# On the utility of Power Spectral Techniques with feature selection Techniques for Effective Mental Task Classification in Non-invasive BCI

Akshansh Gupta[1], Ramesh Kumar Agrawal[1], Jyoti Singh Kirar[1], Javier Andreu-Perez[2,3],
Wei-Ping Ding[4], and Chin-Teng Lin[5],Mukesh Prasad[5]

[1]School of Computer and Systems Sciences, Jawaharlal Nehru University, New Delhi, India
[2]School of Computer Science and Electronic Engineering, University of Essex, United Kingdom
[3]Faculty of Medicine, Imperial College London United Kingdom
[4]School of Computer and Technology, Nantong University, Nantong, China
[5]Centre for Artifical Intelligence, School of Software, FEIT, University of Technology Sydney, Australia

In this paper classification of mental task-root Brain-Computer Interfaces (BCI) is being investigated, as those are a dominant area of investigations in BCI and are of utmost interest as these systems can be augmented life of people having severe disabilities. The BCI model's performance is primarily dependent on the size of the feature vector, which is obtained through multiple channels. In the case of mental task classification, the availability of training samples to features are minimal. Very often, feature selection is used to increase the ratio for the mental task classification by getting rid of irrelevant and superfluous features. This paper proposes an approach to select relevant and non-redundant spectral features for the mental task classification. This can be done by using four very known multivariate feature selection methods viz, Bhattacharya's Distance, Ratio of Scatter Matrices, Linear Regression and Minimum Redundancy & Maximum Relevance. This work also deals with a comparative analysis of multivariate and univariate feature selection for mental task classification. After applying the above-stated method, the findings demonstrate substantial improvements in the performance of the learning model for mental task classification. Moreover, the efficacy of the proposed approach is endorsed by carrying out a robust ranking algorithm and Friedman's statistical test for finding the best combinations and comparing different combinations of power spectral density and feature selection methods.

*Index Terms*—Brain-Computer Interface, Mental Tasks Classification, Feature Extraction, Feature Selection, Power Spectral Density.

## I. INTRODUCTION

A Brain-Computer Interfaces (BCI) [1], [2] is a message transmission framework, through which an individual can communicate for necessities by his or her brain signals, even absence of normal pathway of the computer system and a very effective device for the person with severe motor impairment [3], [4]. It is pragmatic area, which has focused to the design and invent of neuron rooted means to endue solutions for disease prediction , communication and control [5], [6], [7]. On the ground of acquisition of the brain signal BCI is broadly divided in three categories in literature [8], [9], viz, invasive, semi-invasive (electrocorticography (ECoG)) and non-invasive(electroencephalography EEG). Economically nature [10] and calibre to capture brain signals in a non-invasive fashion, EEG is a mostly preferred technique to aquire brain activity for BCI systems [11], [7]. BCI systems can be used as a *Response to mental tasks* system, [12], which is perceived to be more practical for locomotive patients. The basic assumption of this type of system is that mental activities lead to produce task-originated patterns. The BCI system's success depends on the precision of classification assorted mental tasks. These tasks requires extractions of discriminative features from the raw EEG signal to distinguish different mental tasks [13].

In previous studies, the researchers have utilized plenty approaches of feature extraction for better representation of the EEG signal for the classification process in the BCI domain, for example Band Power [14], amplitude values of EEG signals [15], Power Spectral Density (PSD) [16], [17], [18], [19], Autoregressive (AR) and Adaptive Autoregressive (AAR) parameters [20], time-frequency and inverse model-based features [21], [22], [23]. Wavelet Transform (WT) [24], [25] and Empirical Mode Decomposition (EMD) [26], [27], [28], [29], [30], [31], [32] have been used to decompose non-stationary and non-linear EEG signals into smaller frequency components. However, both WT and EMD methods provide low-frequency resolution and may not handle efficiently different overlapping frequency bands [33], [34] present in the EEG. On the other hand, power spectral analysis provides high-frequency resolution. The recording of EEG data occurs from multiple sensors/channels. Hence, the EEG data contains huge number of features but the recording session of the person usually very small in number. That produces, a small number of data samples. Hence, it suffers the curse of dimensionality as the ratio of features and sample is very small [35]. To overcome this problem, reduction of the dimension using feature selection is suggested in literature [36]. In spite of that, no in-depth study has ever been conducted about how to use power spectral features effectively with combination feature selection techniques in BCI the applications.

In this article we provide answers to the following questions:

1) Whether extraction of features using power spectral

techniques helps in mental task classification.

2) Whether the further reduction in dimensionality of features using feature selection approaches improves the classification performance or not.
3) Is multivariate feature selection approach better than univariate feature selection approach?
4) Which conjunction of feature extraction and selection method performs best for mental task classification?

Thus, this present work proposes a procedure of the determination of a compact collection of *relevant and non-redundant features* from the EEG signal in the two-phase approach. The first phase elaborates about the extraction of PSD features from the EEG signal using three different approaches. In the second level, a set of relevant and non-redundant features is sorted by multivariate filter feature selection approach. To investigate the performance of different combinations of PSD method and multivariate feature selection method, experiments are conducted on an open EEG data [7] source. The performance is calculated in terms of classification accuracy and compared with a combination of PSD and a univariate filter feature selection method. In order to rank and compare multiple combinations of power spectral density and feature selection methods Ranking method and Friedman's statistical test were also performed.

The rest of the paper is organized as follows: The Power Spectral Estimation approach has been discussed briefly in Section II.The proposed approach to obtain minimal subset of relevant and non-redundant of the PSD features using multivariate feature selection methods is included in the Section III. The Descriptive information, data and method results are presented in Section IV. In the final, Section V conclusions and future work is discussed.

## II. FEATURE EXTRACTION USING POWER SPECTRAL DENSITY

Power Spectral Density (PSD) is a measure of average power associated with any random sequence [37], which can be catalogued into three categories: (i) Non-parametric, (ii) Parametric and (iii) Subspace. The non-parametric methods are simple to compute and robust. Periodogram based estimation, Bartlett Window, Welch window and Blackman and Tuckey method are examples of this category. However, they do not provide the necessary frequency resolution due to their inability to extrapolate the finite length sequence for data points exceeding the signal length. Another drawback of this approach is spectral leakage [38]. To overcome the drawback of non-parametric methods, parametric method is suggested. The estimation of PSDs values from a given signal in parametric approaches are carried out by assuming that output of the linear system is driven by white noise and then parameters of the system are calculated. Examples are the *Yule-Walker autoregressive (AR) method* [39], the *Burg method* [16], Covariance and modified covariance etc. The commonly used parametric linear system model is the all-pole model which consists of a filter with all zeroes at the

origin and occurs in the z-plane. The output produced by such a filter using white noise as input is an autoregressive (AR) process. Thus, these spectral estimation methods are also sometimes known as *AR methods*. The AR methods tend to aptly describe data spectrum that is "peaky", the data having PSDs value large at certain frequencies, e.g. speech data. Smoother estimates of the PSD are produced by parametric methods than non-parametric methods but are subject to error if the order of model is not chosen correctly. Sub-Space methods are often used when SNR is low. PSDs values are obtained concerning Eigen-decomposition of autocorrelation matrix. For line spectra or spectra having sinusoidal nature Subspace methods are better choice and are also effective in the recognition of sinusoids mixed in noise. However, the subspace methods suffer from the following:The method in all probability does not generate true PSD estimates; it does not store power which is required for processing between the time and frequency domains; and it flunks in getting back the autocorrelation series by computing the inverse Fourier transform of the frequency estimate.

For a given stationary random signal $\mathbf{x}_m$, the PSD $P_{xx}$ is mathematically related to the autocorrelation sequence by Fourier transform, which regarding normalized frequency $f_s$, is given by,

$$P_{xx}(f) = \frac{1}{f_s} \sum_{m=-\infty}^{\infty} R_{xx}(m) e^{-\frac{j2\pi mf}{f_s}} \tag{1}$$

where $f_s$ is the sampling frequency. Fourier transform of the autocorrelation of the signal also gives the PSD. Using the inverse discrete-time Fourier transform from the PSD the correlation sequence can be derived as following:

$$R_{xx} = \int_{-\pi}^{\pi} P_{xx}(\omega) e^{-j\omega m} d\omega = \int_{-f_s/2}^{f_s/2} P_{xx}(f) e^{j2\pi f/f_s} df \tag{2}$$

The average power of the sequence $x_n$ over the entire Nyquist interval is represented by

$$R_{xx}(0) = \int_{-\pi}^{\pi} P_{xx}(\omega) d\omega = \int_{-f_s/2}^{f_s/2} P_{xx}(f) df \tag{3}$$

For a particular frequency band $[\omega_1, \omega_2]$, $(0 \le \omega_1 \le \omega_2 \le \pi)$, the average power of a signal is given by:

$$\overline{P_{[\omega_1,\omega_1]}} = \int_{\omega_1}^{\omega_2} P_{xx}(\omega) d\omega \tag{4}$$

It can be seen from the above expression that $P_{xx}(w)$ represents the power content of a signal in an *extremely small* frequency band, which is why it is known as the power spectral *density*.

