

Sleep Quality Estimation with Adversarial Domain Adaptation: From Laboratory to Real Scenario

Jia-Jun Tong¹, Yun Luo¹, Bo-Qun Ma¹, Wei-Long Zheng¹, Bao-Liang Lu^{1,2,3,*}, Xiao-Qi Song⁴, Shi-Wei Ma⁵

¹Center for Brain-like Computing and Machine Intelligence

Department of Computer Science and Engineering

²Key Laboratory of Shanghai Education Commission for Intelligent Interaction and Cognitive Engineering

³Brain Science and Technology Research Center

Shanghai Jiao Tong University, Shanghai, China

⁴China Railway Lanzhou Group Co., Ltd., Gansu, China

⁵China Academy of Railway Sciences, Beijing, China

Abstract—Previous studies on EEG-based last-night sleep quality estimation mainly focus on evaluation with data from laboratory experiments. However, due to the reality gap constituted with device performance, subject groups, experiment settings and controlled conditions, the models trained solely on laboratory data cannot generalize well to real scenarios. In this work, we investigate the sleep quality estimation for high-speed train drivers as an instance of real-scenario application. Domain adaptation models are adopted to deal with the individual differences across subjects when modeling and testing with the real-scenario data. As it is usually difficult and costly to acquire data and annotate them in real scenarios, the high-quality data in laboratory conditions are used for model trainings. Knowledge from simulation is transferred to reality with domain adaptation methods. A novel approach called Domain Adversarial Neural Network (DANN) is adopted. DANN learns domain independent features through deep networks with an adversarial architecture. The experimental results indicate that DANN outperforms other state-of-the-art methods and achieves 19.55% and 23.50% improvements in terms of accuracy on the cross-subject and cross-scenario tasks, respectively, in comparison with the baseline SVM model.

I. INTRODUCTION

Sleep has always been an important topic, no matter in research field or in daily life. Insufficient sleep may lead to serious impairment in daytime performance, increase the risk of driving and occupational accidents and result in diminished life quality [1]. Akerstedt *et al.* studied the relationship between sleepiness and accidents in transport operation. They concluded that sleepiness causes 15% to 20% of all accidents, surpassing alcohol- or drug-related incidents in all modes of transportation, becoming the largest identifiable and preventable cause of accidents [2]. Besides, driver sleepiness has also become a major threat to the railway safety. Bad sleep quality affects drivers' attention, judgment and execution, which may trigger railway fatalities and gravely damage the security of people's lives and properties [3]. Therefore, a reliable measurement of sleep quality, especially last-night sleep quality, becomes necessary.

Currently, there are subjective approaches judging sleep quality by self-evaluation via questionnaires, interviews and

sleep diary. These self-evaluation approaches usually require the participants to consider their sleep habits of the recent times. Pittsburgh Sleep Quality Index (PSQI) [4] and Epworth Sleep Scale (ESS) [5] are representative ones widely adopted by sleep researchers. PSQI limits the time interval to the last month and ESS refers to the participant's usual way of life in recent times. Taking considerably long time into account makes the measurement more precise and robust in distinguishing patients from healthy people. However, these subjective approaches would fail when evaluating the last-night sleep quality due to the constrains of taking a relatively long period into account.

Various objective approaches have been proposed to precisely measure the sleep quality over a single night. Polysomnography (PSG) [6] is a typical objective method of this kind. The PSG monitors body functions including brain, eye movements, muscle activity, skeletal muscle activation and heart rhythm during sleep. The recorded signals are then analyzed by a doctor who is capable of judging the sleep quality. This approach requires a PSG device attached to the subject during the whole sleeping process, which makes it not an applicable choice in real scenarios. For instance, if PSG is applied to high-speed train driver sleep quality checking, the company would have to purchase a PSG device for every driver, which would be rather expensive, and the sleep quality checking procedure would also be time-consuming and inefficient. Therefore, a cheaper and easier method is in demand.

With the rapid development of wearable electroencephalography (EEG) signal acquiring devices, EEG-based last-night sleep quality measurement is considered as a feasible choice. This approach does not require EEG during the whole sleep process, but only acquires EEG in a short time after the subject wakes up. Ideally, these EEG data would be containing information of the subject's last-night sleep quality. Most EEG-based sleep studies are done in the laboratory environment with precisely controlled sleep duration or sleep deprivation. With different controlled conditions, such as restricting sleep time to be 4, 6 or 8 hours [7], the subjects would be in totally different mental states, which makes it easier to distinguish

*Corresponding author: Bao-Liang Lu (blu@sjtu.edu.cn)

between different sleep quality with EEG signals. These studies achieved promising results on restricted laboratory environments. However, to the best of our knowledge, no studies under real scenarios has been performed yet.

In this paper, we aim to apply the EEG-based last-night sleep quality estimation approach to real scenarios. To acquire real-scenario data, we perform experiments on high-speed train drivers with different last-night sleep qualities at their workplace. We first perform a cross-subject task on real-scenario data. Domain adaptation (DA) methods are adopted to reduce the individual difference across subjects. As the data collection and annotation in real scenarios are usually costly and difficult, using only the real-scenario data to train models is not practicable. We introduce the laboratory data used our in previous work and perform a cross-scenario task in which we train models on laboratory data and test them on real-scenario data. DA methods are also adopted to narrow the reality gap between these two scenarios. Among the adopted DA methods, the Domain Adversarial Neural Network (DANN) method performs the best in terms of classification accuracy with considerable improvement over the other methods.

