

# Supplementary materials for paper: Revealing Common Statistical Behaviors in Heterogeneous Populations

Andrey Zhitnikov   Rotem Mulayoff   Tomer Michaeli

This supplementary document contains:

1. Proofs for Lemmas 1 and 2;
2. An alternative common-covariance estimation algorithm, which can be more efficient on parallel platforms, and an empirical comparison to Alg. 1;
3. Resting-state fMRI experiments on a group of ADHD subjects, and a comparison to graph-Lasso covariance estimation;
4. Analysis of the effect of kernel bandwidth on the pdf estimation in the ABP-PPG experiment.

## 1 Proof of Lemma 1

For convenience we write again the definitions

$$\psi(\mathbf{q}) \triangleq \mathbf{q}^T \boldsymbol{\Sigma}_u \mathbf{q}, \quad (1)$$

$$g_j(\mathbf{q}) \triangleq \mathbf{q}^T \boldsymbol{\Sigma}_{v_j} \mathbf{q}, \quad (2)$$

$$h_m(\mathbf{q}) \triangleq \min_{j \in \{1, \dots, m\}} g_j(\mathbf{q}), \quad (3)$$

and the assumptions of Theorem 1, which are

$$\mathbb{P}(\lambda_{\max}(\boldsymbol{\Sigma}_v) \leq \alpha) = 1 \quad (4)$$

for some  $\alpha > 0$  and

$$\mathbb{P}(\mathbf{q}^T \boldsymbol{\Sigma}_v \mathbf{q} \leq \epsilon) > 0 \quad (5)$$

for every  $\epsilon > 0$  and every unit-norm  $\mathbf{q}$ . Since  $\{g_j(\mathbf{q})\}$  and  $\{h_m(\mathbf{q})\}$  are random functions, we will sometimes explicitly write  $g_j(\omega, \mathbf{q})$  and  $h_m(\omega, \mathbf{q})$ , where  $\omega$  denotes an element from the sample space  $\Omega$ . We begin by demonstrating the first part of the lemma, namely that

$$h_m(\omega, \mathbf{q}) \xrightarrow{\text{a.s.}} 0 \quad (6)$$

for any  $\mathbf{q}$ . To this end, fix  $\mathbf{q}$  and  $\omega \in \Omega$ . By definition, the sequence  $\{h_m(\omega, \mathbf{q})\}_{m=1}^{\infty}$  is monotonically non-increasing and bounded from below, namely  $h_m(\omega, \mathbf{q}) \geq h_{m+1}(\omega, \mathbf{q}) \geq 0$  for every  $m$ . Therefore, this sequence converges for every  $\omega$ . Almost sure convergence implies convergence in probability to the same limit. Therefore, to prove (6), we will show that  $h_m(\omega, \mathbf{q}) \xrightarrow{\mathbb{P}} 0$ , namely that

$$\lim_{m \rightarrow \infty} \mathbb{P}(|h_m(\mathbf{q})| > \epsilon) \rightarrow 0 \quad (7)$$

for every  $\epsilon > 0$ . Note that  $g_j(\mathbf{q}) \geq 0$  for every  $j$ , so that  $h_m(\mathbf{q}) \geq 0$  for every  $m$ , and thus the absolute value in (7) can be omitted. Since the random variables  $\{g_j(\mathbf{q})\}$  are independent and identically distributed, we have that

$$\mathbb{P}(h_m(\mathbf{q}) > \epsilon) = \mathbb{P}(\bigcap_{j=1}^m \{g_j(\mathbf{q}) > \epsilon\}) = \mathbb{P}(\{g_1(\mathbf{q}) > \epsilon\})^m = (1 - \mathbb{P}(g_1(\mathbf{q}) \leq \epsilon))^m \xrightarrow{m \rightarrow \infty} 0, \quad (8)$$

where we used (5). This completes the proof of (6). Next, we prove the second part of the Lemma, namely that

$$h_m(\omega, \mathbf{q}_m) \xrightarrow{\text{a.s.}} 0 \quad (9)$$

for any converging sequence  $\{q_m\}$ . To this end, we will first prove that all the functions  $\{h_m(\mathbf{q})\}$  are Lipschitz w.p. 1. Note from (4) and the definition of  $g_j(\mathbf{q})$  in (2) that each  $g_j(\mathbf{q})$  is Lipschitz w.p. 1. Therefore, all the functions  $\{g_j(\mathbf{q})\}_{j=1}^m$  are simultaneously Lipschitz w.p. 1, with the same Lipschitz constant  $2\alpha$  (a countable union of a.s. events). Fix any  $\omega \in \Omega$  within this a.s. event. The function  $h_1(\omega, \mathbf{y})$  is Lipschitz since it equals  $g_1(\omega, \mathbf{y})$ . Let us show this also holds for  $h_2(\omega, \mathbf{y})$ . For any two vectors  $\mathbf{x}, \mathbf{y}$ , we have that