### A. Welch Method

This method falls under non-parametric approach. For a finite time duration random signal $\mathbf{x}_m$ of $N$ interval length, PSD values are estimated with the help of a periodogram which is the squared modulus of discrete Fourier transform of the signal and is given by

$$P_{\mathbf{xx}}(f) = \frac{1}{N} |\mathbf{x}(f)|^2 \tag{5}$$

Here $f$ corresponds to the frequency of the sequence and $X(f)$ is the Fourier transform of the signal. A periodogram gives asymptotically non biased estimate of power spectrum.

In Welch method, $N$ length signal is divided into $K$ overlapped segments each of length $M$. The $i^{th}$ segment is given by,

$$\mathbf{x}_i(n) = \mathbf{x}(n + iD) \quad (6)$$

Here $n = 0 \dots N-1$, $i = 0 \dots K-1$ and $D$ is overlap segment. For this, a windowed segment periodogram is given by

$$P_{XX}^i(f) = \frac{1}{MU}\left| \sum_{i=0}^{N-1} w(n)\mathbf{x}_i(n)e^{-j2\pi fn} \right|^2 \quad (7)$$

where $w(n)$ is the window function and $U$ is the power of the window function given by,

$$U = \frac{1}{M}\sum_{n=0}^{M-1} w^2(n) \quad (8)$$

The average of $K$ periodograms depicts Welch power spectrum and is given by:

$$P_{XX}^W = \frac{1}{K}\sum_{i=0}^{K-1} P_{XX}^i(f) \quad (9)$$

.

### B. Burg Method

The Burg method [37] is a parametric method of spectral analysis. The PSDs values can be obtained by finding $pth$ order coefficients of an AR process. A $pth$ order real valued AR signal $\mathbf{x}(n)$ (with zero mean) at point $n$ is given by [19].

$$\mathbf{x}(n) = -\sum_{m=1}^{p} a_m x(n-m) + e(n) \quad (10)$$

Here $a_m$ is AR coefficient of $x(n-m)$, $e(n)$ is the error term at point $n$ independent of past terms. Burg algorithm test to find the AR coefficient by applying more data points and minimizes the forward and backward prediction errors in the least squares sense [19], with the AR coefficients constrained to satisfy the Levinson-Durbin recursion. It provides high resolution for short data records. After finding AR coefficients by Burg Algorithm, PSD value $S(f)$ at frequency $f$ is given by:

$$S(f) = \frac{S_e(f)}{|1 + \sum_{i=1}^{p} a_i e^{-j2\pi f i T}|^2} \quad (11)$$

Here $T$ is the sampling period and $S_e(f)$ is spectrum of error sequence which should be flat i.e. independent of frequency. One of foremost concern in AR modelling is the choice of order $p$. To determine $p$, several criterion such as final prediction error (FPE) [40], minimum description length [41], Akaike information criterion (AIC) [42], and autoregressive transfer function [43] are proposed in literature. Among these, AIC is commonly used, which is given by

$$AIC(p) = ln\sigma_{wp}^2 + \frac{2p}{n} \quad (12)$$

where $\sigma_{wp}^2$ is estimated variance in linear prediction error. From Table I, it can be observed that AIC value is minimum

TABLE I
VARIATION OF AIC VALUE FOR A GIVEN ORDER AND A MENTAL TASK.

| Task | Order | | | |
|---|---|---|---|---|
| | 5 | 6 | 7 | 8 |
| Baseline | -1.012 | -1.0117 | -1.0109 | -1.0106 |
| Count | -1.2841 | -1.2851 | -1.2847 | -1.2842 |
| Letter | -1.2574 | -1.259 | -1.2589 | -1.2585 |
| Math | -1.2783 | -1.2772 | -1.2762 | -1.2768 |
| Rot | -1.177 | -1.1768 | -1.176 | -1.1758 |

for order 5 or 6. We have chosen p=6 in our experiments which is also suggested by Kerin & Aunon [7].

### C. Multiple Signal Classification (MUSIC)

Music is an orthogonal subspace decomposition method is based on Pisarenko idea [44] that allows the estimation of low Signal-to-Noise ratio (SNR) frequency components. This method is used to lowers the effect of noise in the analysed signal and finds the optimal frequency resolution in a dynamic signal [45]. Subspace method assumes that any discrete time signal $s[n]$ is representable in the form of $m$ complex sinusoids with a noise $p[n]$ such that

$$s[n] = \sum_{i=1}^{m} \overline{A_i}\, e^{j2\pi f_i} + p[n], \ \ n = 0, 1, 2, \dots, N-1 \quad (13)$$

where $\overline{A_i} = |A_i| e^{\emptyset_i}$ is magnitude of $i^{th}$ complex sinusoid, $m, N, f_i$ and $\emptyset_i$ are frequency signal dimension order, number of sample data, frequency and phase of $i^{th}$ complex sinusoid.

The autocorrelation matrix $\mathbf{R}$ of signal $s[n]$ is given by:

$$\mathbf{R} = \sum_{i=1}^{m} |A_i|^2 p(f_i) p^H(f_i) + \sigma^2 \mathbf{I} \quad (14)$$

where $p(f_i) = \begin{bmatrix} 1 & e^{j2\pi f_i} & e^{j4\pi f_i} & \dots & e^{2\pi(N-1)f_i} \end{bmatrix}^T$ and $\sigma^2$ is variance of white noise signal, H is hermitian transpose and I is the identity matrix. Therefore, it can be observed that $\mathbf{R}$ is a composition of sum of signal and noise autocorrelation matrices such that

$$\mathbf{R} = \mathbf{R_s} + \sigma^2 \mathbf{I} \quad (15)$$

Pisarenko has noticed that variance of noise acts with the smallest eigenvalues of $\mathbf{R}$. The orthogonality of the signal and noise subspace is given as

$$p(f_i)^H v(m+1) = 0, i = 1, 2, \dots, m \quad (16)$$

where $v(m+1)$ is the eigenvector of noise in matrix $\mathbf{R}$ with dimension of $(m+1) \times (m+1)$ The estimation of PSD by Pisarnako is defined as

$$P_{Pisarnako} = \frac{1}{\left| p(f_i)^H v(m+1) \right|^2} \quad (17)$$

PSD estimation by MUSIC gives better performance than Pisarenko due to addition of averaging of extra noise

eigenvectors$(k = m+1, m+2, \ldots, M)$. Estimation of PSDs by MUSIC is given by:

$$P_{MUSIC}(f) = \frac{1}{\sum_{k=m+1}^{M} \left| p(f)^H v_{(k)} \right|^2} \tag{18}$$

Here $p(f)^H v_{(k)} = 0$ for $k = 1, \ldots, m$ using orthogonality of the signal and noise subspace. These PSD values have major peaks at the principal components only. The performance of Music depends on the dimension of the autocorrelation matrix $(M \le N)$

## III. PROPOSED APPROACH-FEATURE SELECTION

The number of PSD values obtained using one of the given three methods from multiple channels would be large, otherwise the number of training samples available is in general relatively small. Hence the method suffers from curse-dimensionality problem [35] .In order to subdue this problem, there is a need to determine a minimal set of pertinent features which can improve classification accuracy of a learning system. This work has proposed an approach to find a minimal subset of relevant feature using multivariate feature selection methods.

Feature selection [46], [36] is one of the widely accepted approach to determine relevant features. In spite of available plenty of research works for the feature selection, not much work has been carried out in the domain of mental task classification. The filter and the wrapper approaches are the two major approaches of feature selection techniques. In filter approach, the step of selecting optimal features set is considered as one of the pre-processing steps of just before applying any machine learning algorithm. This approach adopts only inherent properties of the features and does not consider any virtue of any learning algorithm. Hence, it may not select the optimal feature set for the learning algorithm. Instead, the wrapper approach [46] finds an optimal features subset, which is compatible with the given learning algorithm. The given classifier requires to be trained for each feature of set of the all features separately in the wrapper approach, which makes it more computational costly than filter approach.

Filter approach is further partitioned in two categories on the basis of the way of opting features [36], as Univariate (single feature ranking) and Multivariate (feature subset ranking). Univariate method utilizes a scoring function for measuring relevance of the feature. Implementation of the Univariate method is very simple. In BCI field the researchers, [47], [48], [49], [50] used univariate filter method. The performance of learning model usually improve with the help of reduced set of relevant features obtained by Univariate feature selection method. But it does not captures the correlation among the features. Hence there may be many redundant features in the subset of relevant feature which may take down the performance of learning model. Wrapper method, [51], [52], [7] has been applied to obtain a subset of non-redundant features for the mental task classification. Due to high- dimensionality of feature of EEG data, wrapper approach is not feasible option for mental task classification as it will become more computationally expensive. Hence we have applied both uni-variate

as well as multivariate filter feature selection algorithms. Let us assume we have a data matrix $\mathbf{X}$, of $m$ rows,and $k + 1$ columns, with data sample $\mathbf{x}_i, i = 1, 2, \ldots, m$; containing features set $S = \mathbf{f}_1, \mathbf{f}_2, \ldots \mathbf{f}_k$ and class label $C_1, C_2, \ldots C_n$, where $n \le m$.

