



## Data augmentation for cross-subject EEG features using Siamese neural network



Rongrong Fu<sup>a,\*</sup>, Yaodong Wang<sup>a</sup>, Chengcheng Jia<sup>b</sup>

<sup>a</sup> Measurement Technology and Instrumentation Key Lab of Hebei Province, Department of Electrical Engineering, Yanshan University, Qinhuangdao, China

<sup>b</sup> Department of Electrical, Computer & Biomedical Engineering, Ryerson University, Canada

### ARTICLE INFO

**Keywords:**

Complex motor recognition  
Data augmentation  
Siamese neural network  
Similarity measurement  
Transfer learning

### ABSTRACT

Electroencephalography (EEG) motor intention recognition has been extensively used in robot control, brain rehabilitation and other health care fields. Recently, some algorithms have been proposed based on generative adversarial neural network (GAN) to enhance EEG signal, and have achieved high recognition performance. However, these methods utilize the convolutional kernel method of the GAN, while the optimal convolutional scale of CNN varies from subject to subject. This may lead to the data generated by GAN to lack authenticity and produce data that does not match the ideal situation. Particularly, the performance of data augmentation degrades when the original calibrated EEG is insufficient. To address these issues, we proposed a novel cross-subject Siamese Neural Network (SNN) approach to enhance EEG feature data. Specifically, we used our proposed SNN to construct highly similar extended EEG features of different subjects and successfully improved the performance of motor intention recognition. Then, we design an accurate boundary avoidance task to evaluate the effectiveness of the proposed method. Compared with the traditional experimental paradigm, the coding process of this experiment is more complex, which makes the results more reliable when using the SNN. The extended EEG features display significantly better performance than any other common classifiers in the case of small data size, and it demonstrates that this proposed method can effectively address these issues of existing EEG motor intention recognition methods based on data augmentation and improve the classification performance.

### 1. Introduction

In recent years, brain-computer interface (BCI) [1–5] has been developed as a system for converting mental intention into commands or codes, which allows a direct connection between human brain and external devices. Electroencephalography (EEG) [6–8] is well used due to its ability of detecting neural activity in the brain, and some neural activity can reflect the physiological activity and function of the mind. In contrast to steady-state visual evoked potentials (SSVEP) [9] and event-related potentials (ERP) [10] measured by visual or auditory paradigms, EEG induced by motor imagery [11–13] can provide motor intention in the absence of external stimulation.

Recently, with the widespread attention of deep learning algorithms in pattern recognition applications, it has been explored in the field of BCI [14–17], especially in the classification of motor imagery signals. The reliable and stable performance has a significant improvement after the decade of deep learning algorithms. However, all of these methods rely on the sufficient EEG data to complete the whole calibration cycle to

adapt the system for each user, the EEG signal is also subject to some restrictions. Most of time, there are only hundreds of experimental trials are involved in training set, resulting in a lack of training samples, which limits the performance of motor intention recognition and application of brain-computer interaction.

In terms of deep learning, a good recognition performance can be obtained when data set contains large size of samples for training. McDonnell et al. [18] proposed a shallow neural network to classify large-scale EEG data for motor imagery signals, and obtained great reliability classification effect. However, compared with shallow neural networks, deep neural networks need to train more nodes and require numerous training data to discover the potential of deep neural networks. Therefore, how to solve the problem of achieving high accuracy with limited training trials is an urgent issue, and it is significant to overcome the lack of training trials in motor intention recognition by using data augmentation.

Data augmentation [19–21] is one of the most popular methods that effectively eliminate the drawback of the complex training process, also

\* Corresponding author.

E-mail address: [frr1102@aliyun.com](mailto:frr1102@aliyun.com) (R. Fu).

it provides the generation capability for limited available data. It has been proven that data augmentation can transform or expand existing data to realize the purpose of generating new data [22]. This method has been extensively employed to remove overfitting and improve the effect of classification, and the generated data set has the similar feature distribution in the feature space as the original data set, which can be applied to increase the data size of training samples. In the field of small data image classification, this method has achieved a good classification effect by adding noise, scaling and offsetting to the original image [23]. Correspondingly, it has been applied to generated EEG data in the field of brain-computer interface [24,25]. For example, Krell et al. [26] produced a new EEG signal by offsetting the original EEG data. In Colominas et al. [27], new signals were generated by adding various noises to the original non-stationary signals. The above study proves that data augmentation is beneficial to improve the accuracy and stability of the classification effect in EEG data. At present, some researchers generate EEG data by geometric transformation of existing EEG data. For example, Paris et al. [1] proposed to input Gaussian noise into a system based on Fast Fourier Transform (FFT) in EEG data, so as to endow more features to EEG data and achieve better results in feature extraction and classification. Other researchers focus on generating new EEG data using adversarial neural network. Hartmann et al. [28] proposed a generative adversarial neural network (GAN) for data augmentation, which was trained on the original EEG data and obtained a better classification.

Among the methods of data augmentation using deep learning, admittedly neural network is generally chosen to generate simulated EEG data on the original EEG data [29]. However, either geometric transformation [30] or GAN, due to its network structure, there may be a situation that the generated data lacks authenticity (e.g., the generated EEG signal waveform is highly similar, but there is a negative value in amplitude, which is obviously impossible) [31]. To address this issue, we propose a novel SNN data enhancement method, which is considered as an excellent bridge to measure the similarity of EEG features of CSP. And two groups of EEG features with the highest similarity are combined as the newly generated EEG feature set, to realize the purpose of data augmentation of generating new EEG trials by combining the original EEG trials.

## 2. Materials and methods

### 2.1. Data description

Our data were from Key Lab of Measurement Technology & Instrumentation of Hebei Province. Ten subjects, aged 23–25 years, all of whom had normal hearing, normal/corrected-to-normal visual acuity and no history of any neurological disease, were recruited for this study. All subjects gave written informed consent before experiments. In this study, we designed a “bowl-ball” experimental system which was selected as the experimental paradigm to evoke EEG signals in the subjects. During the experiment, subjects sat in front of the LED screen and controlled the direction of ball movement in the screen with the different hands through the keyboard. The data acquisition equipment is EMOTIV EPOC + 14-Channel Wireless EEG Headset. There are 14 EEG channels named based on the International 10–20 locations, they are: AF3, F7, F3, FC5, T7, P7, O1, O2, P8, T8, FC6, F4, F8, AF4, and the sampling frequency is 128 Hz. The experimental process uses sessions as the basic unit. Each session contains continuously collected EEG data, including 80\*2 trials (There are 80 trials of each type of motor recognition task, and number of motor recognition task classes is 2, which are left hand and right hand.) Finally, the data is truncated and sorted into high-dimensional form Number of samples × Number of channels × Number of Trials × Number of classes(128\*14\*80\*2). With the description of the dataset, we have balanced dataset for 2-class motor recognition task. In the process of acquisition, the experimental environment should be kept quiet and free from unnecessary environmental noise interference. The subjects should try to avoid unnecessary head

movement, and the EEG data collected should be transmitted to the computer through Bluetooth. The experimental procedure is shown in Fig. 1.

### 2.2. Data preprocessing

In this experiment, we collected the original EEG signals generated by 10 subjects operating the task system with different hands under the motor recognition task. This study preprocessed the data, using a notch filter to eliminate 50 Hz power-line interference, and an 8–13 Hz bandpass filter to extract EEG of  $\alpha$  rhythm.