$$g_1(\omega, \mathbf{y}) \geq g_1(\omega, \mathbf{x}) - |g_1(\omega, \mathbf{x}) - g_1(\omega, \mathbf{y})| \geq h_2(\omega, \mathbf{x}) - 2\alpha\|\mathbf{x} - \mathbf{y}\|, \quad (10)$$

where we used (3). Similarly, we have that  $g_2(\omega, \mathbf{y}) \geq h_2(\omega, \mathbf{x}) - 2\alpha\|\mathbf{x} - \mathbf{y}\|$ . Therefore,

$$h_2(\omega, \mathbf{y}) = \min\{g_1(\omega, \mathbf{y}), g_2(\omega, \mathbf{y})\} \geq h_2(\omega, \mathbf{x}) - 2\alpha\|\mathbf{x} - \mathbf{y}\|. \quad (11)$$

This implies that  $h_2(\omega, \mathbf{x}) - h_2(\omega, \mathbf{y}) \leq 2\alpha\|\mathbf{x} - \mathbf{y}\|$ , and by switching the roles of  $\mathbf{x}$  and  $\mathbf{y}$ , we also obtain

$$|h_2(\omega, \mathbf{x}) - h_2(\omega, \mathbf{y})| \leq 2\alpha\|\mathbf{x} - \mathbf{y}\|, \quad (12)$$

demonstrating that  $h_2(\omega, \mathbf{q})$  is Lipschitz with the same constant  $2\alpha$ . In a similar way, we obtain that  $h_m(\omega, \mathbf{q})$  is Lipschitz for every  $m$ . We are now ready to prove (9). We saw that all  $\{h_m(\omega, \mathbf{q})\}$  are Lipschitz a.s. and that  $|h_m(\omega, \mathbf{q})| \rightarrow 0$  a.s. for any  $\mathbf{q}$ . The intersection of these two events is also a.s. Fix any  $\omega$  in this intersection. Then

$$|h_m(\omega, \mathbf{q}_m)| = |h_m(\omega, \mathbf{q}_m) - h_m(\omega, \mathbf{q}^*) + h_m(\omega, \mathbf{q}^*)| \leq 2\alpha\|\mathbf{q}_m - \mathbf{q}^*\| + |h_m(\omega, \mathbf{q}^*)| \xrightarrow{m \rightarrow \infty} 0. \quad (13)$$

This completes the proof.

## 2 Proof of Lemma 2

The functions  $\{\phi(\mathbf{q}) + f_n(\mathbf{q})\}$  are continuous and  $\mathcal{C}$  is a compact set, therefore  $\arg \min_{\mathbf{q} \in \mathcal{C}} \{\phi(\mathbf{q}) + f_n(\mathbf{q})\} \neq \emptyset$ . We have to prove that for every  $\epsilon > 0$  there exists an  $N_\epsilon$  such that

$$\|\mathbf{q}_n - \mathbf{q}^*\| < \epsilon \quad \forall n \geq N_\epsilon. \quad (14)$$

Fix  $\epsilon > 0$  and define the set  $\mathcal{Q} = \{\mathbf{q} \in \mathcal{C} : \|\mathbf{q} - \mathbf{q}^*\| \geq \epsilon\}$ . If the set  $\mathcal{Q}$  is empty, then (14) is obviously satisfied with  $N_\epsilon = 1$  and the proof is complete. If  $\mathcal{Q} \neq \emptyset$ , then  $\phi(\mathbf{q})$  attains a minimum over the set  $\mathcal{Q}$ . This follows from the fact that  $\mathcal{Q}$  is compact and  $\phi(\mathbf{q})$  is continuous over the set  $\mathcal{Q} \subset \mathcal{C}$ . Let

$$\tilde{\mathbf{q}} \in \arg \min_{\mathbf{q} \in \mathcal{Q}} \phi(\mathbf{q}) \quad (15)$$

and denote  $\delta(\epsilon) = \phi(\tilde{\mathbf{q}}) - \phi(\mathbf{q}^*)$ . From

$$\lim_{n \rightarrow \infty} f_n(\mathbf{q}^*) = 0, \quad (16)$$

we know that there exists an  $N_\epsilon$  such that  $0 \leq f_n(\mathbf{q}^*) < \delta(\epsilon)$  for every  $n \geq N_\epsilon$ . Therefore, for every  $n \geq N_\epsilon$ , we have that

$$\phi(\mathbf{q}^*) + f_n(\mathbf{q}^*) < \phi(\mathbf{q}^*) + \delta(\epsilon) = \phi(\tilde{\mathbf{q}}) \leq \underbrace{\min_{\mathbf{q} \in \mathcal{Q}}}_{f_n \geq 0} [\phi(\mathbf{q}) + f_n(\mathbf{q})]. \quad (17)$$

