

Few-shot Class-incremental Learning for EEG-based Emotion Recognition

Tian-Fang Ma¹, Wei-Long Zheng¹, and Bao-Liang Lu^{1,2,3,4,5}

¹ Department of Computer Science and Engineering, Shanghai Jiao Tong University, 800 Dongchuan Rd, Shanghai 200240, China

² RuiJin-Mihoyo Laboratory, RuiJin Hospital, Shanghai Jiao Tong University School of Medicine, 197 Ruijin 2nd Rd, Shanghai 200020, China

³ Key Laboratory of Shanghai Commission for Intelligent Interaction and Cognitive Engineering, Shanghai Jiao Tong University, 800 Dongchuan Rd, Shanghai 200240, China

⁴ Clinical Neuroscience Center, RuiJin Hospital, Shanghai Jiao Tong University School of Medicine, 197 Ruijin 2nd Rd, Shanghai 200020, China

⁵ Brain Science and Technology Research Center, Shanghai Jiao Tong University, 800 Dong Chuan Road, Shanghai 200240, China
{matianfang2676, weilong, bllu}@sjtu.edu.cn

Abstract. Current advanced deep neural networks can greatly improve the performance of emotion recognition tasks in affective Brain-Computer Interfaces (aBCI). Basic human emotions could be induced and electroencephalographic (EEG) signals could be simultaneously recorded. While data of basic common emotions are easier to collect, some complex emotions are low resource in terms of data size and label quality in real life, which would limit the utility of EEG-based emotion recognition models. To enhance the model adaptive capacity of new emotions with few samples, we introduce a few-shot class-incremental deep learning model for emotion recognition. The proposed model consists of a graph convolutional networks (GCN) and a linear classifier. By training the whole network on a base set in a preliminary stage, and fine-tuning the parameters of the linear classifier with very few shots of labeled samples, the model can incrementally learn new types of emotions while preserving knowledge of the old ones. Our experimental results on the SEED-V dataset show that even with very limited new class samples, the fine-tuned pre-trained model could have a fairly good performance on the test set with more emotion classes.

Keywords: EEG · Deep learning · Few-shot class-incremental learning

1 Introduction

Deep learning techniques have largely advanced development and research in brain-computer interface (BCI). As an important branch of BCI, the affective brain-computer interface (aBCI) has also made significant progress in human emotion recognition task [1]. In recent years, EEG-based emotion recognition research has aroused great interest in many interdisciplinary fields from psychology

to engineering, including basic research on emotion theories and applications of aBCIs. In the tasks of aBCI, there are commonly two ways to process the emotions: One way is to map all emotions into a two-dimensional valence-arousal coordinate system [12]. The main challenge is that accurate quantitative labeling is usually difficult to get. The other way is to categorize the emotions into discrete classes [5]. Constantly updated aBCI models are very effective in recognizing basic emotions. Ekman proposed seven basic emotions: fear, anger, joy, sadness, contempt, disgust, and surprise [5]. However, the frequency of occurrence of some new human emotions is relatively low in practice. The sample size of those emotions tends to be small, which makes them difficult to be recognized if the models have not learned the emotion categories in advance.

The problem of identifying new emotional categories can be defined as incremental learning. The ability of incremental learning is to deal with the continuous information flow in the real world and retain, even integrate and optimize old knowledge while absorbing new knowledge. One of the main problems that incremental work at solving to prevent is catastrophic forgetting. Catastrophic forgetting refers to the problem that general machine-learning models have a dramatic drop in the performance on the previously learned tasks [9]. Two typical incremental learning methods are based on regularization and replay respectively. The one using regularization is the learning without forgetting (LwF) algorithm [10]. LwF is a training mode between joint training and fine-tuning training. The model can be updated without using the data from the old task. The other one based on knowledge replay is called Incremental Classifier and Representation Learning (iCaRL) [11]. iCaRL preserves a representative portion of the old data for each old task while training the new data. And it could better remember the characteristics of the data learned from the old task. There are many limitations in the traditional incremental learning model, for example, it is difficult for the model to learn new types of knowledge when the sample size is small, or the model will overfit the new samples when the use of old samples is limited.