The CSP method utilizes spatial filters [32] to minimize the covariance matrix [33] of one class EEG trials and maximize the covariance matrix of another class. In specific, we extract the EEG data of the binary classification problem by the CSP method and obtain the optimal separated spatial distribution of each class. The principle of the algorithm is as follows:

The original EEG data are classified by category, record the left-hand data as  $E_1$ , and record the right-hand data as  $E_2$ .

The covariance matrix of each class of EEG is calculated as:

$$R_i = \frac{E_i E_i^T}{\text{trace}(E_i E_i^T)}, (i = 1, 2), \quad (1)$$

where  $E_i^T$  represents the transpose of  $E_i$ ,  $\text{trace}(\bullet)$  represents the sum of diagonal elements of a matrix.

The corresponding mean covariance matrix of each type of EEG data is calculated as:

$$\bar{R}_i = \frac{1}{T_i} \sum_{k=1}^{T_i} R|Y(k) = i, (i = 1, 2) \quad (2)$$

where  $T_i$  represents the number of samples corresponding to each category and  $Y(k)$  represents the class label of  $k$  experiment.

The spatial filter banks  $W$  are constructed as follows::

$$(\bar{R}_1 + \bar{R}_2)^{-1} \bar{R}_1 = W D W^{-1} \quad (3)$$

where  $D$  is the diagonal matrix composed of each eigenvalue,  $W$  is the matrix composed of eigenvectors.

Given the projection direction of different motor recognition tasks  $P_j$ , the EEG trial  $X$  is projected into the feature space of  $j$ th motor recognition tasks by  $Z_j = P_j X$ , which is treated as  $Z = [Z_1, \dots, Z_C]$  containing the signals filtered by a group of spatial filters. Normalized log-variances features of  $Z$  can be calculated by.

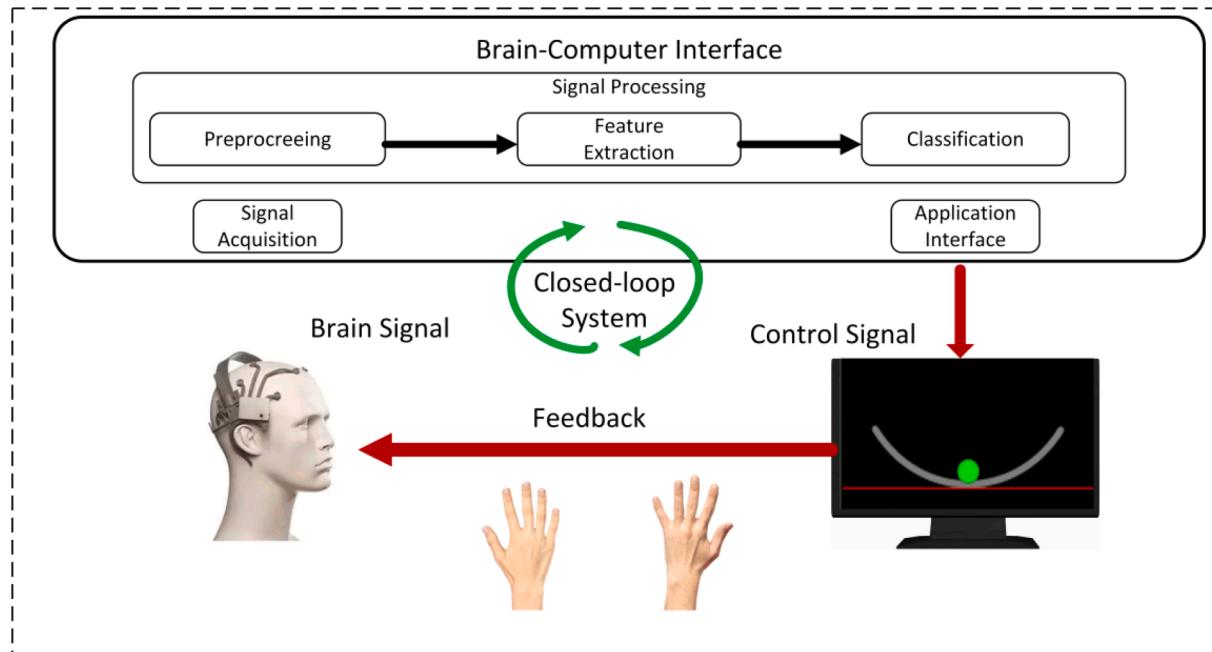
$$f_p = \log\left(\frac{\text{var}(Z_p)}{\sum_{q=1}^M \text{var}(Z_q)}\right), p = 1, \dots, M, \quad (4)$$

where  $\text{var}(Z_p)$  represents the variance of the  $p$ th row vector in  $Z$ ;  $M = C \times m$ , where  $C$  is the number of channels,  $m$  is the number of samples.  $F = \{f_1, \dots, f_M\}$  is a set of normalized log-variance features of EEG signals.

### 2.3. Experimental methods

#### 2.3.1. Cross-subject Siamese neural network

With the CSP method extract EEG features of different subjects, the performance is regularly affected by the lack of authentic EEG data and low feature dimension, which degrades the classification performance. In order to solve this problem, we proposed a novel cross-subject data augmentation method based on the SNN to expand the original EEG data set. The method of expanding the data set to improve the classification performance has been verified in the field of image processing [34]. The prerequisite for finding a suitable EEG dataset for expansion is that the data to be expended has a high similarity to the original EEG data and similar expanded dataset can bring reliable improvement in performance. Therefore, one of the critical issues is how to measure the



**Fig. 1.** Experimental procedure. During the experiment, the subjects controlled the “bowl-ball” system to move from left to right by tapping the keyboard with their left and right hands. The ball was subject to acceleration during the movement and could potentially run out of the bowl. Therefore, subjects were supposed to change the movement of the system using different hands to control the keyboard which was depended on the position of the ball to avoid the ball from escaping from the bowl. If the ball fell out of the bowl, the experiment failed.

similarity of different subjects in the motor recognition task. In machine learning and deep mining algorithms, the input data have similar feature matrices and feature space distributions, however, in the field of brain-computer interface, the feature space distribution of subjects’ EEG is more disparate due to the large individual variability of different subjects. Inspired by You et al. [35] utilization of Siamese neural network to measure image edge similarity, we transfer the SNN model to EEG features to determine the similarity of EEG data from different subjects.

In this paper, we propose a novel cross-subject Siamese Neural Network based on transfer learning for EEG data augmentation. The Siamese Neural Network can be used to judge the similarity of CSP features from different subjects by feeding two groups of EEG features into the SNN composed of two groups of convolutional neural networks (as shown in Fig. 2(a)). The transfer learning based on the SNN can collect two sets of EEG features with the highest similarity (as shown in Fig. 2(b)), and generate a new high-dimensional EEG feature set in this shape which enables to achieve data augmentation for motor recognition tasks. The whole network structure of Siamese neural network consists of four parts: the first part is the input layer, the second part is a Siamese neural network composed of two identical convolutional neural networks, the third part is the measurement layer, the fourth part is full connection layer and softmax output layer.

