

EEG-based Person Authentication System Using A Self-paced Reaching Task

Abstract

Electroencephalography (EEG) has been widely investigated for person authentication because of its advantages of being difficult to fake, impossible to observe or intercept, and unique. In this study, we proposed a new person authentication system based on EEG signals. The self-paced reaching task was applied as it was a natural and common human daily task. Power spectral density (PSD) of delta, theta, alpha and beta bands were extracted from EEG signals of channel C3 as features. After comparing different classification algorithms' effect, a support vector machine was applied as a classifier. Our study showed the superior performance of alpha and beta band PSD features, compared to delta and theta bands. The time course characteristics of alpha and beta bands were also studied. It was showed that alpha band PSD features in $[-2\ 1]$ s and the beta band PSD features in $[-1\ 0]$ s and $[0\ 1]$ s achieved the best results. Moreover, the AUC of the person authentication performance decreased when the combined features of PSD and AR parameters were applied. Overall, the best average accuracy of 78.0% was achieved, which demonstrates the possibility of using a reaching task for person authentication.

Key words: EEG, Person authentication, SVM, PSD, AR

Contents

I.Introduction.....	1
II.Experiment	5
III.Methods	9
1. Features.....	9
1.1. Power Spetral Density	9
1.2. Autoregressive Model	9
2. Classification Models	10
3.2.1 Support Vector Machine	10
3.2.2 Random Forest.....	11
IV.Results	13
1. PSD Analysis.....	13
2. Authentication Performance Comparison of Different Algorithms - SVM/ RF/BP	13
3. Authentication Performance of PSD Features using SVM	14
4.3.1 Authentication Results for PSD Features of Four Frequency Bands	14
4.3.2 Time Course Authentication Performances for Alpha and Beta Band PSD Features	14
4.3.3 Authentication Performances of Combined Alpha and Beta Band PSD Features	15
4.4 Authentication Results for Combined PSD and AR Features using SVM.....	16
V.Conclusions.....	19
References	20
Acknowledgements.....	22

I.Introduction

Associating an identity with an individual is called personal identification. The problem of resolving the identity of a person can be categorized into two fundamentally distinct types of problems with different inherent complexities: (i) authentication and (ii) recognition. Authentication refers to the problems of confirming or denying a person's claimed identity (Am I who I claim I am?). Recognition refers to the problems of establishing a subject's identity (Who I am?) – either from a set of already known identities (closed identification problem) or otherwise (open identification problem).

A number of situations require an identification of a person in our society; accurate identification of a person could defer crime and fraud, streamline business processes and save critical resources. However, the problem of authentication and recognition is very challenging.

Approaches to that problem are to reduce it to the problem of a concrete entity related to the person. Typically, these entities include (i) a person's possession ("something that you possess"), e.g., permit physical access to a building to all persons whose identity could be authenticated by possession of a key; (ii) a person's knowledge of a piece of information ("something that you know"), e.g., permit login access to a system to a person who knows the user-id and a password associated with it. For example, ATMs use a combination of "something that you have" (ATM card) and "something that you know" (PIN) to establish an identity. The problem with the traditional approaches of identification using possession as a means of identity is that the possessions could be lost, stolen, forgotten, or misplaced. Once in control of the identifying possession, any other "unauthorized" person could abuse the privileges of the authorized user. The problem with using knowledge as an identity authentication mechanism is that it is difficult to remember the passwords/PINs; easily recallable passwords/PINs could be easily guessed by the adversaries.

Yet another approach to positive identification has been to reduce the problem of identification to the problem of identifying physical characteristics of the person, namely biometrics. [1]

1. What is Biometrics?

Biometrics, as the measurement of a person's physical features, actions, or behavioral characteristics that are distinguished between individuals, has a long history. Technological developments and intensive research have netted innovative concepts, features, and data processing approaches, and systems automation and implementation have drastically improved security systems.

Biometric systems play a crucial role in current personal identification systems. Based on the specific biometric traits used to extract desired features in a system, biometric identifiers fall into two main classes: physiological traits and behavioral traits. Existing biometric systems analyze humans' physical characteristics and extract features for pattern recognition, such as facial patterns, fingerprints, irises, hand shape, and DNA patterns, to name a few. Among them, an EEG-based biometric system is one interesting possibility.

2. What is Electroencephalography (EEG)?

Electroencephalography (EEG) [2] is an electrophysiological monitoring method to record electrical activity of the brain. It can be performed in two different ways, either as Invasive EEG or Non-invasive EEG. Invasive/Intracranial EEG is recorded directly from the brain through a surgically implanted electrode placed at particular regions of human brain. Non-Invasive/Extracranial EEG is recorded from the surface or cortical layer of the human brain.

The EEG is typically described in terms of (1) rhythmic activity and (2) transients. The rhythmic activity is divided into bands by frequency. To some degree, these frequency bands are a matter of nomenclature (i.e., any rhythmic activity between 8–12 Hz can be described as "alpha"), but these designations arose because rhythmic activity within a certain frequency range was noted to have a certain distribution over the scalp or a certain biological significance.