With the advent of the concept of few-shot learning, newly generated few-shot learning algorithms are designed to learn and generalize from small samples using existing knowledge. Humans can easily build new knowledge from just one or a few examples. However, machine learning algorithms typically require thousands of supervised samples to ensure generalization. As a joint concept, few-shot incremental learning focus on maintaining high performance for base knowledge and good generalization ability for new knowledge with the same model [2, 4, 15].

To make the learning model easily extend to new sets of emotion labels from very few samples, in this paper, we design an EEG-based few-shot class-incremental graph convolutional networks (FSCI-GCN) emotion recognition model. By using samples of the basic classes, the model learns a featured space from the base emotion classes in advance, and continually learns new classes from very few labelled samples by model fine-tuning. There is no limit to the retrieval of old knowledge, we store the extracted feature vectors of the old original data

and lock the model parameter. The final model is effective in recognizing new emotion classes without forgetting the previously learned knowledge.

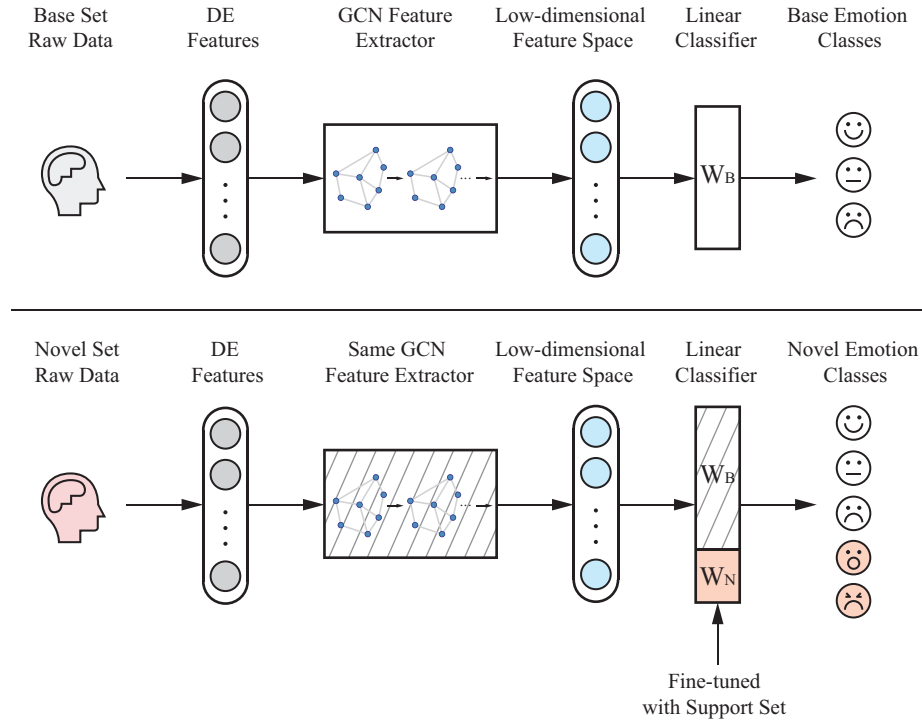


Fig. 1. Illustration of the FSCI-GCN model framework. The top part denotes the pre-training process: the GCN feature extractor and linear classifier of the model are trained with base set. The bottom part denotes the few-shot incremental learning that feature extractor and old weights of the linear classifier are locked and new weights are trained with support set.

2 Methods

2.1 Graph Convolutional Networks and Feature Extractor Pre-training

For an undirected connected graph $\mathcal{G} = (\mathcal{V}, \mathcal{E}, A)$, which consists of a set of nodes \mathcal{V} with $|\mathcal{V}| = n$, a set of edge \mathcal{E} with $|\mathcal{E}| = m$, and the adjacency matrix A . The GCN model is proposed as follow [8]:

$$H^{l+1} = \sigma(\tilde{D}^{-\frac{1}{2}} \tilde{A} \tilde{D}^{-\frac{1}{2}} H^{(l)} W^{(l)}) \quad (1)$$

where $\tilde{A} = A + I$ is the adjacency matrix of the undirected graph (I is the identity matrix). \tilde{D} is the diagonal node degree matrix of \tilde{A} . $W^{(l)}$ is the layer-specific trainable weight matrix. σ denotes the activation function (here used a rectified linear unit). $H^{(l)}$ is the matrix of activations in the l^{th} layer. After a series of graph convolutional layers, we use a max pooling and linear layer to reduce the information in the graph network to a status space.

