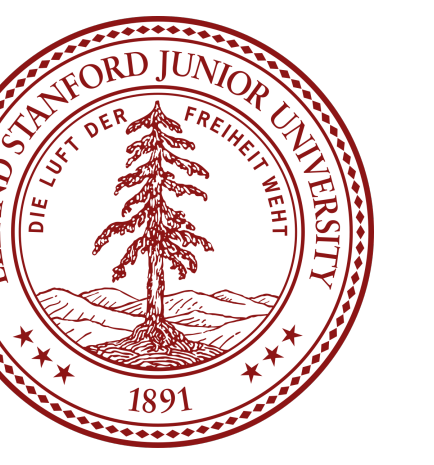




CANONICAL CORRELATION ANALYSIS IN HIGH DIMENSIONS WITH STRUCTURED REGULARIZATION

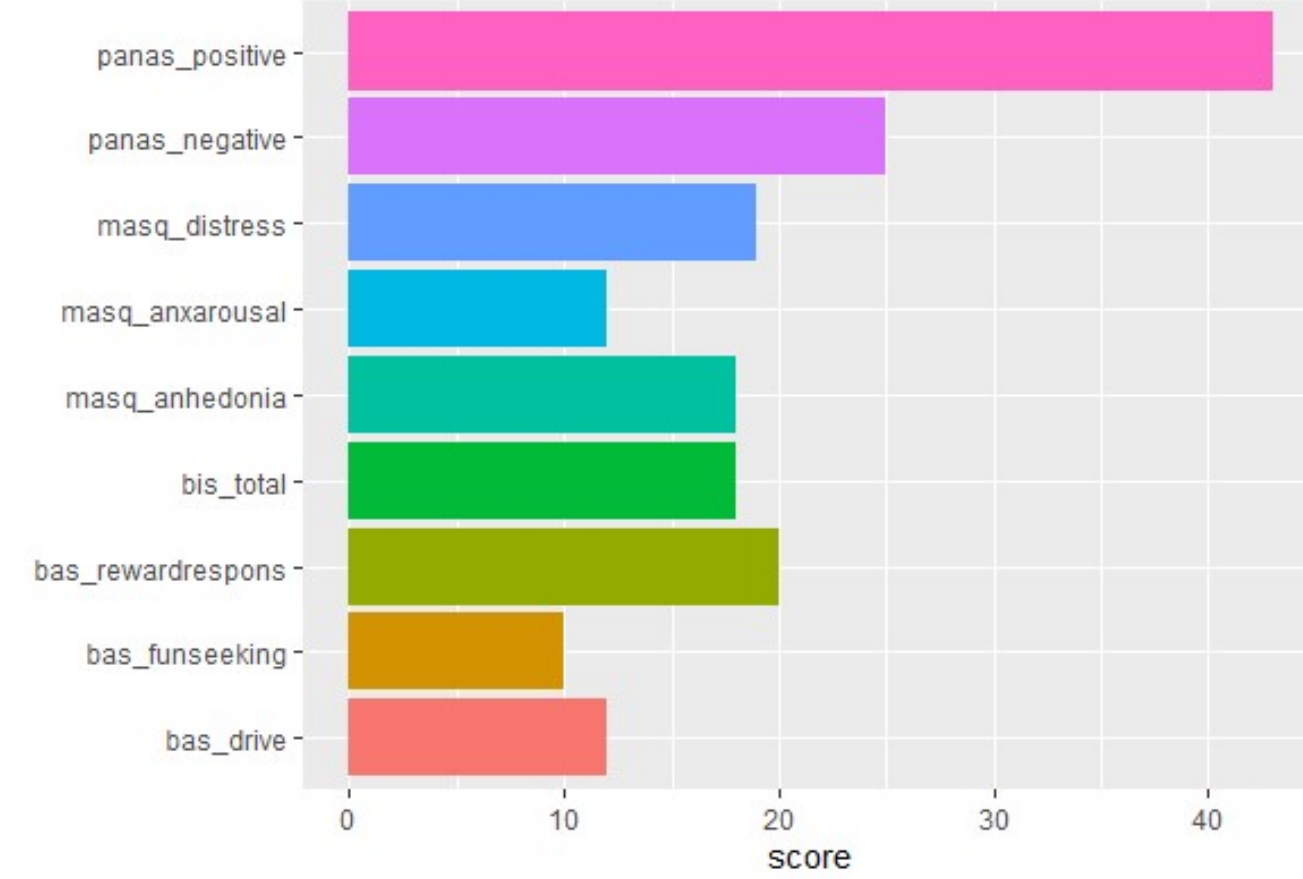
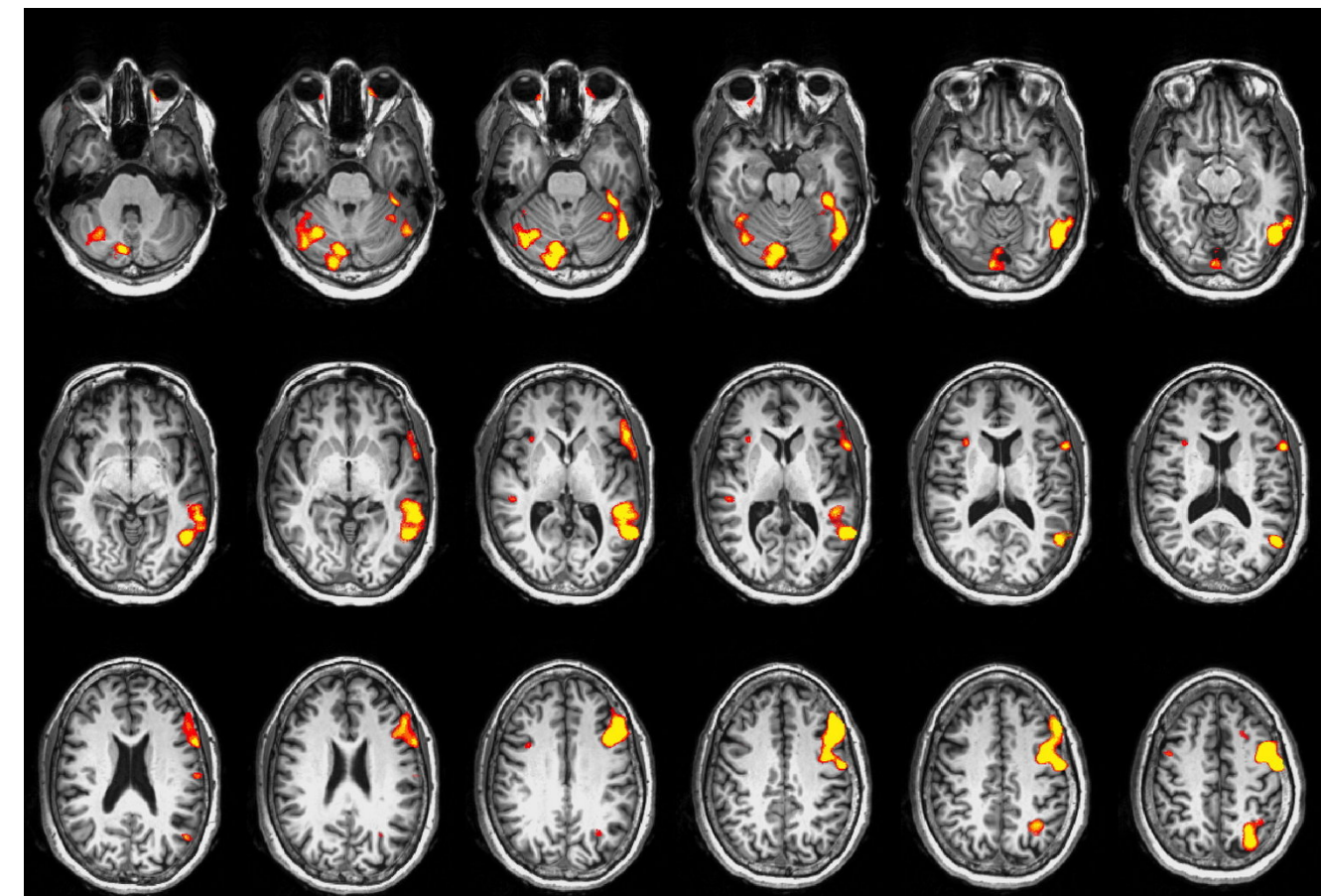


[ELENA TUZHILINA] STANFORD UNIVERSITY, DEPARTMENT OF STATISTICS
JOINT WORK WITH L. TOZZI AND T. HASTIE

MOTIVATION

Canonical correlation analysis (CCA) is a technique for measuring the association between two multivariate data matrices. A regularized modification of canonical correlation analysis (RCCA), imposing an ℓ_2 penalty on the CCA coefficients, is widely used in applications with high-dimensional data. One limitation of such regularization is that it ignores any data structure, which can be ill-suited for some applications. Here we introduce a novel approach that takes the underlying data structure into account. The proposed group regularized canonical correlation analysis (GRCCA), is especially useful when the variables are correlated in groups. We illustrate some computational strategies to avoid excessive computations with regularized CCA in high dimensions. We demonstrate the application of GRCCA method in our motivating application from neuroscience.