**Input layer:** the first three-dimensional feature vectors extract from the EEG data by CSP algorithm and transform into a line map composed of red long dashed lines, green short dashed lines and blue solid lines (as shown in Fig. 3), in which lines of different colors and types represent different dimensions of EEG features. As shown in 2.2 in this paper, we calculate the normalized log-variance features of the EEG signals using the CSP method.  $F = \{f_1, \dots, f_M\}$  is a set of normalized log-variances features of EEG signals, where the dimension of M is 14. Since the 14-dimensional CSP feature graph results in an image that is too complex, it is difficult for the Siamese neural network to compare the similarity of CSP feature graphs between subjects. If we want to compare high-dimensional CSP feature graphs using the Siamese neural network, this requires a large amount of computational resources and training time. Therefore, we choose low-dimensional CSP feature graphs for comparison to validate the effectiveness of our proposed data

augmentation algorithm. Finally, we choose 3-dimensional CSP feature graphs in 14 dimensions to verify the effectiveness of our algorithm. Each input line graph is converted from the  $3 \times 160$  feature matrix of each subject into  $900 \times 900$  line graph, where 3 is the selected first three-dimensional feature vector, 160 is the motor recognition task performed 160 times, and  $900 \times 900$  is the resolution of the line graph.

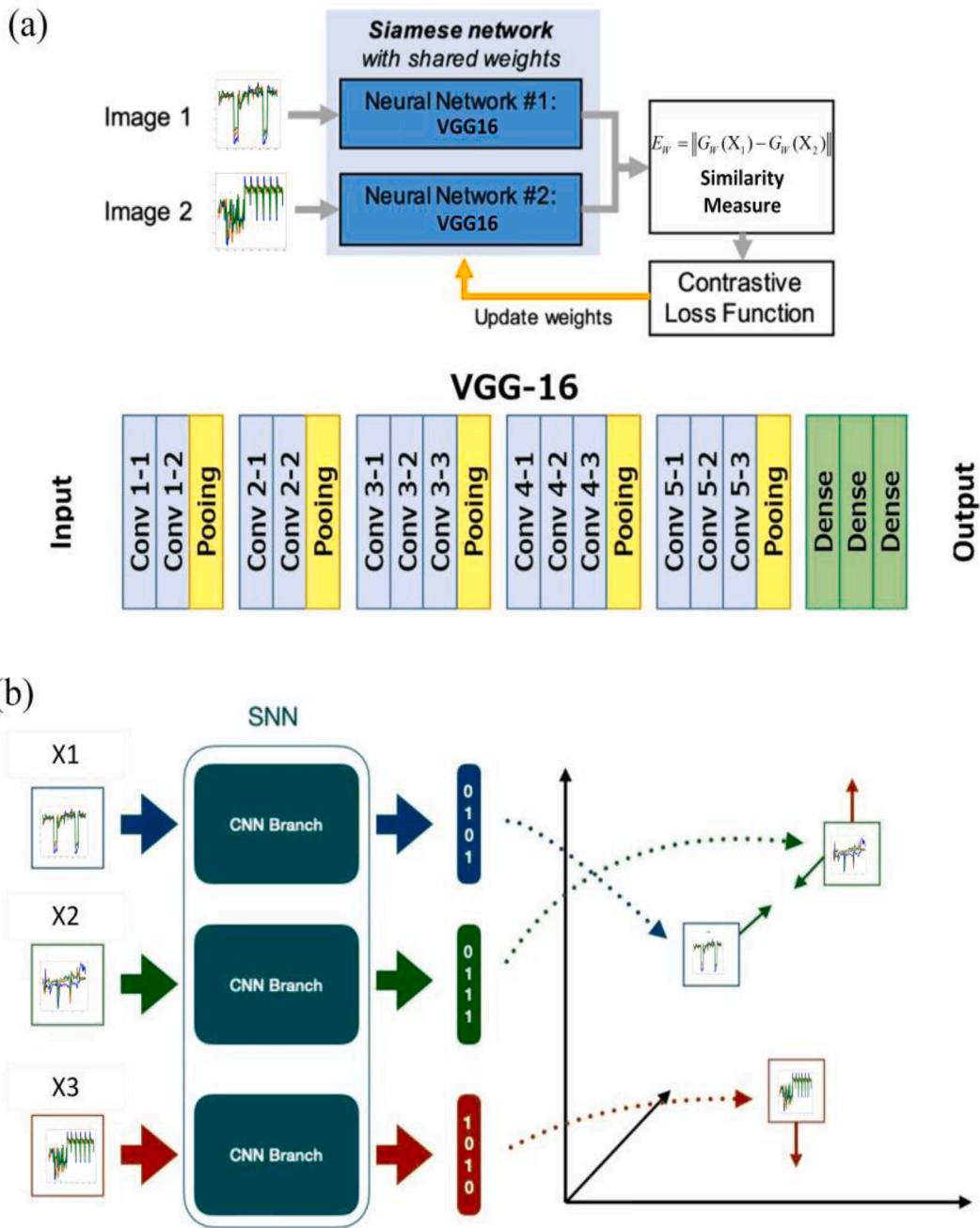
The Siamese Neural Network structure composed of two identical convolutional neural networks: the main role of this network is to extract the two sets of optimal feature vectors from the input feature line graphs of the two subjects. This network structure, that uses two identical convolutional neural networks, facilitates the process of the feature vectors obtained from the feature line graph in the measurement layer. The two convolutional neural networks are VGG16, and the VGG16 is composed of 5 convolution modules and 3 full connection layers. Each convolution module has two or three convolutional cores with a size of  $3 \times 3$  convolution layer and one maximum pool layer. Compared with other models, VGG16 model replaces a large volume product core by several small convolution cores, which avoids the large-scale increase of parameters in network structure, and has more abundant features and stronger learning ability in similarity learning tasks.

**Measurement layer:** two groups of eigenvectors are obtained by convolution neural network, and the two groups of eigenvectors are fed into the measurement layer of Siamese Neural Network to obtain the similarity results.  $E_W = \|G_W(X_1) - G_W(X_2)\|$  is the output of the similarity measurement layer, where  $X_1$  and  $X_2$  is the feature line graph of the network model input,  $G_W(X_1)$  and  $G_W(X_2)$  respectively represents the feature vectors of  $X_1$  and  $X_2$  obtained through the Siamese Neural Network.

**Full connection and SoftMax layer:** the structure of full connection layer–Relu activation function–full connection layer, it can effectively reduce the computational effort when processing large amounts of EEG data.

