



**HAL**  
open science

# t-WDA: A novel Discriminant Analysis applied to EEG classification

Imen Ayadi, Florent Bouchard, Frédéric Pascal

► **To cite this version:**

Imen Ayadi, Florent Bouchard, Frédéric Pascal. t-WDA: A novel Discriminant Analysis applied to EEG classification. 2023 31st European Signal Processing Conference (EUSIPCO), Sep 2023, Helsinki, Finland. 10.23919/eusipco58844.2023.10289799 . hal-04249018

**HAL Id: hal-04249018**

**<https://hal.science/hal-04249018v1>**

Submitted on 19 Oct 2023

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# $t$ -WDA: A novel Discriminant Analysis applied to EEG classification

Imen Ayadi, Florent Bouchard, Frédéric Pascal

Université Paris-Saclay, CNRS, CentraleSupélec, Laboratoire des signaux et systèmes

91190, Gif-sur-Yvette, France

**Abstract**—This paper provides a new classification method of covariance matrices exploiting the  $t$ -Wishart distribution, which generalizes the Wishart distribution. Compared to the Wishart distribution, it is more robust to aberrant covariance matrices and more flexible to distribution mismatch. Following recent developments on this matrix-variate distribution, the proposed classifier is obtained by leveraging the Discriminant Analysis framework and providing original decision rules. The practical interest of our approach is shown thanks to numerical experiments on real data. More precisely, the proposed classifier yields the best results on two standard electroencephalography datasets compared to the best state-of-the-art minimum distance-to-mean (MDM) classifiers.

**Index Terms**—EEG, Covariance matrices,  $t$ -Wishart, Bayesian classification, Discriminant Analysis, BCI.

## I. INTRODUCTION

In signal processing, covariance matrices have recently gained interest in classification problems. Not only do they embed relevant information from signals, but they also fall into the domain of Riemannian geometry. They have proved their merit in radar and image processing [1], biomedical signals analysis, *etc.* For instance, their use has revolutionized the field of Electroencephalography (EEG), which consists in recording brain signals for medical purposes or to help people with motor impairment via brain-computer interfaces (BCI) [2], [3]. The first generation of EEG classification methods uses spatial filters paired with Euclidean classifiers [2]. Such an approach has been outperformed for a decade using spatial covariance matrices as features [4]. There are two main methods for classifying covariance matrices either directly on their native space of symmetric positive definite (SPD) matrices via the Minimum Distance to Mean (MDM) or by projecting them onto a tangent space at a reference SPD matrix and then using a Euclidean classifier for the projected matrices [4].

In order to face the challenges of the intrinsic non-Gaussianity of signals, covariance matrices are often estimated thanks to robust techniques [5] before proceeding to MDM or tangent space projection methods. However, this trick needs to be more efficient in the case of aberrant recordings or mislabelling. Such situations, which often occur when manipulating signals, explain the motivation to consider the outliers at the scale of the covariance matrices and not only at the scale of a single measurement. Accordingly, it is interesting to

design probabilistic classifiers that exploit statistics over SPD [6], particularly heavy-tailed matrix-variate distributions. We would mention that such a strategy was already implemented on matrices (not specifically SPD) [7] where an Expectation-Maximization algorithm is suggested using matrix-variate  $t$ -distributions. Among the distributions handling outliers over SPD matrices, figure the  $t$ -Wishart ones [8], [9]. They are known to generalize the most classical distribution on SPD matrices, the Wishart distribution [10]. The latter is the distribution of sample covariance matrices of random vectors drawn from a multivariate Gaussian distribution. The generalization brought by the  $t$ -Wishart distributions can be similarly interpreted as the extension to the multivariate Gaussian distribution brought by the multivariate  $t$ -distributions [11].

The main contribution of this paper is the development of a new Bayesian classifier,  $t$ -WDA ( $t$ -Wishart Discriminant Analysis), that uses the  $t$ -Wishart distribution for the likelihood of observed data. Moreover, the Discriminant Analysis framework offers a novel interpretation of the most classical versions of MDM. A particular focus is drawn on the merit of  $t$ -WDA for EEG classification.

The paper is organized as follows. Section II provides a brief background about EEG classification, mainly on the methods based on covariance matrices as features. Section III reviews the  $t$ -Wishart distributions. Then, Section IV derives the  $t$ -Wishart Discriminant Analysis classifier. Some numerical experiments in Section V validate its merit for BCI. Finally, concluding remarks and perspectives are drawn in Section VI.

## II. BACKGROUND ON EEG CLASSIFICATION: MDM CLASSIFIER

As the introduction mentions, various methods have been considered to classify EEG signals [2]. The most efficient ones usually rely on covariance matrices. In particular, the reference method in the BCI community is the MDM classifier [4]. This section introduces this EEG classification method, which will be used as a baseline in this work.

