



Published in final edited form as:

*Med Image Anal.* 2014 August ; 18(6): 891–902. doi:10.1016/j.media.2013.10.010.

## Correspondence between fMRI and SNP data by group sparse canonical correlation analysis

Dongdong Lin<sup>a,b</sup>, Vince D. Calhoun<sup>c</sup>, and Yu-Ping Wang<sup>a,b,d,\*</sup>

<sup>a</sup> Biomedical Engineering Department, Tulane University, New Orleans, LA 70118, USA

<sup>b</sup> Center of Genomics and Bioinformatics, Tulane University, New Orleans, LA 70118, USA

<sup>c</sup> The Mind Research Network, Department of Electrical and Computer Engineering, University of New Mexico, Albuquerque, NM 87131, USA

<sup>d</sup> Center for Systems Biomedicine, Shanghai University for Science and Technology, Shanghai, China

### Abstract

Both genetic variants and brain region abnormalities are recognized as important factors for complex diseases (*e.g.*, schizophrenia). In this paper, we investigated the correspondence between single nucleotide polymorphism (SNP) and brain activity measured by functional magnetic resonance imaging (fMRI) to understand how genetic variation influences the brain activity. A group sparse canonical correlation analysis method (group sparse CCA) was developed to explore the correlation between these two datasets which are high dimensional—the number of SNPs/voxels is far greater than the number of samples. Different from the existing sparse CCA methods (sCCA), our approach can exploit structural information in the correlation analysis by introducing group constraints. A simulation study demonstrates that it outperforms the existing sCCA. We applied this method to the real data analysis and identified two pairs of significant canonical variates with average correlations of 0.4527 and 0.4292 respectively, which were used to identify genes and voxels associated with schizophrenia. The selected genes are mostly from 5 schizophrenia (SZ)-related signalling pathways. The brain mappings of the selected voxels also indicate the abnormal brain regions susceptible to schizophrenia. A gene and brain region of interest (ROI) correlation analysis was further performed to confirm the significant correlations between genes and ROIs.

### Keywords

Group sparse CCA; SNP; fMRI; Feature selection; Imaging genetics

## 1. Introduction

Schizophrenia is a complex disease and considered to be caused by the interplay of a number of genetic factors (*e.g.*, change of gene regulation, and alteration of mRNA and SNP) and environmental effects. Genetic factors play an important role in causing schizophrenia disease. People born from a family with a history of schizophrenia have higher risks of schizophrenia than those without a family schizophrenia history. In recent years, many studies have focused on exploring critical genes associated with the schizophrenia. Many potential genetic variants have been reported as possible risk factors such as the G72/G30 gene locus on chromosome 13q (Badner and Gershon, 2002; Abecasis et al., 2004) Gene DISC1 variation (Callicott et al., 2005; Porteous et al., 2006) and copy number variations on gene GRIK3, EFNA5, AKAP5 and CACNG2 (Wilson et al., 2006; Sutrala et al., 2007). In addition to genetic studies, fMRI has also been widely used for the study of schizophrenia because of its capability to identify functional abnormalities within brain regions of schizophrenic patients (Jansma et al., 2004; Li et al., 2007; Meda et al., 2008; Szyck et al., 2009).

Genetic variants and brain region abnormalities are both important markers for the study of schizophrenia. Combining both data can not only contribute to a better understanding of biological mechanisms on brain structure and function but also have the potential to improve the diagnosis and treatments of complex diseases. However, current imaging genetics studies either take brain imaging measurements as endophenotypes to study the associated genetic variants or investigate the effects of a small set of candidate genetic variants on the whole brain measurements (Hamid et al., 2009; Le Cao et al., 2009; Wiley, 2011). It is still challenging to explore the relationship between a large amount of genetic variants and a large number of brain imaging measurements. Therefore, correlative analysis approaches for large-scale multi-modal data analysis are highly demanded.

In this work, we aim to study the effects of multiple SNPs or genes on functional brain activity in schizophrenia. An effective multivariate statistical method is needed. Canonical Correlation Analysis (CCA (Hotelling, 1936)) or Partial Least Squares regression (PLSR (Le Cao et al., 2008)) have been proposed to analyze multi-modal datasets. The CCA aims to maximize the correlation between the linear combinations of variables from two data sets, *e.g.*, a linear combination of SNPs and a linear combination of voxels. However, the method will have the over-fitting issues in analyzing high dimensional data such as SNP and brain imaging data as shown in Fig. 1. Thousands of SNPs with linkage disequilibrium (LD) are detected to reflect the genetic variant at different locus. The number of voxels included in the whole brain fMRI image is also very large (*e.g.*,  $53 \times 63 \times 46$ ). Traditional CCA will perform poorly in such a case due to the multi-collinearity (linear dependence) problem, and thus having computational difficulty (Park-homenko et al., 2009). To address above issue, sparse CCA (sCCA) methods, mostly using the  $l - 1$  norm (CCA- $l_1$ ) or the combination of  $l - 1$  and  $l - 2$  norm (CCA-*elastic net*) penalties, have been developed by introducing the sparse penalties into the traditional CCA model (Waijenborg et al., 2008; Le Cao et al., 2009; Parkhomenko et al., 2009; Witten et al., 2009; Witten and Tibshirani, 2009; Boute and Liu, 2010). Despite the success, they didn't account for group structures within the data in the analysis (*e.g.*, multiple SNPs within the same gene, a group of voxels within the same

region, a group of voxels within the same ROI, etc.), which often exist or are implied by the biological mechanism. For example, SNPs within the same gene have similar functions and act together at the gene or pathway level to affect the brain activity. These SNP effects can be added up to a larger difference (Tyekucheva et al., 2011). Several previous works have shown the benefit of accounting for the group effect of features in the sCCA models (Chen and Liu, 2012; Chen et al., 2013; Lin et al., 2013). However, to our knowledge, little work has been reported to incorporate the group effect into the sCCA model for fMRI and SNP data integration. Motivated by this fact, in this paper, we developed a group sparse CCA model based integration method by imposing the sparse group lasso penalty on the CCA model for the integrative analysis of SNP and fMRI data; please refer to Fig. 1 for illustration. This method has the following advantages: (1) A group of features (voxels/SNPs) will be inspected during the correlation analysis, which can study the joint effects of multiple SNPs on the regions of voxels; (2) feature selection will be performed at both group and single feature level. Irrelevant groups of features as well as single feature within each group can be removed. Our group sparse CCA method can both exploit group information in the correlation analysis while filter out noisy features within the group simultaneously.

The group sparse CCA can estimate the correlation between canonical variates, corresponding to a set of significant SNPs or brain imaging voxels. Based on the estimates, we provided a gene-ROI correlation analysis to further confirm the significance of the correlations between genes and brain functions in ROIs.