### 2.3.2. Transfer learning and data augmentation

Transfer learning is a new research strategy for data set with small sample size [36], which aims to improve the learning effect of the target prediction function in the target domain by using the knowledge of the



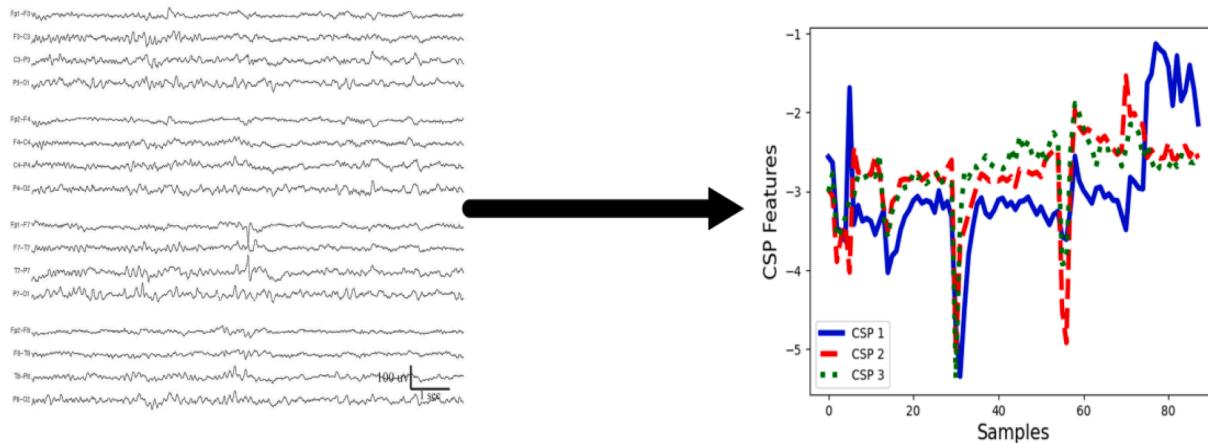
**Fig. 2.** Block diagram of the study protocol. (a) The overall structure diagram of the Siamese Neural Network; (b) The test diagram of the Siamese neural network for three different feature line graphs.

source domain. Transfer learning emphasizes knowledge transfer across similar but different domains, tasks and distributions. In motor recognition task, the individual differences between different subjects are large, which leads to the different distribution of the feature space of the data. In this paper, we investigate the similarity measurement method of Zhan et al. [37] using Siamese neural network for the image edge problem, so as to achieve the similarity measure of EEG data from different subjects.

Data augmentation has been proved to effectively improve the classification performance in neural networks and become one of the common solutions to the problem of data scarcity. In addition, the appearance of GAN provides a new approach and framework for data augmentation. Compared with the traditional augmentation method, this GAN approach works via adversarial training concept and demonstrates more powerful capacity both in feature learning and generation

performance. However, the GAN approach at BCI is not completely satisfactory to researchers, and it still exhibits some problems. Specifically, the GAN method utilizes the convolutional kernel of CNN in the feature learning phase, while the optimal convolutional size of CNN varies from person to person. This may lead to a lack of authenticity in the data generated by GAN, producing data that does not match the ideal situation. In addition, the enhancement effect of GAN on EEG signals is limited when the original calibrated EEG is insufficient.

To address these issues, we propose a novel augmented method based on the Siamese neural network for generating additional EEG to augment the original EEG trials in order to improve the performance of a BCI classifier. In specific, the two EEG features with the highest similarity measured by Siamese neural network are combined into the extended CSP EEG feature set by using data augmentation method, which can effectively improve the classification performance in motor



**Fig. 3.** The original EEG data is processed by CSP to get the front three-dimensional feature line graph.

recognition task. As shown in Fig. 4, the key of this study is to measure the similarity of EEG feature line graphs, so numerous line graphs can be selected as the source domain for training the model, and fine-tuning the model will achieve excellent similarity index of line graphs. Subsequently, the trained model can be transferred to the recognition of EEG feature line graphs to accomplish a great similarity index.

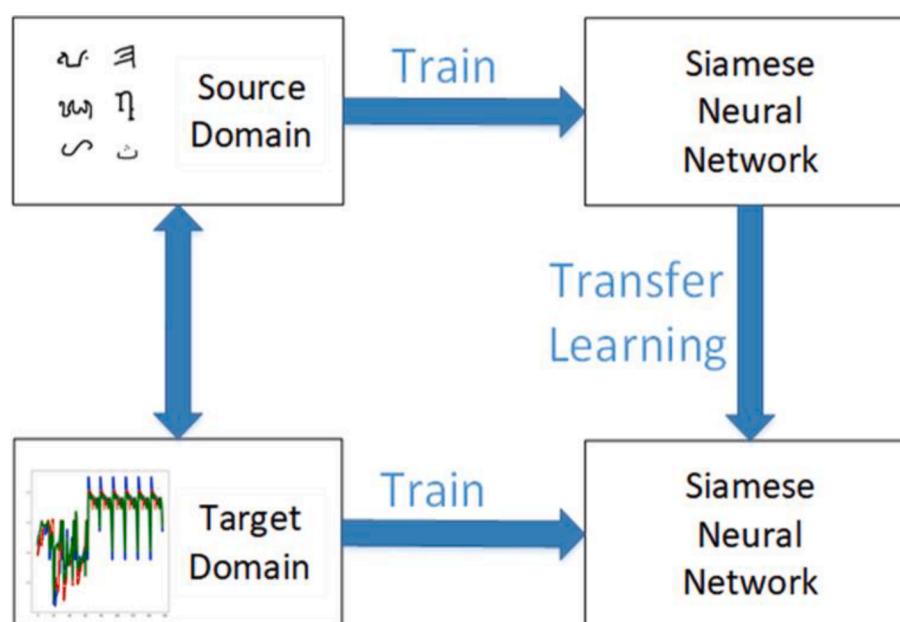
### 3. Results

In this section, we evaluate the augmented performance of the proposed data augmentation method on various machine learning approaches. First, the former three-dimensional CSP features of ten subjects' EEG are extracted and converted into line graphs, as shown in Fig. 5 (a). Taking the EEG features of two different subjects as input, the similarity measurement results of EEG signals of two different subjects can be obtained by Siamese Neural Network. The results of similarity calculation between the EEG features of S01 and S02 and other subjects are shown in Fig. 5 (b). It illustrates that among the ten subjects, S01 and S08 have the highest similarity in three-dimensional feature collection except the experimenter himself, and S02 and S03 have the highest similarity.

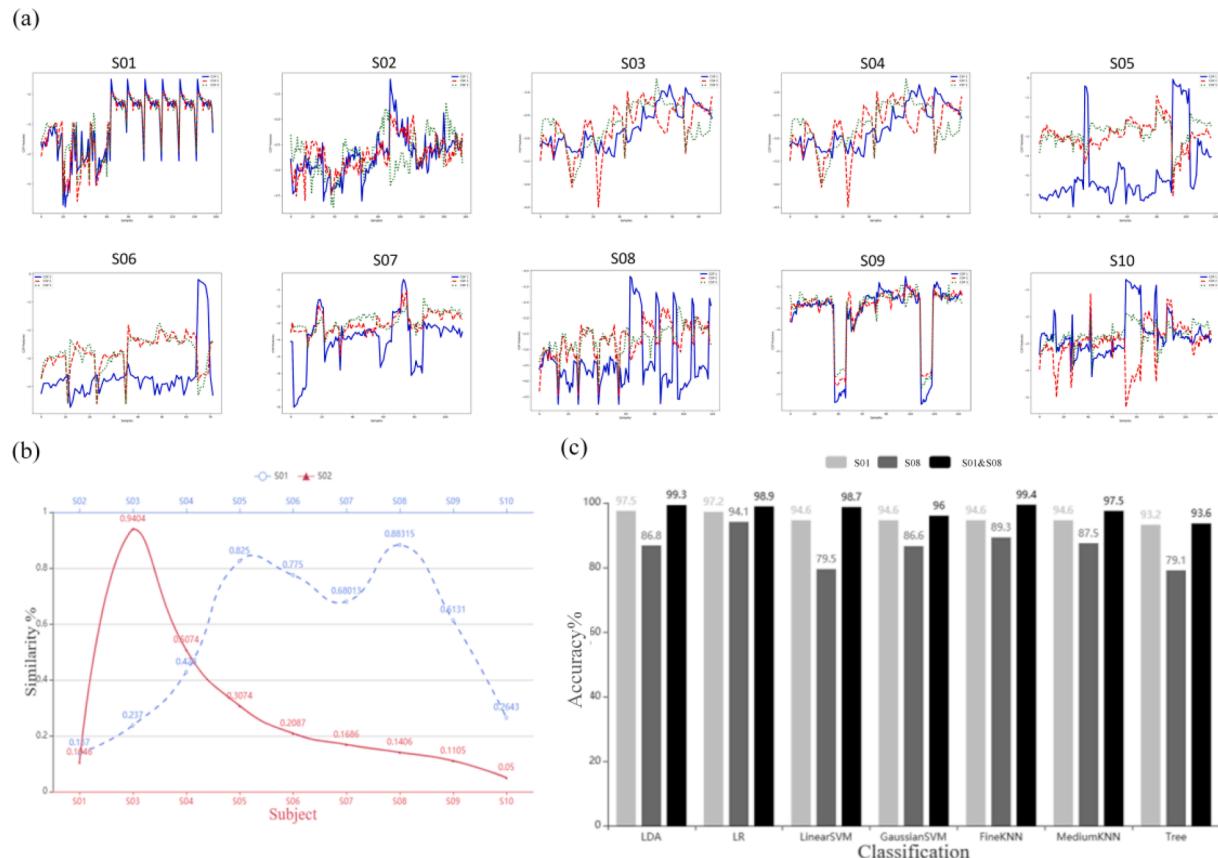
After obtaining the similarity results between different subjects, the