### A. Uni-variate Feature selection

#### 1) Pearson's Correlation

Pearson's correlation coefficient (CORR), [53], [54] is employed to determine linear relationship between two variable. CORR of $i^{th}$ feature vector $(\mathbf{f}_i)$ with the class label vector $(\mathbf{c})$ is given by

$$CORR(\mathbf{f}_i, \mathbf{c}) = \frac{cov(\mathbf{f}_i, \mathbf{c})}{\sigma_{\mathbf{f}_i} \sigma_c} = \frac{E[(\mathbf{f}_i - \mu_i)(\mathbf{c} - \bar{c})]}{\sigma_{\mathbf{f}_i} \sigma_c} \tag{19}$$

where $i = 1, 2, \ldots, k$, $\sigma_{\mathbf{f}_i}$, $\sigma_{\mathbf{c}}$ represent respectively the standard deviations of feature vector $\mathbf{f}_i$ and $\mathbf{c}$. $cov(\mathbf{f}_i, \mathbf{c})$ represents the covariance between $\mathbf{f}_i$ and $\mathbf{c}$, $\mu_i = \frac{1}{k} \sum_{i=1}^{k} X_{ik}$ and $\bar{c} = \frac{1}{k} \sum_{i=1}^{k} c_i$ are the mean of $\mathbf{f}_k$ and $\mathbf{c}$ respectively.

Range of $CORR(\mathbf{f}_i, \mathbf{c})$ falls between -1 & +1. The value nearby to $|1|$, depicts the stronger linear relation among the prescribed variables while zero value implies no correlation between the two variables.

#### 2) Mutual Information

Mutual information [MI] is a feature ranking method on basis of Shannon entropy, which determines relationship between two variables. MI of a feature vector $\mathbf{f}_i$ and the class vector $\mathbf{c}$ can be calculated as [55]:

$$I(\mathbf{f}_i, \mathbf{c}) = \sum P(\mathbf{f}_i, \mathbf{c}) \log \frac{P(\mathbf{f}_i, \mathbf{c})}{P(\mathbf{f}_i) P(\mathbf{c})} \tag{20}$$

where $P(\mathbf{f}_i)$ and $P(\mathbf{c})$ are the marginal probability distribution functions for random variables $\mathbf{f}_i$ and $\mathbf{c}$ respectively and $P(\mathbf{f}_i, \mathbf{c})$ is joint probability distribution. The most extreme estimation of MI demonstrates the higher reliance of the variable on the class label. The advantage of MI is that it can discover even the non-linear dependency between the attribute and the relating class label vector $\mathbf{c}$.

#### 3) Fisher Discriminant Ratio

Fisher Discriminant Ratio (FDR) is a univariate filter feature selection technique which depends on the statistical virtue of the attributes or features. FDR $(\mathbf{f}_i)$ for $i^{th}$ features for two class $C_1$ and $C_2$ can be given as:

$$FDR(\mathbf{f}_i) = \frac{\left(\mu_{1(i)} - \mu_{2(i)}\right)^2}{\sigma_{1(i)}^2 + \sigma_{2(i)}^2} \tag{21}$$

here, $\mu_{1(i)}$ and $\sigma_{1(i)}^2$ are the mean and deviation of the data of class $C_1$ respectively for $i^{th}$ feature.

#### 4) Wilcoxon's Ranksum Test

Wilcoxon Ranksum Test, suggested by [56], is a non-parametric statistical test, accomplishes between data of two classes on the basis of median of the samples having no prior knowledge of probability distribution.

The statistical distinctness $t(\mathbf{f}_i)$ of feature $\mathbf{f}_i$ for known two classes, class $C_1$ and $C_2$ using Wilcoxon's statistics can be defined as [57]:

$$t\left(\mathbf{f}_i\right) = \sum_{l=1}^{N_i} \sum_{m=1}^{N_j} DF((X_{li} - X_{mi}) \le 0) \qquad (22)$$

where $N_i$ and $N_j$ are the number of the data example in class $C_1$ and $C_2$ respectively, $DF$ is the logical discriminative mapping between two classes of data, which defines an estimation of 1 or 0 corresponding to true or false and $X_{li}$, is the expression values of $i^{th}$ feature for $l^{th}$ sample. The value of $t(\mathbf{f}_i)$ lies between zero and $(N_i \times N_j)$. The relevance of the feature can be fined as:

$$R\left(t(\mathbf{f}_i)\right) = \max(t(\mathbf{f}_i), N_i \times N_j - t(\mathbf{f}_i)) \qquad (23)$$

### B. Multivariate Feature Selection

Time-efficient multivariate filter method picks a subset of features, which are relevant to the class label of data and independent from each other. Thus it up dues the limitations of both uni-variate and wrapper approaches. Thus we have opted most widely utilized multivariate filter methods by research community for the dimensionality reduction, are Bhattacharya distance measure [58], Ratio of scatter matrices [59], Linear regression [60] and minimum Redundancy-Maximum Relevance (mRMR) [61]. A brief discussion on the mentioned techniques is given below.

### 1) Bhattacharaya Distance

In literature, Bhattacharya distance is used for find similarity between two continuous or discrete probability distribution. It is a special case of Chernoff distance that provides similarity overlap of the distribution. For multivariate normal probability distribution, Chernoff Distance measure can be written as [62]:

$$J_c = \frac{1}{2}\beta(1-\beta)(\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1)^T[(1-\beta)\boldsymbol{\Sigma}_1 + \beta\boldsymbol{\Sigma}_2]^{-1}(\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1) + \frac{1}{2}\log\frac{|(1-\beta)\boldsymbol{\Sigma}_1 + \beta\boldsymbol{\Sigma}_2|}{|\boldsymbol{\Sigma}_1|^{1-\beta}|\boldsymbol{\Sigma}_2|^{\beta}} \qquad (24)$$

where $\boldsymbol{\mu}_i$ and $\boldsymbol{\Sigma}_i$ are mean vector and covariance matrix for class $C_i$ respectively($i$=1,2).

When $\beta$ is $\frac{1}{2}$ then this distance is called as Bhattacharya distance (BD) [58], which is given as

$$J_B = \frac{1}{8}(\boldsymbol{\mu}_2 - \boldsymbol{\mu}_2)^T(\boldsymbol{\mu}_2 - \boldsymbol{\mu}_2) + \frac{1}{2}\log\frac{(\frac{|\boldsymbol{\Sigma}_1 + \boldsymbol{\Sigma}_2|}{2})}{|\boldsymbol{\Sigma}_1|^{\frac{1}{2}}|\boldsymbol{\Sigma}_2|^{\frac{1}{2}}} \qquad (25)$$

### 2) Ratio of Scatter Matrices

In literature, the trace of ratio of scatter matrices (SR),is a measure of separability, as the trace of a scatter matrix is equal to the sum of the eigenvalues and therefore an indicator of the total variance in the data. How well features cluster around their class mean, as well as, how well they separate the class means. The scatter matrices, within-class scatter matrices,$\mathbf{S}_w$, and between class scatter matrices, $\mathbf{S}_b$, can be defined as

$$\mathbf{S}_w = \sum_{i=1}^{c} P_i E[(\mathbf{x} - \boldsymbol{\mu}_i)^T(\mathbf{x} - \boldsymbol{\mu}_i)] \qquad (26)$$

$$\mathbf{S}_b = \sum_{i=1}^{c} P_i(\boldsymbol{\mu}_i - \mu_{\mathbf{0}})^T(\boldsymbol{\mu}_i - \mu_{\mathbf{0}}) \qquad (27)$$

where $\boldsymbol{\mu}_i$, $P_i$ and $\mu_0$ are mean vector of $i^{th}$ class data, prior probability of $i^{th}$ class data and global mean of data samples respectively.

From the definitions of scatter matrices, the criterion value which has to be maximized, is given as:

$$J_{SR} = \frac{trace(\mathbf{S}_b)}{trace(\mathbf{S}_w)} \qquad (28)$$

When intra cluster distance is very small and the inter cluster distance is very large $J_{SR}$ takes the high value. The main advantage of this criterion that it is not subject any external parameters and assumptions of any probability density function. Also the measure $J_{SR}$ under linear transformation has the advantage of being invariant under linear transformation.