Firstly, the model is trained on a base set with sufficient examples. We jointly train the GCN feature extractor parameters θ and a linear classification layer η by minimizing the following cross-entropy loss with L2 regularization.

$$L(\eta, \theta) = -\frac{1}{N^{(0)}} \sum_{(x,y) \in S^{(0)}} \sum_{c \in C^{(0)}} w_c \log \frac{\exp(\eta_y^\top f_{\hat{\theta}}(x))}{\sum_{c \in C^{(0)}} \exp(\eta_c^\top f_{\hat{\theta}}(x))} y + \alpha(\|\eta\|^2 + \|\theta\|^2) \quad (2)$$

where x, y is pair of input and target, $S^{(0)}$ is the set of all x and y in base class, $N^{(0)}$ is the number of samples in base class, $C^{(0)}$ is the number of classes in base class, w_c is the weight for class c , $f_{\hat{\theta}}$ denotes the GCN-based feature extractor layer which is pre-trained from the base set and θ denotes all the parameters in the feature extractor. α is the hyperparameter of the L2 penalty.

2.2 Few-shot Incremental Learning Step

The learning step follows the notation in the few-shot class-incremental learning (FSCIL) model [15]. Assume a stream of T learning sessions, each session is aligned with a labeled dataset $D^{(0)}, D^{(1)}, \dots, D^{(T)}$. Every dataset $D^{(T)}$ consists of a support set $S^{(T)}$ and a test set (query set) $Q^{(T)}$. Specially, $D^{(0)}$ is referred to the base set and $C^{(0)}$ represents the set of base classes. We assume it contains a large number of examples for every class that existed in $C^{(0)}$. $D^{(1)}$ to $D^{(T)}$ introduce the new classes. For every new dataset $D^{(t)}$, $C^{(t)}$ denote the set of classes expressed in dataset $D^{(t)}$, and $C^{(\leq t)}$ denotes the union set of classes $\bigcup_{j \leq t} C^{(j)}$. In the few-shot incremental learning process, each support set contains only new classes ($C^{(t)} \cap C^{(<t)} = \emptyset$), while each test set evaluates models on a combination of data with the base classes and all classes that have appeared.

The support set contains 5-shot samples for each novel class. Given an incremental session $t < 0$, the linear classifier of the model is fine-tuned so as to perform well in classifying both base classes and novel classes.

Fine-tuning After the preliminary feature extractor training using graph convolutional networks on the base classes, the model is fine-tuned under the loss function $L(\eta)$. We introduce new weight vectors and optimize

$$L(\eta) = L_{CE}(\eta) + \alpha\|\eta\|^2 + \beta R_{ER}^{(t)} + \gamma R_{SR}^{(t)} \quad (3)$$

in which

$$L_{CE}(\eta) = -\frac{1}{N^{(0)}} \sum_{(x,y) \in S^{(0)}} \sum_{c \in C^{(\leq t)}} w_c \log \frac{\exp(\eta_y^\top f_{\hat{\theta}}(x))}{\sum_{c \in C^{(\leq t)}} \exp(\eta_c^\top f_{\hat{\theta}}(x))} y \quad (4)$$

where $R_{ER}^{(t)}$ and $R_{SR}^{(t)}$ denote respectively, the entropy regularization term and the subspace regularization term at session t . Entropy regularization is specifically used to minimize the overlap of class probability distributions of the support set at session t . Subspace regularization minimizes the subspace distance between the new weight vector and the old weight vector of the linear classifier.

To be noted, the denominator in the summation formula of $L_{CE}(\eta)$ is different from (1). Because of the introduction of new labels, the classes change to $C^{(\leq t)}$ instead of $C^{(0)}$.