The rest of the paper is organized as follows. The proposed group sparse CCA model and algorithm are introduced in the section of theory. The group sparse CCA based integration method for SNP and fMRI data is described in the section of method. The validation and comparison of our model with other sCCA models on both simulated and real data analysis are presented in the section of results. The pathway analysis results and limitations of the proposed method are discussed finally.

## 2. Theory

In this section, we first introduces CCA model, based on which the group sparse model is presented. Then we propose a numerical algorithm based on block coordinate descent to solve the model. Finally, we show that the general model we propose can include several existing sCCA models and hence the numerical algorithm can also be applied for their efficient solutions.

### 2.1. Group sparse CCA

We consider two sets of data  $\mathbf{X}$  and  $\mathbf{Y}$  with  $n$  samples, where  $\mathbf{X}$  has  $p$  variables and  $\mathbf{Y}$  has  $q$  variables. Assuming  $p, q \gg n$  and the columns of  $\mathbf{X}$  and  $\mathbf{Y}$  have been standardized such that the mean of each column is zero and the 2-norm is one. The variance matrices and covariance matrices of  $\mathbf{X}$  and  $\mathbf{Y}$  are denoted by  $\sum_{\mathbf{X}\mathbf{X}^T}$ ,  $\sum_{\mathbf{Y}\mathbf{Y}^T}$ ,  $\sum_{\mathbf{X}\mathbf{Y}^T}$  (or  $\sum_{\mathbf{Y}\mathbf{X}^T}$ ) respectively. The CCA model aims to find two loading vectors or projections  $\boldsymbol{\alpha}$ ,  $\boldsymbol{\beta}$  to maximize the correlation between the linear combinations of variables in  $\mathbf{X}$  and  $\mathbf{Y}$ ,  $\boldsymbol{\alpha}^T\mathbf{X}$  and  $\boldsymbol{\beta}^T\mathbf{Y}$  as shown in the following equation:

$$\max_{\alpha, \beta} \alpha^t \sum_{XY} \beta \quad \text{s.t.} \quad \alpha^t \sum_{XX} \alpha = 1, \quad \beta^t \sum_{YY} \beta = 1 \quad (1)$$

There are over-fitting issues in Eq. (1) due to the small samples but high dimensional variables (i.e.,  $p, q \gg n$ ), which also might result in the ill-condition of the matrices  $\sum_{XX}$  and  $\sum_{YY}$ . So, sparse penalties such as the  $l-1$  norm, elastic net, have been imposed on the loading vectors in Eq. (2) in CCA analysis.

$$\min_{\alpha, \beta} -\alpha^t \sum_{XY} \beta + \lambda_1 \Psi(\alpha) + \lambda_2 \Phi(\beta) \quad \text{s.t.} \quad \alpha^t \sum_{XX} \alpha \leq 1, \quad \beta^t \sum_{YY} \beta \leq 1 \quad (2)$$

where  $\Psi(\cdot)$  and  $\Phi(\cdot)$  denote the penalized function on  $\alpha$  and  $\beta$  respectively, which are often taken to be  $l-1$  norm,  $\Psi(\alpha) = \|\alpha\|_1$ ,  $\Phi(\beta) = \|\beta\|_1$ . These sparse penalties can result in a large number of features/variables to zero in  $\alpha$  and  $\beta$ . In addition, the constrained conditions in Eq. (2) are relaxed compared to Eq. (1) so that we can first find the alternative solution in a closed region and then obtain the solution of the optimization problem satisfying

$\alpha^t \sum_{XX} \alpha = 1, \beta^t \sum_{YY} \beta = 1$  (see optimization algorithms in Section 2.2). Similar to conventional CCA, Sparse CCA can also perform multiple pairs of canonical variates extraction iteratively as the following procedure: maximizing the correlation between two data sets by extracting the first pair of sparse loading vectors as well as pair of canonical variates; then removing the effects of the first pair of canonical variates and finding the second sparse loading vectors that maximizes the correlation but is irrelevant to the first pair. This process will not stop until the  $r$ th projection pair is gained ( $r = \text{rank}(X^t Y)$ ).

Group information of variables is expected to be incorporated in the Eq. (2) by changing the penalized function since in many applications, a set of variables often form a group (e.g., SNPs spanning a gene, genes in a pathway). Therefore, we propose a group sparse CCA model by introducing the sparse group lasso penalty into Eq. (2). For simplicity, we will only consider non-overlapping groups in this paper and assume that variables in  $X$  and  $Y$  are partitioned into  $L$  and  $H$  disjoint groups respectively. The following model, namely group sparse CCA (or CCA-sparse group), is proposed to consider group structures existed in the data:

$$\min_{\alpha, \beta} -\alpha^t \sum_{XY} \beta + \lambda_1 \|\alpha\|_G + \tau_1 \|\alpha\|_1 + \lambda_2 \|\beta\|_G + \tau_2 \|\beta\|_1 \quad \text{s.t.} \quad \alpha^t \sum_{XX} \alpha \leq 1, \quad \beta^t \sum_{YY} \beta \leq 1 \quad (3)$$

where  $\|\alpha\|_G = \sum_{l=1}^L \omega_l \|\alpha\|_2$ ,  $\|\beta\|_G = \sum_{h=1}^H \mu_h \|\beta\|_2$  are the group penalties to account for joint effects of features within the same group. This model is more realistic in many cases, e.g., multiple SNPs rather than individual SNPs from the same gene work together as a group to be associated with a disease. The group penalty uses the nondifferentiability of  $\|\alpha_l\|_2$  (or  $\|\beta_h\|_2$ ) at  $\alpha_l = 0$  ( $\beta_h = 0$ ) to set the coefficients of the group to be 0; then the entire group of features will be removed to achieve the group sparsity. While we consider group effects, we can still keep the selection of individual variable/feature. So the  $l-1$  norm penalties on the individual features (i.e.,  $\|\alpha\|_1$  and  $\|\beta\|_1$ ) are imposed.  $\lambda_1$  and  $\lambda_2$  are the

tuning parameters to control the group sparsity while  $\tau_1$  and  $\tau_2$  are used to control individual feature sparsity.  $\omega_l$  and  $\mu_l$  are the weights to adjust for the group size differences. We set them to be  $s_i^{1/2}$ , where  $s_i$  is the  $i$ th group size.

In a particular case, since sparse group lasso penalty is a combination of  $l-1$  norm and group lasso penalty, Eq. (3) can be reduced to the CCA-group model with only group lasso penalty ( $\tau_1 = \tau_2 = 0$ ) and CCA- $l1$  model without group lasso penalty ( $\lambda_1 = \lambda_2 = 0$ ). As discussed above, CCA-group model can tend to select features group by group and keep all features within a selected group while CCA- $l1$  will ignore group effect among features.