II. RELATED WORK

Sleep has been a hot research topic for decades. In the early days, sleep is defined as a human maintaining a specific body gesture and behavioral quiescence with elevated arousal threshold and state reversibility after stimulation [8].

Since the development of technology to amplify and record spontaneous EEG signals, it has been demonstrated that there is a strong correlation between EEG and sleep. Hori *et al.* proposed a nine-stage EEG-based system for identifying successive EEG changes throughout the sleep onset period which represents the greatest advance in understanding the sleep onset process yet achieved [8]. Wolpert *et al.* proposed a sleep stage character system depending on brain waves [9], in which a distinct and typical EEG and physiologic patterns for alert, REM sleep and each stage of NREM has been found. Therefore, researchers in sleep field have always been paying close attention to EEG signals.

In the last decade, the differences and changes of EEG under different sleep qualities have been investigated. Li *et al.* tried to study how lack of sleep influences the event-related potentials under stimulation [10]. Audio stimulation was given to subjects in both cases of sufficient sleeping and lack of sleep and event-related potentials of the subjects were analyzed with the parallel factor analysis method. They found that when subjects are lack of sleep, event-related potentials activities of the subjects tend to appear near to the forehead and Gamma frequency band delay and attenuation would occur. Na *et al.* investigated effects of sleep deprivation on the connection between cerebral hemisphere. They conducted a mutual information analysis on the EEG acquired from normal sleep and sleep deprived subjects and concluded that the connection between cerebral hemisphere was weaker in case of sleep deprivation [11]. Tassi *et al.* focused on the correlation

between partial sleep deprivation and the EEG activity afterward. The spectral analysis applied on the waking EEG during cognitive performance testing shows that alpha activities were increased in both deprived and normal subjects but theta power increased only in the sleep deprived group [12]. The studies mentioned above mainly focus on the comparison between controlled experiment of two cases: the deprived group and the normal group. They usually analyze EEG data in a statistical manner and achieve qualitative conclusions.

In recent years, a minority of researchers take up research into quantitative last-night sleep quality estimation with machine learning methods. Wang *et al.* attempt to measure last-night sleep quality from resting EEG signals [3]. They designed an experiment collecting resting EEG signals from subjects under three discrete sleep conditions: 8 hours sleep, 6 hours sleep, and 4 hours sleep, according to which the EEG data are labeled as three categories. To correctly classify the EEG data, they introduced discriminative graph regularized extreme learning machine together with minimal-redundancy-maximal-relevance feature selection algorithm and achieved a mean accuracy of 83.5% in a leave-one-subject-out validation manner [3]. Zhang *et al.* followed that work and modified the EEG acquisition procedure, which we discuss in the experiment section [7]. They extracted subject independent features with three domain adaptation methods and feed them into SVM to make comparisons. With considerable model performance improvement, they concluded that domain adaptation approaches do have the capability to reduce the differences of EEG data across subjects and sessions [7].

In this work, we focus on last-night sleep quality estimation in real scenarios. In our previous work [3] [7], EEG signals are acquired with wet-electrode devices, which is not applicable to real scenarios due to the laborious preparation process and poor portability of wet-electrode devices. Therefore, we use a dry-electrode device to perform EEG acquiring experiments on high-speed train drivers. We also adapt our experiment settings to the real-scenario conditions. Differences of devices, subject groups, experiment settings and controlled conditions lead to a giant reality gap between laboratory and real-scenario data, making the last-night sleep quality estimation in real scenario a much more challenging task than in laboratory conditions. To address the domain difference problem, we adopt not only traditional DA methods but also a novel DANN method with the capability of learning domain independent features through deep neural networks with an adversarial architecture. DANN outstandingly outperforms other models and turns out to be stable and robust according to our experimental results.

III. METHODS

A. Feature Extraction

EEG signals can be divided into five different frequency bands: Delta (1~4 Hz), Theta (4~8 Hz), Alpha (8~13 Hz), Beta (13~30 Hz) and Gamma (30~50 Hz) bands [13]. The feature extraction procedure converts the time domain signals to the frequency domain and then extracts useful information for the five frequency bands. Differential entropy (DE) features

are commonly used for its reflecting the energy change of EEG signals [14]. We extracted DE features from the 18-channel EEG signals in 5 frequency bands, which add up to a feature vector of length 90. Since the extracted features contain fluctuations caused by noise, we use the linear dynamical system to remove those fluctuations [15].

B. Domain Adaptation models

Domain adaptation helps to transfer knowledge from a source domain to a different but related target domain, even though these domains may have different distributions. Most existing DA methods aim at finding a new feature representation reducing the difference between the distributions of the source and target domains, meanwhile preserving the data properties of the source and target domains. In this work, we would apply DA methods to transfer knowledge from our laboratory experiments to the real-scenario (high-speed train drivers) experiments. Below, we introduce the eight methods adopted in this work.

1) *Transfer Component Analysis (TCA)* [16]: TCA aims to learn a transformation that reduces the distance between the marginal distributions and preserves the important properties of both domains. TCA empirically measures the distance between the source and target distributions by the distance between the empirical means of the two domains in Reproducing Kernel Hilbert Space (RKHS). Besides, to preserve the main properties of both domains, the TCA chooses to maximally preserve their variance.