EEG signals are sinusoidal waves, their amplitude is normally between 0.5 and 100 μ V. After applying a Fourier transform to the row signals and the power spectrum is generated, we have four groups of waves:

Band	Frequency	Location	Normally
Delta	0.5Hz - 3Hz	frontally in adults, posteriorly in children	usually happens during deep sleep
Theta	4Hz - 7Hz	found in locations not related to task at hand	drowsiness in adults and teens
Alpha	8Hz - 13Hz	posterior regions of head, both sides, higher in amplitude on dominant side; central sites (c3-c4) at rest	appears during relaxation without attention and concentration
Beta	14Hz - 31Hz	both sides, symmetrical distribution, most evident frontally	usual working rhythm
Gamma	>32Hz	Somatosensory cortex	Displays during cross-modal sensory processing

Table 1 Comparison of EEG Bands [2]

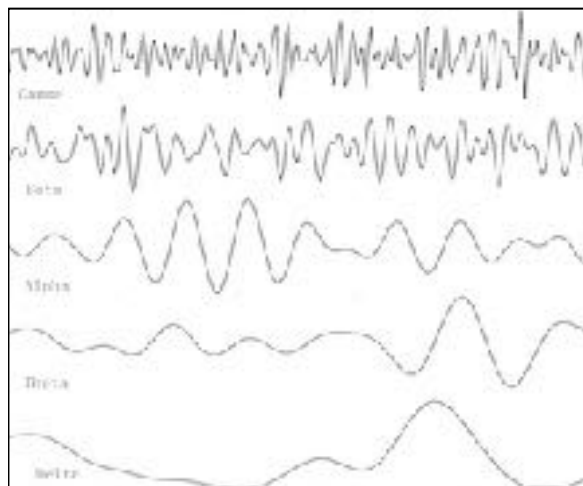


Fig. 1 Brainwaves Graph [2]

3. Related Work

Over the past decades, researchers have been investigated the potential of Electroencephalography (EEG) signals as biometric trait for human authentication purposes [3]. Unlike some conventional biometric traits like finger-print and face recognition techniques, which may be stolen or forged [4][5], EEG signals are difficult to mimic, impossible to intercept, unique and alive person recording required.

In EEG-based biometrics research field, various human tasks, discriminative features and classification methods have been studied. Palaniapan et al. [6] proposed a person authentication technique using gamma band visual evoked potentials from ten subjects when the individual is seeing a picture. Phuoc Nguyen et al. [7] investigated EEG-based person authentication of 9 subjects during a motor imagery task and obtained a success rate of 94%. Pham et al. [8] investigated the possibility of emotional states for person authentication using AR parameters and power spectral density (PSD) features.

4. Motivation and Objects

However, the stability and operability of the tasks used currently are still open to discussion. Various tasks should be further investigated for EEG-based person authentication as the tasks applied now are far from natural and versatile.

In this paper, a practicable person authentication system based a self-paced reaching task, which is a common and natural human daily task, was designed. PSD features of delta, theta, alpha and beta bands were extracted as subject-specific features as they were widely used [9]. Moreover, the time course characteristic of EEG signals of the self-paced reaching task for person authentication were investigated.

5. Organization

The rest of the paper is organized as follows. Section II explains the proposed self-paced reaching task as well as the data collection and preprocessing details. Section III describes the EEG features and the authentication methods used. Results and conclusions are presented separately in Section IV and V.

II.Experiment

1. Subjects

Thirty healthy volunteers participated in this experiment in the school of Data and Computer Science, Sun Yat-Sen University, China. All subjects had normal or corrected to normal vision and reported no history of neurological disorder. The procedures conformed to the Declaration of Helsinki. After detailed explanation of the procedures, all subjects signed a written informed consent. Data of nine subjects was removed from analysis because of their wrong experiment behaviors.

2. Experiment Design

Subjects were seated in front of a dell XPS laptop with a 13-inch touch screen. The opening angle of the touch screen was 135 degrees. The distance between the subject and the laptop was adjusted individually to allow comfortable reach to the touch screen.

The experiment interface was designed using Java. The JFrame framework was applied, which mainly called Swing, a common development toolkit for developing application user interfaces. Swing has many kinds of listeners in the event handling, which are suitable for different scenarios such as ActionListener, AdjustmentListener, ItemListener, etc.

As the experiment requires displaying a specific graph cyclically, each loop starts in a timed manner, and then the subject's touch and release events are monitored. The number of loops is internally initialized as required by the experiment. In addition, a number of colored ball's colors and appearance positions are randomly generated in advance of the experiment in order to improve the efficiency of the program operation. Subject's press and release time during the experiment and the balls' color and position on the screen will be recorded in a file as a reference for subsequent analysis (See Fig. 2).

In our experiment, one trial contained a holding window, a touching window and a relax window. A gray screen was showed to indicate the beginning of each trial. After two seconds, five colorful balls with one black ball were simultaneously presented on the screen. The five colorful balls were randomly selected from a pool of ten balls with different colors and the position of each colorful ball was randomly generated. Subjects were required to hold the black ball using their forefinger for at least three seconds. Then subjects self-paced released

the black ball and touched each colorful ball in which order the subjects preferred to. The black ball disappeared from the screen when it was released, and the colorful ball changed to be gray once be touched. A time window of three seconds was designed for subjects to relax after subjects touched all the colorful balls.