## 2.2. Iterative optimization algorithm

To solve Eq. (3), we rewrite Eq. (3) based on the singular vector decomposition of matrix  $\mathbf{K}$  and impose sparse penalties on vectors  $\mathbf{u}$ ,  $\mathbf{v}$

$$\min_{\mathbf{u}, \mathbf{v}} \|\mathbf{K} - d\mathbf{u}\mathbf{v}^t\|_F^2 + \lambda_1 \|\mathbf{u}\|_G + \tau_1 \|\mathbf{u}\|_1 + \lambda_2 \|\mathbf{v}\|_G + \tau_2 \|\mathbf{v}\|_1 \quad s.t. \quad \|\mathbf{u}\|_2^2 \leq 1, \quad \|\mathbf{v}\|_2^2 \leq 1 \quad (4)$$

where  $\mathbf{K} = \sum_{XX}^{-1/2} \sum_{XY} \sum_{YY}^{-1/2} = \sum_i d_i \mathbf{u}_i \mathbf{v}_i$ , and  $d_i$  is the positive square root of the  $i$ th eigenvalue of  $\mathbf{K}^t \mathbf{K}$ .  $\mathbf{u}_i$ ,  $\mathbf{v}_i$  are the  $i$ th eigenvectors of  $\mathbf{K}^t \mathbf{K}$  corresponding to  $d_i$ .  $\|\mathbf{u}\|_2^2 = 1$ ,  $\|\mathbf{v}\|_2^2 = 1$  should be satisfied when the solution of the optimization problem is obtained. The loading vectors can then be derived by

$$\boldsymbol{\alpha}_i = \sum_{XX}^{-1/2} \mathbf{u}_i, \quad \text{and} \quad \boldsymbol{\beta}_i = \sum_{YY}^{-1/2} \mathbf{v}_i \quad (5)$$

The matrices  $\sum_{XX}$  and  $\sum_{YY}$  in Eq. (5) might be ill-conditioned because of the high dimensionality of data. We adopt Witten and Tibshirani's (2009) method by replacing the covariance matrices with identity matrices  $\mathbf{I}$  and hence penalizing the vectors  $\mathbf{u}$ ,  $\mathbf{v}$  instead of the loading vectors  $\boldsymbol{\alpha}$ ,  $\boldsymbol{\beta}$ .

Since the problem in Eq. (4) is still not a convex optimization with respect to  $\mathbf{u}$  and  $\mathbf{v}$ , for simplicity, we decouple the problem into two simple biconvex optimizations by fixing  $\mathbf{u}$  and  $\mathbf{v}$  alternatively. An iterative algorithm is then derived to solve the problem, as shown in Table 1.

Taking the solution of (a) for example, one can find the solution with the following Lagrange form formulation:

$$\min_{\mathbf{u}} -tr(\mathbf{K}\mathbf{v}\mathbf{u}^t) + \lambda_1 \sum_{l=1}^L \omega_l \|\mathbf{u}_l\|_2 + \tau_1 \|\mathbf{u}\|_1 + \Delta(\|\mathbf{u}\|_2^2 - 1) \quad (6)$$

where  $\Delta$  is the parameter to make  $\|\mathbf{u}\|_2^2 = \sum_{l=1}^L \|\mathbf{u}_l\|_2^2 = 1$ . This is converted into a simple optimization problem with separable objection function and sparse group lasso penalty. A

block coordinate decent algorithm has been developed to solve this problem efficiently as shown in Table 2.

The coordinate decent algorithm has been shown to be effective in solving generalized linear regression models (Wu et al., 2009; Friedman et al., 2010a, 2010b), especially for underdetermined system. It estimates parameter one by one by fixing the other parameters unchanged. Similarly, instead of estimating parameter individually, block coordinate decent algorithm will estimate a block of parameters each time while fixing the other blocks of parameters. In data set  $\mathbf{X}$ , each group  $k = 1, 2, \dots, L$ , will be inspected one by one. If a group is selected, we will then select each feature in the group by the coordinate decent with the soft-thresholding. Since the optimization is convex, the optimal solution of Eq. (6) is determined by a sub-gradient equation and will converge to a global minimum.

The CCA-group model is a special case of group sparse CCA when  $\tau_1 = 0$ ; the soft-threshold operator in Step (2) (Table 2) need to be changed to the simpler form:

$$\|S((\mathbf{KV})^{(k)}, \tau_1)\|_2 = \|(\mathbf{KV})^{(k)}\|_2 \quad (7)$$

The general sparse CCA formula Eq. (3) that we propose can include a variety of sparse CCA models used before (Le Cao et al., 2009; Witten and Tibshirani, 2009), e.g., CCA- $l_1$  (Eq. (8)) and CCA-elastic net (Eq. (9)) models. The coordinate decent algorithm can also be applied to solve both models.

$$\Psi(\mathbf{u}) = \|\mathbf{u}\|_1, \quad \Phi(\mathbf{v}) = \|\mathbf{v}\|_1 \quad (8)$$

$$\Psi(\mathbf{u}) = (1 - \delta_1) \|\mathbf{u}\|_2^2 + \delta_1 \|\mathbf{u}\|_1, \quad \Phi(\mathbf{v}) = (1 - \delta_2) \|\mathbf{v}\|_2^2 + \delta_2 \|\mathbf{v}\|_1 \quad (9)$$

where  $\delta_1, \delta_2$  are the parameters, controlling the trade-off between  $l_2$  norm and  $l_1$  norm penalties.

### 3. Method

We applied group sparse CCA to investigate the association of functional brain regions with genetic variations as shown in Fig. 1. Components extracted from fMRI represent brain regions expressing the functional difference in different subjects. Components from SNP data are linear combinations of SNPs from different genes that may have associations with the disease. After preprocessing, the collected SNPs and ROI-based voxels are both still high dimensional with a large number of features compared to the number of samples. We then used the group sparse CCA to estimate two group loading vectors  $\mathbf{u}$  and  $\mathbf{v}$ , from which a pair of canonical variates is obtained. The loading vectors for each component reflect the effect size of the features on the correlation. Then  $\mathbf{u}$  and  $\mathbf{v}$  were also used to perform gene-ROI correlation analysis to identify the significantly correlated genes and ROIs.