2) *Information-Theoretical Learning (ITL)* [17]: ITL assumes that in a latent feature space induced by a linear transformation $\mathbf{L} \in \mathbb{R}^{d \times D}$, data in source and target domains are clustered according to their labels and the clusters from two domains that possess the same label are geometrically close. Then a k -nearest neighbors model is used to estimate the label of each sample in the target domain. The negated mutual information between the target data and the estimated label is used to approximate the classification error, which is to be minimized. Another k -nearest neighbors model is used to estimate the domain label of each sample in both domains. To express the desideratum that clusters possessing the same label are geometrically close, the negated mutual information between samples and domain labels is maximized.

3) *Geodesic Flow Kernel (GFK)* [18]: GFK parameterizes the process of the source domain smoothly changing to the target domain with a geodesic flow $\Phi(t)$, where $t \in [0, 1]$. $\Phi(t)$ smoothly changes from the basis of subspace for the source domain to that of the target domain when t gradually changes from 0 to 1. GFK then expands the original feature \mathbf{x} by projecting it onto all these subspaces, which gives us an infinite-dimensional feature vector \mathbf{z}^∞ . Then the inner product in this infinite dimensional space would be

$$\langle \mathbf{z}_i^\infty, \mathbf{z}_j^\infty \rangle = \int_0^1 (\Phi(t)^T \mathbf{x}_i)^T (\Phi(t)^T \mathbf{x}_j) = \mathbf{x}_i^T \mathbf{G} \mathbf{x}_j. \quad (1)$$

This is a ‘‘kernel trick’’, and the matrix \mathbf{G} has a closed form solution.

4) *Joint Distribution Adaptation (JDA)* [19]: JDA hopes to find projection matrix \mathbf{A} that adapts the joint distributions and the conditional distribution of features \mathbf{x} and labels \mathbf{y} simultaneously. Since there is no labeled data in the target domain, it uses a classifier f to generate pseudo target label as an approximation of actual label. In order to achieve more accurate approximation, an iterative pseudo label refinement strategy is introduced to refine the transformation and the classifier.

5) *Subspace Alignment (SA)* [20]: SA represents the source and target domains as subspaces described by eigenvectors. It then seeks a mapping function that aligns the source and target subspaces as a domain adaptation solution. Using PCA, it selects d eigenvectors for each domain as their corresponding source (\mathbf{X}_S) and target (\mathbf{X}_T) subspaces. Then, a linear transformation matrix \mathbf{M} is used to align \mathbf{X}_S with \mathbf{X}_T . \mathbf{M} is learned by minimizing the Bregman matrix divergence.

6) *Transfer Joint Matching (TJM)* [21]: TJM proposes to adapt domains with a transformation T that (a) minimize the *Maximum Mean Discrepancy* distance between the source and target domains, and (b) reweight the source instances through a structured sparsity. TJM works in RKHS to match both first- and high-order statistics. As there is a row-sparsity regularizer in the optimization objective, TJM solves the problem by an iterative approach.

7) *Maximum Independence Domain Adaptation (MIDA)* [22]: MIDA shares the same intuition with TCA. It aims at minimizing the domain difference while preserving the data properties with transformation Φ . With kernel trick, the kernel \mathbf{K} matrix can be easily induced. A projection matrix $\tilde{\mathbf{W}}$ is used to transform the kernel map. MIDA assumes that $\tilde{\mathbf{W}} = \Phi(\mathbf{X})\mathbf{W}$. MIDA measures the dependence between projected samples and domain features by Hilbert-Schmidt independence criterion (HISC). Besides, MIDA preserves the properties of data by keeping the variance in the transformed feature space.

8) *Domain Adversarial Neural Network (DANN)* [23]: DANN is a domain adaptation approach with deep architectures that can be trained with labeled source domain data and unlabeled target domain data. Its adaptation behavior is achieved by augmenting a normal feed-forward model with a few standard layers that work as a domain discriminator and a gradient reversal layer that makes the model be trained procedure in an adversarial manner. Figure 1 depicts the architecture we adopt in this work.

DANN estimates the dissimilarity of domains in feature representation f by looking at the loss of domain classifier G_d . Therefore, at training time, DANN seeks the parameters θ_f that maximize the loss of G_d and meanwhile, seeks for the parameters θ_d that minimize the loss of G_d . By introducing a gradient reversal layer (GRL) $R_\lambda(\cdot)$ between the feature f and the domain classifier G_d , one can define GRL as a pseudo-function by its forward- and back-propagation behaviors:

$$R_\lambda(f) = f \quad (2)$$

$$\frac{\partial R_\lambda(f)}{\partial f} = -\lambda I, \quad (3)$$

where λ is a tradeoff parameter. The objective function can then be optimized with statistic gradient descent (SGD) and also other optimization methods based on SGD.

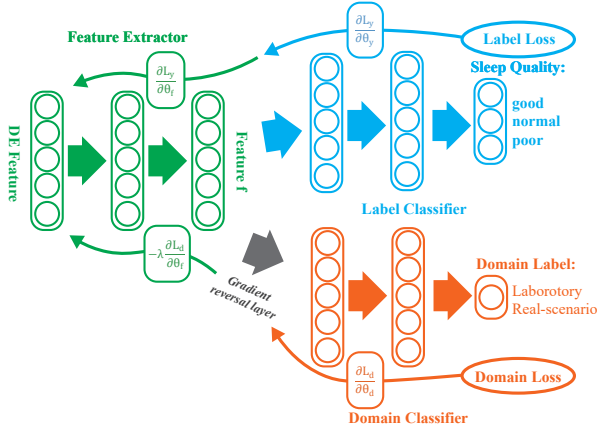


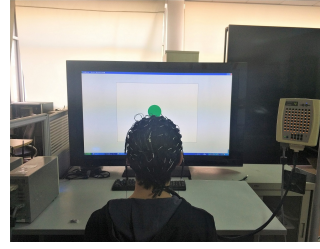
Fig. 1. The DANN model with three components: a feature extractor G_f , a label predictor G_y and a domain classifier G_d . The feature extractor maps the input x to a feature vector $f \in \mathbb{R}^m$ with all of its layer parameters denoted as θ_f , i.e. $f = G_f(x; \theta_f)$. The label predictor maps f to a predicted label \hat{y} with all of its layer parameters denoted as θ_y . The domain classifier maps f to a predicted domain mark \hat{d} with parameters denoted as θ_d .