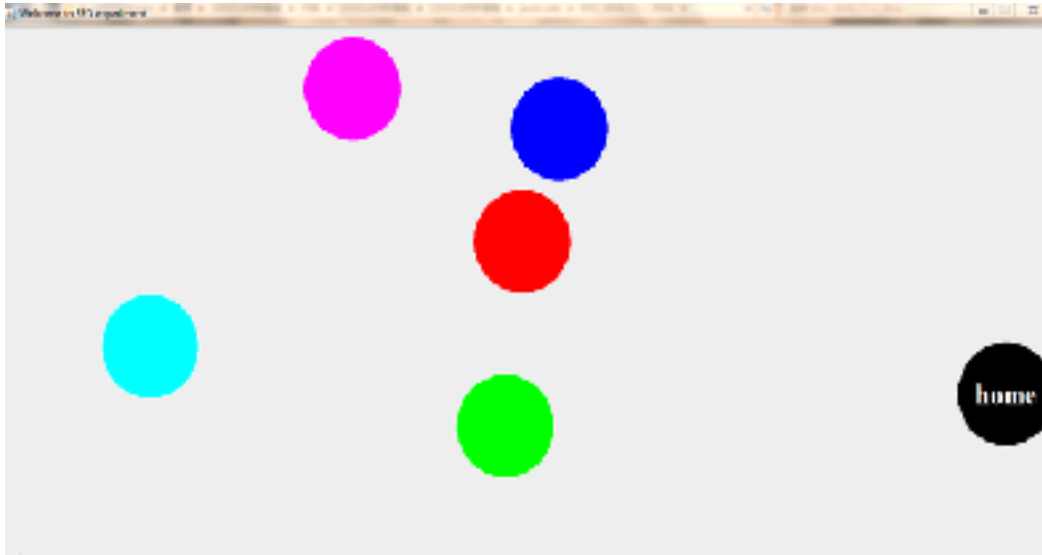


Fig. 2 Experiment Interface

For each subject, the whole experiment contained ten blocks and each block had 10 trials, i.e. a total of 100 trials were finished (See Fig. 3).



Fig. 3 A subject in the experiment

3. EEG Data Collection and Preprocessing

EEG signals were acquired at 500Hz using a BrainAmp DC amplifier (Brain Product) with 64 active Ag/AgCl EEG electrodes when subjects performed the experiment. The 64

electrodes were mounted in a head-cap (actiCap, Brain Product) according to the extended 10/20-System. The reference electrode was placed on FCz and the ground on Afz (See Fig. 4).

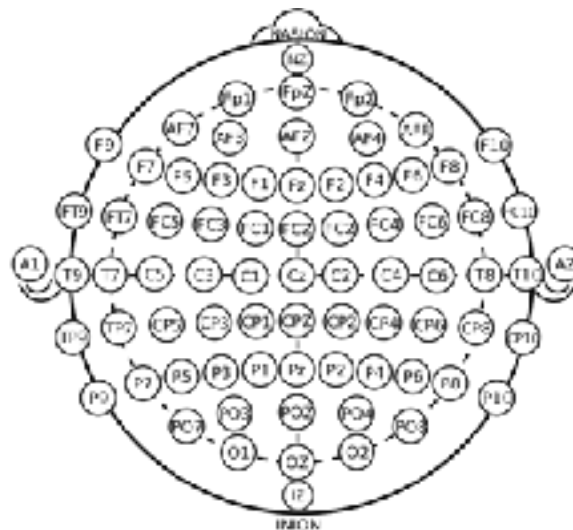


Fig. 4 10/20 System (EEG)

The EEG data were preprocessed using the open source Matlab toolbox EEGLAB [10]. The continuous EEG data were first band-pass filtered in the frequency range of 1-30 Hz. A common average reference was computed as the average voltage amplitude of the EEG data from all EEG channels. EEG signals were segmented into epochs of $[-2 \ 1]$ s time-locked to the onset of the movement. The $[-2 \ 0]$ s was defined as movement intention period and $[0 \ 1]$ s as movement execution period. Each epoch was baseline corrected to the EEG data in time window $[-500 \ 0]$ ms time-locked to the onset of the five colorful circles and two black balls (See Fig. 5).

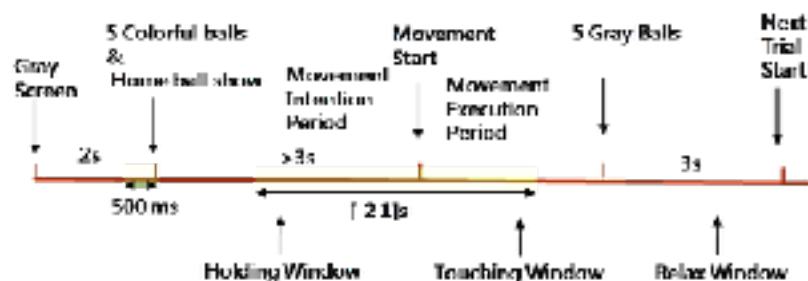


Fig. 5 Timeline of the self-paced reaching experiment

EEG data in the opposite motor area were selected for analysis. A single channel, C3, was

selected in this study as all the subjects in this study were right-handed. To improve the poor spatial resolution of EEG data, a large Laplacian filter was then applied, i.e. the average of EEG data of C1, C5, CP1, CP3, CP5, FC1, FC3 and FC5 was subtracted from the EEG data of C3.