### 3.1. Group sparse CCA for fMRI and SNP data analysis

We represented fMRI data collected from participants as a set of spatial voxels. These voxels were divided into 116 ROIs based on the aal (automated anatomical labelling) template (Tzourio-Mazoyer et al., 2002). These ROIs were assumed to be spatially independent but the voxels within each ROI may correlate with each other. These voxels were grouped by ROIs so that we can perform the whole brain analysis (Ng and Abugharbieh, 2011). For SNP data, we extracted those SNPs from preselected 74 reported SZ-risk genes. These SNPs were grouped at gene level (Liu et al., 2013). Hence the  $\mathbf{X}$  and  $\mathbf{Y}$  matrices can be constructed as follows:

$$\begin{aligned} \mathbf{X} &= \left[ \mathbf{G}_X^1, \mathbf{G}_X^2, \dots, \mathbf{G}_X^L \right], & \mathbf{G}_X^k &\in R^{n \times l_k}, \quad k=1, 2, \dots, L \\ \mathbf{Y} &= \left[ \mathbf{G}_Y^1, \mathbf{G}_Y^2, \dots, \mathbf{G}_Y^H \right], & \mathbf{G}_Y^t &\in R^{n \times h_t}, \quad t=1, 2, \dots, H. \end{aligned}$$

where  $L = 74$  indicates the number of genes in SNP data;  $H = 116$  is the number of ROIs used in fMRI data, and  $n$  is the number of samples.  $l_k$  and  $h_t$  are the number of SNPs and voxels contained in the  $k$ th gene and  $t$ th ROI respectively. Under this data structure, two sparse loading vectors  $\mathbf{u}$  and  $\mathbf{v}$  (or  $\boldsymbol{\alpha}$  and  $\boldsymbol{\beta}$ ) were also grouped in the same way. This way the original data were projected into a low dimensional space.

Four tuning parameters ( $\lambda_1; \tau_1; \lambda_2; \tau_2$ ) were used in the model to control the group sparsity ( $\lambda_1; \lambda_2$ ) and individual feature sparsity ( $\tau_1, \tau_2$ ). A cross-validation can be used to select the optimal parameters but is time-consuming. Therefore, we divided the cross-validation into two steps: (1) using the CCA-group method to select optimal  $\lambda_1; \lambda_2$  by considering only the group penalty; (2) selecting the optimal  $\tau_1, \tau_2$  based on the selected  $\lambda_1; \lambda_2$ : In summary, the procedure of performing group sparse CCA on the data analysis is as follows:

1. Decompose matrix  $\mathbf{K}$  in Eq. (5) using SVD to initialize the loading vectors and normalize the vectors with the  $l - 2$  norm. Use the two-step cross validation to obtain the optimal tuning parameters.
2. Perform the sub-optimization analysis with respect to each loading vector ( $\mathbf{u}$  or  $\mathbf{v}$ ) in each modality alternatively by the block coordinate decent algorithm. A stopping criterion is assessed for each modality.
3. Two loading vectors are updated and the stopping criteria for both loading vectors are assessed. If they are satisfied, the solution is obtained, go to Step 4. Otherwise, go back to Step 2.
4. Test the significance of the predict correlation ( $corr_j$ ) by permutation. If the correlation is significant, go to Step 5. Otherwise, there is no significant correlation between the current two data sets.
5. Perform gene-ROI analysis to explore the pair-wise correlation between each gene and ROI.

6. Calculate the remaining matrix and repeat step (1–4) to extract the next pair of loading vectors until  $r$  pairs of loading vectors are obtained or the extracted predict correlation ( $corr_j$ ) is not significant.

### 3.2. Tuning parameters selection

The k-fold cross validation was recommended by Waaijenborg and Zwinderman (2009) and Parkhomenko et al. (2009) for parameter selection. Parkhomenko et al. (2009) adopted a criterion that maximizes the mean absolute canonical correlation value of the testing set as shown in the following equation:

$$\Delta_{corr1} = \frac{1}{k} \sum_{i=1}^k |cor(\mathbf{X}_i \hat{\mathbf{u}}^{-i}, \mathbf{Y}_i \hat{\mathbf{v}}^{-i})| \quad (10)$$

Waaijenborg et al. (2008) considered the mean difference between the canonical correlations of the training and testing subsets as in the following equation:

$$\Delta_{corr2} = \frac{1}{k} \sum_{i=1}^k |cor(\mathbf{X}_{-i} \hat{\mathbf{u}}^{-i}, \mathbf{Y}_{-i} \hat{\mathbf{v}}^{-i}) - cor(\mathbf{X}_i \hat{\mathbf{u}}^{-i}, \mathbf{Y}_i \hat{\mathbf{v}}^{-i})| \quad (11)$$

This criterion determines the number of variables which tend to have the same correlations in both training and testing subsets. It is sensitive to the correlation sign change that if the correlation of testing subset changes sign it would be penalized more than when the sign would not change.

Witten and Tibshirani (2009) proposed another permutation based method for optimizing parameters. The whole data set (including both training and testing data sets) will be used to compute the non-permuted correlation  $d_0$  and then  $\mathbf{X}$  is permuted with  $T$  times to calculate the permuted correlation  $d_{i=1, \dots, T}$ . The parameters having the largest  $p$ -value as in Eq. (12) will be selected.

$$p_{permute} = \left( d_0 - \frac{1}{T} \sum_i d_i \right) / std(d_{i=1, \dots, T}) \quad (12)$$

Here,  $\hat{\mathbf{u}}^{-i}$  and  $\hat{\mathbf{v}}^{-i}$  are estimated loading vectors from training data set.  $\mathbf{X}_{-i}$  and  $\mathbf{Y}_i$  are training data set in which subset  $j$  is deleted.  $\mathbf{X}_i$  and  $\mathbf{Y}_i$  are the testing data set. Based on these three criteria, we performed a simulation study using our proposed group sparse CCA model and other sCCA models to choose the optimal criterion (the details are described in Section 4.1).

### 3.3. Gene-ROI correlation analysis

After the parameters were determined, we obtained the estimates of vectors  $\mathbf{u}$ ,  $\mathbf{v}$  as well as the corresponding canonical variates. We used these estimates to measure the strength and significance of the correlation between genes and ROIs. The null hypothesis of no correlation between a gene and ROI can be written as  $H_0 : \rho_{11} = \rho_{12} = \dots = \rho_{ck} = \rho_{s_i s_j} = 0$  versus the alternative hypothesis  $H_A : \exists c; k > 0; \rho_{ck} \neq 0$ , where  $s_i$  and  $s_j$  are the number of SNPs and voxels in Gene $_i$  and ROI $_j$  respectively;  $\rho_{ck}$  is the correlation between the  $c$ th SNP



and the  $k$ th voxel. To test the hypothesis, we first calculated the pair-wise correlation between the SNPs from Gene $_i$  and voxels from ROI $_j$ . SNPs and voxels were indicated by non-zero coefficients in  $\mathbf{u}$  and  $\mathbf{v}$  respectively. Then, to avoid the potential bias due to varying gene or ROI size, we averaged the correlation of each SNP-voxel pair as the test statistics as follows.

$$\rho_{ij} = \frac{1}{S_i S_j} \sum_{c=1}^{s_i} \sum_{k=1}^{s_j} |\rho_{ck}| \quad (13)$$

By comparing our observed statistic  $\rho_{ij}$  with the null statistics  $\rho_{ij}^0$  with  $T$  times permutations of the samples, we can evaluate the significance of the correlation by

$$p \text{ value} = \left\{ \rho_b^0 \geq \rho_{ij}; b=1, 2, \dots, T \right\} / T \quad (14)$$

This gene-ROI correlation analysis can also be applied to the gene-gene correlation and ROI-ROI correlation studies.