EEG feature sets of the two subjects with the highest similarity are selected and combined into a new dataset as the augmented EEG feature set for these two subjects. Since the EEG signals of S01 and S08 have the highest similarity, the EEG features of these two subjects are adopted to merge and expand. We evaluate the classification performance of the augmented feature set on various machine learning methods. To ensure the validity of this method, we evaluate the augmented performance of limited EEG data using 5-fold cross-validation. As shown in Fig. 0.5 (c), augmentation of EEG data from designed experiment with the proposed method significantly improved the classification accuracy from 90% to 96%. Note that since the EEG trials collected are not particularly adequate, the obtained accuracy may not be stable. Nevertheless, the augmentation with the proposed method improves the classification performance. In other words, it shows that our proposed EEG augmentation method does not depend on the size of the EEG trials.

Considering that the improvement of classification accuracy may be due to the increase of feature set dimension, we further verify the reliability of the augmentation performance for the proposed method. Therefore, when we verify this method, we randomly merge a subject's feature set with the feature set of S02. For example, we select S07 to merge a new data set of S02 and S07. The five-fold cross-validation is performed on the merged dataset of S02&S07 using the above



**Fig. 4.** Transfer learning process.



**Fig. 5.** Evaluation of the augmented performance for the SNN method. (a) Line graphs of the former three-dimensional CSP features for different subjects. The long red dashed line, the short green dashed line and the solid blue line represent the different EEG features, respectively. (b) The similarity between S01, S02 and other subjects is calculated by Siamese Neural Network. (c) Cross-validation results for multiple classification methods between S01 and S08.

classification methods. As shown in Table 1, the comparison results between the data sets of S02&S07 with LDA, SVM, KNN and logistic regression and the data sets of S02&S03 show that the average classification accuracy of the proposed method is improved by about 4%.

According to Fig. 5 and Table 1, experiments show that the feature set with the highest similarity can be obtained by using the Siamese Neural Network algorithm. Compared with the original feature set, the classification accuracy of the expanded feature set in a variety of classification methods is improved by about 6%. Considering the improvement of the classification effect caused by the increased dimensions of the data set, compared with the feature set of the random subject, the classification performance of the feature set combined by this algorithm has improved by about 4%.

To further verify that this method can improve the classification performance in a variety of machine learning methods, receiver operating feature curve (ROC) is selected to verify the reliability of this method in SVM, LDA, QDA and logistic regression methods. The

horizontal axis represents false positive rate, the vertical axis represents true positive rate, and the area under the ROC curve represents AUC value. Fig. 6 (a) shows the effect of the feature set of subject S01 compared with the augmented feature set of subject S01. The ROC curve of the augmented feature set in machine learning algorithms such as QDA, LDA, SVM and logistic regression is closer to the upper left corner, the area under the curve is larger, the AUC value is larger, and the classification performance is better. Fig. 6(b) shows the effect of the feature set of subject S02 compared with the augmented feature set of subject S02. The ROC curve of the augmented feature set is closer to the upper left corner, with larger area under the curve, larger AUC value and better classification performance. Fig. 6 can verify the improvement of the classification accuracy and reliability of the classification performance of this data augmentation method on multiple machine learning classification methods of multiple subjects.

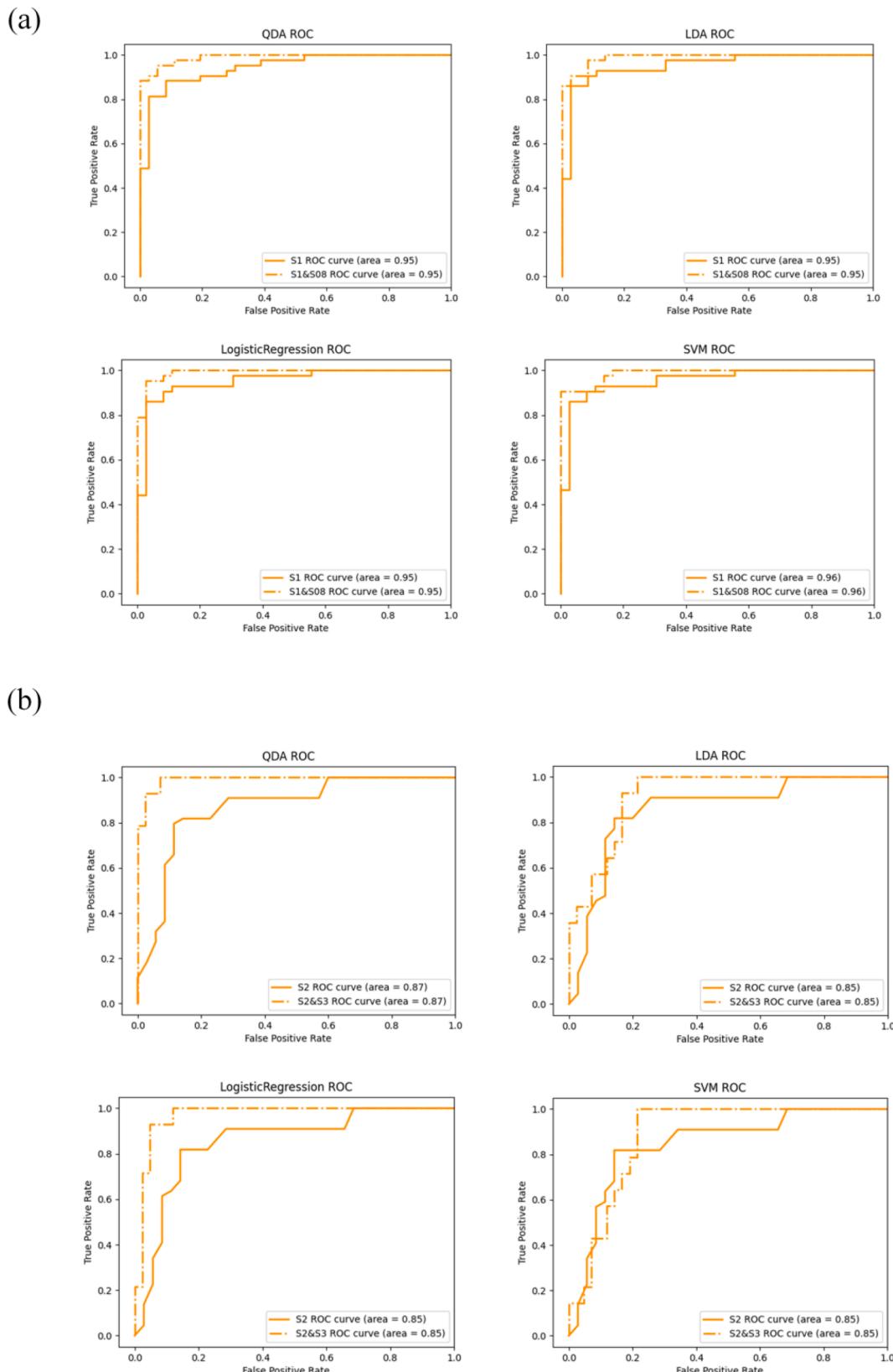
We randomly select two subjects S05 and S06 among ten subjects and verify the effectiveness of this method in each of the two subjects. As shown in Fig. 7, the results of similarity measures show that subject S05 has the highest similarity with subject S01 and subject S06 has the highest similarity with subject S07. Since the EEG signals of S05 and S01 have the highest similarity, the EEG features of these two subjects are adopted to merge and expand. S06 and S07 are also adopted to merge and expand. To ensure the validity of this method, we evaluate the augmented performance of limited EEG data using 5-fold cross-validation. As shown in Table 2, augmented S05 and S06 have improved performance on a variety of classification methods.