DATA



Brain activations $X \in \mathbb{R}^{n \times p}$: magnetic resonance imaging obtained during a gambling task.

Behavioral scores $Y \in \mathbb{R}^{n \times q}$: self-reports assessing various aspects of reward-related behaviors.

$n = 153$ participants; $p = 90,368$ greyordinates; $q = 9$ scores

KERNEL TRICK

Idea: find a linear transformation

$$V = \begin{matrix} p \times n \\ \boxed{} \end{matrix} \quad R = XV = \begin{matrix} n \times n \\ \boxed{} \end{matrix}$$

such that RCCA for (X, Y) is equivalent to RCCA for (R, Y)

Step-by-step procedure:

- $X = UDV^T = \begin{matrix} n \times n \\ \boxed{} \end{matrix} \begin{matrix} n \times n \\ \boxed{} \end{matrix} \begin{matrix} n \times p \\ \boxed{} \end{matrix}$
- set $R = XV = UD$ and solve RCCA problem for $(R, Y) \implies$ get α_R, β_R
- recover coefficients $\alpha_X = V\alpha_R$
- variables stay the same $R\alpha_R = X\alpha_X$

GROUP STRUCTURE

Motivation: brain features come in groups (aka brain regions). How to take into account the group structure?



GROUP RCCA

GRCCA optimization problem:

$$\begin{aligned} &\text{maximize } \alpha^T \Sigma_{XY} \beta \\ &\text{w.r.t. } \alpha \in \mathbb{R}^p, \beta \in \mathbb{R}^q \\ &\text{s.t. } \alpha^T \Sigma_{XX} \alpha = 1 \\ &\quad \beta^T \Sigma_{YY} \beta = 1 \\ &\quad \sum_{k=1}^K \|\alpha_k - \bar{\alpha}_k\|^2 \leq t_1 \\ &\quad \sum_{k=1}^K p_k \bar{\alpha}_k^2 \leq s_1 \end{aligned}$$

Modified correlation coefficient:

$$\frac{\alpha^T \Sigma_{XY} \beta}{\sqrt{\alpha^T (\Sigma_{XX} + K(\lambda_1, \mu_1)) \alpha} \sqrt{\beta^T \Sigma_{YY} \beta}}$$

where $K(\lambda_1, \mu_1) = \lambda_1(I - C) + \mu_1 C$

$$\text{and } C = \begin{bmatrix} \frac{11^T}{p_1} & 0 & \dots & 0 \\ 0 & \frac{11^T}{p_2} & \dots & 0 \\ \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & \frac{11^T}{p_K} \end{bmatrix}$$

Lemma

GRCCA for (X, Y) is equivalent to RCCA for (\tilde{X}, Y) where

$$\tilde{X} = (X_1 - \bar{X}_1, \dots, X_K - \bar{X}_K, \sqrt{\frac{p_1 \lambda_1}{\mu_1}} \bar{X}_1, \dots, \sqrt{\frac{p_K \lambda_1}{\mu_1}} \bar{X}_K)$$

CANONICAL CORRELATION ANALYSIS

Goal: given two random vectors $x = (x_1, \dots, x_p)$ and $y = (y_1, \dots, y_q)$

$$\text{maximize } \text{cor}(\alpha^T x, \beta^T y) \text{ w.r.t. } \alpha, \beta$$

- canonical coefficients α and β
- canonical variates $\alpha^T x$ and $\beta^T y$
- canonical correlation $\text{cor}(\alpha^T x, \beta^T y)$

Correlation coefficient:

$$\rho(\alpha, \beta) = \text{cor}(\alpha^T x, \beta^T y) \approx \frac{\alpha^T \Sigma_{XY} \beta}{\sqrt{\alpha^T \Sigma_{XX} \alpha} \sqrt{\beta^T \Sigma_{YY} \beta}}$$