### 3.4. Statistical evaluations

Although our CCA-group sparse method aims to obtain a higher correlation with fewer features selected, in our real data both number of voxels and SNPs are much larger than that of subjects. In order to test the significance of the correlation obtained by the CCA-sparse group method, we took the permutation-based testing. We first obtained  $\mathbf{u}$  and  $\mathbf{v}$  based on the optimal parameters, and then permuted the order of subjects in SNP data randomly while keeping the subject order in fMRI data unchanged. The correlation for each permutation can be calculated by the same  $\mathbf{u}$  and  $\mathbf{v}$ , which is expected to estimate the null distribution of the correlation. The  $p$ -value can be estimated by large number of permutations. The lower the  $p$ -value, the more significant the detected correlation is.

## 4. Results

### 4.1. Simulation

To assess the performance of the proposed group sparse CCA method, we first simulated two correlated data sets and then we compared group sparse CCA with the other penalized CCA methods such as the CCA-group and CCA- $l_1$  on these simulated data.

Two data sets of SNP data  $\mathbf{X}$  with  $p$  SNPs and fMRI data  $\mathbf{Y}$  with  $q$  voxels were simulated. To correlate the SNPs with the voxels, a latent model similar to (Parkhomenko et al., 2009) was used. We first set a latent variable  $\Upsilon = \{\gamma_i | i = 1, \dots, n\}$  with normal distribution

$N(0, \sigma_\gamma^2)$  to have the similar effect on the associated SNPs and voxels in the two data sets.

Then  $\mathbf{X}$  and  $\mathbf{Y}$  data set were simulated from a multivariate normal distribution by  $\mathbf{X}^T \mathbf{X}$  and

$\text{diag} \left( \sum_{x \in X} \right)$  respectively, where the vectors  $x_i \in R^p$ ,  $y_i \in R^q$  are the observations of the  $i$ th sample in  $\mathbf{X}$  and  $\mathbf{Y}$ .  $\boldsymbol{\alpha} = [\alpha_1, \dots, \alpha_j, \dots, \alpha_p]$ ,  $\boldsymbol{\beta} = [\beta_1, \dots, \beta_k, \dots, \beta_q]$ ,  $\alpha_j = 0$ ,  $\beta_k = 0$ , if  $x_j$ ,  $y_k$  are the correlated variables; otherwise, the variables would be considered as noise with zero means ( $\alpha_j = 0$ ,  $\beta_k = 0$ ).  $\gamma_i$  is the  $i$ th observation of  $\Upsilon$ .  $\sigma_\gamma^2$ ,  $\sigma_e^2$ , are the variances of  $\Upsilon$  and noise

variable, which were used to control the correlation between SNPs and voxels.  $\sum_{p \times p}$  and  $\sum_{q \times q}$  were the variance–covariance matrices of each data set. They were used to simulate the group effect within each data set. For each data set, we set the correlation between the correlated variables within the same group as from the uniform distribution  $U(0.1, 0.3)$ .

The samples of two data sets were simulated from the multivariate normal distribution. Then SNP data  $\mathbf{X}$  will be converted into categorical variables with three levels (e.g.,  $-1, 0, 1$ ) by specifying the minor allele frequency (MAF)  $p$  from the uniform distribution  $\text{Unif}(0.2, 0.4)$ . The probability of an observation being from a particular level was set to be  $p^2, p(1-p)$  and  $(1-p)^2$  respectively. For each variable, the expected number of observations in each level was  $n_1 = p^2n, n_2 = p(1-p)n, n_3 = (1-p)^2n$ . We sorted the observations in each variable and assign value 1 to those  $n_1$  observations with the highest value,  $-1$  to those  $n_3$  observations with lowest value and 0 in between (Simon and Tibshirani, 2012). To simulate the fMRI data  $\mathbf{Y}$  with a brain structure similar to that of real data, we referred to the standard aal template with 116 brain ROIs containing 54277 voxels. The group structure of voxels is based on the ROIs. Voxels correlated with SNPs were generated by latent variable model.

Irrelevant voxels within each ROI were simulated from  $y_t \sim N(\mu_0^t, \sum_0^t)$  where  $\mu_0^t$  and  $\sum_0^t$  are mean and covariance matrix estimated from the  $t$ th ROI in real data.

First, we compared the three parameter optimization criteria by applying them to the simulation study. Two data sets  $\mathbf{X}$  with 200 SNPs and  $\mathbf{Y}$  with 200 voxels were simulated. Both data sets were divided into 20 groups with group size 10. 20 SNPs from 4 groups in set  $\mathbf{X}$  and 30 voxels from 5 groups in set  $\mathbf{Y}$  were set to be correlated. Standard deviation  $\sigma_y = 1$  and  $\sigma_e = 0.3$ . The sample size was 100 and 50 simulations were replicated. The total true positive (TTP) and total discordance (TD: the total false positive plus the total false negative) were used for the comparison. Fig. 2 shows the box-plot of TTP and TD by using different criteria. It can be seen that, by all the methods,  $corr1$  tends to give higher TTP but also introduces higher TD; the use of  $p_{permute}$  can cause the decrease of TD rapidly at the cost of low TTP; and  $corr2$  is a trade-off between these two criteria by keeping a high TTP and decreasing TD to a low level. Hence the last one is chosen as our parameter selection criterion.

Second, we compared the performance of three methods in the simulation. SNP data  $\mathbf{X}$  with 400 SNPs are divided into 20 groups evenly. 4 groups were randomly selected and 15 SNPs randomly selected from these 4 groups were generated to be correlated with voxels, as shown in Fig. 3(a). For the sake of simplicity, 5% voxels randomly from the first 4 regions with totally 214 voxels were associated with SNPs (Fig. 3(e)). Fig. 3(b–d, f–h) shows the results of recovered loading vectors  $\mathbf{u}$  and  $\mathbf{v}$  by CCA- $l1$ , CCA-group and CCA-sparse group methods. A 5-fold cross-validation was used to select the optimal parameters. From Table 3, it can be seen that the CCA-sparse group method can better estimate true  $\mathbf{u}$  and  $\mathbf{v}$  than other two methods. For the SNP data, all three methods can identify true SNPs (60) while CCA- $l1$  has more false positives (36, shown in Fig. 3(b)) than those of CCA-group (20) and CCA-sparse group (4). For the fMRI data, CCA- $l1$  also misses out more true voxels (29 out of 214) in Fig. 3(f) when selecting  $\mathbf{v}$ . The CCA-group, however, can better recover all the

groups with true features. Nevertheless, it selects all features from the group, which results in very high false positive (4097, shown in Fig. 3(g)). CCA-sparse group can not only identify the group structure to find more true variables than CCA-*l1* but also remove those noisy features in the group, leading to the least false positive (113 false voxels versus 912 by CCA-*l1* and 4097 by CCA-group).