III. Methods

1. Features

1.1. Power Spectral Density

Power spectral density of a time series is a positive real function of the frequency components. The PSD is defined as the discrete time fourier transform of the covariance sequence:

$$\Phi(\omega) = \sum_{k=-\infty}^{\infty} r(k)e^{-i\omega k} \quad (3.1)$$

where the auto covariance sequence $r(k)$ is defined as

$$r(k) = E \left\{ s(k)s^*(t-k) \right\} \quad (3.2)$$

and $s(t)$ is the discrete time signal $\{s(t); t = 0, \pm 1, \pm 2, \dots\}$, assumed to be a sequence of random variables with zero mean.

In this study, the Welch's method using periodogram was used for estimating the power of a time series at different frequencies. To estimate the PSD accurately, a Hamming window with 50% overlap was applied in this study.

1.2. Autoregressive Model

(1) Autoregressive Model Parameters

Autogressive model is a linear prediction formulas that best describe the signal generation system. Each sample $s(n)$ in an AR model is considered to be linearly related, with respect to a number of its previous samples:

$$s(n) = - \sum_{k=1}^p a_k s(n-k) + x(n) \quad (3.3)$$

where a_k , $k = 1, 2, \dots, p$ are the linear parameters, n denotes the discrete sample time and $x(n)$ is the noise input. AR has been also widely applied in EEG-based person authentication [11].

(2) AR Model Order Selection

An important issue in AR modeling is the selection of the appropriate model order. An underdetermined order may result in a smoothed spectral estimate with a poor resolution, whereas an overdetermined order could cause spurious peaks in the spectral

estimate and lead to spectral line splitting. Commonly used criteria for model order selection include Akaike Information Criterion (AIC), Schwarz-Bayes Criterion (SBC), also known as the Bayesian Information Criterion (BIC), Akaike's Final Prediction Error Criterion (FPE), and Hannan-Quinn Criterion (HQ). In brief, each criterion is a sum of two terms: one is the prediction error of the model, and the second term is the number of freely estimated parameters in the model, which increases with increasing model order. We need to minimize both terms by selecting an appropriate model order to gain estimation accuracy and efficiency.

According to the finite sample theory and Levinson-Durbin recursion analysis, the order selection depends on the reflection coefficients as the selection criteria, since the decrease in the residual variance is related to reflection coefficients. Model order overestimation could be determined from the significant decay of reflection coefficient value towards zero. In this thesis, we used AIC to determine the order of the model.

$$\text{AIC}(r) = (N - M)\log\sigma_e^2 + 2r \quad (3.4)$$

$$\sigma_e^2 = \frac{1}{N - M} \sum_{t=M+1}^N e_t^2 \quad (3.5)$$

N is the length of the data record, M is the maximal order employed in the model, $(N-M)$ is the number of the data samples used for calculating the likelihood function, r is the number of independent parameter presented in the model, and the optimal r is the minimum of $\text{AIC}(r)$.

2. Classification Models

3.2.1 Support Vector Machine

Support vector machine (SVM) constructs the optimal hyperplane

$$f(x) = \omega^T \phi(x) + b \quad (3.6)$$

to separate the training data into two classes in the training phase. SVM maximizes the hyperplane margin and minimizes the cost of error by solving the following optimization problem:

$$\min \frac{1}{2} \|\omega\|^2 + C \sum_{i=1}^l \delta_i \quad (3.7)$$

subject to $y_i [\omega^T \phi(x_i) + b] \geq 1 - \delta_i, \delta_i \geq 0, i = 1, 2, \dots, l$.

In the test phase, an SVM is used by computing the sign of

$$f(x) = \sum_i^{N_s} \alpha_i y_i \phi(s_i)^T \phi(x) + b = \sum_i^{N_s} \alpha_i y_i K(s_i, x) + b \quad (3.8)$$

where $K = \phi(x_i)^T \phi(x_j)$, s_i are support vectors, N_s is the number of support vectors. ϕ is the mapping function.

In this study, radial basis function (RBF) kernel and C-support vector classification (C-SVC) algorithm were used. A total of 21 classifiers were constructed to separate each subject from the others. A 10-fold cross validation was applied to calculate the authentication accuracy and AUC.

3.2.2 Random Forest

Random forest classifier consists of many individual classification trees (in this work the number of trees was 25), where each tree is a classifier by itself that is given a certain weight for its classification output. The classification outputs from all trees is used to determine the overall classification output which is done by choosing the mode (the output with most votes) of all trees classification output. The steps for building the random forests classifier can be summarized as follows:

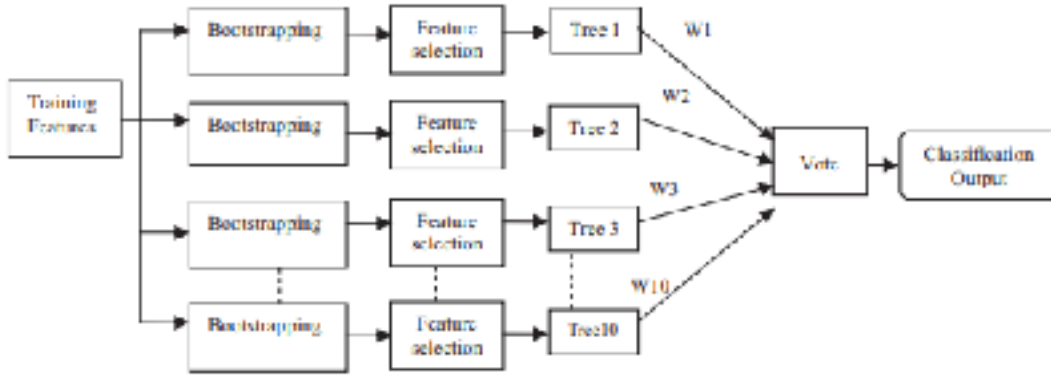


Fig. 6 Structure of Random Forest Classifier [12]