### 3) Linear Regression

Linear regression is a statistical approach, which determines casual link of an independent variable upon a dependent variable. The class label of the data is recognized as the target dependent variable and the feature that affect the target is known as independent variable. There may be many features which can affect the class of the data, thus in such case multiple regression analysis would be more appropriate. A multiple regression model with $k$ independent features $\mathbf{f_1}, \mathbf{f_2}, \ldots, \mathbf{f_k}$ and a class variable $y$ can be written as [60];

$$y_i = \beta_0 + \beta_1 X_{i1} + ... + \beta_k X_{ik} + \zeta_i, i = 1, 2, ..., n \qquad (29)$$

where $\beta_0, \beta_1, ..., \beta_k$ is set of fixed values calculated by the class label $y$ and observed values of $\mathbf{X}$ and $\zeta_i$ is the error term. The sum of squared error (SSE) is given by

$$SSE = \sum_{i=1}^{n}(y_i - y_i^p)^2 \qquad (30)$$

where $y_i$ and $y_i^p$ are observed and predicated values respectively. The lower value of SSE depicts preferable regression model. The total sum of squares (SSTO) can be calculated as:

$$SSTO = \sum_{i=1}^{n}(y_i - \bar{y})^2 \qquad (31)$$

where $\bar{y}$ is the mean value of $y_i, i = 1, 2, ..., n$. The criterion function $J_{LR}$ is given as:

$$J_{LR} = R^2 = 1 - \frac{SSE}{SSTO} \qquad (32)$$

The value of $J_{LR}$ lies between 0 and 1. The feature for which the value of $J_{LR}$ is higher is selected.

### 4) minimum Redundancy-Maximum Relevance

minimum-redundancy maximum-relevance (mRMR) [60] is based on mutual information to discover a subset of features that have minimum redundancy among themselves and maximum relevance with the class labels. mRMR uses mutual information $I(\mathbf{f}_i, \mathbf{f}_l)$ as a measure of similarity between two feature vector $\mathbf{f}_i$ and $\mathbf{f}_l$, which is given as pursues:

$$I(\mathbf{f}_i, \mathbf{f}_l) = \sum_{k,l} p(\mathbf{f}_k, \mathbf{f}_l) \log(\frac{p(\mathbf{f}_i, \mathbf{f}_l)}{p(\mathbf{f}_i)p(\mathbf{f}_l)}) \qquad (33)$$

where $p(\mathbf{f}_i), p(\mathbf{f}_l)$ are the marginal probabilities of $k^{th}$ and $l^{th}$ features respectively and $p(\mathbf{f}_i, \mathbf{f}_l)$ is selected joint probability density. The relevance between the set of features S and the target class label vector $\mathbf{c}$, denoted by $REL$, is expressed as:

$$REL = \frac{1}{|S|} \sum_{\mathbf{f}_i \in S} I(\mathbf{f}_i, \mathbf{c}) \qquad (34)$$

The average redundancy among features in the set $S$, denoted by $RED$, is defined as:

$$RED = \frac{1}{|S|^2} \sum_{\mathbf{f}_i, \mathbf{f}_l \in S} I(\mathbf{f}_i, \mathbf{f}_l) \qquad (35)$$

where $S$ denotes the subset of features and $|S|$ denotes the number of features in set $S$. Minimum redundancy and maximum relevance is measured by:

$$J_{MID} = max(f_i)[REL - RED] =$$
$$max(f_i) \left[ \frac{1}{|S|} \sum_{\mathbf{f}_i \in S} I(\mathbf{f}_i, c) - \frac{1}{|S|^2} \sum_{\mathbf{f}_i, \mathbf{f}_l \in S} I(\mathbf{f}_i, \mathbf{f}_l) \right] \qquad (36)$$

Clearly, the maximum values of $J_{MID}$ can be achieved with minimum redundancy among features and maximum relevance with target vector.

## IV. RESULTS AND DISCUSSION

### A. Data

For the simulation of our proposed framework, we have utilized a freely available Mental Task Classification data-set which has been used first time in the work of(Keirn and Aunon, 1990). Seven subjects (person) participated in the recording of this EEG dataset, but we did not utilize of Subject 4 due to incomplete information. In the baseline task (relax: B) each subject was instructed to carry out five different mental tasks ; the mental Letter Composing task (L); the Non trivial Mathematical task (M); the Visualizing Counting of numbers written on a blackboard task (C), and the Geometric Figure Rotation task (R). Each recording consists of the five trials of each of above said five mental tasks. EEG signals are recorded from C3, C4, P3, P4, O1 and O2 electrode position with A1 and A2 as the reference electrode as shown in Figure 1. Each trial is recorded for 10 seconds duration recorded with the sampling rate of 250 per second, which resulted in 2500 samples points per trial.

Overall flow of the proposed approach for mental task classification is shown in Figure 2. The proposed approach consists four steps: segmentation feature extraction, feature selection and classification; to distinguish two different mental tasks. The main contribution of the work is employment of filter feature selection algorithm to enhance performance of learning algorithm for the classification of the mental tasks.

### B. Feature Formation

For feature vector formulation, each trial data is pre-processed by decomposing into half-second segments, gen-erating 20 fragments per trial for each subject. The extraction of features has been carried out from each signal using three
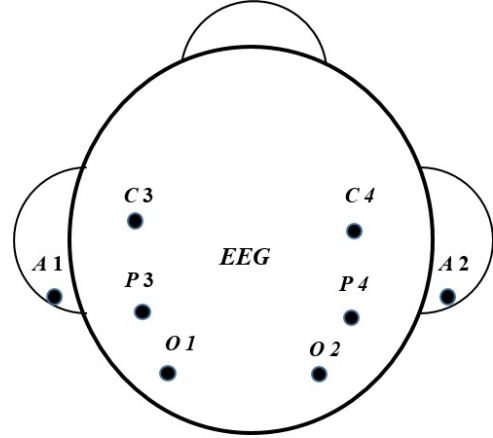


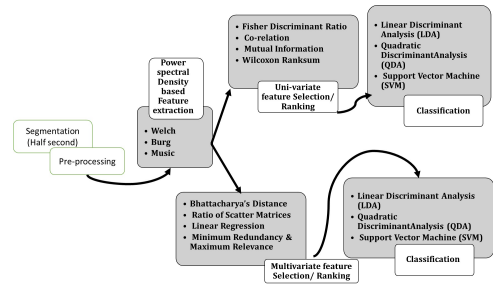Fig. 1. Electrode placement of EEG recording adapted from [19].



Fig. 2. Flow Diagram of the proposed approach for mental task classification.

different power spectral density approaches such as Welch, Burg, and MUSIC separately. A total of 52 PSD values are obtained from each channel. Combining PSD values of all six channels, each signal is represented regarding 312 PSD values. PSD values obtained for different tasks using Burg (parametric approach) for all six channels are shown in Figure 3, which shows that the extracting features from Burg PSD approach are effective in distinguishing different mental tasks. It can be also observed that PSD values at some frequency values differ considerably among different mental tasks (e.g. Frequency range of 6-9 Hz for channel C3, 6-13 Hz for channel C4, 6-13 Hz for channel P3, 6-16 Hz for channel P4, 6-9 Hz for channel O1 and 16-19 Hz for channel O2). This difference in PSD values can help in distinguishing different mental tasks. While PSD values at some frequency values take similar values (e.g. Frequency values above 15 Hz for C3, above 17 Hz for channel C4, above 13 Hz for channel O1, above 30 Hz for channel O2, above 20 Hz for channel P3 and above 22 Hz for channel P4) and cannot help in distinguishing different mental tasks. Similar observations are also noted for Welch and MUSIC methods. This suggests that all features (PSD values) are not relevant for mental task classification.