The MDM classifier can be seen as a supervised  $k$ -means classification method: it aims to minimize the dispersion inside classes according to a given distance or divergence. In EEG, signals  $\mathbf{X} \in \mathbb{R}^{p \times n}$  ( $p$ , number of electrodes;  $n$ , number of samples) are supposed centered, *i.e.*,  $\mathbb{E}[\mathbf{X}] = \mathbf{0}$ .

<sup>1</sup>This work is partially supported by a public grant overseen by the French National Research Agency (ANR) through the program UDOPIA, project funded by the ANR-20-THIA-0013-01<sup>1</sup>.

The covariance matrix of  $\mathbf{X}$  is estimated with the sample covariance matrix (SCM), defined as

$$\mathbf{C} = \frac{1}{n} \mathbf{X} \mathbf{X}^T. \quad (1)$$

By construction,  $\mathbf{C} \in \mathcal{S}_p^{++}$ , the set of  $p \times p$  SPD matrices. In the training step, a barycenter  $\hat{\Sigma}_k$  is computed for each class  $k$ . Given covariance matrices  $\{\mathbf{C}_i\}_{i=1}^{N_k}$  in class  $k$ , it is

$$\hat{\Sigma}_k = \operatorname{argmin}_{\Sigma \in \mathcal{S}_p^{++}} \sum_{i=1}^{N_k} d(\mathbf{C}_i, \Sigma), \quad (2)$$

where  $d$  is a divergence or the square of a distance on  $\mathcal{S}_p^{++}$ . In the testing step, a covariance matrix  $\mathbf{C}$  is assigned to the class  $y$  whose barycenter  $\hat{\Sigma}_y$  is the closest. Formally, the decision rule is:

$$y = \operatorname{argmin}_{k \in [1, K]} d(\mathbf{C}, \hat{\Sigma}_k). \quad (3)$$

Several divergences/distances have been considered; see e.g., [3] for a review. The original choice, which is also one of the most natural since it corresponds to the Fisher distance of the multivariate Gaussian distribution [12], is the natural Riemannian – or affine-invariant – distance defined as

$$\delta_{\mathbb{R}}^2(\mathbf{A}, \mathbf{B}) = \|\log(\mathbf{B}^{-1} \mathbf{A})\|_2^2, \quad (4)$$

where  $\|\cdot\|_2$  denotes the Frobenius norm and  $\log(\cdot)$  is the matrix logarithm. The corresponding barycenter of  $\{\mathbf{C}_i\}_{i=1}^N$  is the so-called geometric mean, which is the unique solution to the fixed-point equation

$$\hat{\Sigma} = \hat{\Sigma} \exp\left(\frac{1}{N} \sum_{i=1}^N \log(\hat{\Sigma}^{-1} \mathbf{C}_i)\right). \quad (5)$$

where  $\exp(\cdot)$  is the matrix exponential. As no closed-form expression is known for  $p > 2$ , it is computed using a recursive algorithm based on (5) (see, e.g., [4]).

Another choice of particular interest for  $d$  is the Kullback-Leibler divergence between two centered multivariate Gaussian distributions [13], defined as

$$d_{\text{KL}}(\mathbf{A}, \mathbf{B}) = \operatorname{tr}(\mathbf{A} \mathbf{B}^{-1}) - \log |\mathbf{A} \mathbf{B}^{-1}| - p, \quad (6)$$

where  $\operatorname{tr}(\cdot)$  and  $|\cdot|$  denote the trace and determinant operators, respectively. Since the Kullback-Leibler divergence is not symmetric, one can obtain two barycenters. Here, we use the “right” barycenter, i.e., we choose  $\mathbf{A} = \mathbf{C}_i$  and  $\mathbf{B} = \Sigma$  as in (2) (see [14] for a review on the barycenters of this Kullback-Leibler divergence). The corresponding right barycenter of  $\{\mathbf{C}_i\}_{i=1}^N$  is simply the arithmetic mean

$$\hat{\Sigma} = \frac{1}{N} \sum_{i=1}^N \mathbf{C}_i. \quad (7)$$

The MDM classifiers associated with (4) and (6) are denoted RMDM and rKLMDM, respectively. In Section IV, it is shown that they can be seen as Discriminant Analysis classifiers under specific assumptions.

### III. WISHART AND $t$ -WISHART DISTRIBUTIONS

In this section, we present the Wishart and  $t$ -Wishart distributions. After providing definitions, we give estimation methods of their parameter, which lies in  $\mathcal{S}_p^{++}$ . While the maximum likelihood estimator (MLE) of the Wishart distribution is known in closed form, one has to employ an iterative algorithm to obtain the  $t$ -Wishart’s one. In this work, we propose a method relying on Riemannian optimization that is very similar to the one in [15]. However, it is a bit simpler as it relies on the usual affine-invariant metric on  $\mathcal{S}_p^{++}$  rather than on the Fisher information metric of the distribution.