Building a random tree starts at the top of the tree with all the training dataset (in-bag data set). The first step involves selecting an attribute (feature) at the root node and then splitting the training data into subsets for every possible value of the attribute. This makes a branch for each possible value of the attribute. The splitting and the selection of the root node are done based on the information gain of the splitting of the attribute (the attribute with the highest information gain is selected at the root node). The information gain (IG) of splitting the training data set (Y) into subsets (Y_i) can be defined as:

$$IG = - \sum \frac{|Y_i|}{|Y|} E(Y_i) \quad (3.9)$$

$$E(Y_i) = - \sum_{j=1}^N p_j \log_2(p_j) \quad (3.10)$$

$|.$ is the size of the set; N is the number of the class to be classified ($N = 2$).

The node is split if the information gain is positive; otherwise the node is not split and becomes a leaf node in the training subset. The process is repeated recursively at each branch node using the subset that reaches the branch and the remaining attributes. The selection of the next attribute is also done based on the highest information gain of the remaining attributes. The process of splitting (tree growing) continues until all attributes are selected. The most occurring class in the training subset that reached that node is the classification output.

3.2.3 Backpropagation Neural Networks

Backpropagation is one of the self-learning methods of ANN to give desired answers. ANN is a parallel computing system emulating the ability of the biological neural network by interconnecting many artificial neurons. A two-layered network with an input layer and an output layer is adopted in this study. The input layer consists of 50 neurons and the output layer consists of one neuron; the layers are interconnected by sets of correlation weights. The neurons receive inputs from the initial inputs or the interconnections and produce outputs by transformation using an adequate nonlinear transfer function. A common transfer function is the sigmoid function expressed by $f(x) = (1 + e^{-x})^{-1}$ and in this study, we have used log-sigmoid function. In the learning process, the interconnection weights are adjusted using an error convergence technique to obtain a desired output for a given input. The learning constant (0.01) and the training iteration (100) have been used in the network structure [13].

IV. Results

1. PSD Analysis

Fig. 7 and 8 show the PSDs of two example subjects in the experiment. The x-axis is the time in second and the y-axis is the frequency in Hertz. The spectrogram shows that the PSD distributions vary between subjects during the movement intention and execution periods.

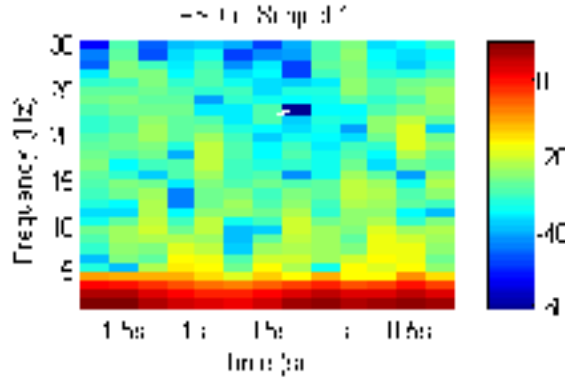


Fig. 7 Time-frequency representation for Subject 1. Time 0s represents the onset of movement.

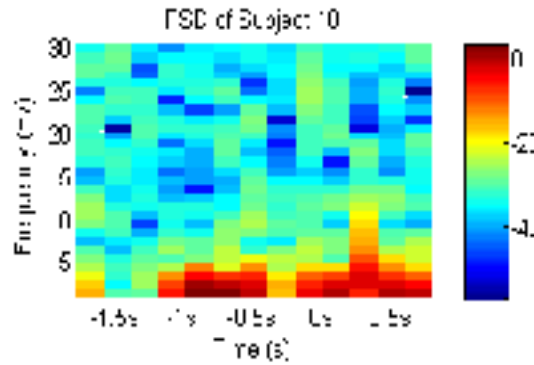


Fig. 8 Time-frequency representation for Subject 10. Time 0s represents the onset of movement.

2. Authentication Performance Comparison of Different Algorithms - SVM/RF/BP

Fig. 9 shows the authentication performance of different machine learning algorithms – Support Vector Machine, Random Forest and Backpropagation Neural Networks. Although the accuracy of Random Forest and Backpropagation Neural Network methods are higher than that of Support Vector Machine, their true positive rates are far lower than that of the later. Therefore we have adopted the Support Vector Machine as the classifier of the authentication system.