Entropy Regularization We use the entropy regularizer for the support set fine-tuning process. The approach of entropy regularization was introduced as a semi-supervised learning method [6], and later used as a few-shot learning baseline for image recognition [3]. The regularizer minimizes a low Shannon Entropy H . In our case, the transductive fine-tuning solves for minimizing the following loss:

$$R_{ER}^{(t)} = \frac{1}{N} \sum_{(x,y) \in S^{(t)}} \mathbb{H}(p_\eta(\cdot|x)) \quad (5)$$

In which N is the number of samples of each new class. $p_\eta(\cdot|x)$ is the distribution new class samples. \mathbb{H} denotes the Shannon Entropy. Minimizing the Shannon Entropy allows the fine-tuned model to predict a high probability of the support sets being classified into their right labels.

Subspace regularization Multiple previous works showed that constraining parameters for related tasks lie on the same manifold or the same linear subspace [7]. The potential feature space shared by all classes is useful for class increments [13]. Regularizing the subspaces spanned by all base class weight vectors encourages the classification of new categories to rely on semantics rather than pseudo-features, in other words, making the feature space of the new category to be consistent with the subfeature space of the existing task to the greatest extent [2].

Given a parameter for an incremental class η_c and base class parameters $\eta_{j \in C^{(0)}}$, we first compute the subspace target m_c for each class. The subspace regularizer is defined by η_c and m_c :

$$R_{SR}^{(t)}(\eta) = \sum_{c \in C^{(t)}} \|\eta_c - m_c\|^2 \quad (6)$$

where m_c is the projection of η_c onto the space spanned by $\eta_{j \in C^{(0)}}$:

$$m_c = P_{C^{(0)}}^T \eta_c \quad (7)$$

Let $P^{C^{(0)}}$ contain the orthogonal basis vectors of each subspace spanned by the initial set of base weights $\eta_{j \in C^{(0)}}$, which can be computed by using the QR decomposition:

$$[P^{C^{(0)}} \ Q'] \begin{bmatrix} R \\ \mathbf{0} \end{bmatrix} = \eta_{C^{(0)}}^T \quad (8)$$

Subspace regularization does not assume that the data of all labels are available at the beginning. In the learning process, tasks arrive in an incremental way and predictions can be made over all categories that have been learned so far.

3 Experiment Setup

3.1 The SEED-V Dataset

The SEED-V dataset is one of EEG datasets used for emotion recognition from the SEED series (SJTU Emotion EEG Dataset)¹. The original SEED dataset contains EEG data of 12 subjects with 3 labeled basic emotions which are positive, negative, and neutral. The SEED-V dataset included fear and disgust as the fourth and fifth emotions and collected EEG data and eye movement data from another 16 subjects (6 males and 10 females) [17]. A total number of 24 video clips are used for the stimulation of five categories of emotion: happy, neutral, sad, fear, and disgust. Sixteen subject participants are recruited for the experiment. Each participant is required to watch the video clips in 3 sessions (24 clips randomly placed for each session). In each session, the video clips of every emotion label occurred the same number of times. The 45 video clips in a session are placed in 3-fold order (15 clips each), with one emotion for each category in a fold, for the convenience of cross-validation [16].

3.2 Feature Extraction

In the SEED-V dataset, the original EEG signals are recorded by the ESI NeuroScan System with 62 electrode channels at a sampling rate of 1000Hz. For pre-processing, the raw EEG signals of all participants are applied to a band-pass filter between 1 and 75 Hz to reduce the influence of artifacts and drift. Then the filtered EEG signals are down-sampled from 1000 Hz to 200 Hz to reduce the computational complexity. Both power spectral density (PSD) and differential entropy (DE) features are extracted from the 200 Hz down-sampled signal. Both features are computed within a 4-second non-overlapping Hanning window in five frequency bands: delta (1-4 Hz), theta (4-8 Hz), beta(14-31 Hz), and gamma(31-50 Hz) for each channel. The total dimension of each EEG feature in a sample is 310. The linear dynamic system algorithm was used for feature smoothing [14].

Preliminary works showed that using the DE features of all five frequency bands is the most effective predictor of emotion [17–19]. Thus, we use DE features of all five frequency bands (a total dimension of 310 features) in both the pre-training and fine-tuning process of our FSCI-GCN model.

