

## Motivation

The main goal is to develop a robust & nonparametric outlier detection method to identify artifact-contaminated fMRI scans.

- fMRI data contain **artifacts** from hardware and physiological sources, which **reduce** the signal-to-noise ratio (**SNR**) and violate common statistical assumptions.
- “**Scrubbing**” or removal of the artifactual volumes in fMRI scans is needed to improve the quality of the statistical results.
- We develop a **robust outlier detection approach** that considers the **high dimensional, auto-correlated, non-Gaussian** aspects of fMRI data.

## Background

Distance-based outlier detection approaches:

- Existing methods using distances to detect outliers: Mahalanobis Distance (MD) and Robust Distance (RD).
- Since MD suffers from “masking” and “swamping” effects, RD became a popular method to identify outliers.

**MCD-based robust distance [1]:**

- This RD approach uses the minimum covariance determinant (MCD) estimates of mean and covariance which are calculated by dividing data into two subsets: included set and excluded set.
- Included set contains observations located very close to the center of the data which are unlikely to contain outliers so that mean and covariance estimates are obtained from this subset robustly.
- MCD-based RD of observation is calculated analogously to MD but using the MCD estimates of  $\mu$  &  $\Sigma$ .

$$MD(x_i) = \sqrt{(x_i - \hat{\mu})^T \hat{\Sigma}^{-1} (x_i - \hat{\mu})}$$

$$RD(x_i) = \sqrt{(x_i - \hat{\mu}_{MCD})^T \hat{\Sigma}_{MCD}^{-1} (x_i - \hat{\mu}_{MCD})}$$

for  $MVNormal(\mu, \Sigma)$  with  $N=100$  &  $p=2$

for 4 ICs selected fMRI (see Step0 in Methodology)

$\mu = (0,0)^T$  and  $\Sigma = \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1 \end{bmatrix}$

$\hat{\mu}_{MCD} = (0.03, 0.42)^T$  and  $\hat{\Sigma}_{MCD} = \begin{bmatrix} 1.01 & 0.82 \\ 0.82 & 0.75 \end{bmatrix}$

**Distribution of MCD distances:**

- Hardin & Rocke (2005) derived an F-distribution of MCD distances of excluded observations obtained from independently, identically, and Normally distributed data (H&R's result) [2].
- Unfortunately, the empirical distribution of MCD distance has often been observed to deviate from the theoretical distribution. Previous work employing MCD distance attempted to improve the distributional fit of the theoretical distribution, for example using median matching [3].

**How does violation of the independence assumption affect H&R's results?**

$$i.i.d \text{ Normal} \rightarrow FPR = 1.3\%$$

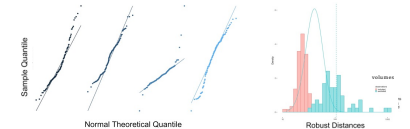
$$AR(1) \text{ } \phi = 0.3 \rightarrow FPR = 1.3\%$$

$$AR(1) \text{ } \phi = 0.6 \rightarrow FPR = 1.5\%$$

The independence violation has a minor effect on the validity of the theoretical F distribution.

**How does violation of the normality assumption affect H&R's result in fMRI data?**

- The Autism Brain Imaging Data Exchange (ABIDE) contains subjects with autism spectrum disorder (ASD) and typically developing subjects (ABIDE [4]) which includes 1112 subjects from 17 international sites (Repetition Time ranges 2-3 sec).
- We select a single slice of a single session from a typically developing subject that is relatively free of artifacts for visualization.
- The QQ-plots of randomly selected components clearly show the deviation from the Normality.



The scaled F-distribution failed to fit the distribution of the RD's excluded observations which shows the effect of violation of the Normality assumption.

**How can the identity assumption hold?**

Mean and variance detrending can be applied to meet the identity assumption of H&R's result.

**fMRI data violate the independence and Normality assumptions of H&R's result preventing us from using their theoretical distribution to detect outliers. We develop a nonparametric robust outlier detection method appropriate for fMRI data.**

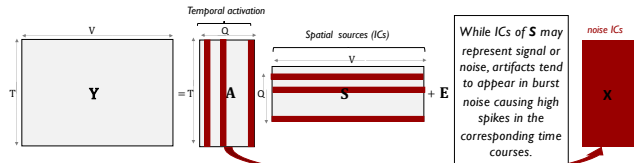
## Methodology

The goal is to estimate the upper quantile of MCD distances in outlier-free data (threshold for outliers) using a non-parametric bootstrap approach.