IV. EXPERIMENTS

Since we aim at transferring knowledge from the laboratory to the real scenario, it is necessary to have data from both scenarios. The laboratory data used in this work are acquired in our previous studies of sleep quality estimation [7], and in this work, we collected the real-scenario data from high-speed train drivers. Here, for comparison between these two types of experiments and easy understanding of how different the real-scenario data is from the laboratory ones, we will describe both of them in details. The data are preprocessed after acquisition and DE features are extracted from the preprocessed data. We then perform the cross-subject and cross-scenario tasks on the extracted features. A baseline SVM model and eight DA methods are then adopted in both tasks.

A. Laboratory Sleep Experiment

The Laboratory experiment aims at collecting resting EEG after sleep of different quality. Each experiment consists of two parts: sleep and EEG acquisition. 10 subjects (six males and four females, age range: 21-26, mean: 23.57, std: 1.62) were recruited for the experiments. At the night before EEG acquisition, the subject was required to sleep for a specified amount of time: 4, 6 or 8 hours corresponding to a good, normal or poor sleep quality [3]. To ensure that the subject slept for the time as required, a smart band was attached to them monitoring their wrist motions, which can rarely be detected when human fall asleep [7]. After the subject woke up, the EEG acquisition procedure started in less than an hour and lasted for 30 minutes. The EEG signals were recorded with a 62-channel wet electrode cap using the ESI NeuroScan system. Electrodes on the cap are placed according to the



(a) Laboratory environment while acquiring EEG for student subjects



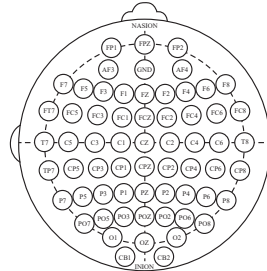
(b) Real-scenario experiment environment while acquiring EEG for high-speed train drivers



(c) Wet electrodes EEG cap



(d) Dry electrodes DSI-24 headset



we performed our experiment on the high-speed train drivers right before they start to work. 70 subjects (70 males, age range: 25-49, mean: 37.10, std: 4.69) are recruited. At the night before EEG acquisition, the sleep time of the subjects were not controlled. Before the experiments began, the subjects were required to fill out a questionnaire about their sleep last night. The questions investigated sleep duration, dreams and current feelings. In this work, we use the sleep duration as an index of subjects' sleep quality, as shown in Table I.

The EEG acquisition procedure is performed 30 minutes right before the subjects went to work. During EEG recording, the subject was required to stare at a green dot on the screen for 1 minute and close for another 1 minute alternately for 3 times. During the whole session, the subjects were required to count silently to keep a peaceful and concentrated state. The EEG signals were recorded with an 18-channel dry electrode cap using the DSI-24 device at a sampling rate of 300 Hz. Electrodes on the cap are placed according to the international 10-20 system. The experiment environment and dry electrode DSI-24 device are shown in the right column of Figure 2.

C. Differences between experiments

In order to keep the nature of the real scenario, instead of imposing restrictions on subjects, we adapt our experiment to the limitations in the real world, thus making the real-scenario experiments vary from the laboratory ones in many different ways. In order to keep our experiment short to fit the time limitation, we shrink our actual EEG recording time from 30 minutes to 6 minutes. In laboratory experiments, 64-channel wet electrode devices are used to record the EEG signals. However, these devices require laborious setup processes, therefore do not meet the demand of frequent acquisition of EEG signals in the real scenario. This demand is well addressed with the dry electrode device DSI-24 headset, which is more convenient to use. Moreover, we do not control the subjects' sleep time and experiment start time, so that the mental states of subjects during the experiments are close to those in their daily lives. For the ease of understanding, we summarized the difference between experiments in different scenarios as Table II.

D. Preprocessing

Some of the hardware induced difference can be eliminated by preprocessing methods. Preprocessed data in both scenarios would be sharing the same format and we can apply same subsequent processing procedures to them.

1) *Electrodes Selection:* In the laboratory, a wet electrodes cap with 62 channels was used and in the real scenario, a dry electrodes headset with 18 channels was used. The former one's electrodes are placed according to the higher-resolution international 10-20 system and the later according to the normal international 10-20 system. As a matter of fact, extra electrodes in higher resolution 10-20 system are added using the 10% division, which fills in intermediate sites halfway between those of the existing 10-20 system. In other words, if we remove the extra electrodes in the higher resolution

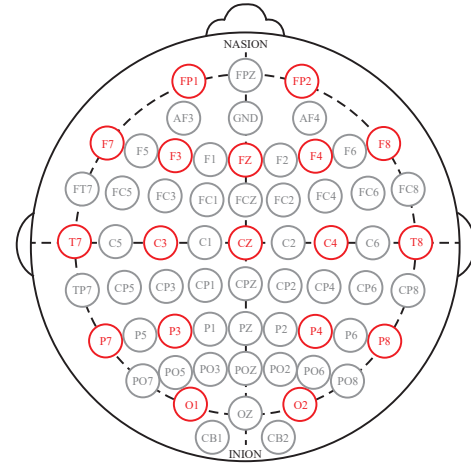


Fig. 3. Electrodes selection on laboratory device. Those electrodes marked red are selected ones matching electrodes on DSI-24 headset. The CMF electrode on DSI-24 headset dose not record any signal, so the corresponding PZ electrode is not marked.

system, the electrodes left would match perfectly with the normal system. We selected the corresponding electrodes from the laboratory device, as shown in Figure 3.

