

# SlopeCCA and gslopeCCA: sorted L1 penalized canonical correlation analysis

Alexej Gossmann and Yu-Ping Wang

Multiscale Bioimaging and Bioinformatics Laboratory, Tulane University, New Orleans, LA



Tulane University

## Introduction

- ▶ Canonical correlation analysis (Hotelling, 1936), abbr. CCA, is a classical statistical technique, which can be used to make sense of the cross-correlation of two sets of measurements collected on the same set of samples (e.g., two genomic assays for the same set of cancer cells, or fMRI imaging and DNA sequencing data for the same mental illness patients).
- ▶ Sparse CCA (e.g., Parkhomenko et. al., 2009, or Witten et. al., 2009) extends the classical CCA to high-dimensional data by imposing a sparsity assumption on the CCA solution.
- ▶ Determination of the sparsity level, or the model tuning parameters, remains a challenging problem.
- ▶ We propose a definition of false discovery rate (FDR) for CCA, and (group) sorted  $\ell_1$  penalized CCA methods (slopeCCA and gslopeCCA). We show that slopeCCA and gslopeCCA are adaptive to the unknown sparsity level and keep the (group) FDR below a user-specified target threshold.

## Assumed data structure

- ▶ Let  $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^p$  be independent  $\mathcal{N}(\mathbf{0}, \Sigma_X)$ .
- ▶ Let  $\mathbf{y}_1, \dots, \mathbf{y}_n \in \mathbb{R}^q$  be independent  $\mathcal{N}(\mathbf{0}, \Sigma_Y)$ .
- ▶ Assume that  $\text{Cov}(\mathbf{x}_k, \mathbf{y}_k) = \Sigma_{XY} \in \mathbb{R}^{p \times q}$  for all  $k \in \{1, \dots, n\}$ , and that  $\text{Cov}(\mathbf{x}_k, \mathbf{y}_j) = \mathbf{0}$  whenever  $k \neq j$ .
- ▶ Define random matrices  $\mathbf{X} \in \mathbb{R}^{n \times p}$  and  $\mathbf{Y} \in \mathbb{R}^{n \times q}$

$$\mathbf{X} := \begin{bmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_n^T \end{bmatrix}, \quad \mathbf{Y} := \begin{bmatrix} \mathbf{y}_1^T \\ \mathbf{y}_2^T \\ \vdots \\ \mathbf{y}_n^T \end{bmatrix}.$$

- ▶ We can think of  $\mathbf{X}$  and  $\mathbf{Y}$  as two datasets containing respectively  $p$  and  $q$  features for the same  $n$  independent samples, where the features in the two datasets are cross-correlated with cross-covariance matrix  $\Sigma_{XY}$ .

## Canonical correlation analysis (CCA)

- ▶ CCA chooses  $\mathbf{u} \in \mathbb{R}^p$  and  $\mathbf{v} \in \mathbb{R}^q$  to maximize the sample correlation between  $\mathbf{X}\mathbf{u}$  and  $\mathbf{Y}\mathbf{v}$ ; i.e., the CCA optimization problem is

$$\arg \max_{\mathbf{u} \in \mathbb{R}^p, \mathbf{v} \in \mathbb{R}^q} \widehat{\text{Cov}}(\mathbf{X}\mathbf{u}, \mathbf{Y}\mathbf{v}) = \arg \max_{\mathbf{u} \in \mathbb{R}^p, \mathbf{v} \in \mathbb{R}^q} \frac{1}{n} \mathbf{u}^T \mathbf{X}^T \mathbf{Y} \mathbf{v},$$

subject to

$$\widehat{\text{Var}}(\mathbf{X}\mathbf{u}) = \frac{1}{n-1} \mathbf{u}^T \mathbf{X}^T \mathbf{X} \mathbf{u} = 1, \quad \widehat{\text{Var}}(\mathbf{Y}\mathbf{v}) = 1.$$

The solution is called first pair of canonical vectors.

- ▶ Subsequent pairs of canonical vectors are restricted to be uncorrelated with the previous ones.
- ▶ The problem is degenerate if  $n \leq \max\{p, q\}$ .