### C. Application of Uni-variate Feature Selection

To determine relevant features that can distinguish different mental tasks, four different univariate methods: Correlation (Cor), Fisher Discriminant Ratio (FDR), Mutual Informa-tion (MI) and Wilcoxon's Rank Sum Test (Ranksum) are investigated in our experiment. FDR score corresponding to
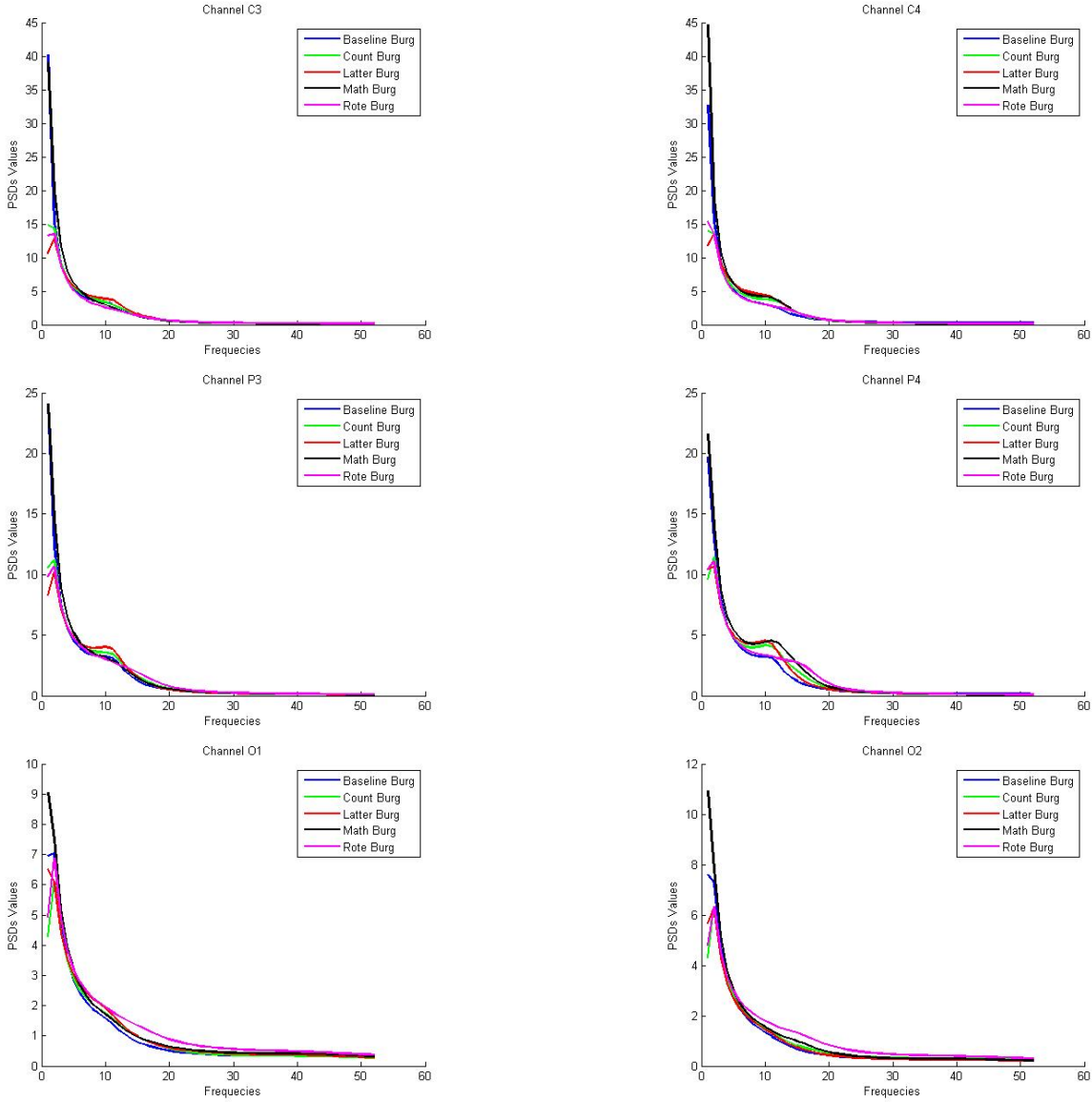
Fig. 3.   Comparison of features of different mental tasks using Burg method.

features obtained from each of the three PSD approaches to distinguish Baseline task from Count Task is shown from Figure 4 to Figure 6. It can be seen from these figures that FDR score corresponding to few features is very high and very less for others. This suggests that some features are more relevant than others. Similar observations are also noted for other univariate methods and other pairs of tasks. For all univariate feature selection methods, the top 25 -ranked features are incrementally added to develop the decision model using forward feature selection approach. Comparison of different methods is reported in terms of maximum average classification accuracy for top features of 10 runs of 10 cross-validations. We have used three well-known classifiers: Linear Discriminant Analysis (LDA), Quadratic Discriminant Analysis (QDA) and Support Vector Machine (SVM) in our experiments. Figure 7 shows a comparison of all combinations of three PSD approaches and four univariate methods with each of the three PSD approaches without any feature selection

in terms of average classification accuracy (over six subjects for all combination of tasks). We can observe the following from Figure 7

- In general, the classification accuracy of all the three PSD approaches improves with the use of univariate feature selection method with all three classifiers.
- Among all combinations of PSD approaches, univariate methods, and classifiers, the maximum classification accuracy is achieved with the combination of Burg, FDR, and SVM.
- Among four univariate feature selection methods, maximum classification accuracy is achieved with FDR.

### D. Application of Multivariate Feature Selection

Figure 8 shows a color map of correlation values among top 20 relevant features obtained using the combination of FDR and Burg method to distinguish Baseline task from Count Task. It can be noted that some of the correlation
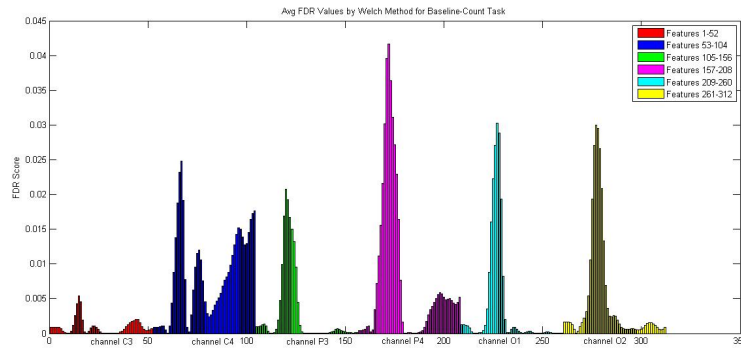
Fig. 4. Fisher Discriminant Ratio score for a pair of Baseline task and Count Task for features extracted using Welch.
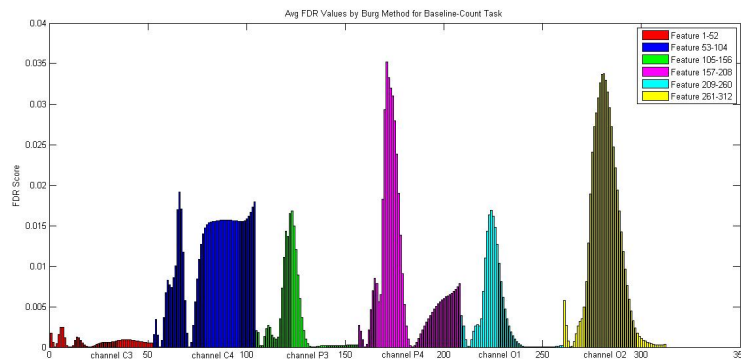


Fig. 5. Fisher Discriminant Ratio score for a pair of Baseline task and Count Task for features extracted using Burg.
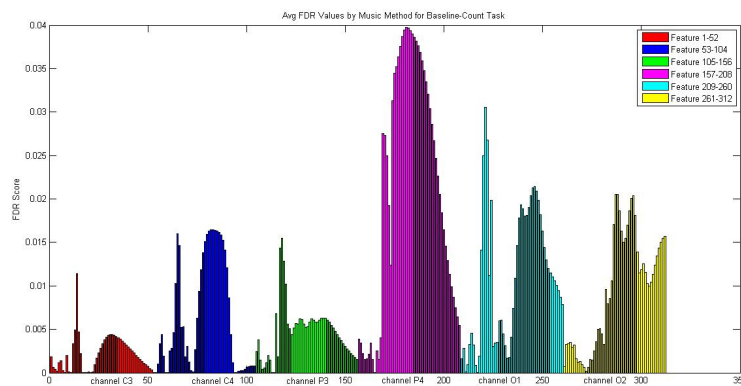


Fig. 6. Fisher Discriminant Ratio score for a pair of Baseline task and Count Task for features extracted using MUSIC.

values take a high value which depicts that such features are correlated (redundant) among themselves. Similar observations are also noted for other combinations of PSD approaches and univariate methods for another pair of tasks. This observation suggests the need to determine a subset of relevant and non-redundant features to further improve the performance of mental task classification. For this, we used four well known multivariate methods: linear regression (LR), Bhattacharya distance (BD), Scatter Ratio (SR), Minimum Redundancy-Maximum Relevance (mRMR) to obtain minimal subset of non-redundant and relevant features using forward feature

selection approach. Figure 9 shows a comparison of all combinations of three PSD approaches and four multivariate methods with the combination of PSD approaches and FDR (best performing univariate method) in terms of average classification accuracy. We can observe the following from Figure 9

- Among all combinations of PSD approaches, multivariate feature selection methods and classifiers, the maximum classification accuracy is achieved with the combination of Burg, LR, and LDA.
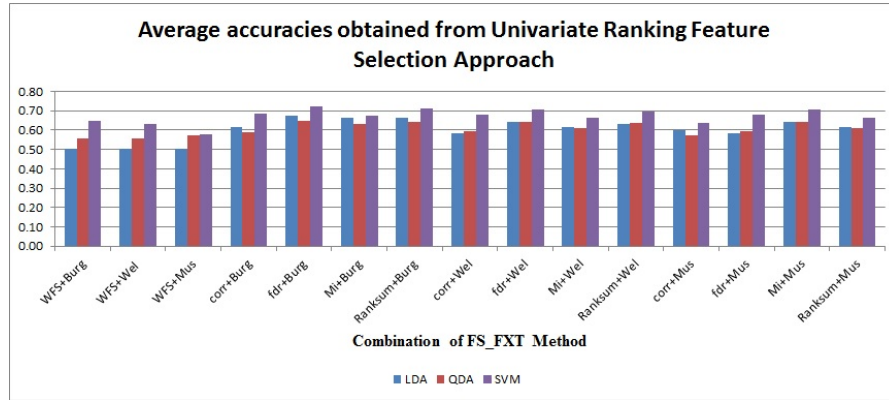- The performance of all combination of PSD approaches