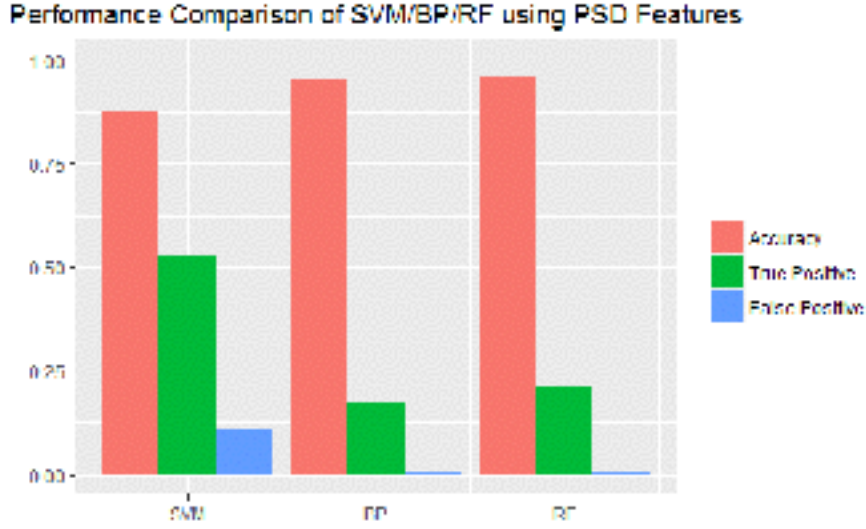


Fig. 9 Performance Comparison of SVM/BP/RF

3. Authentication Performance of PSD Features using SVM

4.3.1 Authentication Results for PSD Features of Four Frequency Bands

Fig. 10 presents the authentication results of SVM using PSD features of four frequency bands, i.e. delta (1-3Hz), theta (4-7Hz), alpha (8-13Hz) and beta (14-30Hz) for each subject. As the bandwidth is much larger in beta band (14-30Hz), the beta band power spectral densities were averaged in each frequency range (14-17) Hz, (18-21) Hz, (22-25) Hz and (26-30) Hz to derive 4 features for each subject. It shows that the alpha and beta bands achieve better authentication performances than delta and theta bands.

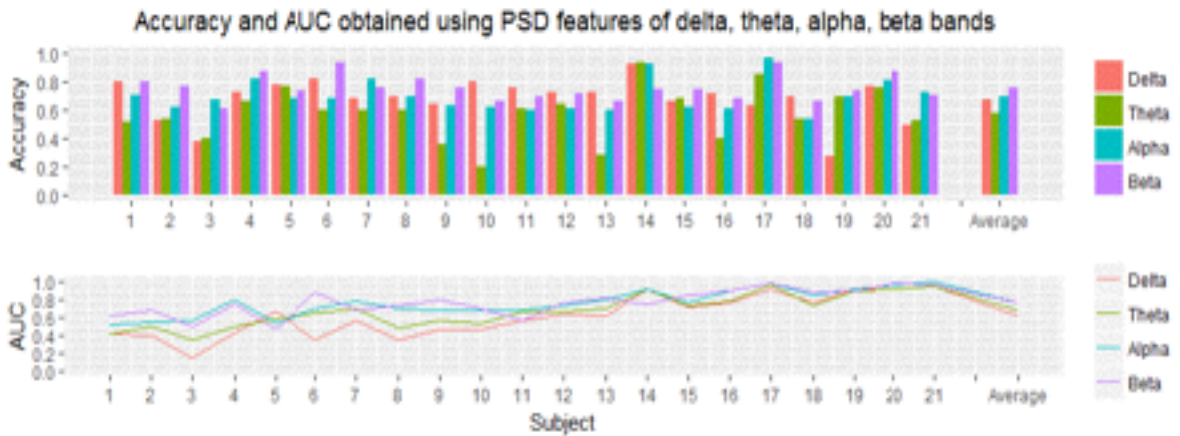


Fig. 10 Accuracy and AUC obtained using PSD features of delta, theta, alpha, beta bands.

4.3.2 Time Course Authentication Performances for Alpha and Beta Band PSD Features

To investigate the time source authentication performances, each epoch of [-2 1] s was

further segmented into three segments ($[-2 \ -1]$ s, $[-1 \ 0]$ s, $[0 \ 1]$ s). Person authentication based on PSD features in each time segment was conducted for alpha and beta bands separately. Fig. 11 and 12 present the average authentication results of SVM using alpha and beta band PSD features in each time segment across all subjects. The x-axis represents each time segment and the y-axis presents the accuracy and AUC separately. It is suggested that alpha band PSD features in the time period of $[-2 \ 1]$ s and the beta band PSD features in the time period of $[-1 \ 0]$ s and $[0 \ 1]$ s achieve a better result than in other time segments.

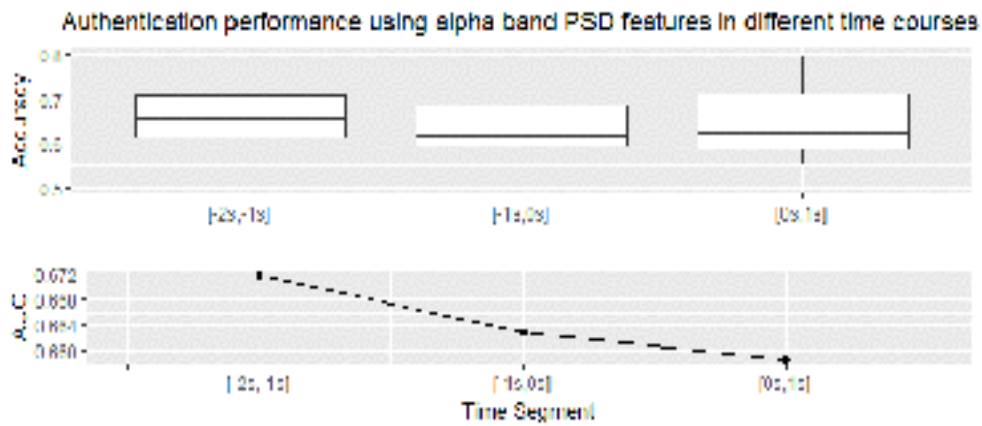


Fig. 11 Authentication performance using alpha band PSD features in different time courses.