## False discovery rate (FDR) for CCA

- ▶ For  $i \in \{1, 2, \dots, p\}$ , we call the  $i$ th entry of the estimate  $\hat{\mathbf{u}}$  a *false discovery*, if  $\hat{u}_i \neq 0$  but  $(\Sigma_{XY})_{i,j} = 0$  for all  $j \in \{1, 2, \dots, q\}$ .
- ▶ Let  $V_{\hat{\mathbf{u}}}$  denote the number of false discoveries in  $\hat{\mathbf{u}}$ .
- ▶ Let  $R_{\hat{\mathbf{u}}}$  be the total number of non-zero entries  $\hat{\mathbf{u}}$ .
- ▶ Analogously, define  $R_{\hat{\mathbf{v}}}$  and  $V_{\hat{\mathbf{v}}}$ .

### Definition (FDR)

Define the false discovery rate in  $\mathbf{u}$  by

$$\text{FDR}(\hat{\mathbf{u}}) := \mathbb{E} \left( \frac{V_{\hat{\mathbf{u}}}}{\max\{R_{\hat{\mathbf{u}}}, 1\}} \right),$$

and analogously define  $\text{FDR}(\hat{\mathbf{v}})$ .

- ▶ **gFDR** is a generalization of **FDR** for the case when variables form groups (Brzyski et. al., 2016).

## Sorted L1 norm (Bogdan et. al., 2015)

The sorted  $\ell_1$  norm in  $p$  dimensions is given by,

$$J_{\lambda}(\mathbf{u}) := \sum_{i=1}^p \lambda_i |\mathbf{u}|_{(i)},$$

where  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p \geq 0$  is a non-increasing regularizing sequence, and  $|\mathbf{u}|_{(1)} \geq |\mathbf{u}|_{(2)} \geq \dots \geq |\mathbf{u}|_{(p)}$  is the order statistic of the magnitudes of the vector  $\mathbf{u} \in \mathbb{R}^p$  (absolute values ranked in non-increasing order).

## SlopeCCA and gslopeCCA

Inspired by SLOPE (Bogdan et. al., 2015) and Group SLOPE (Brzyski et. al., 2016, Gossmann et. al., 2015), we propose two novel sparse CCA methods with the goal of FDR and gFDR control.

### Definition (slopeCCA: sorted L1 penalized CCA)

$$\text{minimize}_{\mathbf{u} \in \mathbb{R}^p, \mathbf{v} \in \mathbb{R}^q} \left\{ -\mathbf{u}^T \mathbf{X}^T \mathbf{Y} \mathbf{v} + \sqrt{n} J_{\lambda^u}(\mathbf{u}) + \sqrt{n} J_{\lambda^v}(\mathbf{v}) \right\},$$

subject to  $\|\mathbf{u}\|_2 \leq 1, \|\mathbf{v}\|_2 \leq 1$ .