The probability density function (pdf) of a random matrix  $\mathbf{S} = \mathbf{X} \mathbf{X}^T \in \mathcal{S}_p^{++}$ , with  $\mathbf{X} \in \mathbb{R}^{p \times n}$  ( $n \geq p$ ), following the Wishart distribution  $\mathcal{W}(n, \Sigma)$  with center  $\Sigma \in \mathcal{S}_p^{++}$  is, up to a normalization factor,

$$f^{\mathcal{W}}(\mathbf{S}|\Sigma) \propto |\Sigma|^{-\frac{n}{2}} |\mathbf{S}|^{\frac{n-p-1}{2}} \exp\left(-\frac{1}{2} \operatorname{tr}(\Sigma^{-1} \mathbf{S})\right). \quad (8)$$

As for the usual multivariate  $t$ -distribution, the  $t$ -Wishart distribution is obtained by replacing the exponential in (8) with another density generator. The pdf of  $\mathbf{S} = \mathbf{X} \mathbf{X}^T$  following the  $t$ -Wishart distribution  $t\text{-}\mathcal{W}(n, \Sigma, \nu)$  with center  $\Sigma$  and degree of freedom (d.o.f.)  $\nu > 0$  is, up to a normalization factor,

$$f^{t\text{-}\mathcal{W}}(\mathbf{S}|\Sigma) \propto |\Sigma|^{-\frac{n}{2}} |\mathbf{S}|^{\frac{n-p-1}{2}} \left(1 + \frac{\operatorname{tr}(\Sigma^{-1} \mathbf{S})}{\nu}\right)^{-\frac{\nu+n p}{2}}. \quad (9)$$

Notice that, compared to the Wishart distribution, the  $t$ -Wishart distribution adds a dependence on the columns of  $\mathbf{X}$ . The independence in the case of the Wishart distribution is a direct consequence of the properties of the exponential function. Thus, substituting it with another density generator cancels this property.

Given independent and identically distributed (i.i.d.) samples  $\{\mathbf{S}_i\}_{i=1}^N$ , the log-likelihood  $\mathcal{L}$  associated to the distribution with pdf  $f$  is  $\mathcal{L}(\Sigma) = \sum_i \log f(\mathbf{S}_i|\Sigma)$ . The log-likelihoods of the Wishart and  $t$ -Wishart distributions are denoted  $\mathcal{L}^{\mathcal{W}}$  and  $\mathcal{L}^{t\text{-}\mathcal{W}}$ , respectively. To obtain the MLE, one needs to solve the optimization problem  $\hat{\Sigma} = \operatorname{argmax}_{\Sigma} \mathcal{L}(\Sigma)$ . In the case of the Wishart distribution, it is

$$\hat{\Sigma}^{\mathcal{W}} = \frac{1}{nN} \sum_{i=1}^N \mathbf{S}_i. \quad (10)$$

For the  $t$ -Wishart distribution, we must develop an iterative algorithm. Here, we perform a Riemannian gradient descent [16] on  $\mathcal{S}_p^{++}$  to obtain the MLE  $\hat{\Sigma}^{t\text{-}\mathcal{W}}$ . Given a sequence of iterates  $\{\hat{\Sigma}^{(j)}\}$ , the next iterate is

$$\hat{\Sigma}^{(j+1)} = R_{\hat{\Sigma}^{(j)}}\left(t_j \nabla \mathcal{L}^{t\text{-}\mathcal{W}}\left(\hat{\Sigma}^{(j)}\right)\right), \quad (11)$$

where  $\nabla \mathcal{L}^{t\text{-}\mathcal{W}}$  is the Riemannian gradient of  $\mathcal{L}^{t\text{-}\mathcal{W}}$  in  $\mathcal{S}_p^{++}$ ,  $R_{\cdot}(\cdot)$  is a retraction on  $\mathcal{S}_p^{++}$  (mapping from tangent spaces back onto the manifold), and  $t_j$  is a stepsize computed through a linesearch [16]. To define the Riemannian gradient, one has to choose a Riemannian metric on  $\mathcal{S}_p^{++}$ . Here, we select the usual affine-invariant metric. Given  $\hat{\Sigma} \in \mathcal{S}_p^{++}$  and