2) *Resetting Reference:* As shown in Table II, the EEG data recorded with different devices have inconsistent references. The real-scenario data use the mean of A1 and A2 electrodes as references. However, the laboratory data use a specific REF electrode as the reference. Besides, the wet-electrode cap we used does not have any of the A1 or A2 electrodes, and also the DSI-24 headset does not have a REF electrode. To address the inconsistency, we decided to use the mean signal of all 18 electrodes as reference [7].

E. Modeling

We extract DE features for EEG signals in both scenarios with time window set to 1s. For the laboratory data, we collect EEG signals for 10 subjects with 3 experiments per subject and each experiment lasts about 30 minutes. Therefore, we have 58172 laboratory samples. For the real-scenario data, we collected EEG signals for 70 subjects with 1 experiment per subject and the valid time of each experiment is 3 minutes exactly. Therefore, we have 12600 real-scenario samples.

We mainly perform two modeling tasks in this work: (a) training and testing on real-world data, and (b) training on laboratory data but testing on real-world data. While performing the first task, the model performance is evaluated in a 5-fold cross validation manner. The real-world data is split into 5 parts, each containing 14 subjects and no subject's data should appear in multiple validation parts. Therefore, in the following discussion we refer to the first task as the **Cross-Subject** one and the second as the **Cross-Scenario** one. For the ease of comparison between two tasks, we split the real-scenario data into 5 parts in the same way as the cross-subject task while evaluating on the cross-scenario task.

After the feature extraction, We use the LIBLINEAR [24] SVM with default parameters as the baseline model on both

TABLE II
DIFFERENCES BETWEEN LABORATORY AND REAL-SCENARIO EXPERIMENTS

Categories	Specifications	Laboratory	Real-scenario
Subjects	Age	21-26	25-49
	Sex-distribution	4 males, 6females	70 males
	Occupation	Students	High-speed train drivers
	Experiments per subject	3	1
Controlled conditions	Sleep time	4,6,8 hours	Not controlled
	Monitoring sleep	Yes	No
	Head Cleaning	Yes	No
	Experiment Starts	≤ 1 hour after wake up	1-10 hour after wake up
Experiment settings	Experiment duration	30 min	6 min
	Task	Eyes open	Eyes open and close
	Environment	Quite and isolated	Noisy
Devices	Electrods type	Wet electrodes	Dry electrodes
	Resolution	62 channels	18 channels
	Reference	REF(Between CPZ and PZ)	A1 and A2 (On the ears)
	Common Mode Follower	0	1 (CMF)
	Sample rate	1000 Hz	300 Hz

tasks. In the cross-subject task, while performing a 5-fold cross-validation, we view the fold held out as the target domain and the other 4 folds together as the source domain. In the cross-scenario task, we consider the whole real-world dataset as the target domain and the laboratory dataset as the source domain. TCA, ITL, GFK, JDA, SA, TJM and MIDA are run on all data as dimensionality reduction step and then SVM classifier is trained on transformed labeled source domain data and tested on the transformed target domain data. To boost the performance of these models, we pick the best subspace dimension parameter from the range $\{10, \dots, 80\}$. The DANN model was trained with labeled source domain data and unlabeled target domain data, and the target domain labels can only be seen in evaluation steps. The tradeoff parameter λ of the GRL in DANN is set to 0.01. In the following discussions, we refer to the DA models except for DANN as the **baseline DA models** and refer to the baseline model as the **baseline SVM model**.

V. RESULTS AND DISCUSSION

We first compare DANN with baseline DA methods and also the baseline SVM model in terms of classification accuracy. The classification accuracies of cross-subject and cross-scenario models are summarized as Table III and IV.

DANN achieves much better performance than baseline DA models on both tasks. The average classification accuracy of DANN on cross-subject task is 73.72%. The performance improved 6.97% compared to MIDA (best baseline DA model), and 19.55% compared to the baseline SVM model. On the cross-scenario task, the accuracy of DANN model is 65.29%, which is 7.83% and 23.50% higher than the MIDA (best baseline DA model) and the baseline SVM model, respectively. We also observed that the standard deviation of DANN test accuracy across the 5 folds is lower than all the baseline DA models. Especially on the cross-scenario task, the lowest standard deviation of baseline DA models is 10.6%, while

DANN achieves a standard deviation of 3.80%. This suggests that DANN is more stable than baseline DA models.