Fig. 7.  Comparison of different combination of univariate methods and PSD methods in terms of classification accuracy.
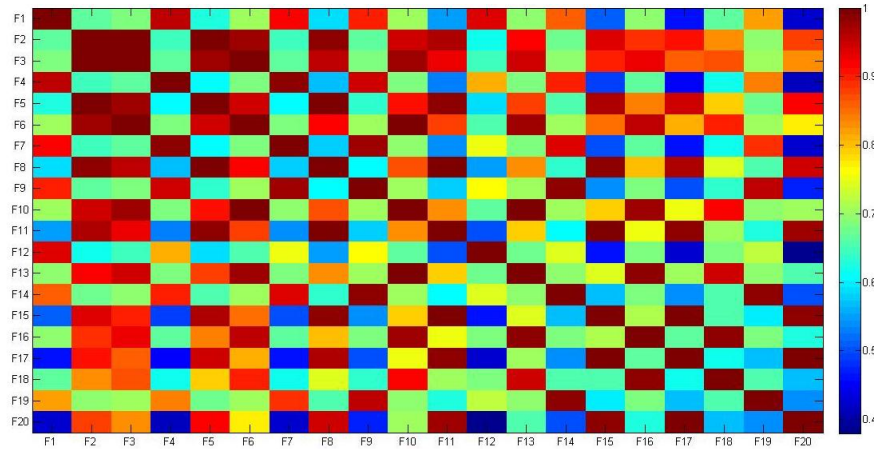


Fig. 8.  Colormap of Correlation values for top 20 PSD features obtained using combination of FDR .
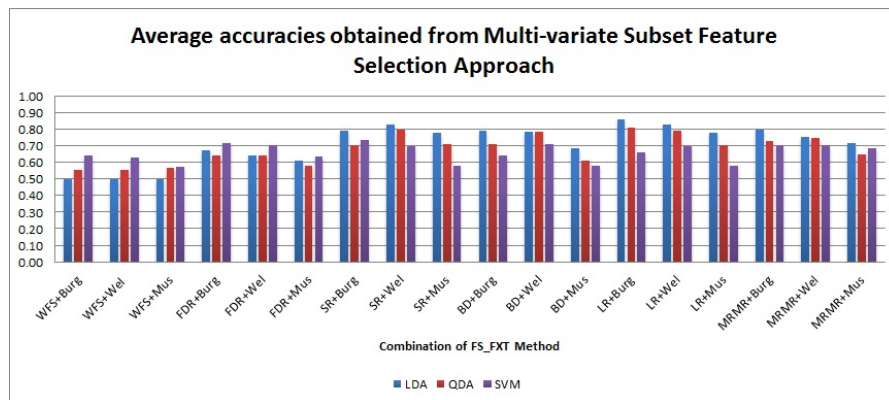


Fig. 9.  Comparison of all combinations of three PSD approaches and four multivariate methods with combination of PSD approaches and FDR in terms of average classification accuracy.

and multivariate methods is better in comparison to the combination of PSD approaches and FDR for LDA and QDA in terms of classification accuracy.

- The performance of MUSIC is worst among three PSD approaches with univariate as well as multivariate feature selection methods.

### E. The Rankings of Respective Combinations of Feature Extraction and Selection Methods

To investigate the relational rank of both univariate and multivariate methods feature selection techniques in combination with a feature extraction method, we have utilized the robust ranking approach [63], on the ground of percentage gain in classification accuracy with respect to without applying any feature selection method [64].

Figure 10 shows twenty-four combinations of FS-FXT methods which are the feature selection and extraction methods. These methods are compared on the basis of percentage gain in accuracy of the different combination of selection and extraction methods and their corresponding ranks. From the Figure 10, we can observe that the combination of multivariate feature selection with all three feature extraction is ranked better in comparison to the combination of univariate feature selection and all three feature extraction methods except one combination (BD-MUSIC). Among all combination of selection and extraction methods, the combination of LR and Burg is best, whereas the team of MUSIC and Ranksum performs the worst.

### F. Friedman statistical test

We have applied a non-parametric statistical test known as Friedman test in order to compare the statistically significant difference evolving in various combination of the feature selection and the PSD methods. From Table II, it can be noted that almost (11 out of 12) all combinations of multivariate feature selection with PSD methods obtained better rank than the combination of univariate feature selection method and PSD methods. The SEL-EXT pair performance is also examined with respect to a control method, i.e., the one that emerges with the lowest rank (combination of LR and Burg). In the comparison of the control method with other 23 combinations of feature selection and feature extraction method, adjusted p-values [65] we computed in order to take into account the error accumulated and provide the correct correlation. A set of post-hoc procedures to provide correct correlation is defined in the literature. The adjusted p-values in the method are computed in order to find whether the control method shows any statistical difference when compared with the remaining methods. For pair-wise comparisons, the widely used post hoc methods to obtain adjusted p-values are [65]: Holm, Hochberg and Hommel procedures. Table III illustrates adjusted p-values for the Hommel procedure. The values in Table III represents the p-value when the pair-wise comparison with control method(Burg+LR) is conducted. The bold values suggest the significant difference observed from the control method (Burg+LR) with the combinations at the significance level of 0.05. This demonstrates that combination

#### TABLE II
AVERAGE RANKING USING FRIEDMAN'S STATISTICAL TEST.

| Different Combination | Ranking |
|---|---|
| LR + Burg | **1.85** |
| SR + Welch | 2.2 |
| LR + Welch | 2.5 |
| BD + Welch | 3.84 |
| mRMR + Burg | 5.45 |
| SR + Burg | 5.85 |
| mRMR + Welch | 6.3 |
| BD + Burg | 8.2 |
| SR + MUSIC | 10.25 |
| LR + MUSIC | 10.79 |
| mRMR + MUSIC | 11.2 |
| CORR + Welch | 11.6 |
| FDR + Burg | 12.65 |
| FDR + MUSIC | 14.15 |
| CORR + MUSIC | 15.3 |
| MI + Welch | 15.85 |
| CORR + Burg | 18.1 |
| RANKSUM + MUSIC | 18.7 |
| BD + MUSIC | 19 |
| MI + Burg | 19.15 |
| FDR + Welch | 20.35 |
| RANKSUM + Burg | 20.9 |
| MI + MUSIC | 22.15 |
| RANKSUM + Welch | 23.649 |

#### TABLE III
ADJUSTED p-VALUES FOR THE HOMMEL PROCEDURE.

| Different Combination | unadjusted p | p Homm |
|---|---|---|
| Ranksum + Welch | **5.43E-12** | **1.25E-10** |
| MI + MUSIC | **1.37E-10** | **3.01E-09** |
| Ranksum + Burg | **1.70E-09** | **3.57E-08** |
| FDR + Welch | **4.91E-09** | **9.82E-08** |
| MI + Burg | **4.48E-08** | **8.07E-07** |
| BD + MUSIC | **5.85E-08** | **9.95E-07** |
| Ranksum + MUSIC | **9.91E-08** | **1.68E-06** |
| CORR + Burg | **2.77E-07** | **4.43E-06** |
| MI + Welch | **9.55E-06** | **1.43E-04** |
| CORR + MUSIC | **2.11E-05** | **2.95E-04** |
| FDR + MUSIC | **1.00E-04** | **0.001305** |
| FDR + Burg | **6.37E-04** | **0.007647** |
| CORR + Welch | **0.0020477** | **0.020477** |
| mRMR + MUSIC | **0.0031092** | **0.027983** |
| LR + MUSIC | **0.0046513** | **0.037211** |
| SR + MUSIC | **0.0079** | 0.0632 |
| BD + Burg | **0.0446384** | 0.312469 |
| mRMR + Welch | 0.1593641 | 0.637456 |
| SR + Burg | 0.2059032 | 0.823613 |
| mRMR + Burg | 0.2549452 | 0.91187 |
| BD + Welch | 0.5270893 | 0.91187 |
| LR + Welch | 0.837144 | 0.91187 |
| SR + Welch | 0.9118703 | 0.91187 |

of Burg with LR performs significantly better than all combinations of univariate method and feature extraction methods. It also performs significantly better than few combinations of multivariate method and feature extraction method.