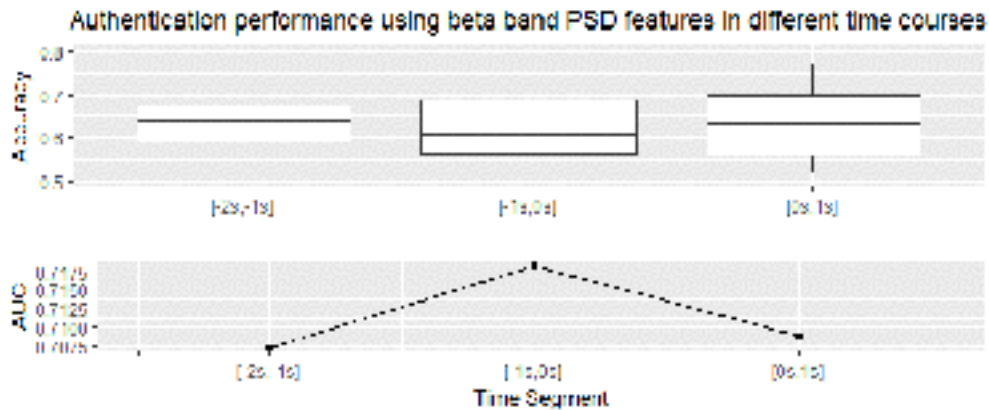


Fig. 12 Authentication performance using beta band PSD features in different time courses.

4.3.3 Authentication Performances of Combined Alpha and Beta Band PSD Features

Based on the previous results, the alpha band PSD features in the time period of $[-2 \ 1]$ s

and the beta band PSD features in the time period of $[-1\ 1]$ s were combined. Thus, the dimension of the combined features for each trial is 14. It is showed that the combined features achieved a better authentication performance (Fig.13).

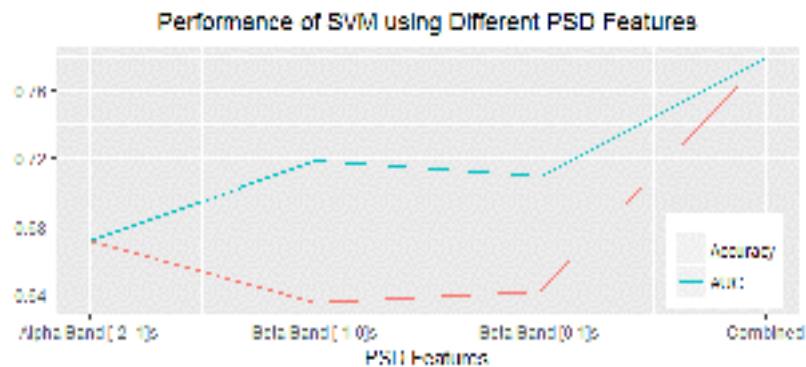


Fig. 13 Authentication performance using different PSD features.

4.4 Authentication Results for Combined PSD and AR Features using SVM

The AR model order of 6 was selected based on Akaike's information criterion(AIC) (Fig. 14). It can be seen that, compared to AR parameters, PSD features gives better performance.

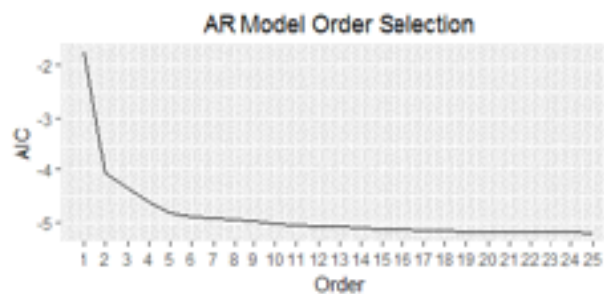


Fig. 14 AR model order selection

When combining the features of AR and PSD, the AUC of the person authentication result decreases despite the accuracy increases a bit (Fig. 15).

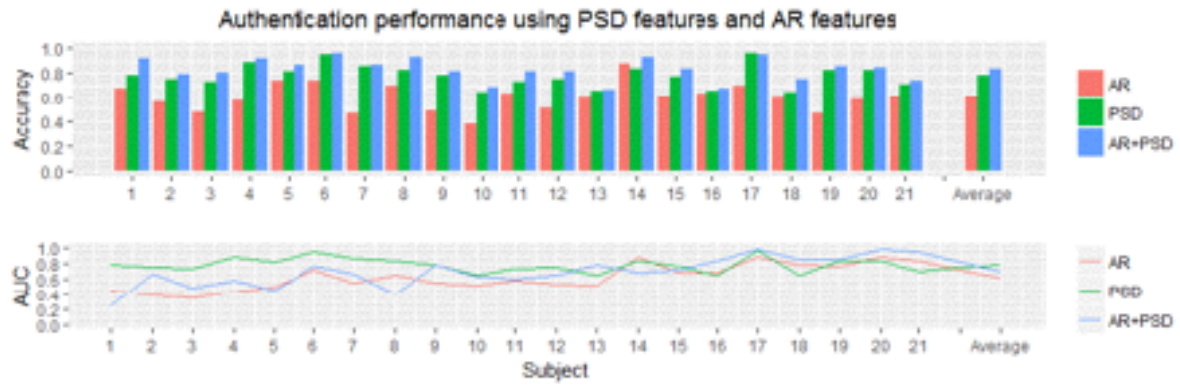


Fig. 15 Authentication performance using PSD features and AR features.