TABLE III
ACCURACY (%) OF CROSS-SUBJECT TASK

Method	Fold 1	Fold 2	Fold 3	Fold 4	Fold 5	Mean \pm Std
SVM	59.57	48.41	57.42	55.71	49.74	54.17 \pm 4.87
TCA	58.37	72.98	63.77	69.48	53.85	63.69 \pm 7.82
ITL	66.07	66.79	57.14	76.43	56.71	64.63 \pm 8.13
GFK	59.41	66.79	53.73	58.45	58.49	59.37 \pm 4.70
JDA	59.41	59.52	76.11	60.44	55.91	62.28 \pm 7.92
SA	77.42	67.38	60.36	70.68	46.47	64.46 \pm 11.7
TJM	65.91	74.72	58.73	62.22	59.37	64.19 \pm 6.53
MIDA	71.35	69.76	66.91	69.21	56.51	66.75 \pm 5.94
DANN	67.50	74.64	71.23	76.98	78.25	73.72\pm4.38

TABLE IV
ACCURACY (%) OF CROSS-SCENARIO TASK

Method	Fold 1	Fold 2	Fold 3	Fold 4	Fold 5	Mean \pm Std
SVM	43.86	40.35	42.18	38.98	43.59	41.79 \pm 2.10
TCA	54.60	61.19	38.89	67.34	56.87	55.60 \pm 10.6
ITL	53.69	63.10	24.68	60.71	49.41	50.32 \pm 13.7
GFK	49.64	62.46	40.99	70.95	49.25	54.66 \pm 11.9
JDA	47.70	61.43	35.71	66.11	66.94	55.58 \pm 13.5
SA	42.74	62.54	37.50	61.23	54.17	51.64 \pm 11.1
TJM	41.67	60.28	38.97	59.88	60.60	52.28 \pm 10.9
MIDA	56.47	66.35	38.69	68.93	56.87	57.46 \pm 11.9
DANN	69.10	64.68	60.79	62.6296	69.25	65.29\pm3.80

We can also see that all the transfer learning methods achieved a better classification accuracy than the baseline SVM model on both tasks. With domain adaptation, the mean accuracy is improved at least 5.20% and 9.85% on cross-subject and cross-scenario, respectively. We can see that the DA methods do eliminate the difference between the source and target domain. However, the major limitation of baseline DA models is that they use predefined functionalities to measure the domain difference and transform the data in a

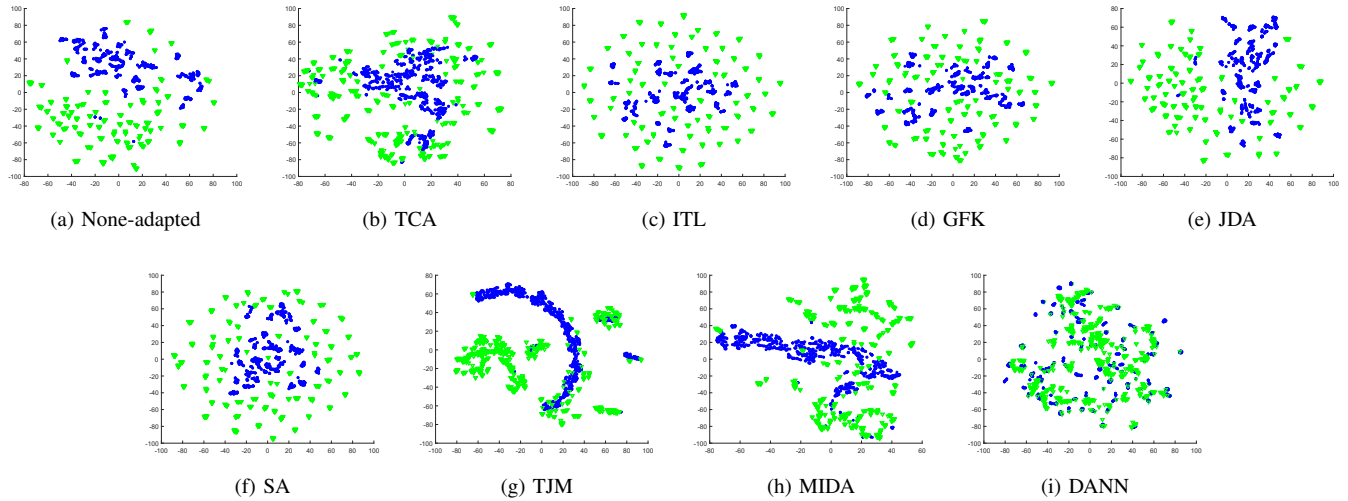


Fig. 4. The effects of adaptation on cross-scenario task (Best viewed in color). The figure shows t-SNE [25] visualizations of features. (a) shows features with no adaptation; (b)-(h) show features after transformation by baseline DA methods; (i) shows the feature transformed by the feature extractor of DANN. Blue dots correspond to the source domain examples, while green triangles represent the target domain ones. The source and target domain features after transformation are much closer in DANN than baseline DA models