tangent vectors  $\xi, \eta \in \mathcal{S}_p$  (set of  $p \times p$  symmetric matrices), it is  $\langle \xi, \eta \rangle_{\hat{\Sigma}} = \text{tr}(\hat{\Sigma}^{-1} \xi \hat{\Sigma}^{-1} \eta)$ . The Riemannian gradient  $\nabla \mathcal{L}^{t\text{-}\mathcal{W}}(\hat{\Sigma})$  at  $\hat{\Sigma}$  is then defined as the only tangent vector such that, for all  $\xi \in \mathcal{S}_p$ ,  $\langle \nabla \mathcal{L}^{t\text{-}\mathcal{W}}(\hat{\Sigma}), \xi \rangle_{\hat{\Sigma}} = D \mathcal{L}^{t\text{-}\mathcal{W}}(\hat{\Sigma})[\xi]$ , where  $D$  denotes the directional derivative. One can show

$$\nabla \mathcal{L}^{t\text{-}\mathcal{W}}(\hat{\Sigma}) = \frac{1}{2} \sum_{i=1}^N \frac{\nu + np}{\nu + \text{tr}(\hat{\Sigma}^{-1} \mathbf{S}_i)} \mathbf{S}_i - \frac{nN}{2} \hat{\Sigma}. \quad (12)$$

Finally, the retraction we choose is the second-order approximation of geodesics on  $\mathcal{S}_p^{++}$  (a generalization of straight lines to manifolds), which is arguably the best choice from a numerical perspective [17]. Given  $\hat{\Sigma}$  and  $\xi$ , it is

$$R_{\hat{\Sigma}}(\xi) = \hat{\Sigma} + \xi + \frac{1}{2} \xi \hat{\Sigma}^{-1} \xi. \quad (13)$$

#### IV. DISCRIMINANT ANALYSIS CLASSIFIERS

In this section, a novel Bayesian classifier is proposed for signal classification. Inspired by Linear/Quadratic Discriminant Analysis (LDA/QDA) in the multivariate case [18], the introduced method generalizes the concept of Discriminant Analysis to the matrix-variate case. It exploits Wishart and  $t$ -Wishart distributions to model the likelihood of observed samples knowing their associated labels.

##### A. Wishart and $t$ -Wishart Discriminant Analysis classifiers

Rather than considering SCMs,  $\mathbf{S} = \mathbf{X}\mathbf{X}^T$  is used. The following proposition introduces the  $t$ -Wishart Discriminant Analysis, denoted  $t\text{-}\mathcal{WDA}$ .

**Proposition 1** ( $t\text{-}\mathcal{WDA}$ ). *The decision rule of the  $t$ -Wishart Discriminant Analysis, for a testing covariance matrix  $\mathbf{S}$ , is given by*

$$\hat{y}(\mathbf{S}) = \underset{k \in [1, K]}{\text{argmax}} \delta_k^{t\text{-}\mathcal{W}}(\mathbf{S}), \quad (14)$$

with the discriminant function

$$\delta_k^{t\text{-}\mathcal{W}}(\mathbf{S}) = \log(\hat{\pi}_k) - \frac{n}{2} \log |\hat{\Sigma}_k| - \frac{\nu + np}{2} \log \left( 1 + \frac{\text{tr}(\hat{\Sigma}_k^{-1} \mathbf{S})}{\nu} \right), \quad (15)$$

where  $\hat{\pi}_k$  is the proportion of the class  $k$  in the training set;  $\hat{\Sigma}_k$  is the estimated barycenter of training covariance matrices of the class  $k$ , computed thanks to (11), and  $\nu$  is the degree of freedom of the model<sup>1</sup>.

*Proof.* Being a probabilistic classifier, the  $t\text{-}\mathcal{WDA}$  has the following discrimination rule

$$\hat{y}(\mathbf{S}) = \underset{k \in [1, K]}{\text{argmax}} P(y = k | \mathbf{S}). \quad (16)$$

Using Bayes formula,  $P(y = k | \mathbf{S}) \propto \pi_k f_k(\mathbf{S})$ , where  $\pi_k$  is the prior probability of the class  $k$  and  $f_k(\mathbf{S})$  is the conditional probability density function observed on the class  $k$ .

$t\text{-}\mathcal{WDA}$  approaches the classification problem by assuming that  $\mathbf{S} | y = k$  follows the  $t$ -Wishart distribution  $t\text{-}\mathcal{W}(n, \Sigma_k, \nu)$

where  $\Sigma_k$  is considered as the center of the class  $k$ . Hence, by plugging in the pdf defined in (9) and applying the logarithm function, (16) becomes

$$\hat{y}(\mathbf{S}) = \underset{k \in [1, K]}{\text{argmax}} \log(\pi_k f^{t\text{-}\mathcal{W}}(\mathbf{S} | \Sigma_k)). \quad (17)$$

where  $f^{t\text{-}\mathcal{W}}(\mathbf{S} | \Sigma_k)$  is defined in (9).