### V. CONCLUSION

In this paper, we have examined the performance of the combination of three different PSD approaches, with four well-known uni-variates as well as four very popular multivariates, filter feature selection methods. The experimental findings demonstrate that the multivariate feature selection
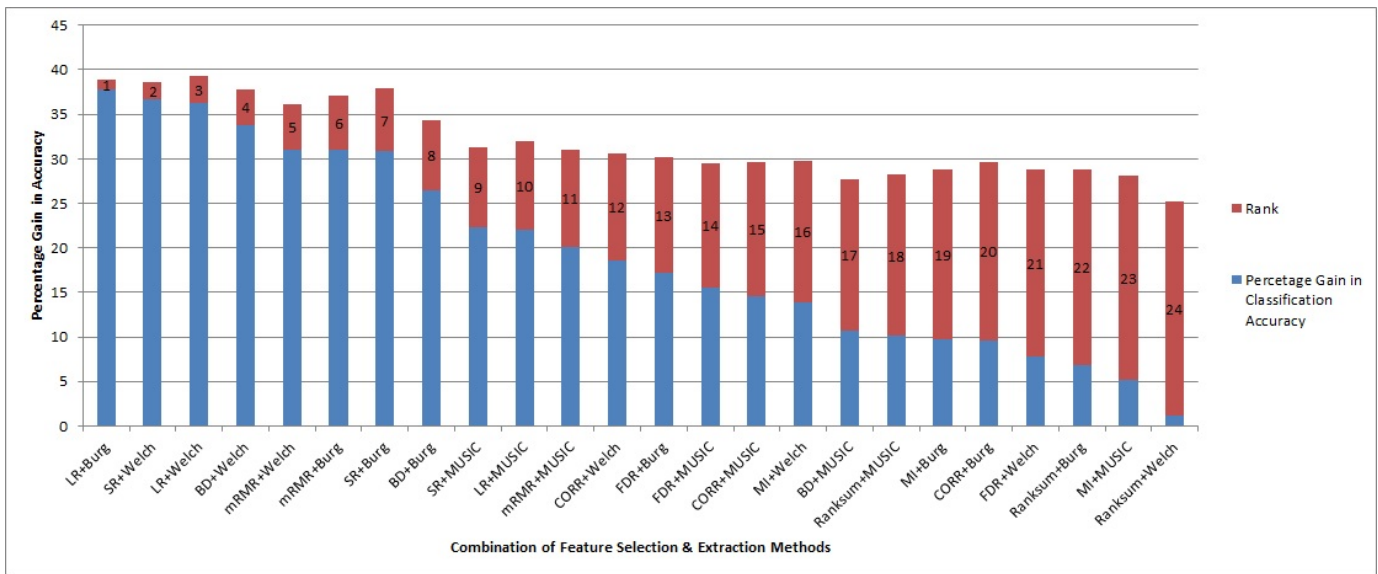
Fig. 10. Ranking of different combinations of Feature Extraction and Selection methods

algorithms endue more distinguishable feature set for the mental task classification, compared with univariate feature selection approach. The outcome determined features for the mental task classification by a minimal subset of relevant and non-redundant features. Experimental results demonstrate significant improvement in classification accuracy utilizing the selected feature selection methods. It is observed that the performance of multivariate filter feature selection methods is, in general, better than univariate filter feature selection methods. The combination of Burg method, LR and Linear Discriminant Analysis(LDA) achieved maximum classification accuracy among all other combinations.

Statistical tests also endorsed that the performance of the combination of Burg and the linear regression is notably different from the majority of the combinations. It has also been observed that for mental task classification multivariate feature selection approach works better than univariate feature selection approach in most of the cases with the conjunction of power spectral density approach.

In the future, we would like to extract spectral density of different brain frequency separately. Since the comparisons and investigations have been done on binary mental task classification, we would, therefore, like to extend this approach for multi-class mental tasks classification.

### REFERENCES

[1] B. Graimann, B. Allison, and G. Pfurtscheller, "Brain–computer interfaces: A gentle introduction," in *Brain-Computer Interfaces*. Springer, 2010, pp. 1–27.

[2] M. Rahman, M. Chowdhury, and S. Fattah, "An efficient scheme for mental task classification utilizing reflection coefficients obtained from autocorrelation function of eeg signal," *Brain informatics*, vol. 5, no. 1, p. 1, 2018.

[3] V. Gandhi, G. Prasad, D. Coyle, L. Behera, and T. M. McGinnity, "Eeg-based mobile robot control through an adaptive brain–robot interface," *IEEE Transactions on Systems, Man, and Cybernetics: Systems*, vol. 44, no. 9, pp. 1278–1285, 2014.

[4] X. Zhao, Y. Chu, J. Han, and Z. Zhang, "Ssvep based brain-computer interface controlled functional electrical stimulation system for upper extremity rehabilitation," *IEEE Transactions on Systems, Man, and Cybernetics: Systems*, vol. 46, no. 7, pp. 947–956, 2016.

[5] C. W. Anderson, E. A. Stolz, and S. Shamsunder, "Multivariate autoregressive models for classification of spontaneous electroencephalographic signals during mental tasks," *Biomedical Engineering, IEEE Transactions on*, vol. 45, no. 3, pp. 277–286, 1998.

[6] F. Babiloni, F. Cincotti, L. Lazzarini, J. Millan, J. Mourino, M. Varsta, J. Heikkonen, L. Bianchi, and M. Marciani, "Linear classification of low-resolution eeg patterns produced by imagined hand movements," *Rehabilitation Engineering, IEEE Transactions on*, vol. 8, no. 2, pp. 186–188, 2000.

[7] Z. A. Keirn and J. I. Aunon, "A new mode of communication between man and his surroundings," *Biomedical Engineering, IEEE Transactions on*, vol. 37, no. 12, pp. 1209–1214, 1990.

[8] A. Kübler, *Brain Computer Communication: Development of a Brain Computer Interface for Locked-in Patients on the Basis of the Psychophysiological Self-regulation Training of Slow Cortical Potentials (SCP)*. Schwäbische Verlags-Gesellschaft, 2000.

[9] G. Schalk, "Brain–computer symbiosis," *Journal of neural engineering*, vol. 5, no. 1, p. P1, 2008.

[10] W. He, Y. Zhao, H. Tang, C. Sun, and W. Fu, "A wireless bci and bmi system for wearable robots," *IEEE Transactions on Systems, Man, and Cybernetics: Systems*, vol. 46, no. 7, pp. 936–946, 2016.

[11] F. Akram, H.-S. Han, and T.-S. Kim, "A p300-based brain computer interface system for words typing," *Computers in biology and medicine*, vol. 45, pp. 118–125, 2014.

[12] A. Bashashati, M. Fatourechi, R. K. Ward, and G. E. Birch, "A survey of signal processing algorithms in brain–computer interfaces based on electrical brain signals," *Journal of Neural engineering*, vol. 4, no. 2, p. R32, 2007.

[13] Y. Zhang, B. Liu, X. Ji, and D. Huang, "Classification of eeg signals based on autoregressive model and wavelet packet decomposition," *Neural Processing Letters*, vol. 45, no. 2, pp. 365–378, 2017.

[14] G. Pfurtscheller, C. Neuper, D. Flotzinger, and M. Pregenzer, "Eeg-based discrimination between imagination of right and left hand movement," *Electroencephalography and clinical Neurophysiology*, vol. 103, no. 6, pp. 642–651, 1997.

[15] M. Kaper, P. Meinicke, U. Grossekathoefer, T. Lingner, and H. Ritter, "Bci competition 2003-data set iib: support vector machines for the p300 speller paradigm," *Biomedical Engineering, IEEE Transactions on*, vol. 51, no. 6, pp. 1073–1076, 2004.

[16] S. Chiappa, N. Donckers, S. Bengio, and F. Vrins, "Hmm and iohmm modeling of eeg rhythms for asynchronous bci systems." in *ESANN*, 2004, pp. 193–204.

[17] N. N. Neshov, A. H. Manolova, I. R. Draganov, K. T. Tonschev, and O. L. Boumbarov, "Classification of mental tasks from eeg signals using spectral analysis, pca and svm," *Cybernetics and Information Technologies*, vol. 18, no. 1, pp. 81–92, 2018.

[18] M. M. Moore, "Real-world applications for brain-computer interface technology," *Neural Systems and Rehabilitation Engineering, IEEE Transactions on*, vol. 11, no. 2, pp. 162–165, 2003.

[19] R. Palaniappan, R. Paramesran, S. Nishida, and N. Saiwaki, "A new brain-computer interface design using fuzzy artmap," *Neural Systems and Rehabilitation Engineering, IEEE Transactions on*, vol. 10, no. 3, pp. 140–148, 2002.

[20] W. D. Penny, S. J. Roberts, E. A. Curran, M. J. Stokes *et al.*, "Eeg-based communication: a pattern recognition approach," *IEEE Transactions on Rehabilitation Engineering*, vol. 8, no. 2, pp. 214–215, 2000.

[21] L. Qin, L. Ding, and B. He, "Motor imagery classification by means of source analysis for brain–computer interface applications," *Journal of Neural Engineering*, vol. 1, no. 3, p. 135, 2004.

[22] B. Kamousi, Z. Liu, and B. He, "Classification of motor imagery tasks for brain-computer interface applications by means of two equivalent dipoles analysis," *Neural Systems and Rehabilitation Engineering, IEEE Transactions on*, vol. 13, no. 2, pp. 166–171, 2005.

[23] M. Congedo, F. Lotte, and A. Lécuyer, "Classification of movement intention by spatially filtered electromagnetic inverse solutions," *Physics in medicine and biology*, vol. 51, no. 8, p. 1971, 2006.