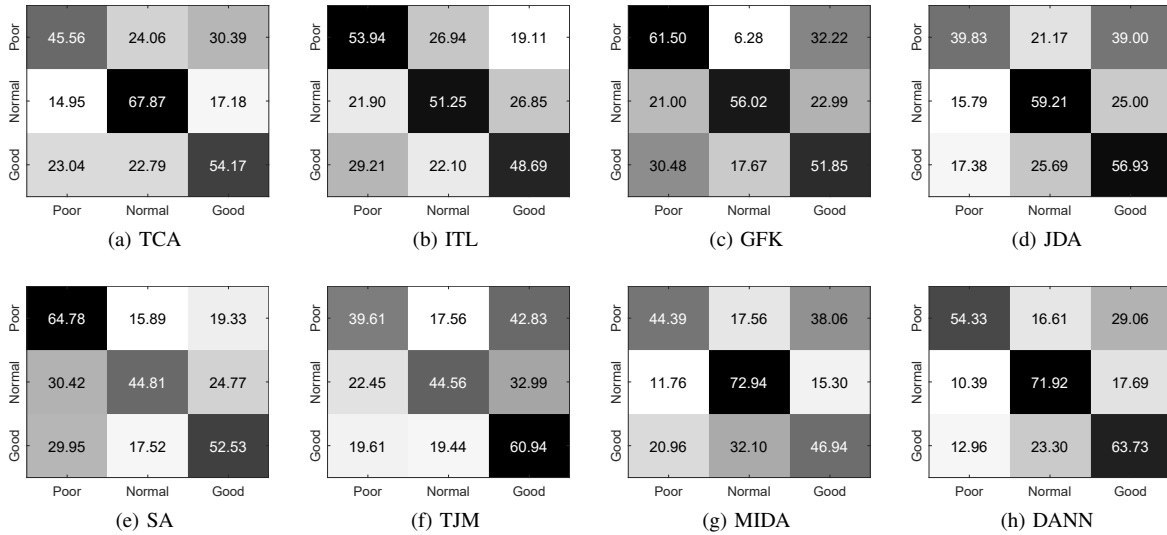


Fig. 5. The confusion matrix of DA models on cross-scenario task. The vertical axis represents the labels marked by sleep duration according to Table I; The horizontal axis represents the predicted label of each DA model. The number on each block represents the percentage of samples with that the true label (vertical ones) being classified to the predicted label (horizontal ones).

predetermined way, which is rather rigid compared to the DANN method. DANN represents both the domain difference and the transformation with multilayer neural networks and can be trained in an adversarial manner with SGD based methods, which makes it easier to reduce the domain discrepancy between source and target domains. Figure 4 shows how DANN succeeded at aligning feature distributions on the cross-scenario task. Another drawback of the insufficient adaptation of baseline DA methods is that it may cause negative transfer in some components of the target domain. We can observe this phenomenon in Table IV, where negative transfer occurs in Fold 3 on TCA, ITL, JDA, SA and TJM models, and yet DANN shows no negative transfer in all 5 Folds.

It is noteworthy that the real-scenario data is severely unbalanced, with 1800 poor-sleep-quality samples, 4320 normal-sleep-quality samples and 6480 good-sleep-quality samples. Thus the classifiers may benefit from biased predictions. For a simple example, if one classifies all samples to be of good sleep quality, it would achieve over 50% classification accuracy. With that in mind, we visualize the confusion matrix of each model on cross-scenario task to see whether the models have made biased predictions (see Figure 5). We observed that TJM and JDA do suffer from the biased prediction problem. TJM achieved a classification accuracy of 60.94% on good-sleep-quality samples while its accuracy drops sharply to 39.61% on poor-sleep quality-samples, which is close to

a random prediction result. In contrast, DANN's prediction preserves a 54.33% classification accuracy on poor-sleep-quality samples.

Lastly, we compare the performance of models on different tasks. All models perform with a relatively lower classification accuracy on the cross-scenario task compared to the cross-subject task, which is predictable as the former is more challenging. As we have illustrated in Table II, the cross-scenario task is dealing with the reality gap constituted with differences of four categories. The intersubject variation that the cross-subject task focuses on is merely one of these categories. Therefore, the baseline SVM model performs with a low accuracy (41.79%) on cross-scenario task. In spite of the difficulties, DANN achieves remarkable improvement over baseline DA methods, which shows that DANN performs more robustly on difficult domain adaptation tasks than baseline DA models.

VI. CONCLUSION

We have studied last-night sleep quality estimation in real scenarios. To acquire the real-scenario data, we have performed EEG acquisition for high-speed train drivers without changing their work schedules. In order to meet the demand of practical application, we have introduced a dry electrode device to collect the EEG signals. To estimate the sleep quality of the drivers, we have evaluated two tasks: the cross-subject one and the cross-scenario one. Eight DA methods have been adopted to both tasks to reduce the domain discrepancy between source and target domains. In this work, we have adopted a novel DANN approach, which is based on neural networks with an adversarial architecture. The results have shown that DANN approach considerably outperforms baseline models on both tasks in terms of classification accuracy. Moreover, we have also analyzed the models from various perspectives and concluded that the DANN approach is also stable and robust on the challenging cross-scenario task.

ACKNOWLEDGMENT

This work was supported in part by grants from the National Key Research and Development Program of China (Grant No. 2017YFB1002501), the National Natural Science Foundation of China (Grant No. 61673266), the Major Basic Research Program of Shanghai Science and Technology Committee (Grant No. 15JC1400103), ZBYY-MOE Joint Funding (Grant No. 6141A02022604), the Technology Research and Development Program of China Railway Corporation (Grant No. 2016Z003-B), and the Fundamental Research Funds for the Central Universities.