Finally, we evaluated the performance of three models with respect to different noise levels, which affected the values of correlation between two data sets. 200 SNPs from 20 genes were simulated in data set  $X$ . The group size was 10. 50 voxels from 5 ROIs were randomly selected to be correlated with 20 SNPs in  $X$ . The sample size was 150. According to the correlation estimate in (Parkhomenko et al., 2009), we simulated different level of correlations between two data sets by changing the standard deviation of the noise variable  $\sigma_e$  from 0.1 to 1 with interval 0.1. From the results in Fig. 4, we can see that when the true correlation increases with the decrease of noise variance, more true variables can be recovered with less total discordance by all three methods. The CCA-group model can recover the most correlated variables but has the highest total discordance. When  $\sigma_e$  is larger than 0.6 (the estimated highest correlation is 0.39), the number of TTP and TD is relatively unchanged. When  $\sigma_e$  decreases, more true variables can be identified by all methods and the TD of CCA-*l1* and CCA-sparse group decreases rapidly.

## 4.2. Experimental results

In this study, subject recruitment and data collection were conducted by The Mind Clinical Imaging Consortium (MCIC). Two types of data (SNP and fMRI) were collected from 208 subjects including 92 schizophrenia patients (age:  $34 \pm 11$ , 22 females) and 116 healthy controls (age:  $32 \pm 11$ , 44 females). All of them were provided written informed consents. Healthy participants were free of any medical, neurological or psychiatric illnesses and had no history of substance abuse. By the clinical interview of patients for DSM IV-TR Disorders (22) or the Comprehensive Assessment of Symptoms and History, patients met criteria for DSM-IV-TR schizophrenia (23). Antipsychotic history was collected as part of the psychiatric assessment.

**4.2.1. fMRI data collection and preprocessing**—fMRI data were collected during a sensor motor task, a block-design motor response to auditory stimulation. The images were acquired on a Siemens3T Trio Scanner and 1.5 T Sonata with echo-planar imaging (EPI) sequences taking parameters (TR = 2000 ms, TE = 30 ms (3.0 T)/40 ms (1.5 T), field of view = 22 cm, slice thickness = 4 mm, 1 mm skip, 27 slices, acquisition matrix =  $64 \times 64$ , flip angle =  $90^\circ$ ). Data were pre-processed with SPM5 (<http://www.fil.ion.ucl.ac.uk/spm>) and were realigned, spatially normalized and resliced to  $3 \times 3 \times 3$  mm, smoothed with a  $10 \times 10 \times 10$  mm<sup>3</sup> Gaussian kernel, and analyzed by multiple regression considering the stimulus and their temporal derivatives plus an intercept term as repressors. Finally the stimulus-on versus stimulus-off contrast images were extracted with  $53 \times 63 \times 46$  voxels and all the voxels with missing measurements were excluded. 116 ROIs were extracted based on the aal brain atlas, which resulted in 41236 voxels left for analysis.

**4.2.2. Genotyping and preprocessing**—A blood sample was obtained for each participant and DNA was extracted. Genotyping for all participants was performed at the Mind Research Network using the Illumina Infinium HumanOmni1-Quad assay covering 1,140,419 SNP loci. Bead Studio was used to make the final genotype calls. PLINK software package (<http://pngu.mgh.harvard.edu/~purcell/plink>) was used to perform a series of standard quality control procedures, resulting in the final dataset spanning 777,635 SNP loci. Each SNP was categorized into three clusters based on their genotype and was represented with discrete numbers: 0 for ‘BB’ (no minor allele), 1 for ‘AB’ (one minor allele) and 2 for ‘AA’ (two minor alleles). SNP with >20% missing data were deleted and missing data were further imputed. SNPs with minor allele frequency <1% were removed. To reflect the influence of genetic variation on brain behavior, SNPs included in top 75 schizophrenia genes listed on the SZGene database (<http://www.szgene.org/>) were selected for the analysis. This procedure yielded 3082 SNPs, which were annotated with 74 genes. There was no SNP found in the remaining one gene.

**4.2.3. Correlation analysis**—The fMRI voxels were grouped based on ROIs while SNPs were grouped by genes. Then we applied our CCA-sparse group method to the analysis of the correlation between two data sets. 208 subjects were randomly divided into two subsets: 150 subjects for training and the remaining 58 ones for testing. In training data set, we fit three models: (i) CCA-group lasso, i.e., using group lasso penalty on genes and ROIs, (ii) CCA- $l_1$ , i.e., only imposing  $l_1$  norm on the effects of all SNPs and voxels, and (iii) CCA-sparse group model, i.e., using both group-level and single feature level regularizations. The optimal parameters were obtained from training data by 5-fold cross validation. The models were estimated as well as the features were selected from the training data using the optimal parameters. Then, these estimated models were applied to test data to predict the correlation between two data sets.

To have a stable feature selection, we performed random sampling from 208 subjects repeatedly for  $B$  times, selected the same proportion of subjects for training and test data sets, and fitted three models on each sub-sample. We assumed that those SNPs and voxels selected more frequently are more valuable for exploring the correlation between two data sets. For each sub-sample,  $\hat{\mathbf{u}}_b, \hat{\mathbf{v}}_b$  were estimated,  $b = 1, 2, \dots, B$ . A measure of feature importance can be computed by frequency of their appearance defined as the selection probability in (15). A set of features with high value will be selected by a cut-off threshold.

$$p_j = \frac{1}{B} \sum_{b=1}^B I(\hat{t}_j^b \neq 0) \quad (15)$$

where  $\hat{t}_j^b$  is taken to be  $\hat{u}_j^b$  (or  $\hat{v}_j^b$ ), i.e., the loading coefficient of sCCA model corresponding to the  $j$ th SNP (or voxel), and  $I(\cdot)$  is indicator function.