Concerning the prior distribution of the class  $k$ ,  $\pi_k$  is simply estimated by the proportion of training samples belonging to class  $k$ , denoted as  $\hat{\pi}_k$ . Notice that  $\hat{\pi}_k$  is a consistent estimate of  $\pi_k$ , a direct consequence of the strong law of large numbers.

The centers  $\{\Sigma_k\}_{k=1}^K$  are also estimated during the training step. A natural estimation choice is the MLE of  $t\text{-}\mathcal{W}(n, \Sigma, \nu)$ , denoted  $\{\hat{\Sigma}_k\}_{k=1}^K$ . They are computed as described in Section III. We assume that the MLE of the center for the  $t$ -Wishart distribution is consistent<sup>2</sup>. As a result of the continuous mapping theorem and the Slutsky lemma,  $\log(\hat{\pi}_k f^{t\text{-}\mathcal{W}}(\mathbf{S} | \hat{\Sigma}_k))$  converges in probability to  $\log(\pi_k f^{t\text{-}\mathcal{W}}(\mathbf{S} | \Sigma_k))$  when the size of the training set tends to  $+\infty$ . The term depending only on  $\mathbf{S}$  is then neglected due to the argmax operator, which concludes the proof.  $\square$

**Remark 1.** *As mentioned in Proposition 1, the d.o.f.  $\nu$  is a hyperparameter; estimating it during the training step is beyond the scope of this paper. A possible way to tune it is to perform a grid search. Furthermore, we would highlight that in the actual model, a unique d.o.f. is chosen for all classes for simplicity. It would be more realistic to estimate different d.o.f. for each class to improve discrimination, especially when the centers of classes are close.*

Similarly, the Wishart Discriminant Analysis classifier denoted  $\mathcal{WDA}$ , is derived in the following proposition.

**Proposition 2** ( $\mathcal{WDA}$ ). *The decision rule of the Wishart Discriminant Analysis, for a testing covariance matrix  $\mathbf{S}$ , is given by*

$$\hat{y}(\mathbf{S}) = \underset{k \in [1, K]}{\text{argmax}} \delta_k^{\mathcal{W}}(\mathbf{S}), \quad (18)$$

with the discriminant function

$$\delta_k^{\mathcal{W}}(\mathbf{S}) = \log(\hat{\pi}_k) - \frac{n}{2} \log |\hat{\Sigma}_k| - \frac{1}{2} \text{tr}(\hat{\Sigma}_k^{-1} \mathbf{S}).$$

where  $\hat{\pi}_k$  is the proportion of the class  $k$  in the training set, and  $\hat{\Sigma}_k$  is the estimated barycenter of training covariance matrices of the class  $k$ , computed thanks to Eq. (10).

*Proof.* The proof follows the same steps as for the  $t\text{-}\mathcal{WDA}$ .  $\square$

The discriminant function of the  $t\text{-}\mathcal{WDA}$  approximates the  $\mathcal{WDA}$ 's one when the dof  $\nu$  tends to  $+\infty$ .

**Remark 2.** *Let us examine the decision boundary between two classes  $k$  and  $l$  for  $t\text{-}\mathcal{WDA}$  and  $\mathcal{WDA}$ . We recall that it corresponds to the region  $\{\mathbf{S} \in \mathcal{S}_p^{++} : \delta_k(\mathbf{S}) = \delta_l(\mathbf{S})\}$ , where  $\delta_k := \delta_k^{t\text{-}\mathcal{W}}$  for  $t\text{-}\mathcal{WDA}$  and  $\delta_k := \delta_k^{\mathcal{W}}$  for  $\mathcal{WDA}$ .*

<sup>1</sup>considered here as a hyperparameter

<sup>2</sup>Due to the lack of space, this will be proved in a forthcoming paper

The decision boundary is linear in  $\mathbf{S}$  since it is a solution to an equation involving the trace of  $\mathbf{M}_{kl}\mathbf{S}$ , where  $\mathbf{M}_{kl}$  is a symmetric matrix that depends on the model parameters. This makes the proposed classifier very close to LDA. However, finding decision boundaries for  $\mathcal{WDA}/t\text{-}\mathcal{WDA}$  is more complex due the constraint  $\mathbf{S} \in \mathcal{S}_p^{++}$ .

### B. Interpretation of RMDM and rKLMDM in the framework of Discriminant Analysis

As mentioned in Section II, rKLMDM can be related to the proposed  $\mathcal{WDA}$  classifier. In fact, if the training set is balanced, rKLMDM corresponds exactly to  $\mathcal{WDA}$ . First, one can observe that the MLE for the center of the i.i.d. samples  $\{\mathbf{S}_i\}_{i=1}^N$  following the Wishart distribution is the arithmetic mean of  $\{\mathbf{C}_i\}_{i=1}^N$  (with the notation  $\mathbf{S}_i = n\mathbf{C}_i$ ). Therefore, rKLMDM and  $\mathcal{WDA}$  share the same barycenters of classes. Furthermore, one can write

$$n \log |\hat{\Sigma}_k| + \text{tr}(\hat{\Sigma}_k^{-1} \mathbf{S}) = n(d_{\text{KL}}(\mathbf{C}|\hat{\Sigma}_k) + p + \log |\mathbf{C}|).$$

Hence, rKLMDM and  $\mathcal{WDA}$  have the same decision rule.