CCA optimization problem:

$$\begin{aligned} &\text{maximize } \alpha^T \Sigma_{XY} \beta \\ &\text{w.r.t. } \alpha \in \mathbb{R}^p, \beta \in \mathbb{R}^q \\ &\text{s.t. } \alpha^T \Sigma_{XX} \alpha = 1 \\ &\quad \beta^T \Sigma_{YY} \beta = 1 \end{aligned}$$

Solution: via SVD of $\Sigma_{XX}^{-\frac{1}{2}} \Sigma_{XY} \Sigma_{YY}^{-\frac{1}{2}}$

REGULARIZATION

Motivation: CCA doesn't work for $p > n$

Modified correlation coefficient:

$$\frac{\alpha^T \Sigma_{XY} \beta}{\sqrt{\alpha^T (\Sigma_{XX} + \lambda_1 I) \alpha} \sqrt{\beta^T \Sigma_{YY} \beta}}$$

Shrinkage property:

$$\begin{aligned} &\text{maximize } \alpha^T \Sigma_{XY} \beta \\ &\text{w.r.t. } \alpha \in \mathbb{R}^p, \beta \in \mathbb{R}^q \\ &\text{s.t. } \alpha^T \Sigma_{XX} \alpha = 1 \\ &\quad \beta^T \Sigma_{YY} \beta = 1 \\ &\quad \|\alpha\| \leq t_1 \end{aligned}$$

Solution: via SVD of $(\Sigma_{XX} + \lambda_1 I)^{-\frac{1}{2}} \Sigma_{XY} \Sigma_{YY}^{-\frac{1}{2}}$

CCA LIBRARY

```
library(CCA)
rcc(X = activation, Y = behavior, lambda1 = 10, lambda2 = 0)
```