### Definition (gslopeCCA: group sorted L1 penalized CCA)

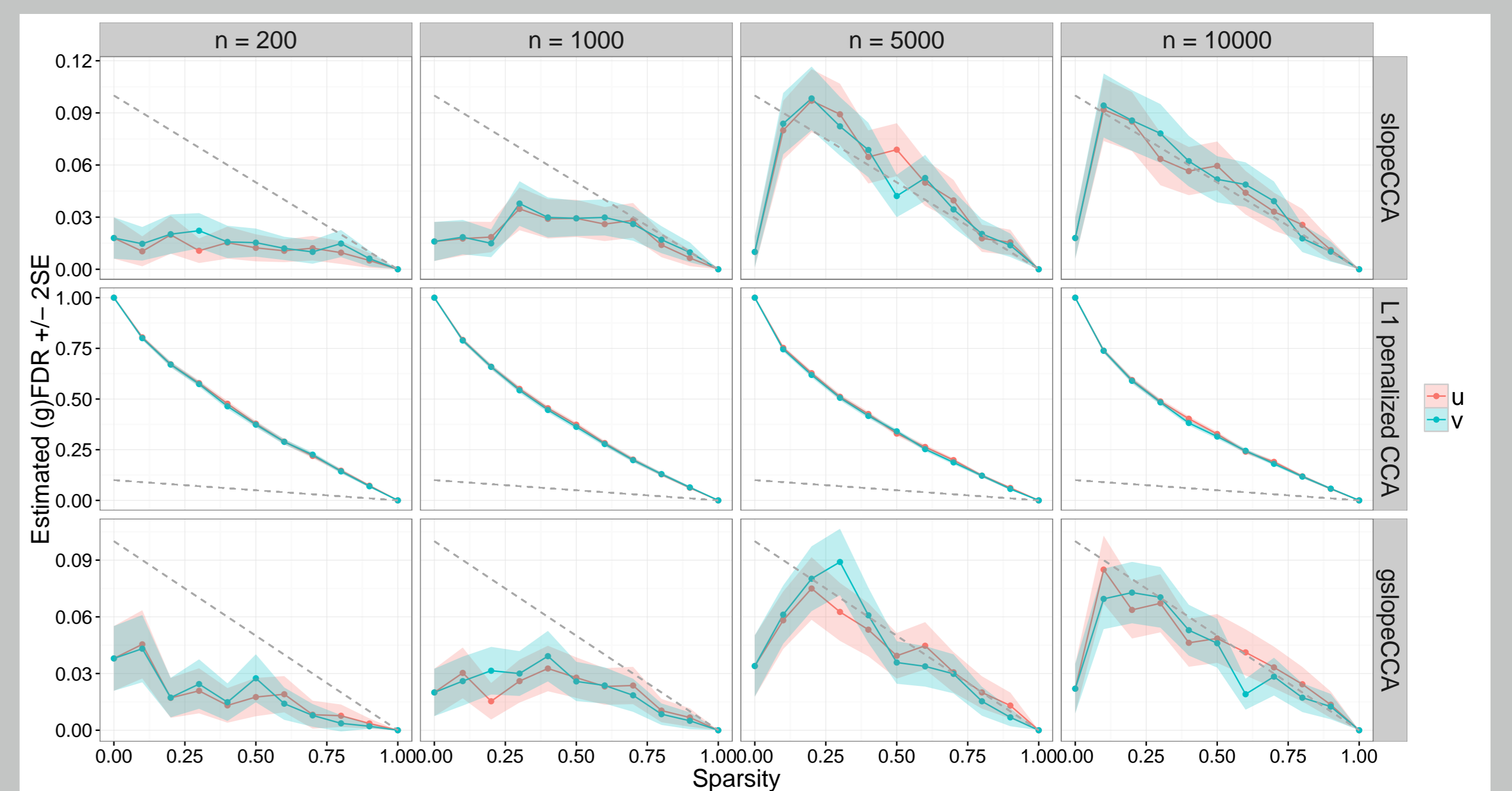
$$\text{minimize}_{\mathbf{u} \in \mathbb{R}^p, \mathbf{v} \in \mathbb{R}^q} \left\{ -\mathbf{u}^T \mathbf{X}^T \mathbf{Y} \mathbf{v} + \sqrt{n} J_{\lambda^u}(\|\mathbf{u}_1\|_2, \dots, \|\mathbf{u}_m\|_2)^T \right. \\ \left. + \sqrt{n} J_{\lambda^v}(\|\mathbf{v}_1\|_2, \|\mathbf{v}_2\|_2, \dots, \|\mathbf{v}_m\|_2)^T \right\},$$

subject to  $\|\mathbf{u}\|_2 \leq 1, \|\mathbf{v}\|_2 \leq 1$ .

Higher-order pairs of canonical vectors can be found by applying slopeCCA (or gslopeCCA) to a residual matrix, obtained from  $\mathbf{X}^T \mathbf{Y}$  and the previously found canonical pairs. The slopeCCA and gslopeCCA optimization problems are biconvex and can be solved by alternating minimization algorithms.