Similarly, under the condition of a balanced training set, RMDM can be seen as the Discriminant Analysis classifier associated with the Riemannian Gaussian distribution on  $\mathcal{S}_p^{++}$  defined in [19]. We recall that the pdf of the Riemannian Gaussian distribution  $\mathcal{G}(\Sigma, \sigma^2)$  with center  $\Sigma$  and dispersion parameter  $\sigma > 0$  is, up to a normalization factor,  $f(\mathbf{C}|\Sigma) \propto \exp(-\frac{1}{2\sigma^2} \delta_{\text{R}}^2(\Sigma, \mathbf{C}))$ . Moreover, its MLE corresponds to the geometric mean of the samples. The observed likelihood associated with RMDM in a Discriminant Analysis framework is obtained with the model  $\mathbf{C}|y = k \sim \mathcal{G}(\Sigma_k, \sigma^2)$  (or equivalently,  $\mathbf{S}|y = k \sim \mathcal{G}(n\Sigma_k, \sigma^2)$ ). The barycenters of classes in the MDM classifier are exactly the MLEs for centers of Riemannian Gaussian distributions. Under the assumption of equal proportions  $\pi_k$ 's, the decision rule for the MDM given by (3) with the affine-invariant Riemannian distance (4) is exactly recovered.

## V. NUMERICAL EXPERIMENTS

This section aims to show the practical interest of the proposed  $t\text{-}\mathcal{WDA}$  classifier. To do so, numerical experiments are conducted on two different real EEG datasets. The datasets, which are available on the MOABB platform<sup>3</sup>, correspond to two different BCI paradigms. The first one [3] contains steady-state visually evoked potentials (SSVEP) while the second one [20] concerns motor imagery (MI). In both cases, the classification performances of RMDM,  $\mathcal{WDA}$  (or rKLMDM) and  $t\text{-}\mathcal{WDA}$  are compared.

### A. SSVEP dataset

SSVEP signals are natural responses to repetitive visual stimuli at specific frequencies, *i.e.*, the visual cortex synchronizes with the stimuli. In this work, we employ the dataset used in [3], where 12 subjects are asked to look at light emitting diodes (LEDs) blinking at three different frequencies:

17, 13 and 21 Hz. EEG signals are acquired on eight electrodes located around the visual cortex and the sampling rate is set to 256 Hz. For each subject, the recordings of two to five sessions are available. A session contains 32 labeled trials of five seconds each. Data are equally divided into four classes: a class per stimulus frequency and a resting class where no light is blinking.

Before proceeding to the classification, a preprocessing step is needed. This allows us to better extract the information at the frequencies of interest in the data. Following [3], each EEG trial  $\mathbf{X}$  becomes

$$\tilde{\mathbf{X}} = \begin{bmatrix} \mathbf{X}^{(17)} \\ \mathbf{X}^{(13)} \\ \mathbf{X}^{(21)} \end{bmatrix}, \quad (19)$$

where  $\mathbf{X}^{(f)}$  is the result of backward-forward filtering of  $\mathbf{X}$  with a band-pass region  $[f - \Delta f, f + \Delta f]$ , where  $\Delta f = 0.5$  Hz. Therefore,  $\tilde{\mathbf{X}} \in \mathbb{R}^{p \times n}$  with  $p = 3 \times 8 = 24$  and  $n = 5 \times 256 = 1280$ . Transformed signals  $\tilde{\mathbf{X}}$  are centered and their covariances are computed.

A within-session classification is run: for each subject, a session is randomly divided into a training set composed of 20 trials and a testing set of 12 trials such that all classes are equally represented in the training set. A Monte-Carlo approach of 100 permutations is performed in each session. For the  $t\text{-}\mathcal{WDA}$ , the degree of freedom  $\nu$  is set to 10. The means and standard deviations of test accuracies for each subject are reported in Figure 1.

One can observe that the RMDM yields the best results for two subjects (2 and 4);  $\mathcal{WDA}$  is the most accurate for three subjects (6, 8, and 9); and  $t\text{-}\mathcal{WDA}$  features the best performance for the remaining seven subjects. Averaging over all the subjects,  $t\text{-}\mathcal{WDA}$  has a 3.11% gain over RMDM and a 1.19% gain over  $\mathcal{WDA}$ . Thus, our proposed  $t\text{-}\mathcal{WDA}$  classifier appears advantageous on these SSVEP data.