[24] S. G. Mallat, "A theory for multiresolution signal decomposition: the wavelet representation," *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, vol. 11, no. 7, pp. 674–693, 1989.

[25] A. Gupta, R. Agrawal, and B. Kaur, "A three phase approach for mental task classification using eeg," in *Proceedings of the International Conference on Advances in Computing, Communications and Informatics*. ACM, 2012, pp. 898–904.

[26] N. E. Huang, Z. Shen, S. R. Long, M. C. Wu, H. H. Shih, Q. Zheng, N.-C. Yen, C. C. Tung, and H. H. Liu, "The empirical mode decomposition and the hilbert spectrum for nonlinear and non-stationary time series analysis," *Proceedings of the Royal Society of London. Series A: Mathematical, Physical and Engineering Sciences*, vol. 454, no. 1971, pp. 903–995, 1998.

[27] P. F. Diez, E. Laciar, V. Mut, E. Avila, and A. Torres, "Classification of mental tasks using different spectral estimation methods," *Biomedical Engineering, CAB de Mello, Ed. InTech*, pp. 287–306, 2009.

[28] A. Gupta and R. Agrawal, "Relevant feature selection from eeg signal for mental task classification," in *Advances in Knowledge Discovery and Data Mining*. Springer, 2012, pp. 431–442.

[29] V. Bajaj and R. B. Pachori, "Classification of seizure and nonseizure eeg signals using empirical mode decomposition," *IEEE Transactions on Information Technology in Biomedicine*, vol. 16, no. 6, pp. 1135–1142, 2012.

[30] T. M. Rutkowski, D. P. Mandic, A. Cichocki, and A. W. Przybyszewski, "Emd approach to multichannel eeg data—the amplitude and phase components clustering analysis," *Journal of Circuits, Systems, and Computers*, vol. 19, no. 01, pp. 215–229, 2010.

[31] A. S. Fine, D. P. Nicholls, and D. J. Mogul, "Assessing instantaneous synchrony of nonlinear nonstationary oscillators in the brain," *Journal of neuroscience methods*, vol. 186, no. 1, pp. 42–51, 2010.

[32] D. Mylonas, C. Siettos, I. Evdokimidis, A. Papanicolaou, and N. Smyrnis, "Modular patterns of phase desynchronization networks during a simple visuomotor task," *Brain topography*, vol. 29, no. 1, pp. 118–129, 2016.

[33] B. M. Battista, C. Knapp, T. McGee, and V. Goebel, "Application of the empirical mode decomposition and hilbert-huang transform to seismic reflection data," *Geophysics*, vol. 72, no. 2, pp. H29–H37, 2007.

[34] S. Adamczak, W. Makieła, and K. Stępień, "Investigating advantages and disadvantages of the analysis of a geometrical surface structure with the use of fourier and wavelet transform," *Metrology and Measurement Systems*, vol. 17, no. 2, pp. 233–244, 2010.

[35] R. Bellman, R. E. Bellman, R. E. Bellman, and R. E. Bellman, *Adaptive control processes: a guided tour*. Princeton University Press Princeton, 1961, vol. 4.

[36] I. Guyon and A. Elisseeff, "An introduction to variable and feature selection," *The Journal of Machine Learning Research*, vol. 3, pp. 1157–1182, 2003.

[37] P. Stoica and R. L. Moses, *Spectral analysis of signals*. Pearson/Prentice Hall Upper Saddle River, NJ, 2005.

[38] J. G. Proakis, *Digital signal processing: principles algorithms and applications*. Pearson Education India, 2001.

[39] G. Pfurtscheller, C. Neuper, A. Schlogl, and K. Lugger, "Separability of eeg signals recorded during right and left motor imagery using adaptive autoregressive parameters," *Rehabilitation Engineering, IEEE Transactions on*, vol. 6, no. 3, pp. 316–325, 1998.

[40] H. Akaike, "Fitting autoregressive models for prediction," *Annals of the institute of Statistical Mathematics*, vol. 21, no. 1, pp. 243–247, 1969.

[41] J. Rissanen, "A universal prior for integers and estimation by minimum description length," *The Annals of statistics*, pp. 416–431, 1983.

[42] H. Akaike, "A new look at the statistical model identification," *Automatic Control, IEEE Transactions on*, vol. 19, no. 6, pp. 716–723, 1974.

[43] E. Parzen, "On consistent estimates of the spectrum of a stationary time series," *The Annals of Mathematical Statistics*, pp. 329–348, 1957.

[44] S. H. Kia, H. Henao, and G.-A. Capolino, "A high-resolution frequency estimation method for three-phase induction machine fault detection," *Industrial Electronics, IEEE Transactions on*, vol. 54, no. 4, pp. 2305–2314, 2007.

[45] E. D. Übeyli, "Implementing eigenvector methods/probabilistic neural networks for analysis of eeg signals," *Neural networks*, vol. 21, no. 9, pp. 1410–1417, 2008.

[46] R. Kohavi and G. H. John, "Wrappers for feature subset selection," *Artificial intelligence*, vol. 97, no. 1, pp. 273–324, 1997.

[47] I. Koprinska, "Feature selection for brain-computer interfaces," in *New Frontiers in Applied Data Mining*. Springer, 2010, pp. 106–117.

[48] G. Rodriguez-Bermudez, P. J. Garcia-Laencina, and J. Roca-Dorda, "Efficient automatic selection and combination of eeg features in least squares classifiers for motor imagery brain–computer interfaces," *International journal of neural systems*, vol. 23, no. 04, 2013.

[49] C. Guerrero-Mosquera, M. Verleysen, and A. N. Vazquez, "Eeg feature selection using mutual information and support vector machine: A comparative analysis," in *Engineering in Medicine and Biology Society (EMBC), 2010 Annual International Conference of the IEEE*. IEEE, 2010, pp. 4946–4949.

[50] M. Murugappan, N. Ramachandran, Y. Sazali *et al.*, "Classification of human emotion from eeg using discrete wavelet transform," *Journal of Biomedical Science and Engineering*, vol. 3, no. 04, p. 390, 2010.

[51] S. Bhattacharyya, A. Sengupta, T. Chakraborti, A. Konar, and D. Tibarewala, "Automatic feature selection of motor imagery eeg signals using differential evolution and learning automata," *Medical & biological engineering & computing*, vol. 52, no. 2, pp. 131–139, 2014.

[52] N. S. Dias, M. Kamrunnahar, P. M. Mendes, S. Schiff, and J. H. Correia, "Feature selection on movement imagery discrimination and attention detection," *Medical & biological engineering & computing*, vol. 48, no. 4, pp. 331–341, 2010.

[53] K. Pearson, "Notes on the history of correlation," *Biometrika*, pp. 25–45, 1920.

[54] S. Dowdy, S. Wearden, and D. Chilko, *Statistics for research*. John Wiley & Sons, 2011, vol. 512.

[55] C. E. Shannon and W. Weaver, "The mathematical theory of communication (urbana, il," 1949.

[56] F. Wilcoxon, "Individual comparisons by ranking methods," *Biometrics bulletin*, pp. 80–83, 1945.

[57] S. Li, X. Wu, and M. Tan, "Gene selection using hybrid particle swarm optimization and genetic algorithm," *Soft Computing*, vol. 12, no. 11, pp. 1039–1048, 2008.

[58] A. Bhattacharyya, "On a measure of divergence between two multinomial populations," *Sankhyā: The Indian Journal of Statistics*, pp. 401–406, 1946.

[59] P. A. Devijver and J. Kittler, *Pattern recognition: A statistical approach*. Prentice-Hall London, 1982, vol. 761.

[60] H.-S. Park, S.-H. Yoo, and S.-B. Cho, "Forward selection method with regression analysis for optimal gene selection in cancer classification," *International Journal of Computer Mathematics*, vol. 84, no. 5, pp. 653–667, 2007.

[61] H. Peng, F. Long, and C. Ding, "Feature selection based on mutual information criteria of max-dependency, max-relevance, and min-redundancy," *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, vol. 27, no. 8, pp. 1226–1238, 2005.

[62] H. Chernoff, "A measure of asymptotic efficiency for tests of a hypothesis based on the sum of observations," *The Annals of Mathematical Statistics*, pp. 493–507, 1952.

[63] R. Adhikari and R. Agrawal, "Performance evaluation of weights selection schemes for linear combination of multiple forecasts," *Artificial Intelligence Review*, pp. 1–20, 2012.

[64] A. Gupta and D. Kumar, "Fuzzy clustering-based feature extraction method for mental task classification," *Brain informatics*, vol. 4, no. 2, pp. 135–145, 2017.

[65] J. Derrac, S. García, D. Molina, and F. Herrera, "A practical tutorial on the use of nonparametric statistical tests as a methodology for comparing evolutionary and swarm intelligence algorithms," *Swarm and Evolutionary Computation*, vol. 1, no. 1, pp. 3–18, 2011.