This demonstrates that for every  $n > N_\epsilon$  there exists no point at distance larger than  $\epsilon$  from  $\mathbf{q}^*$ , at which the function  $\phi(\mathbf{q}) + f_n(\mathbf{q})$  attains a lower value than at  $\mathbf{q}^*$ . In other words, the minimum of  $\phi(\mathbf{q}) + f_n(\mathbf{q})$  must be attained at a point  $\mathbf{q}_n$ , which satisfies (14). This completes the proof.

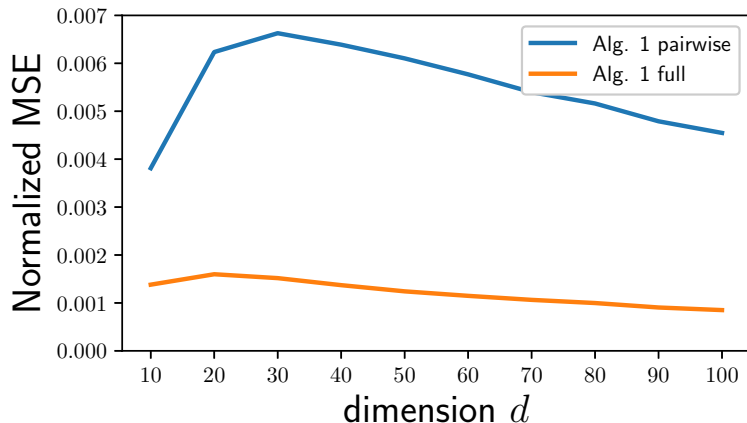


Figure 1: **Common covariance estimation with the full and pairwise versions of Alg. 1 at high SNR.** We used 100 noisy covariance matrices to estimate the underlying common covariance, using Alg. 1 in the full and pairwise form. In this simulation,  $\text{SNR} = 1$ . Approximation of the MSE is based on 100 trials.

### 3 Accelerated covariance estimation

To speed up the estimation of the common covariance on parallel platforms, one can break the task of estimating the full  $d \times d$  covariance matrix  $\Sigma_{\mathbf{u}}$  into  $d \times (d - 1)/2$  sub-tasks of estimating all  $2 \times 2$  sub-covariance-matrices of  $\Sigma_{\mathbf{u}}$ . Those sub-tasks can be solved in parallel, while averaging the  $d - 1$  different estimates obtained for the diagonal entries. This procedure is not guaranteed to yield a positive semidefinite estimate  $\hat{\Sigma}_{\mathbf{u}}$ , so that negative eigenvalues have to be truncated. However, this method enjoys the same asymptotic guarantees as the direct estimation of the full covariance.

To study the effect of this approach on the estimation accuracy for a finite number of subjects, we next compare between Alg. 1 and this pairwise version in a simulation. We take the common covariance to be the identity matrix  $\Sigma_{\mathbf{u}} = \mathbf{I}$  and generate the noise covariance according to equation (27) of the main text, where now  $\Lambda_j = \text{diag}\{\beta_1^j \dots \beta_d^j\}$  with

$$\beta_k^j \sim U[0, b] \quad \forall j = 1, \dots, m \quad \forall k = 1, \dots, d. \quad (18)$$

For the rotation matrices, we first draw a matrix  $\tilde{\mathbf{M}}_j$  with iid entries uniformly distributed on the interval  $[-50, 50]$ . Then, the matrix  $\mathbf{M}_j$  is taken to be the  $\mathbf{Q}$  matrix from the  $\mathbf{QR}$  decomposition of  $\tilde{\mathbf{M}}_j$ . Figures 1 and 2 depict the MSE attained by both algorithms (normalized by the number of entries in the estimated matrix,  $d^2$ ), as a function of the dimension  $d$ , for two different SNR levels. As can be seen, in all the tested settings, Alg. 1 in its full form gives more accurate results. We conclude that the potential improvement in running time offered by this approach, typically comes at the cost of reduced estimation accuracy.