### B. MI dataset

MI is a mental process during which a subject mentally simulates a physical action, *i.e.*, he/she imagines moving their right or left hand, feet, tongue, *etc.* In this work, we consider the BNCI2014001 dataset [20]. It contains the EEG recordings from 9 subjects. Signals are acquired via 22 electrodes with a sampling rate of 250 Hz. For each subject, two sessions composed of 288 trials are available. Only the four last seconds of the trial are taken into account. Thus, for each trial, we have  $\mathbf{X} \in \mathbb{R}^{p \times n}$  with  $p = 22$  and  $n = 4 \times 250 = 1000$ .

A within-session classification is run: cross-validation of stratified 5-folds is considered on each session, dividing it into 228 training samples and 60 testing samples with a balanced representation of the 4 classes. For the  $t\text{-}\mathcal{WDA}$ , the degree of freedom  $\nu$  is set to 10. The means and standard deviations of test accuracies for each subject are reported in Figure 2.

One can observe that  $\mathcal{WDA}$  never features the best performance on this dataset. RMDM is the most accurate for two subjects (1 and 5) and has a similar accuracy for subject 2 compared to  $t\text{-}\mathcal{WDA}$ . Finally,  $t\text{-}\mathcal{WDA}$  yields the best results

<sup>3</sup><https://github.com/NeuroTechX/moabb>

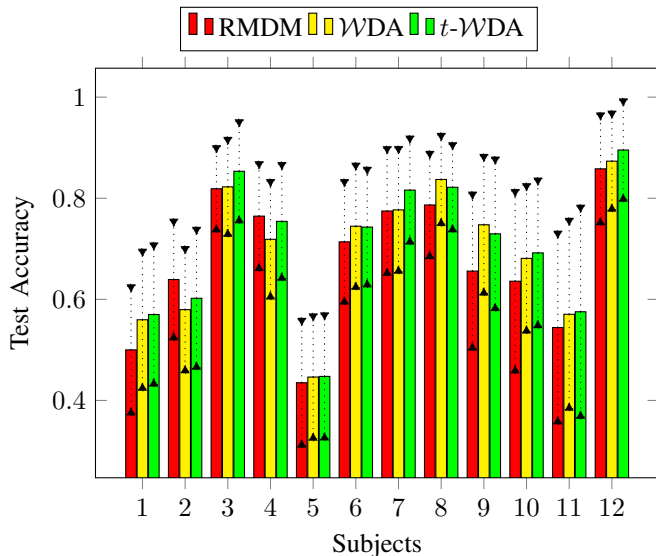


Fig. 1. Comparison of classifiers on the SSVEP dataset [3]. Standard deviations are plotted in dotted lines. The degree of freedom for  $t$ -WDA is  $\nu = 10$ .

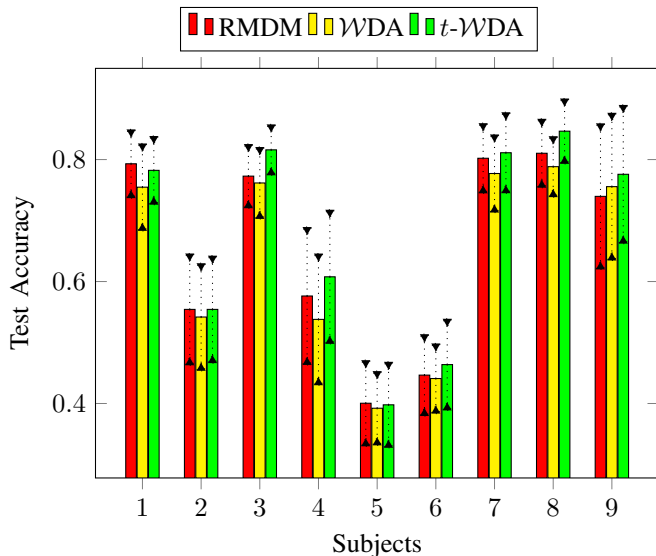


Fig. 2. Comparison of classifiers on the MI dataset. Standard deviations are plotted in dotted lines. The degree of freedom for  $t$ -WDA is  $\nu = 10$ .

for the five remaining subjects. Averaging over all the subjects,  $t$ -WDA has a 1.78% gain over RMDM and a 3.39% gain over WDA. Again,  $t$ -WDA appears to be the most powerful classifier on this MI dataset compared to RMDM and WDA.

## VI. CONCLUSION AND PERSPECTIVES

In this paper, we propose an original classification method that exploits the  $t$ -Wishart distribution over the set of matrices  $S_p^{++}$ . The classifier is obtained by leveraging the Discriminant Analysis framework. On two standard EEG BCI datasets, our proposed method outperforms best state-of-the-art MDM methods.