#### 4. Discussions

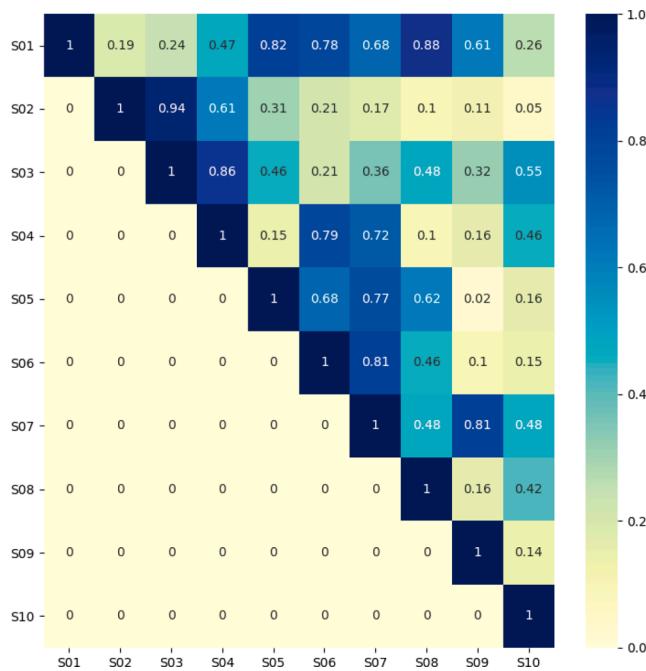
In this paper, we propose a new data augmentation method to

**Table 1**  
Cross-validation results of S02 and S03 for various classification methods.

Classification	S02	S03	S02&S03	S02&S07
LDA	89.2%	79.7%	93.7%	92.9%
Logistic Regression	88.6%	76.6%	99.4%	90.2%
Linear SVM	91.1%	78.5%	94.9%	95.5%
Fine Gaussian SVM	86.7%	83.5%	94.3%	94.6%
Medium Gaussian SVM	92.4%	81.6%	98.7%	94.6%
Coarse Gaussian SVM	91.8%	81.0%	98.1%	94.6%
Cubic SVM	91.1%	77.8%	98.1%	92.9%
Fine KNN	91.1%	86.7%	98.1%	96.4%
Medium KNN	91.8%	76.6%	96.2%	93.8%
Coarse KNN	91.1%	50.0%	61.4%	49.1%



**Fig. 6.** The ROC curves of different subjects in different machine learning algorithms. (a) ROC curves of S01 and S01&S08 in different machine learning algorithms, (b) ROC curves of S02 and S02&S03 in different machine learning algorithms.



**Fig. 7.** Results of similarity measures for CSP features across different subjects. (S = subject) Each square represents the result of the similarity measure between two subjects.

**Table 2**  
Cross-validation results of S05 and S06 for various classification methods.

Classification	S05	S01	S05&S01	S06	S07	S06&S07
LDA	99.2%	88.4%	97.5%	92.4%	95.0%	99.2%
Logistic Regression	97.5%	90.2%	97.5%	96.6%	99.2%	99.2%
Linear SVM	91.1%	94.6%	99.2%	97.5%	96.7%	100%
Fine Gaussian SVM	86.4%	94.6%	92.4%	87.3%	91.7%	85.6%
Medium Gaussian SVM	99.2%	94.6%	100%	97.5%	95.8%	97.5%
Coarse Gaussian SVM	96.6%	94.6%	98.3%	94.1%	87.5%	96.6%
Cubic SVM	96.6%	94.6%	99.2%	97.5%	98.3%	100%
Fine KNN	97.5%	94.6%	100%	96.6%	98.3%	99.2%
Medium KNN	96.6%	94.6%	97.5%	94.1%	95.8%	98.3%
Coarse KNN	49.2%	49.1%	49.2%	49.2%	50.0%	49.2%

improve the accuracy of motor intention recognition task under the case of limited training samples. As shown in Fig. 7, we use Siamese Neural Network to measure the similarity of the three-dimensional EEG feature line graphs with the EEG signal, and to generate new high-dimensional EEG features. The Siamese Neural Network is a simple measure to find the difference between different subjects and easy to understand and compute. According to the results in Table 1, compared with results given by the original EEG features, the accuracy of classification has improved for a variety of classification methods. This means that similarity measures of EEG features reflect the differences in feature mapping across different subjects. This difference may be caused by the individual differences of different subjects' coding methods in the brain for the same motor intention recognition task. The Siamese Neural Network can be seen as a bridge to study the individual differences between different subjects. We attribute the innovation of this method mentioned to the following three aspects.

- (1) Previous researchers have generally focused on the augmentation effect of the data augmentation method and ignored the problem of the authenticity of the generated data. Our method merges the

highly similar EEG of different subjects as the expanded EEG data, so that our expanded data are from the real data and solve the problem of lack of authenticity of the generated data.

- (2) In addition, our proposed data augmentation method establishes a subject-to-subject connection compared to GAN and other network structures. This approach effectively reduces the possibility of overfitting of the expanded data. Since overfitting is often caused by the lack of data diversity, we use the multi-subject EEG data as the source of the expanded data to increase the diversity of the expanded data while achieving the purpose of data augmentation.
- (3) Our proposed data augmentation method has a simple structure and can migrate to accomplish the data augmentation requirements on a wide range of problems. In the experiments, for the training process of Siamese Neural Network, the amount of EEG data is too small to meet the training requirements of Siamese Neural Network. We use transfer learning to migrate the weight model learned by Siamese Neural Network from other dataset to the problem of EEG feature set similarity measurement.