**Step 0: Dimension reduction & selection**

**Goal:** To identify a low dimensional subspace primarily related to artifactual signals.

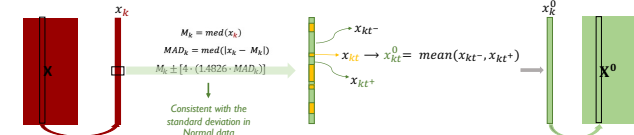
We adopt the “ICA projection scrubbing” method of [4] using kurtosis to identify independent components (ICs) representing burst noise.



**Step 1: Univariate outlier imputation**

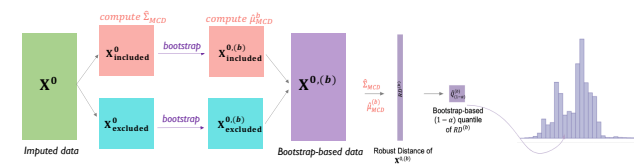
**Goal:** To eliminate the effect of outliers on the upper quantile of the distribution of MCD distances.

We use a robust empirical rule & impute based on neighboring time points.



**Step 2: Bootstrap Procedure**

**Goal:** To estimate the  $(1-\alpha)$  quantile of the distribution of MCD distances in outlier-free data.



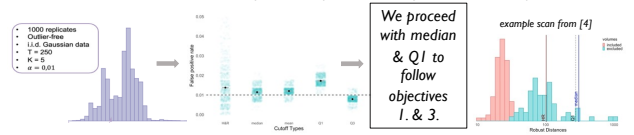
**Step 3: Multivariate Outlier Identification**

**Goal:** To determine a threshold value for outlier detection from the bootstrap-based quantiles

$\hat{Q}_{1-\alpha}^{(b)}, b = 1, \dots, B$ .

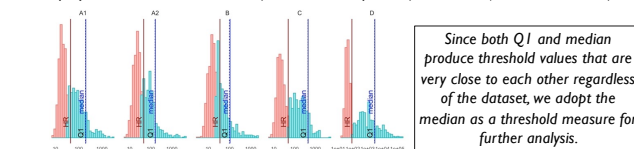
Possible objectives when summarizing across bootstrap samples

- Average false positive rate (FPR) of  $\alpha \rightarrow$  use mean or median
- FPR of at most  $\alpha$  with some statistical certainty  $\rightarrow$  use upper quantile
- FPR of at least  $\alpha$  with some statistical certainty to maintain power for detecting outlier  $\rightarrow$  use lower quantile



## Application on real datasets

We employ Human Connectome Project with 4 unique subjects' scans (A1, A2, B, C, and D).



