demonstrates an extraordinary effect on training MLP models with weight matrices initialized with a uniform distribution in the range [0, 1]. The experiment shows that the BP algorithm can become stuck in local optimal points, and RSRA-LSRA, while capable of escaping local optimal solutions, cannot reach its full potential when the hidden node number is limited—in other words, when the subspace is limited. The combination of RSRA-LSRA and the BP algorithm leverages both algorithms to achieve optimal performance for MLP models with a relatively smaller number of hidden nodes.

In the future, we aim to apply this approach to different network architectures, such as CNN, LSTM, and GNN. Simultaneously, we plan to reduce the time complexity and space complexity of the SRA. We may explore alternative methods to find the rotation matrix for the weight matrix.

## References

[1]  A. Krizhevsky, I. Sutskever, and G. E. Hinton. "Imagenet classification with deep convolutional neural networks". In: *Communications of the ACM* 60.6 (2017), pp. 84–90.

[2]  A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin. "Attention is all you need". In: *Advances in neural information processing systems* 30 (2017).

[3]  V. Mnih, K. Kavukcuoglu, D. Silver, A. Graves, I. Antonoglou, D. Wierstra, and M. Riedmiller. "Playing atari with deep reinforcement learning". In: *arXiv preprint arXiv:1312.5602* (2013).

[4]  R. Miotto, F. Wang, S. Wang, X. Jiang, and J. T. Dudley. "Deep learning for healthcare: review, opportunities and challenges". In: *Briefings in bioinformatics* 19.6 (2018), pp. 1236–1246.

[5]  A. M. Ozbayoglu, M. U. Gudelek, and O. B. Sezer. "Deep learning for financial applications: A survey". In: *Applied Soft Computing* 93 (2020), p. 106384.

[6]  U. Seiffert. "Multiple layer perceptron training using genetic algorithms." In: *ESANN*. Citeseer. 2001, pp. 159–164.

[7]  Y. LeCun, Y. Bengio, et al. "Convolutional networks for images, speech, and time series". In: *The handbook of brain theory and neural networks* 3361.10 (1995), p. 1995.

[8]  L. R. Medsker and L. Jain. "Recurrent neural networks". In: *Design and Applications* 5 (2001), pp. 64–67.

[9]  D. E. Rumelhart, G. E. Hinton, and R. J. Williams. *Learning internal representations by error propagation*. Tech. rep. California Univ San Diego La Jolla Inst for Cognitive Science, 1985.

[10] I. Livieris and P Pintelas. "A survey on algorithms for training artificial neural networks". In: *University of Patras, Department of Mathematics, Educational Software Development Laboratory, University of Patras, GR-265* 4 (2008).

[11] W.-D. K. Ma, J. Lewis, and W. B. Kleijn. "The HSIC bottleneck: Deep learning without back-propagation". In: *Proceedings of the AAAI conference on artificial intelligence*. Vol. 34. 04. 2020, pp. 5085–5092.

[12] A. Choromanska, B. Cowen, S. Kumaravel, R. Luss, M. Rigotti, I. Rish, P. Diachille, V. Gurev, B. Kingsbury, R. Tejwani, et al. "Beyond backprop: Online alternating minimization with auxiliary variables". In: *International Conference on Machine Learning*. PMLR. 2019, pp. 1193–1202.

[13] J. C. Spall. "A one-measurement form of simultaneous perturbation stochastic approximation". In: *Automatica* 33.1 (1997), pp. 109–112.

[14] D. Krotov and J. J. Hopfield. "Dense associative memory for pattern recognition". In: *Advances in neural information processing systems* 29 (2016).

[15] M. Moonen, P. Van Dooren, and J. Vandewalle. "A singular value decomposition updating algorithm for subspace tracking". In: *SIAM Journal on Matrix Analysis and Applications* 13.4 (1992), pp. 1015–1038.

[16] B. Yang. "Projection approximation subspace tracking". In: *IEEE Transactions on Signal Processing* 43.1 (1995), pp. 95–107. DOI: 10.1109/78.365290.

[17] J. A. Lee, M. Verleysen, et al. *Nonlinear dimensionality reduction*. Vol. 1. Springer, 2007.

[18] S. Cheng, Y. Wang, H. Huang, D. Liu, H. Fan, and S. Liu. "Nbnet: Noise basis learning for image denoising with subspace projection". In: *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 2021, pp. 4896–4906.

[19] K. Fukui, N. Sogi, T. Kobayashi, J.-H. Xue, and A. Maki. "Discriminant feature extraction by generalized difference subspace". In: *IEEE Transactions on Pattern Analysis and Machine Intelligence* 45.2 (2022), pp. 1618–1635.

[20] M. Wax and A. Adler. "Detection of the number of signals in uniform arrays by invariant-signal-subspace matching". In: *IEEE Transactions on Signal Processing* 70 (2022), pp. 1270–1281.

[21] K. Adhikari, R. J. Vaccaro, and D. D. Sartori. "Shift invariant sparse arrays and their optimal signal and noise subspaces". In: *Signal Processing* 198 (2022), p. 108579.

[22] R. Roy, A. Paulraj, and T. Kailath. "ESPRIT–A subspace rotation approach to estimation of parameters of cisoids in noise". In: *IEEE transactions on acoustics, speech, and signal processing* 34.5 (1986), pp. 1340–1342.