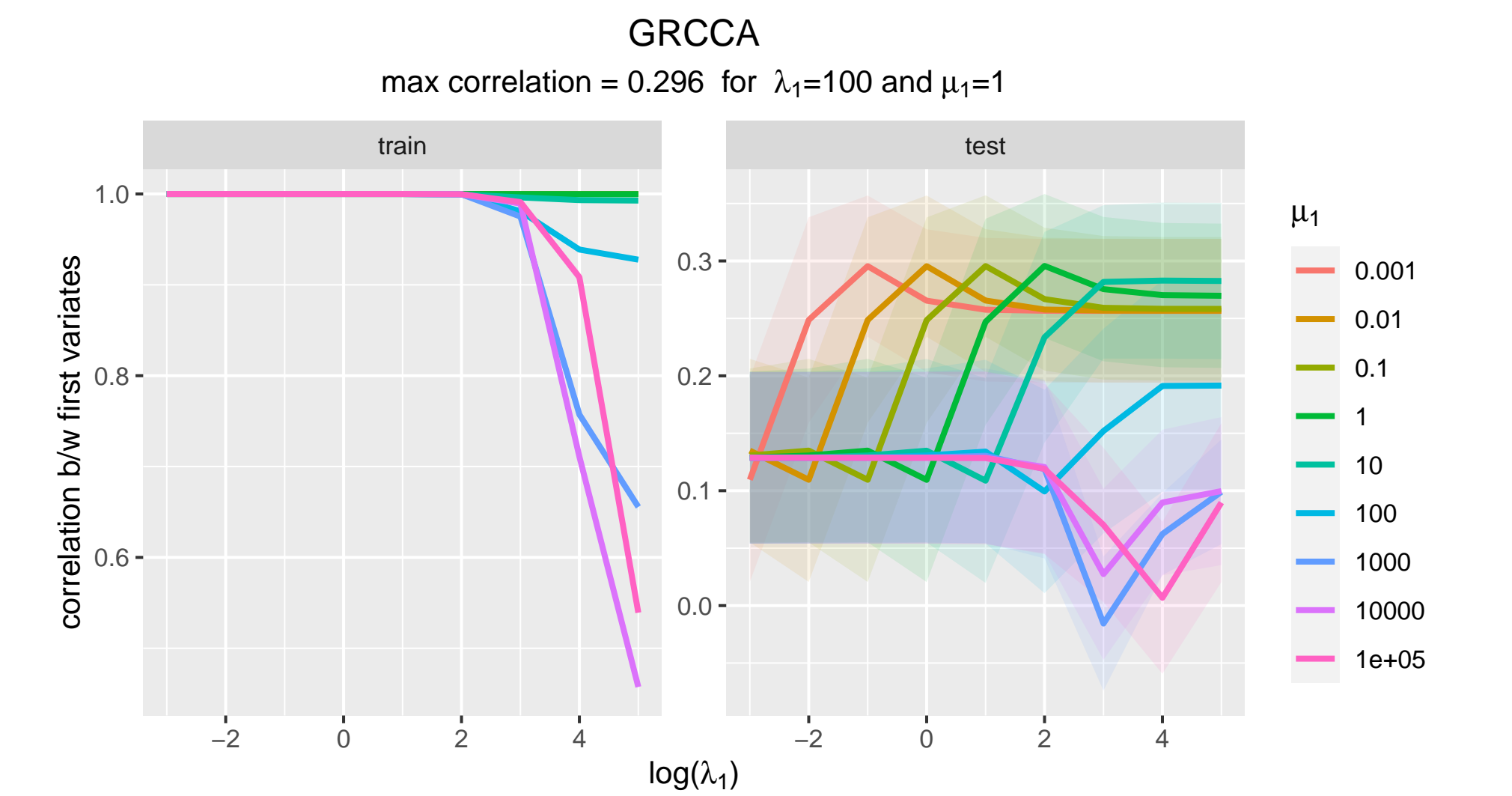
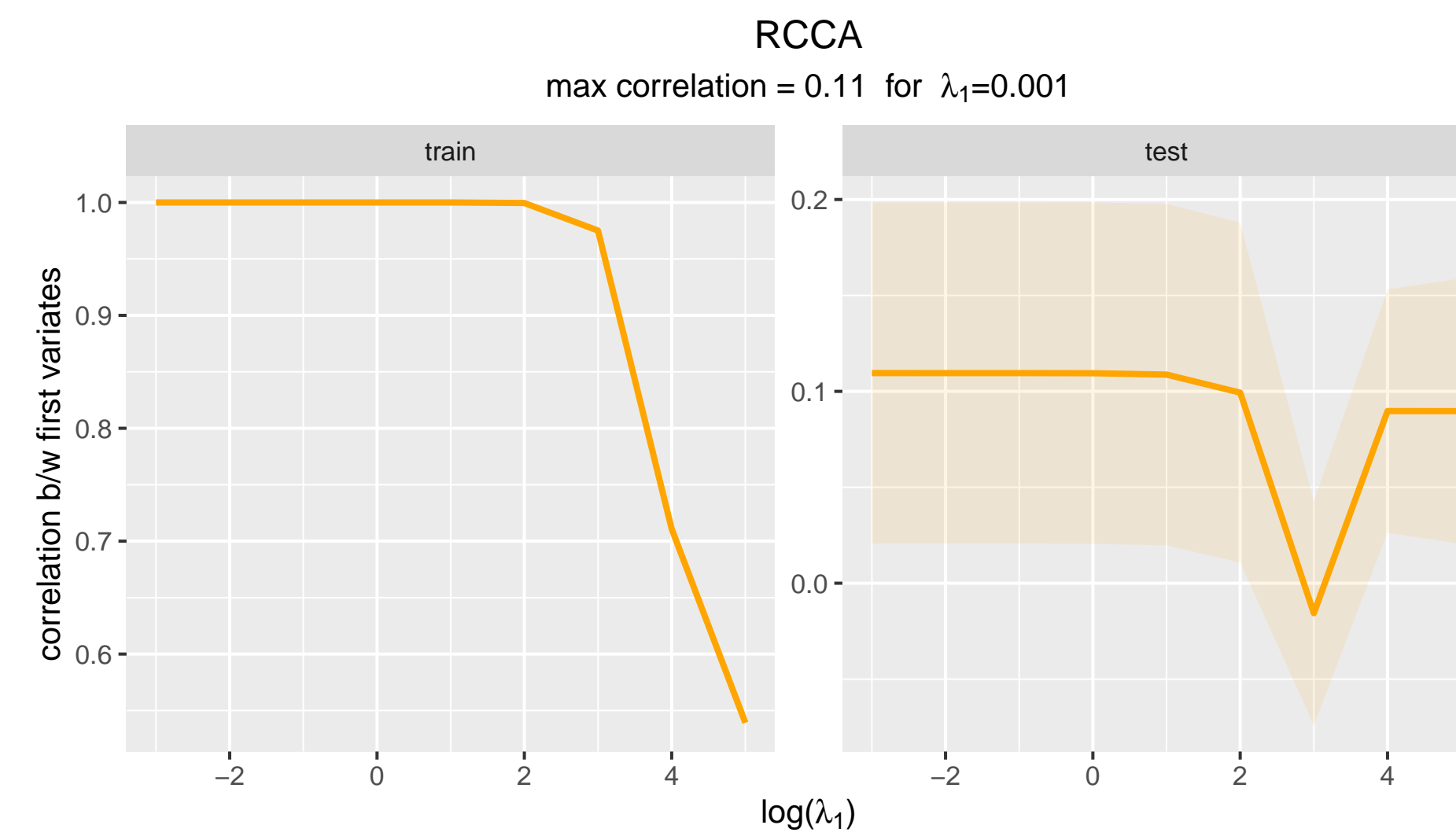
Error: cannot allocate vector of size 62.1 Gb
Traceback:

- rcc(X = activation, Y = behavior, lambda1 = 10, lambda2 = 0)
- var(X, na.rm = TRUE, use = "pairwise")

Problem:

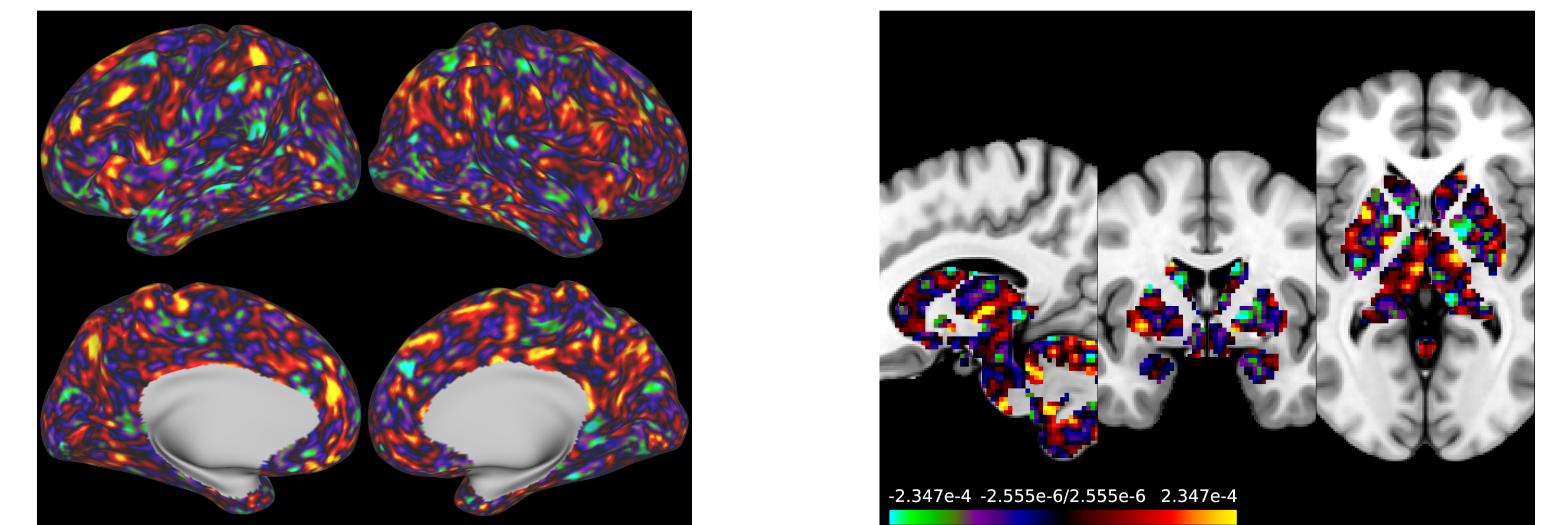
$\Sigma_{XX} \in \mathbb{R}^{p \times p}$
 $\Sigma_{XY} \in \mathbb{R}^{p \times q}$
are too large when $p \approx 90K$