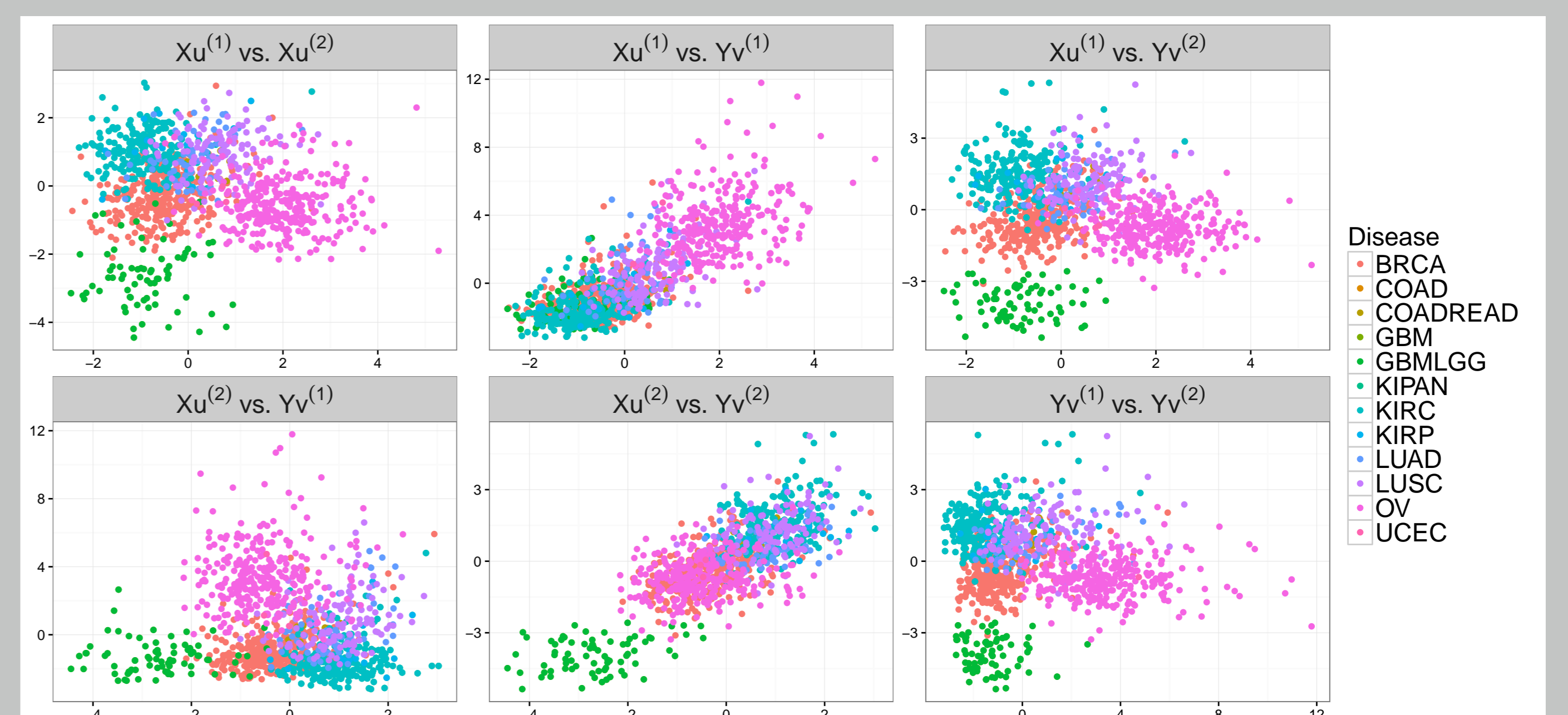
## Asymptotic FDR properties of slopeCCA and gslopeCCA

- ▶  $\mathbf{X} \in \mathbb{R}^{n \times p}$  and  $\mathbf{Y} \in \mathbb{R}^{n \times q}$  are Gaussian with  $p = q = 300$ , varying  $n$ .
- ▶  $\Sigma_X, \Sigma_Y$  are *diagonal* (slopeCCA and PMA) or *block-diagonal* (gslopeCCA).
- ▶ 11 evenly spaced sparsity levels between 0 (i.e.,  $\mathbf{X}$  is uncorrelated with  $\mathbf{Y}$ ) and 1 (i.e., every feature of  $\mathbf{X}$  is correlated to some feature of  $\mathbf{Y}$ , and vice versa) were considered, where  $(\Sigma_{XY})_{ij} \in \{0, 0.5\}$ .
- ▶ For comparison the sparse CCA of the widely used R package PMA (Witten et. al., 2009) is applied to the same data, with the tuning parameters determined by a permutation based approach (via `PMA::CCA.permute`).
- ▶ Dashed line represents the theoretical bound on FDR and gFDR for slopeCCA and gslopeCCA with target FDR of 0.1 (Theorem not shown).



## Application to TCGA data

- ▶ We apply gslopeCCA to methylation and mRNA data for 12 diseases from The Cancer Genome Atlas (TCGA);  $\mathbf{X} \in \mathbb{R}^{1436 \times 24981}$ ,  $\mathbf{Y} \in \mathbb{R}^{1436 \times 20255}$ .
- ▶ Canonical variates reveal differences between cancer types.
- ▶ A drawback: the obtained canonical vectors are not as sparse as desired, and therefore the results are difficult to interpret.



## Acknowledgements

Our work is partially supported by NIH R01 GM109068, R01 MH104680, R01 MH107354.

Created with L<sup>A</sup>T<sub>E</sub>X Beamer poster <http://www-16.informatik.uni-aachen.de/~dreuw/latexbeamerposter.php>