In this paper, the Siamese Neural Network is used to find the feature set of two subjects with the smallest individual differences between different subjects, and the feature set of two different subjects is combined as a new feature set. The expanded feature set not only adds a lot of training samples for the training model, that increases the accuracy and reliability of the model classification, but eliminates the influence of individual differences on the model classification performance, so as to achieve the purpose of data augmentation. Considering cross-subject similarity measures in data augmentation is an intuitive counterbalance to the overfitting problem that occurs in some models, such as GAN.

- (4) Compared with other data augmentation methods using generative neural network, our new high-dimensional EEG feature sets are all from the original low dimensional EEG feature sets, which solve the problem that there is a certain deviation in the authenticity between the simulated EEG signal and the real EEG signal, providing a good dimensionality reduction effect. Since this method chooses to use raw EEG data to generate similar EEG data, the generated EEG data is slightly less similar to the raw EEG data compared to the GAN generated EEG data.

In addition, we choose CSP to extract EEG features, but this data augmentation method proposed in this paper can also be used in other feature extraction methods such as time-frequency domain, wavelet transform and Hilbert transform, and we think that it is feasible to use this method to generate high channel EEG data on the original low channel EEG data. Therefore, we believe that this method not only solves the problem of insufficient training samples of EEG data, but also provides a new idea for mapping low channel EEG data to high channel EEG data, and makes it possible to realize brain computer interaction system on low channel EEG data.

## 5. Conclusions

In the motor intention recognition task, the classification accuracy will be affected when the training samples are limited. Data augmentation is usually used to solve the problem of insufficient training samples in machine learning, and GAN is widely used to generate new EEG data in deep learning. In this paper, we propose a data augmentation method based on Siamese Neural Network for CSP EEG features, and utilize this method to improve the classification accuracy of motor intention signals. During the experiment, we collected EEG experimental data of precise motion control for a complex boundary avoidance task. The experimental results indicate that the accuracy and stability of classification performance can be improved by using this data

augmentation method in a variety of classification methods for motor intention signals. The cross-validation method proves that this method can improve the accuracy and stability of classification performance by 6% on average compared with the current existing classification methods. In particular, this paper further wide the application scope of data augmentation and shows that similarity measurement can provide a new idea for data augmentation.

## CRediT authorship contribution statement

**Rongrong Fu:** Conceptualization, Methodology, Resources, Writing – review & editing. **Yaodong Wang:** Software, Formal analysis, Writing – original draft, Validation, Investigation. **Chengcheng Jia:** Formal analysis, Writing – original draft, Validation, Investigation.

## Declaration of Competing Interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

## Acknowledgments

This work was supported by the National Natural Science Foundation of China [grant numbers 62073282, 61973262, 61806174]; Natural Science Foundation of Hebei Province [grant number E2018203433]; the Central Guidance on Local Science and Technology Development Fund of Hebei Province [grant number 206Z0301G], China Postdoctoral Science Foundation [grant number 2016M600193]; and Hebei Province Funding Project for Returned Overseas Scholar [grant number CL201727].