3.3 Evaluation Details

For the SEED-V dataset, we use 3-fold cross-validation. Due to the fact that few-shot samples have high randomness, we design a secondary 3-fold cross-validation based on the session term (as shown in Fig. 2). The EEG data of the

¹ <https://bcmi.sjtu.edu.cn/home/seed/seed-v.html>

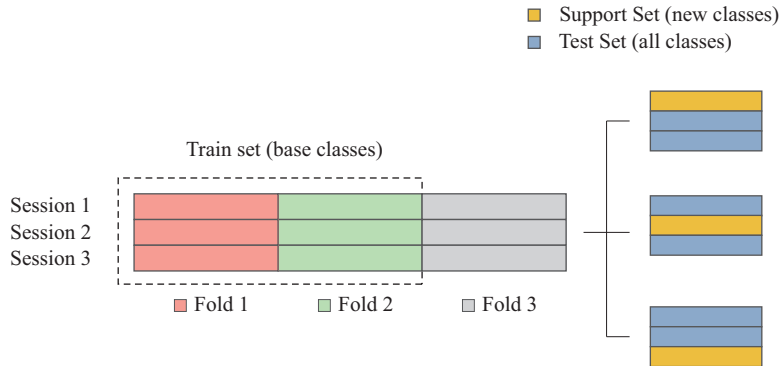


Fig. 2. Cross-validation partitioning of the SEED-V dataset.

three base classes (happy, neutral, and sad) in fold 1 and fold 2 are considered as the base set. In fold 3, EEG data with two new labels (fear and disgust) belongs to the support set. All five emotions including both base classes and novel classes in fold 3, session 2 and 3 form the test set. By parity of reasoning, fold 3 session 2 will be the support set and the other two sessions would form the test set, the same goes for fold 3 session 3. Each class in the training set has one shot or five shots. The one-shot and five-shot data are selected from the support set under a uniform distribution. Within the first fold (fold 1 and 2 are the base set), a secondary cross-validation yields three pairs of support set and test set. All three folds are used for hyperparameter selection and average accuracy estimation of the primary fold for each subject.

4 Experiment Results

4.1 Single-class Increment Result

Table 1 and Table 2 show the basic model performance on the base set and the experiment result when one single emotion class is entered into the FSCI-GCN model respectively. For the references: the support vector machine (SVM) baseline denotes the basic linear partition accuracy of the SEED-V dataset in multidimensional space. The GCN model and FSCI-GCN model use exactly the same network structure. The SVM and GCN model result in the table denotes the overall accuracy rate both models could get when using the complete data of new class (novel class) and train base class and new class together. We use the iCaRL model as the incremental learning baseline [11]. From the result, the accuracy of GCN model for the classification of three types of emotions (happy, neutral, sadness) in SEED-V reaches 81.19%. The accuracy of the FSCI-GCN model (5-shot) reaches 62.25%.

When training with a full-shot support set, the performances of both the iCaRL model and FSCI-GCN model increase. The 4-class full-shot accuracy rate

of the FSCI-GCN model is higher than the iCaRL model baseline regardless of shot numbers.

Table 1. Performance of different models on the 3-class base set.

Model (Base-class)	KNN	SVM	MLP	GCN
Mean	0.5890	0.6558	0.7631	0.8119
Std.	0.2095	0.2117	0.1572	0.1544

Table 2. Performance of different models on the 4-class test set.

Model (4-class)	Mean	Std.
iCaRL Baseline [11] (5-shot)	0.5357	0.2043
iCaRL Baseline (Full-shot)	0.5882	0.1635
FSCI-GCN (5-shot)	0.6225	0.1788
FSCI-GCN (Full-shot)	0.6876	0.1189
SVM	0.6310	0.1704
GCN	0.7598	0.1415

4.2 Multiple-class Increment Result

Table 3 shows the experiment result of an increment of multiple classes. The support set consists of few-shot samples of fear and disgust emotion. And the test set consists of all five emotion classes. The accuracy of the 5-class GCN model is 67.96%, much lower than the 4-class accuracy. The FSCI-GCN model has exactly the same parameters in the feature extractor layer and base weight and bias as the GCN (base) model. After training with only 5-shot samples of each novel classes, the FSCI-GCN model can notably recognize new classes and the overall accuracy can reach 51.81%.