### 4 Comparison of common covariance matrices of Control subjects and subjects with ADHD

We repeated the analysis of Sec. 5.3 in the main text on a group of 141 subjects diagnosed with ADHD. In addition to our estimator and the Riemann mean and Euclidean mean estimators, we also show results obtained with the graph Lasso estimator (Friedman et al., 2008) applied to all data points from all the subjects. This estimator assumes a sparse precision matrix and selects its  $L_1$  regularization weight through cross-validation. As can be seen in Figs. 3 and 4 below, our estimator consistently detects higher correlations within known networks (around the main diagonal).

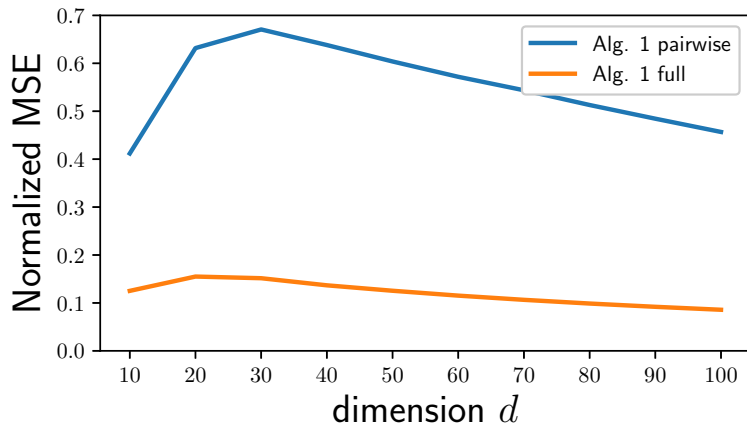


Figure 2: **Common covariance estimation with the full and pairwise versions of Alg. 1 at low SNR.** We used 100 noisy covariance matrices to estimate the underlying common covariance, using Alg. 1 in the full and pairwise form. In this simulation,  $\text{SNR} = 0.1$ . Approximation of the MSE is based on 100 trials.

## 5 The effect of the kernel bandwidth in naive KDE and in our common density estimator

In Sec. 5.4 of the main text, we demonstrated the ability of our common pdf estimator to reveal delicate structures that are not seen with naive KDE applied on all data points from all patients. To show that this phenomenon is not a matter of appropriate selection of bandwidth in the KDE, we now repeat this experiment with several bandwidths, ranging from ones leading to over-smoothness to ones leading to under-smoothness. As can be seen in Fig. 5 below, there exists no bandwidth with which the vertical trace is seen in the KDE estimate. At the same time, this trace is seen with all bandwidths in our pdf estimate.

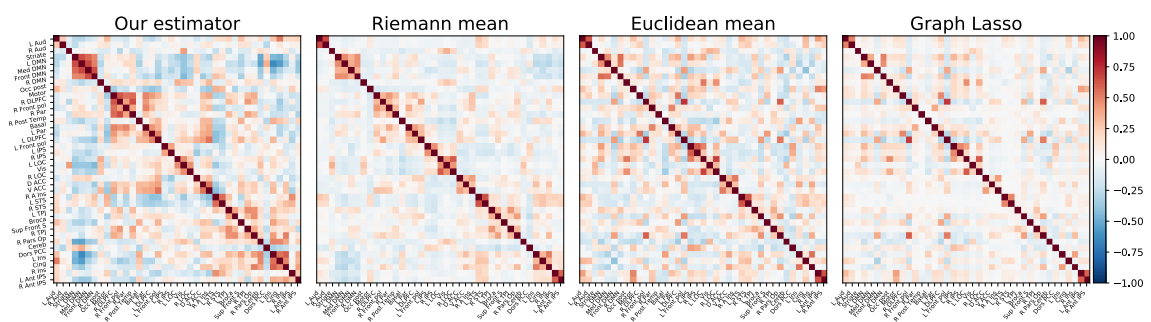


Figure 3: Comparison of group level correlation matrices of control subjects (same experiment as in the main text).

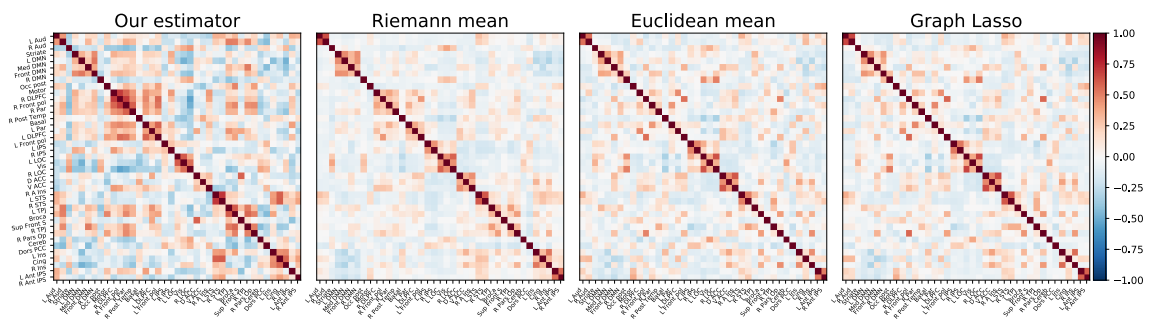


Figure 4: Comparison of group level correlation matrices of ADHD subjects.