By the algorithm in Section 2.2, we can derive several pairs of canonical variates by iteratively implementing the CCA decomposition. Here, we only show the first two pairs of canonical variates.  $B = 50$  replicates were performed. To compare the performance, the test correlations  $\text{cor}(\mathbf{X}_{test}\hat{\mathbf{u}}, \mathbf{Y}_{test}\hat{\mathbf{v}})$  based on the estimated loading vectors  $(\hat{\mathbf{u}}, \hat{\mathbf{v}})$  in training data

by three methods were calculated. The numbers of selected features were also used for comparison. Fig. 5(a) shows that CCA-group method chose the largest number of SNPs and voxels compared to the other two methods while CCA-sparse group method selected the least features due to the double sparsity constraints on loading vectors. As shown in Table 4, CCA-sparse group generally performs better than the other two methods by obtaining higher predict correlations with comparable low variance. The low variance of these correlations demonstrates the robustness of the estimation. The correlation variance by CCA-sparse group method is also comparable with that of the other two methods. In addition, the small training and test correlation difference ( $r_{corr} = 0.146, 0.1331, 0.1471$  for three methods respectively) indicates the advantage of using criterion Eq. (11) for tuning parameters selection, i.e., the stability of the criterion regardless of the method used.

We further discuss the selected features. Choosing SNPs and voxels by ranking the selection probabilities ( $p_j$  in Eq. (14)) with a cut-off threshold 0.3, we summarized those top ranked SNPs and voxels in two pairs of variates by three methods in Fig. 5(b and c). There are a high proportion of overlapped SNPs and voxels selected by these methods. Table 5 lists top ranked SNPs and corresponding genes in two pairs of variates by CCA-sparse group method. The first pair contained 51 SNPs from 16 genes correlated with 756 voxels from 36 ROIs with the average correlation 0.4527,  $p < 0.001$ . 30 SNPs from genes ERBB4 and MAGI2 and 212 voxels from ROIs 7 and 8 were selected by all methods. The other pair had 19 SNPs from 9 genes correlated with 1558 voxels from 29 ROIs, and the average correlation was 0.4292,  $p < 0.001$ . 3 SNPs from MAGI2 and 251 voxels from ROIs 47, 51 were overlapped by all three methods. Table 6 listed the top ranked brain ROIs in the first and second pair of variates respectively. Fig. 6 shows the brain mapping of genomic correlated ROIs. Several brain regions have been reported to be susceptible to schizophrenia from other neuroimaging studies (Shenton et al., 2001; Kumari et al., 2002; Onitsuka et al., 2004; Torrey, 2007; Bellani et al., 2010; Pinault, 2011). They include superior, middle, and medial frontal gyrus, inferior parietal lobule, superior temporal gyrus, thalamus, parahippocampal gyrus, cingulate gyrus, and lingual gyrus, which provided additional evidences that these disease relevant brain regions may be affected by the correlated genomic variations.

**4.2.4. Gene-ROI correlation analysis**—Based on the two pairs of canonical variates, we further explored and verified the correlation between genes and ROIs. We found that some regions of brain might be correlated with a set of genes. For each gene-ROI, gene-gene and ROI-ROI correlation, 10000 permutations (as mentioned in Method section) were performed to test the significance. In Fig. 7(a), the absolute values of significant gene-ROI correlation are  $0.1862 \pm 0.0317$  (mean  $\pm$  SD,  $p < 0.005$ ). As shown in Fig. 8(a), Gene ERBB4, NRG1, MAGI2 and GABRG2 show correlations with some ROIs such as 3, 4, 7, 8 (the index of ROI is defined by the aal template (Tzourio-Mazoyer et al., 2002)) with the correlation ( $\rho = 0.1822, 0.1483, 0.1335, 0.1782$ ,  $p < 0.004$ ) respectively. These ROIs mostly consist of superior, middle, medial front gyrus and pre-central gyrus located at frontal lobe which contains primary motor cortex and has been suspected to have abnormal changes in schizophrenia patients (Honey et al., 2005; Kiehl et al., 2005). ROIs 75, 76, 77, 11 and 13(not shown) are also found to be correlated with gene ERBB4, NRG1, and MAGI2 with

the correlation value of 0.1822, 0.1649 and 0.1541 respectively. These ROIs are mainly located at thalamus (right) which plays a critical role in coordinating the pass of information between brain regions. Many studies show the association between dysfunction of thalamus with schizophrenia (Kiehl and Liddle, 2001; Clinton and Meador-Woodruff, 2004; Sui et al., 2011). These three genes are also correlated with each other, which may have the similar effects on the ROIs. In addition, ROIs 91, 92, 99, 100, 103, 111, 112 are correlated with many genes (*e.g.* GRIN2B, CHL1, ERBB4, NRG1, FOXP2, MAGI1, GABRG2, MAGI2). These ROIs are located at declive of cerebellum and culmen. Their relationship with schizophrenia is not clear yet but several previous publications have reported significant difference of these regions between normal control and schizophrenia patients (Kim et al., 2009) and there are models of schizophrenia which discuss importance of the cerebellum (Andreasen and Pierson, 2008).

In the second pair of canonical variates (Fig. 7(b)), the absolute value of significant gene-ROI correlation is  $0.1713 \pm 0.0308$  (mean  $\pm$  SD,  $p < 0.005$ ). There are two patterns apparent in the figure. The first pattern corresponds to ROIs 40, 43, 44, 45, 46, 47, 48, 49, 50, 51, 52, 53, 54, 55, 56, 67, 68, which are located at occipital lobe including lingual gyrus, posterior cingulate, precuneus, parahippocampagyrus and superior, middle and inferior occipital gyrus as shown in Fig. 8(b). Four genes ERBB4, CHL1, GRM3 and MAGI2 were significantly correlated with these ROIs at 0.1776, 0.1356, 0.1762 and 0.1429,  $p < 0.005$  respectively. These critical regions have been widely studied and are shown to have potential relationship with schizophrenia (Kiehl et al., 2005; Fransson and Marrelec, 2008; Kim et al., 2009). The second pattern are ROIs 1, 24, 57 from precentral gyrus at parietal lobe corresponding to genes SNPA29, MTHFR, SLC1A and NRG1 with the correlation of 0.2251, 0.1899, 0.1816 and 0.21 respectively.

## 5. Discussion and conclusion

In this paper, we proposed a novel method to explore the relationship between genomic data and fMRI brain imaging data by considering the group effects of the variables in the data. We introduced the group sparse CCA method and the numerical implementation based on the regularized SVD and block coordinate decent algorithm. The performance of group sparse CCA model was compared with other sCCA models in the simulation study, showing that our group sparse CCA method could better recover the true correlations with lower false positive and total discordance. Then we applied the method to correlation analysis between the SNP data and fMRI imaging data. Two pairs of canonical variates with significant correlations were identified. There are 5 pathways implied by these identified genes, which may involve with the biological processes related to schizophrenia. The SZ-risk genes correlated with brain regions have also been reported to be susceptible to schizophrenia, which further validates the results of our method.