V.Conclusions

In this study, a practicable EEG based biometric authentication system is designed. The self-paced reaching task is applied as it is a natural and common human daily task, compared to visual evoked tasks, motor imagery tasks, emotional tasks, etc. The best average accuracy of our person authentication is 78.0%. And the average AUC is 0.751. Our experiment validates the usage of self-paced reaching task for EEG person authentication. Further investigation is necessary to optimize the performance of the proposed method in larger populations.

References

- [1] Jannatul Ferdous, “Assessing Self-similarity and Cross-similarity Between EEG Patterns for Biometrical Applications,” Unpublished MS dissertation, Lamar University, 2016.
- [2] W. Khalifa, A. Salem, M. Roushdy, K. Revett, “A Survey of EEG Based User Authentication Schemes,” 8th International Conference on INFormatics and Systems (INFOS2012), 2012.
- [3] J. Thorpe, P. Van Oorschot, and A. Somayaji, “Pass-thoughts: Authenticating with our minds,” Proceedings of the New Security Paradigms Workshop, NSPW, 2005.
- [4] J. Maatta, A. Hadid, and M. Pietikainen, “Face spoofing detection from single images using micro-texture analysis,” Biometrics (IJCB), 2011 International Joint Conference on, pp. 1–7, 2011.
- [5] Charles Arthur, “Iphone 5s fingerprint sensor hacked by germany’s chaos computer club”, 2014.
- [6] R. Palaniappan, P. Raveendran, “Individual identification technique using visual evoked potential signals”, Electronics Letters, vol. 38, pp. 1634-1635, 2002.
- [7] Phuoc Nguyen, Dat Tran, Xu Huang, and Wanli Ma, “Motor Imagery EEG-Based Person Verification”, International Work-Conference on Artificial Neural Networks (IWANN), Advances in Computational Intelligence, pp. 430-438, 2013.
- [8] T. Pham, W. Ma, D. Tran, D. S. Tran, and D. Phung, “A study on the stability of EEG signals for user authentication,” Neural Engineering (NER), 7th International IEEE/EMBS Conference on, 2015, pp. 122-125, 2015.
- [9] Salah. Altahat, Michael.Wagner and Elisa.Martinez-Marroquin, “Robust electroencephalogram channel set for person authentication, Acoustics,” Speech and Signal Processing (ICASSP), 2015 IEEE International Conference on, 2015.
- [10] Delorme Arnaud, Scott Makeig, “EEGLAB: an open source toolbox for analysis of single-trial EEG dynamics including independent component analysis,” J. Neurosci. Methods, vol. 134, no. 1, pp. 9–21, 2004.
- [11] R. Palaniappan, P. Raveendran, S. Nishida and N. Saiwaki, “Autoregressive spectral analysis and model order selection criteria for EEG signals,” Proceedings of

TENCON 2000, vol. 2, pp.126–129, 2000.

[12] Luay Fraiwan, Khaldon Lweesy, Natheer Khasawneh, Heinrich Wenz, Hartmut Dickhaus, “Automated sleep stage identification system based on time–frequency analysis of a single EEG channel and random forest classifier,” *Computer methods and programs in biomedicine*, pp. 11-19, 2012.

[13] Ching-Piao Tsai, Tsong-Lin Lee, “Backpropagation neural network in tidal-level forecasting,” *Journal of waterway, port, coastal and ocean engineering*, pp. 195-202, 1999.

Acknowledgements

I would like to express my gratitude to all those who have offered me invaluable help during my undergraduate study in the completion of this present thesis.

My deepest gratitude goes first and foremost to Dr. Yang Lingling, my supervisor, for her valuable suggestions and consistent instruction to this thesis. She has devoted lots of time to listening to me and helping me work out many problems. Without her painstaking efforts in revising and polishing my drafts, the completion of the present thesis would not have been possible.

Besides my supervisor, I am greatly indebted to many other teachers. Prof. Lu Songsong's encouragement and unwavering support has sustained me through frustration and depression. Special thanks should go to express my heartfelt gratitude to Prof. Wang Xueqin and Prof. Lin Chenshun, who have brought me into the world of statistics and data science. I feel grateful to have the opportunities to work with Prof. Zeng Yan and Dr. Huang Zhihong, from whose research I benefited greatly. My sincere thanks are also given to Prof. Huang Guanhua, Prof. Chen Chi-rung, Prof. Hsu Long for their care during my exchange study at National Chiao Tung University.

Last but not least, my thanks would go to my beloved family for their loving considerations and great confidence in me. I also owe my sincere gratitude to my friends and my fellow classmates in both School of Geography and Planning and School of Mathematics at Sun Yat-sen University for their company all through these years.