Table 3. Performance of different models on the 5-class test set.

Model (4-class)	Mean	Std.
iCaRL Baseline [11] (5-shot)	0.4201	0.1854
iCaRL Baseline (Full-shot)	0.4635	0.1611
FSCI-GCN (5-shot)	0.5181	0.1506
FSCI-GCN (Full-shot)	0.5763	0.0907
SVM	0.5940	0.1538
GCN	0.6796	0.1271

Fig. 3 shows the confusion matrix of FSCI-GCN model (5-shot) on the 4-class and 5-class test set. Comparing the confusion matrix of the GCN (base)

model on base classes, and the FSCI-GCN model on all classes of novel set, the FSCI-GCN model significantly improves the recognition rate of new classes. The confusion matrix in Fig. 3 shows that, while the 43% of the new emotion 'fear' are correctly recognized, the FSCI-GCN model also mistakenly recognize about 17% of the old emotions as the new emotion 'fear'. That means if there are no new emotions in the test set, the model accuracy would even decrease. Due to the small number of new samples, the model compromises the recognition rate of old categories in order to improve the recognition ability of new categories.

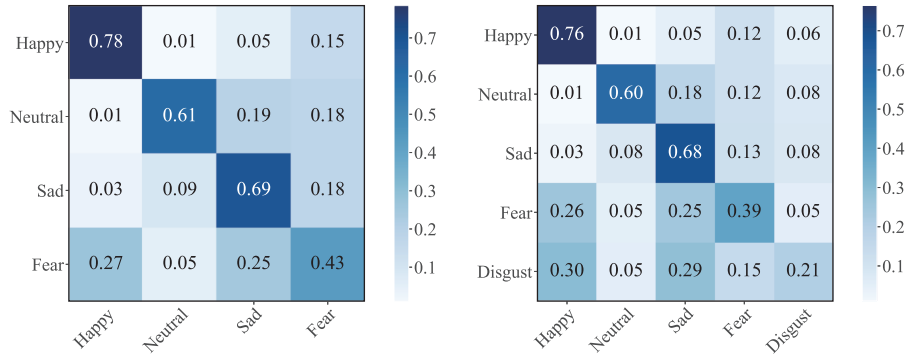


Fig. 3. Confusion matrices of the FSCI-GCN model on the test set.

4.3 Discussion

The capacity of the few-shot incremental learning model to improve recognition rate depends on base model and subjects. In general, the average accuracy across all subjects is significantly improved. However, for subjects with low data quality in which the base model can not distinguish the base classes well, the performance of new model barely improves. Also, comparing results from single-class increment versus multiple-class increment of the SEED-V dataset, the performance of the FSCI-GCN model declines as the number of novel classes increases. With one single novel class, the model is easier to distinguish novel classes from old ones from the feature space. This is consistent with the conclusion in image recognition. In addition, the FSCI-GCN model also bears some limitations. The model is not built under a zero-shot condition. If the model is not trained with new emotion data, it could not identify emotions that are different from the old categories. This is room for improvement in the future.

5 Conclusion

In this paper, we have proposed a few-shot incremental GCN-based model for EEG emotion recognition. For EEG emotion recognition models that have been

trained to recognize basic emotion labels, the model framework expands to a set of new weights that can be fine-tuned. By adopting entropy regularization and subspace regularization on the training process of the fine-tuned linear classifier, the model can balance the old training samples and the new ones, and make predictions on new labels while avoiding the catastrophic forgetting of old knowledge. To reduce the impact of randomness of small samples, we have applied a secondary three-fold cross-validation for the partition of the support set and test set. The test result on both datasets shows that the model can significantly increase the recognition rate of new samples.

Acknowledgements. This work was supported in part by grants from the National Natural Science Foundation of China (No. 61976135), MOST 2030 Brain Project (No. 2022ZD0208500), Shanghai Municipal Science and Technology Major Project (No. 2021SHZDZX), SJTU Global Strategic Partnership Fund (2021 SJTUHKUST), Shanghai Marine Equipment Foresight Technology Research Institute 2022 Fund (No. GC3270001/012), and GuangCi Professorship Program of RuiJin Hospital Shanghai Jiao Tong University School of Medicine.