REFERENCES

- [1] D. J. Buysse, C. F. Reynolds, T. H. Monk, S. R. Berman, and D. J. Kupfer, "The pittsburgh sleep quality index: A new instrument for psychiatric practice and research," *Psychiatry Research*, vol. 28, no. 2, pp. 193–213, 1989.
- [2] T. Akerstedt, "Consensus statement: fatigue and accidents in transport operations," *Journal of Sleep Research*, vol. 9, no. 4, pp. 395–395, 2000.
- [3] L.-L. Wang, W.-L. Zheng, H.-W. Ma, and B.-L. Lu, "Measuring sleep quality from EEG with machine learning approaches," in *the IEEE International Joint Conference on Neural Networks*, 2016, pp. 905–912.
- [4] D. J. Buysse, C. F. Reynolds, T. H. Monk, S. R. Berman, and D. J. Kupfer, "The pittsburgh sleep quality index: a new instrument for psychiatric practice and research," *Psychiatry Research*, vol. 28, no. 2, pp. 193–213, 1989.
- [5] M. W. Johns, "A new method for measuring daytime sleepiness: the epworth sleepiness scale," *Sleep*, vol. 14, no. 6, pp. 540–545, 1991.
- [6] Y.-K. Teng, L.-C. Chiang, K.-H. Lue, S.-W. Chang, L. Wang, S.-P. Lee, H. Ting, and S.-D. Lee, "Poor sleep quality measured by polysomnography in non-obese asthmatic children with or without moderate to severe obstructive sleep apnea," *Sleep Medicine*, vol. 15, no. 9, pp. 1062–1067, 2014.
- [7] X.-Z. Zhang, W.-L. Zheng, and B.-L. Lu, "EEG-based sleep quality evaluation with deep transfer learning," in *International Conference on Neural Information Processing*. Springer, 2017, pp. 543–552.
- [8] R. D. Ogilvie, "The process of falling asleep," *Sleep Medicine Reviews*, vol. 5, no. 3, pp. 247–270, 2001.
- [9] E. A. Wolpert, "A manual of standardized terminology, techniques and scoring system for sleep stages of human subjects," *Archives of General Psychiatry*, vol. 20, no. 2, pp. 246–247, 1969.
- [10] N. Li, Y. Wang, M. Wang, and H. Liu, "Effects of sleep deprivation on gamma oscillation of waking human EEG," *Progress in Natural Science*, vol. 18, no. 12, pp. 1533–1537, 2008.
- [11] S. H. Na, S.-H. Jin, and S. Y. Kim, "The effects of total sleep deprivation on brain functional organization: mutual information analysis of waking human EEG," *International Journal of Psychophysiology*, vol. 62, no. 2, pp. 238–242, 2006.
- [12] P. Tassi, A. Bonnefond, O. Engasser, A. Hoeft, R. Eschenlauer, and A. Muzet, "EEG spectral power and cognitive performance during sleep inertia: the effect of normal sleep duration and partial sleep deprivation," *Physiology & Behavior*, vol. 87, no. 1, pp. 177–184, 2006.
- [13] K. Mann, P. Bäcker, and J. Röschke, "Dynamical properties of the sleep EEG in different frequency bands," *International Journal of Neuroscience*, vol. 73, no. 3-4, pp. 161–169, 1993.
- [14] L.-C. Shi, Y.-Y. Jiao, and B.-L. Lu, "Differential entropy feature for EEG-based vigilance estimation," in *35th Annual International Conference of the IEEE Engineering in Medicine and Biology Society*, 2013, pp. 6627–6630.
- [15] X.-W. Wang, D. Nie, and B.-L. Lu, "Emotional state classification from EEG data using machine learning approach," *Neurocomputing*, vol. 129, pp. 94–106, 2014.
- [16] S. J. Pan, I. W. Tsang, J. T. Kwok, and Q. Yang, "Domain adaptation via transfer component analysis," *IEEE Transactions on Neural Networks*, vol. 22, no. 2, p. 199, 2011.
- [17] G. Niu, B. Dai, M. Yamada, and M. Sugiyama, "Information-theoretic semi-supervised metric learning via entropy regularization," in *29th International Conference on Machine Learning*, July 2012, pp. 89–96.
- [18] B. Gong, Y. Shi, F. Sha, and K. Grauman, "Geodesic flow kernel for unsupervised domain adaptation," in *the IEEE Conference on Computer Vision and Pattern Recognition*, June 2012, pp. 2066–2073.
- [19] M. Long, J. Wang, G. Ding, J. Sun, and P. S. Yu, "Transfer feature learning with joint distribution adaptation," in *the IEEE International Conference on Computer Vision*, 2013, pp. 2200–2207.
- [20] B. Fernando, A. Habrard, M. Sebban, and T. Tuytelaars, "Unsupervised visual domain adaptation using subspace alignment," in *the IEEE International Conference on Computer Vision*, 2013, pp. 2960–2967.
- [21] M. Long, J. Wang, G. Ding, J. Sun, and P. S. Yu, "Transfer joint matching for unsupervised domain adaptation," in *the IEEE Conference on Computer Vision and Pattern Recognition*, 2014, pp. 1410–1417.
- [22] K. Yan, L. Kou, and D. Zhang, "Learning domain-invariant subspace using domain features and independence maximization," *IEEE Transactions on Cybernetics*, vol. 48, no. 1, pp. 288–299, Jan 2018.
- [23] Y. Ganin and V. Lempitsky, "Unsupervised domain adaptation by backpropagation," in *the International Conference on Machine Learning*, 2015, pp. 1180–1189.
- [24] R.-E. Fan, K.-W. Chang, C.-J. Hsieh, X.-R. Wang, and C.-J. Lin, "LIBLINEAR: A library for large linear classification," *Journal of Machine Learning Research*, vol. 9, no. Aug, pp. 1871–1874, 2008.
- [25] D. Vazquez, A. M. Lopez, J. Marin, D. Ponsa, and D. Geronimo, "Virtual and real world adaptation for pedestrian detection," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 36, no. 4, pp. 797–809, 2014.