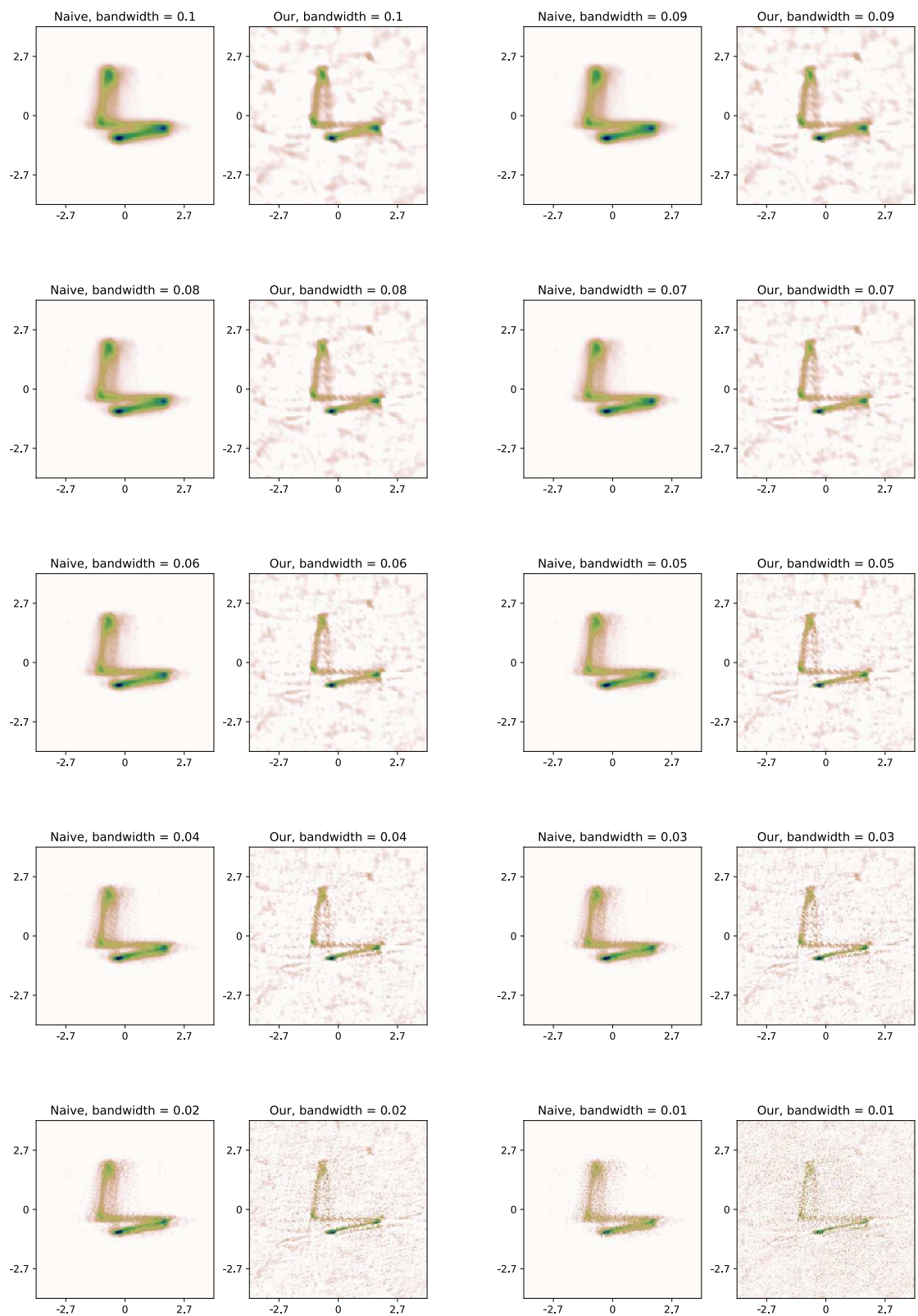


Figure 5: Common probability density estimation of PPG and ABP with a bandwidths ranging from 0.1 (top left) to 0.01 (bottom right).