References

1. Alarcao, S.M. and Fonseca, M.J.: Emotions recognition using EEG signals: A survey. *IEEE Transactions on Affective Computing*, 10(3), pp.374-393 (2017)
2. Akyürek, A.F., Akyürek, E., Wijaya, D. and Andreas, J.: Subspace Regularizers for Few-Shot Class Incremental Learning. *arXiv preprint arXiv:2110.07059* (2021)
3. Dhillon, G.S., Chaudhari, P., Ravichandran, A. and Soatto, S.: A baseline for few-shot image classification. *arXiv preprint arXiv:1909.02729* (2019)
4. Dong, S., Hong, X., Tao, X., Chang, X., Wei, X. and Gong, Y.: May. Few-shot class-incremental learning via relation knowledge distillation. In: *Proceedings of the AAAI Conference on Artificial Intelligence* Vol. 35, No. 2, pp. 1255-1263 (2021)
5. Ekman, P.: An argument for basic emotions. *Cognition & Emotion*, 6(3-4), pp.169-200 (1992)
6. Grandvalet, Y. and Bengio, Y.: Semi-supervised learning by entropy minimization. *Advances in Neural Information Processing Systems*, 17 (2004)
7. Jacob, L., Vert, J.P. and Bach, F.: Clustered multi-task learning: A convex formulation. *Advances in Neural Information Processing Systems*, 21 (2008)
8. Kipf, T.N. and Welling, M.: Semi-supervised classification with graph convolutional networks. *arXiv preprint arXiv:1609.02907* (2016)
9. Kirkpatrick, J., Pascanu, R., Rabinowitz, N., Veness, J., Desjardins, G., Rusu, A.A., Milan, K., Quan, J., Ramalho, T., Grabska-Barwinska, A. and Hassabis, D.: Overcoming catastrophic forgetting in neural networks. In: *Proceedings of the National Academy of Sciences*, 114(13), pp.3521-3526 (2017)
10. Li, Z. and Hoiem, D.: Learning without forgetting. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 40(12), pp.2935-2947 (2017)
11. Rebuffi, S.A., Kolesnikov, A., Sperl, G. and Lampert, C.H.: icarl: Incremental classifier and representation learning. In: *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pp. 2001-2010 (2017)
12. Russell, J.A.: A circumplex model of affect. *Journal of Personality and Social Psychology*, 39(6), p.1161 (1980)

13. Schonfeld, E., Ebrahimi, S., Sinha, S., Darrell, T. and Akata, Z.: Generalized zero-shot learning via aligned variational autoencoders. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops, pp. 54-57 (2019)
14. Shi, L.C. and Lu, B.L.: August. Off-line and on-line vigilance estimation based on linear dynamical system and manifold learning. In 2010 Annual International Conference of the IEEE Engineering in Medicine and Biology, pp. 6587-6590. IEEE. (2010)
15. Tao, X., Hong, X., Chang, X., Dong, S., Wei, X. and Gong, Y.: Few-shot class-incremental learning. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 12183-12192 (2020)
16. Zhao, L.M., Li, R., Zheng, W.L. and Lu, B.L.: March. Classification of five emotions from EEG and eye movement signals: complementary representation properties. In 2019 9th International IEEE/EMBS Conference on Neural Engineering (NER) (pp. 611-614). IEEE. (2019)
17. Zheng, W.L. and Lu, B.L.: Investigating critical frequency bands and channels for EEG-based emotion recognition with deep neural networks. IEEE Transactions on Autonomous Mental Development, 7(3), pp.162-175 (2015)
18. Zheng, W.L., Zhu, J.Y. and Lu, B.L.: Identifying stable patterns over time for emotion recognition from EEG. IEEE Transactions on Affective Computing, 10(3), pp.417-429 (2017)
19. Zheng, W.L., Liu, W., Lu, Y., Lu, B.L. and Cichocki, A.: Emotionmeter: A multi-modal framework for recognizing human emotions. IEEE Transactions on Cybernetics, 49(3), pp.1110-1122 (2018)