This work also yields several perspectives. As of now, the degree of freedom of the  $t$ -Wishart distribution is a hyperparameter. Instead, it should be estimated during the training step and different degrees of freedom for the different classes should be considered. Other Elliptical Wishart distributions could also be investigated.

## REFERENCES

- [1] O. Tuzel, F. Porikli, and P. Meer, "Pedestrian detection via classification on Riemannian manifolds," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 30, no. 10, pp. 1713–1727, 2008.
- [2] F. Lotte, L. Bougrain, A. Cichocki, M. Clerc, M. Congedo, A. Rakotomamonjy, and F. Yger, "A review of classification algorithms for eeg-based brain–computer interfaces: a 10 year update," *Journal of Neural Engineering*, vol. 15, no. 3, 2018.
- [3] S. Chevallier, E. Kalunga, Q. Barthélemy, and E. Monacelli, "Review of Riemannian distances and divergences, applied to SSVEP-based BCI," *Neuroinformatics*, vol. 19, 01 2021.
- [4] A. Barachant, S. Bonnet, M. Congedo, and C. Jutten, "Multiclass brain–computer interface classification by Riemannian geometry," *IEEE Transactions on Biomedical Engineering*, vol. 59, no. 4, pp. 920–928, 2012.
- [5] R. A. Maronna, R. D. Martín, V. J. Yohai, and M. Salibián-Barrera, *Robust statistics: theory and methods (with R)*. John Wiley & Sons, 2019.
- [6] A. K. Gupta and D. K. Nagar, *Matrix variate distributions*. Chapman and Hall/CRC, 2018.
- [7] G. Z. Thompson, R. Maitra, W. Q. Meeker, and A. F. Bastawros, "Classification with the matrix-variate  $t$  distribution," *Journal of Computational and Graphical Statistics*, vol. 29, no. 3, pp. 668–674, 2020.
- [8] K.-T. Fang and T. W. Anderson, *Statistical inference in elliptically contoured and related distributions*. Allerton Press, 1990.
- [9] F. J. Caro-Lopera, G. González-Farías, and N. Balakrishnan, "On generalized Wishart distributions-I: Likelihood ratio test for homogeneity of covariance matrices," *Sankhya A*, vol. 76, no. 2, pp. 179–194, 2014.
- [10] J. Wishart, "The generalized product moment distribution in samples from a normal multivariate population," *Biometrika*, vol. 20A, no. 1-2, pp. 32–52, 12 1928.
- [11] E. Ollila, D. E. Tyler, V. Koivunen, and H. V. Poor, "Complex elliptically symmetric distributions: Survey, new results and applications," *IEEE Transactions on signal processing*, vol. 60, no. 11, pp. 5597–5625, 2012.
- [12] L. T. Skovgaard, "A Riemannian geometry of the multivariate normal model," *Scandinavian journal of statistics*, pp. 211–223, 1984.
- [13] S. Kullback, "Information theory and statistics," *New York: Dover*, 1968.
- [14] Z. Chebbi and M. Moakher, "Means of hermitian positive-definite matrices based on the log-determinant  $\alpha$ -divergence function," *Linear Algebra and its Applications*, vol. 436, no. 7, pp. 1872–1889, 2012.
- [15] I. Ayadi, F. Bouchard, and F. Pascal, "Elliptical Wishart distribution: maximum likelihood estimator from information geometry," in *International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2023.
- [16] P.-A. Absil, R. Mahony, and R. Sepulchre, *Optimization Algorithms on Matrix Manifolds*. Princeton, NJ, USA: Princeton University Press, 2008.
- [17] B. Jeuris, R. Vandebril, and B. Vandereycken, "A survey and comparison of contemporary algorithms for computing the matrix geometric mean," *Electronic Transactions on Numerical Analysis*, vol. 39, pp. 379–402, 2012.
- [18] B. Ghoghj and M. Crowley, "Linear and quadratic discriminant analysis: Tutorial," 2019. [Online]. Available: <https://arxiv.org/abs/1906.02590>
- [19] S. Said, L. Bombrun, Y. Berthoumiou, and J. H. Manton, "Riemannian gaussian distributions on the space of symmetric positive definite matrices," *IEEE Transactions on Information Theory*, vol. 63, no. 4, pp. 2153–2170, 2017.
- [20] M. Tangermann, K.-R. Müller, A. Aertsen, N. Birbaumer, C. Brunner, C. Brunner, R. Leeb, C. Mehring, K. J. Miller, G. Mueller-Putz *et al.*, "Review of the BCI competition IV," *Frontiers in neuroscience*, p. 55, 2012.