CROSS VALIDATION RESULTS

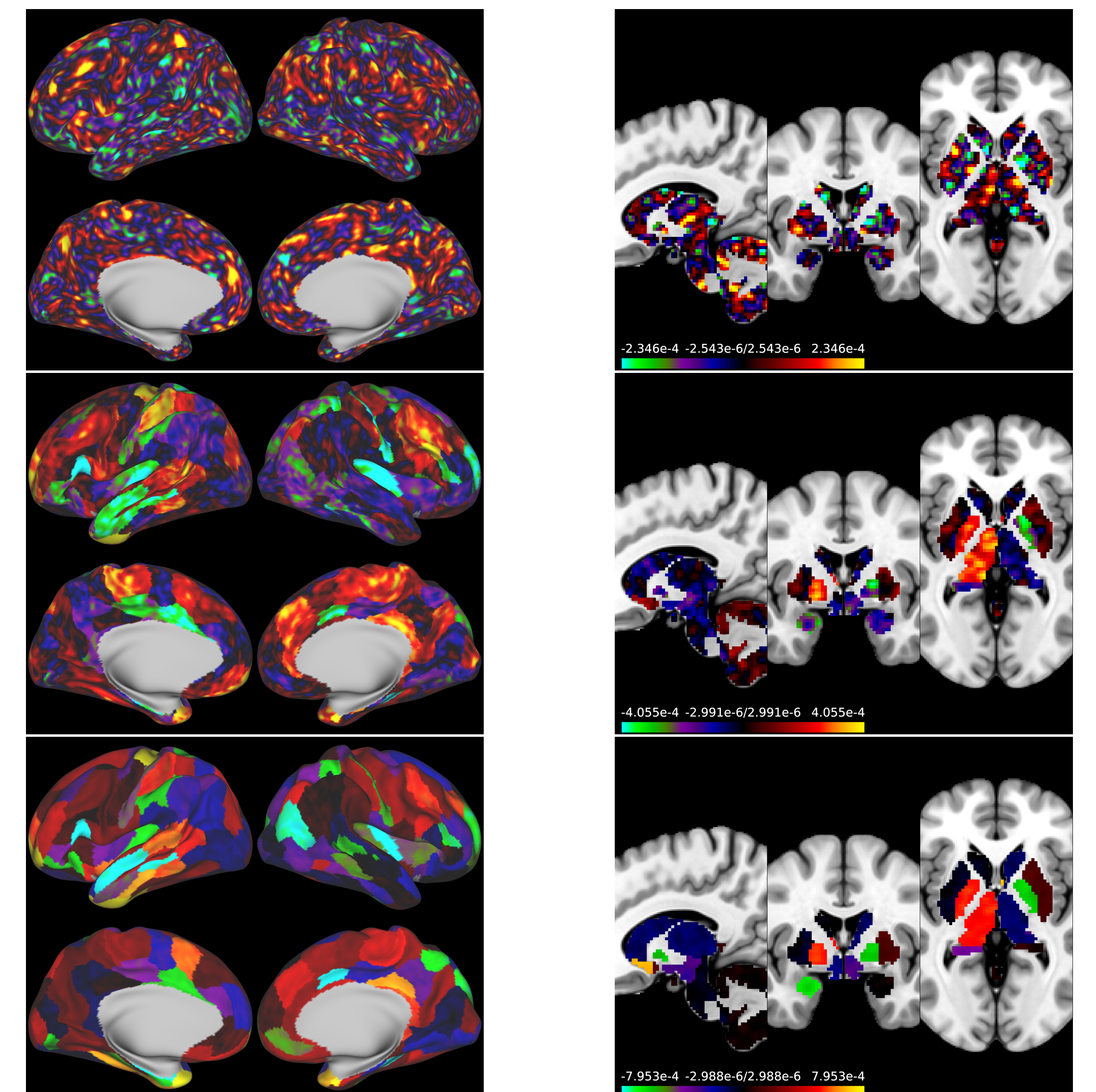


CANONICAL COEFFICIENTS

RCCA ($\lambda_1 = 0.001$)



GRCCA ($\lambda_1 = 1, 10, 100$ and $\mu_1 = 1$)



REFERENCES

- Hotelling. Relations between two sets of variables. *Biometrika*, 28, 1936.
- Gonzalez et al. CCA: An R Package to Extend Canonical Correlation Analysis. *Journal of Statistical Software*, 23(12), 2008.
- Leurgans et al. Canonical Correlation Analysis when the Data are Curves. *Journal of the Royal Statistical Society*, 55(3), 1993.