## References

- [1] A. Paris, G.K. Atia, A. Vosoughi, S.A. Berman, A new statistical model of electroencephalogram noise spectra for real-time brain-computer interfaces, *IEEE Trans. Biomed. Eng.* 64 (2017) 1688–1700, <https://doi.org/10.1109/TBME.2016.2606595>.
- [2] R. Mane, T. Chouhan, C. Guan, BCI for stroke rehabilitation: motor and beyond, *J. Neural Eng.* 17 (2020) 041001. 10.1088/1741-2552/aba162.
- [3] U. Chaudhary, N. Mrachacz-Kersting, N. Birbaumer, J. Taylor, D. Farina, Neuropsychological and neurophysiological aspects of brain-computer-interface (BCI) control in paralysis, *J. Physiol.* 9 (2020).
- [4] S. Perdikis, L. Tonin, S. Saeedi, C. Schneider, J.R. del Millán, The Cybathlon BCI race: successful longitudinal mutual learning with two tetraplegic users, *PLoS Biol.* 16 (2018), e2003787, <https://doi.org/10.1371/journal.pbio.2003787>.
- [5] S. Saha, K.A. Mamun, K. Ahmed, R. Mostafa, A. Khandoker, S. Darvishi, M. Baumert, Progress in Brain Computer Interfaces: Challenges and Trends, (n.d.) 20.
- [6] M.H. Lee, O.Y. Kwon, Y.J. Kim, H.K. Kim, Y.E. Lee, J. Williamson, S. Fazli, S.-W. Lee, EEG dataset and OpenBMI toolbox for three BCI paradigms: an investigation into BCI illiteracy, *GigaScience* 8 (2019) giz002, <https://doi.org/10.1093/gigascience/giz002>.
- [7] D.J. McFarland, J.R. Wolpaw, EEG-based brain–computer interfaces, *Curr. Opin. Biomed. Eng.* 4 (2017) 194–200, <https://doi.org/10.1016/j.cobme.2017.11.004>.
- [8] M. Del Pozo-Banos, J.B. Alonso, J.R. Ticay-Rivas, C.M. Travieso, Electroencephalogram subject identification: a review, *Expert Syst. Appl.* 41 (2014) 6537–6554, <https://doi.org/10.1016/j.eswa.2014.05.013>.
- [9] C. Guger, G. Edlinger, W. Harkam, I. Niedermayer, G. Pfurtscheller, How many people are able to operate an EEG-based brain-computer interface (BCI), *IEEE Trans. Neural Syst. Rehabil. Eng.* 11 (2003) 145–147, <https://doi.org/10.1109/TNSRE.2003.814481>.
- [10] S. Park, H.S. Cha, J. Kwon, H. Kim, C.H. Im, Development of an Online Home Appliance Control System Using Augmented Reality and an SSVEP-Based Brain-Computer Interface, (n.d.) 2.
- [11] E.R. Paitel, M.R. Samii, K.A. Nielson, A systematic review of cognitive event-related potentials in mild cognitive impairment and Alzheimer's disease, *Behav. Brain Res.* 396 (2021), 112904, <https://doi.org/10.1016/j.bbr.2020.112904>.
- [12] V. K., D. A., M. J., S. M., A. A., S.A. Iraj, A novel method of motor imagery classification using eeg signal, *Artificial Intelligence in Medicine*. 103 (2020) 101787. 10.1016/j.artmed.2019.101787.
- [13] J. Luo, X. Gao, X. Zhu, B. Wang, N. Lu, J. Wang, Motor imagery EEG classification based on ensemble support vector learning, *Comput. Methods Programs Biomed.* 193 (2020), 105464, <https://doi.org/10.1016/j.cmpb.2020.105464>.
- [14] A. Al-Saeq, S.A. Dawwd, J.M. Abdul-Jabbar, Deep learning for motor imagery EEG-based classification: a review, *Biomed. Signal Process. Control* 63 (2021), 102172, <https://doi.org/10.1016/j.bspc.2020.102172>.
- [15] A. Craik, Y. He, J.L. Contreras-Vidal, Deep learning for electroencephalogram (EEG) classification tasks: a review, *J. Neural Eng.* 16 (2019) 031001. 10.1088/1741-2552/ab0ab5.
- [16] H. Dose, J.S. Möller, H.K. Iversen, S. Puthusserypady, An end-to-end deep learning approach to MI-EEG signal classification for BCIs, *Expert Syst. Appl.* 114 (2018) 532–542, <https://doi.org/10.1016/j.eswa.2018.08.031>.
- [17] K.M. Tsioris, V.C. Pezoulas, M. Zervakis, S. Konitsiotis, D.D. Koutsouris, D. I. Fotiadis, A long short-term memory deep learning network for the prediction of epileptic seizures using EEG signals, *Comput. Biol. Med.* 99 (2018) 24–37, <https://doi.org/10.1016/j.combiomed.2018.05.019>.
- [18] M.D. McDonnell, T. Vladusich, Enhanced image classification with a fast-learning shallow convolutional neural network, *Int. Joint Conf. Neural Networks (IJCNN) 2015* (2015) 1–7, <https://doi.org/10.1109/IJCNN.2015.7280796>.
- [19] A. Sharma, D.B. Jayagopi, Towards efficient unconstrained handwriting recognition using Dilated Temporal Convolution Network, *Expert Syst. Appl.* 164 (2021), 114004, <https://doi.org/10.1016/j.eswa.2020.114004>.
- [20] F. Fahimi, S. Dosen, K.K. Ang, N. Mrachacz-Kersting, C. Guan, Generative adversarial networks-based data augmentation for brain-computer interface, *IEEE Trans. Neural Netw. Learn. Syst.* 32 (2021) 4039–4051, <https://doi.org/10.1109/TNNLS.2020.3016666>.
- [21] Y. Luo, B.L. Lu, EEG Data Augmentation for Emotion Recognition Using a Conditional Wasserstein GAN, in: 2018 40th Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC), IEEE, Honolulu, HI, 2018: pp. 2535–2538. 10.1109/EMBC.2018.8512865.
- [22] J. Lemley, P. Corcoran, Deep learning for consumer devices and services 4—A review of learnable data augmentation strategies for improved training of deep neural networks, *IEEE Consum. Electron. Mag.* 9 (2020) 55–63, <https://doi.org/10.1109/MCE.2019.2959075>.
- [23] P. Chlap, H. Min, N. Vandenberg, J. Dowling, L. Holloway, A. Haworth, A review of medical image data augmentation techniques for deep learning applications, *J. Med. Imag. Radiat. Oncol.* 65 (2021) 545–563, <https://doi.org/10.1111/1754-9485.13261>.
- [24] F. Wang, S. Zhong, J. Peng, J. Jiang, Y. Liu, Data augmentation for EEG-based emotion recognition with deep convolutional neural networks, in: K. Schöeffmann, T.H. Chalidabhongse, C.W. Ngo, S. Aramvith, N.E. O'Connor, Y.-S. Ho, M. Gabouj, A. Elgammal (Eds.), MultiMedia Modeling, Springer International Publishing, Cham, 2018, pp. 82–93, [https://doi.org/10.1007/978-3-319-73600-6\\_8](https://doi.org/10.1007/978-3-319-73600-6_8).
- [25] X.R. Zhang, M.Y. Lei, Y. Li, An amplitudes-perturbation data augmentation method in convolutional neural networks for EEG decoding, in: in: 2018 5th International Conference on Information, Cybernetics, and Computational Social Systems (ICCSS), 2018, pp. 231–235, <https://doi.org/10.1109/ICCSS.2018.8572304>.
- [26] M.M. Krell, A. Seeland, S.K. Kim, Data Augmentation for Brain-Computer Interfaces: Analysis on Event-Related Potentials Data, *ArXiv:1801.02730* [Cs, q-Bio]. (2018). <http://arxiv.org/abs/1801.02730>.
- [27] M.A. Colominas, G. Schlotthauer, M.E. Torres, Improved complete ensemble EMD: A suitable tool for biomedical signal processing, *Biomed. Signal Process. Control* 14 (2014) 19–29, <https://doi.org/10.1016/j.bspc.2014.06.009>.
- [28] K.G. Hartmann, R.T. Schirrmeister, T. Ball, EEG-GAN: Generative adversarial networks for electroencephalographic (EEG) brain signals, *ArXiv:1806.01875* [Cs, Eess, q-Bio, Stat]. (2018). <http://arxiv.org/abs/1806.01875> (accessed September 18, 2021).
- [29] G. Dai, J. Zhou, J. Huang, N. Wang, HS-CNN: a CNN with hybrid convolution scale for EEG motor imagery classification, *J. Neural Eng.* 17 (2020), 016025, <https://doi.org/10.1088/1741-2552/ab405f>.
- [30] K. Zhang, G. Xu, Z. Han, K. Ma, X. Zheng, L. Chen, N. Duan, S. Zhang, Data augmentation for motor imagery signal classification based on a hybrid neural network, *Sensors* 20 (2020) 4485, <https://doi.org/10.3390/s20164485>.
- [31] F. Lotte, Signal processing approaches to minimize or suppress calibration time in oscillatory activity-based brain-computer interfaces, *Proc. IEEE* 103 (2015) 871–890, <https://doi.org/10.1109/JPROC.2015.2404941>.
- [32] K.K. Ang, Z.Y. Chin, H. Zhang, C. Guan, Filter Bank Common Spatial Pattern (FBCSP) in Brain-Computer Interface, in: 2008 IEEE International Joint Conference on Neural Networks (IEEE World Congress on Computational Intelligence), 2008: pp. 2390–2397. 10.1109/IJCNN.2008.4634130.
- [33] H. Lu, H.L. Eng, C. Guan, K.N. Venetsanopoulos, Regularized common spatial pattern with aggregation for EEG classification in small-sample setting, *IEEE Trans. Biomed. Eng.* 57 (2010) 2936–2946, <https://doi.org/10.1109/TBME.2010.2082540>.
- [34] X. Xuan, B. Peng, W. Wang, J. Dong, On the generalization of GAN image forensics, in: Z. Sun, R. He, J. Feng, S. Shan, Z. Guo (Eds.), Biometric Recognition, Springer International Publishing, Cham, 2019, pp. 134–141, [https://doi.org/10.1007/978-3-030-31456-9\\_15](https://doi.org/10.1007/978-3-030-31456-9_15).
- [35] W. You, H. Zhang, X. Zhao, A siamese CNN for image steganalysis, *IEEE Trans. Inform. Forensic Secur.* 16 (2021) 291–306, <https://doi.org/10.1109/TIFS.2020.3013204>.
- [36] L. Huang, Y. Chen, Dual-path siamese CNN for hyperspectral image classification with limited training samples, *IEEE Geosci. Remote Sensing Lett.* 18 (2021) 518–522, <https://doi.org/10.1109/LGRS.2020.2979604>.
- [37] Y. Zhan, K. Fu, M. Yan, X. Sun, H. Wang, X. Qiu, Change detection based on deep siamese convolutional network for optical aerial images, *IEEE Geosci. Remote Sensing Lett.* 14 (2017) 1845–1849, <https://doi.org/10.1109/LGRS.2017.2738149>.