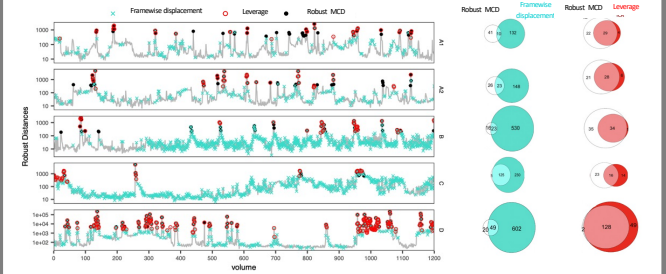
## References

- A preprint version of this poster: <https://github.com/fatma1990/RobustOutlierDetection/blob/main/RobustOutlierDetection.pdf>
- Rousseeuw, Peter J., and Bert C Van Zomeren. 1990. "Unmasking Multivariate Outliers and Leverage Points." *Journal of the American Statistical Association* 85 (411): 633-39.
- Hardin, Johanna, and David M Rocke. 2005. "The Distribution of Robust Distances." *Journal of Computational and Graphical Statistics* 14 (4): 928-46.
- Mejia, A.F., Nabel, M.B., Boyan, A., Caffo, B., and Lindquist, M.A. 2017. PCA leverage: outlier detection for high-dimensional functional magnetic resonance imaging data. *Biostatistics*, 18(3), pp.521-536.
- Di Martino, A., Yan, C.G., Li, Q., Denio, E., Castellanos, F.X., Alberts, K., Anderson, J.S., Assaf, M., Bookheimer, S.Y., DiPrete, M., and Dean, B. 2014. The autism brain imaging data exchange towards a large-scale evaluation of the intrinsic brain architecture in autism. *Molecular psychiatry*, 19(6), pp.659-667.
- Pham, D., McDonald, D., Ding, L., Nabel, M.B., and Mejia, A. 2021. Less is more: balancing noise reduction and data retention in fMRI with projection scrubbing. *arXiv preprint arXiv:2108.00319v2*.
- Van Essen, D.C., Smith, S.M., Barch, D.M., Behrens, T.E., Yacoub, E., Ugurbil, K., and VU-Min HCP Consortium. 2013. The VU-Min human connectome project: an overview. *Neuroimage*, 80, pp.62-79.
- Power, J.D., Mitra, A., Laumann, T.O., Snyder, A.Z., Schlaggar, B.L., and Petersen, S.E. 2014. Methods to detect, characterize, and remove motion artifact in resting state fMRI. *Neuroimage*, 84, pp.320-341.

## Results: Comparison with other scrubbing methods

We compare our method with the following alternative scrubbing approaches.

- A statistical approach: ICA projection scrubbing with leverage [5]**
  - Uses the same dimension reduction & selection approach but uses leverage as a non-robust distance measure, which can lead to masking and/or swamping effects.
- A hardware approach: Motion scrubbing based on framewise displacement (FD) [7]**
  - No principled or universally accepted threshold value
  - Low sensitivity (not data-driven, can only identify motion-based artifacts)
  - Low specificity (tends to be over-aggressive due to the use of the low FD threshold in practice)



**Conclusions from the comparison with Leverage**

- Leverage fails to detect some volumes with high RD in A1, A2, B, and C.
- Artifactual volumes are almost doubled by using Robust MCD instead of Leverage in A1, A2, B, and C.
- However, our MCD distance approach seems to be inefficient to detect artifacts in D which may be due to strong deviations from Normality in the univariate outlier imputation step (see Future Work, 2)

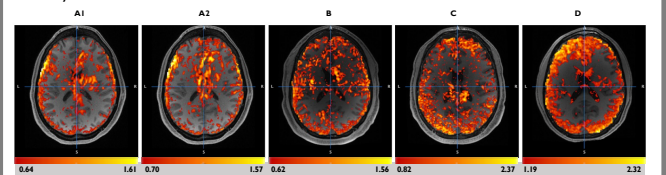
**Conclusions from the comparison with Framewise Displacement (FD)**

- Many volumes flagged artifacts with motion scrubbing which do not exhibit large MCD distances.
- Some volumes with low motion have large MCD distance which is possibly due to lagged effects of motion or other sources of artifacts.

## Results: Spatial patterns of artifacts

Unlike existing scrubbing methods, ICA projection allows us to localize the source of burst noise artifacts. We display the spatial patterns of artifactual volumes identified by our method.

- Recall that  $X$  is obtained by selecting a subset of the detrended columns of  $A$ . Let  $A^*$  and  $S^*$  contain only the selected rows and columns of  $A$  and  $S$ .
- For each outlying volume  $t$ , an image of artifact intensity can be obtained by multiplying the  $t$  row of  $A^*$  by  $S^*$ .
- Since each of these volumes cannot be displayed due to space constraints, we only visualize the average across all outlying volumes, giving an overall measure of artifact intensity at every voxel of the brain.



- The images show a ring of intensity on the outer edge of the brain, which is often caused by head movements. It is possible that these head movements cause mis-localizing signals, which cause artifacts in the brain volume.
- We see similar spatial patterns for subject A across two different scans.

## Future work

- Developing a multiple imputation approach to account for uncertainty
- Develop a Monte Carlo-based approach using transformations to Normality as an alternative to bootstrap
- Validating on simulated fMRI data and test-retest fMRI data