We identified two pairs of canonical variates with significant correlations and further verified those linked components between gene-ROI with significant pair-wise correlation. Those identified genes are from the preselected 74 genes from  $\alpha$  database which were reported to be susceptible to SZ. Some of these genes may not directly show group difference in our fMRI data. Therefore, in this section, we focus on those group-

discriminating genetic factors associated with motor response task-related brain function disruption in schizophrenia. We use the difference of minor allele frequency (MAF) between case and control group to evaluate the effects of the SNPs in the gene correlated with ROI component. Significantly higher MAF in case group indicates the positive discriminating effect while lower MAF indicates the negative discriminating effect. Pathway analysis of these important genes is performed through Ingenuity Pathway Analysis (IPA: Ingenuity Systems, <http://www.ingenuity.com>). We found those selected genes involving in 5 pathways as shown in Fig. 9.

1. Neuregulin signalling pathway: NRG1 and ERBB4. These two genes have significant effects in both pairs of canonical variates. It has been reported that the NRG1-ERBB4 modulates some plausible neurobiological mechanisms, i.e., neuronal plasticity in human brain which may be altered in SZ (Buonanno, 2010). In our results, SNPs 'rs10090544', 'rs16878394', 'rs7843384' (NRG1) and 'rs16847732', 'rs16847769', 'rs2008506' in ERBB4 exhibit higher MAFs in patients; 'rs2466063' (NRG1) and the MAFs in 'rs11903508', 'rs16846111', 'rs16846352' in ERBB4 are lower in patients indicating negative effect.
2. Glutamate receptor signalling: GRIN2B, GRM3 and SLC1A2. These genes are involved in encoding ionotropic glutamate receptors and histamine receptors which may regulate the neurotransmitter transmission associated with SZ (Bishop et al., 2005). Several researches have supported the possibility of GRIN2B conferring the susceptibility to schizophrenia (Ohtsuki et al., 2001; Qin et al., 2005). GRM3 may be involved in the pathophysiology of schizophrenia, and its associated cognitive impairment, especially at prefrontal and hippocampal regions (Egan et al., 2004). In our results, 'rs2299218', 'rs802432' in GRM3 and 'rs2284425' in GRIN2B have lower MAFs in patients showing negative effects while 'rs3026164' (GRIN2B), 'rs802425' (GRM3) and 'rs11033095' (SLC1A2) have positive effects on the disease.
3. GABA receptor signalling: GABRB2 and GABRG2. GABA neuronal dysfunction has been found to associate with cognitive impairment of schizophrenia. We identified a SNP 'rs153303' (GABRB2) with lower MAF in schizophrenia, showing negative effect while 'rs211037' (GABRG2) has higher MAF in patients with a possible opposite influence.
4. Calcium signalling pathway: HTR3A, ERBB4 and PPP3CC. Recent biochemical research supports that schizophrenia may be produced by alterations in various intracellular molecules as long as these alterations lead to abnormal functioning of some central intracellular regulatory pathways and one of unifying elements of molecular changes in schizophrenia is their association with potential altered Ca<sup>2+</sup> signalling (Lidow, 2003). 'rs1150219' (HTR3A) and 'rs7010861' (PPP3CC) both show higher MAFs in patients in the results.
5. Tight junction: MAGI1, MAGI2. These genes attend the coding for MAGI proteins which may relate to SZ by influencing the development and communication between nerve cells (Karlsson et al., 2012). In our results, 'rs9311944', 'rs6809559' (MAGI1), and 'rs10230275', 'rs10229284', 'rs1330490'

(MAGI2) present higher MAFs in patients while 'rs9821646', 'rs2061937' (MAGI1), 'rs10230752', 'rs10953782' (MAGI2) exhibit lower MAFs in patients.

The tuning parameters selection criteria used in this work are also important. The performance of Waaijenborg et al.'s criterion is better than those of other two criteria in terms of feature selection accuracy, which is consistent with the conclusion from (Waaijenborg and Zwinderman, 2009). However, these parameter selection criteria are not directly used for the purpose of feature selection; a better criterion can be developed. In addition, the selection of tuning parameters is via cross-validation, relying on large training data, which is often not available. Alternatively, Bayesian information criterion and Akaike information criterion (Posada and Buckley, 2004) can be used for the optimal parameter selection.

## Acknowledgments

This work is partially supported by both NSF and NIH. It is also supported by Shanghai Eastern Scholarship Program.

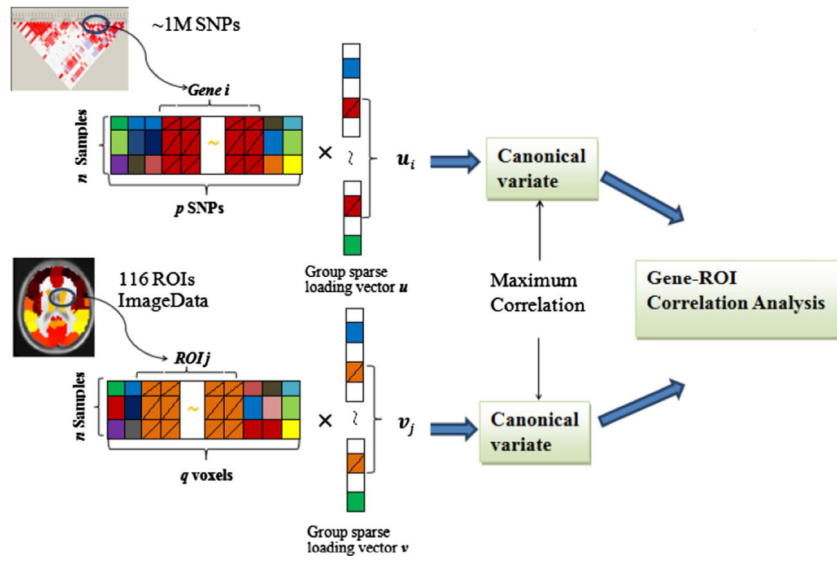
## References

- Abecasis GR, Burt RA, et al. Genomewide scan in families with schizophrenia from the founder population of Afrikaners reveals evidence for linkage and uniparental disomy on chromosome 1. *Am. J. Hum. Genet.* 2004; 74(3):403–417. [PubMed: 14750073]
- Andreasen NC, Pierson R. The role of the cerebellum in schizophrenia. *Biol. Psychiatr.* 2008; 64(2): 81–88.
- Badner JA, Gershon ES. Meta-analysis of whole-genome linkage scans of bipolar disorder and schizophrenia. *Mol. Psychiatr.* 2002; 7(4):405–411.
- Bellani M, Ferro A, et al. The potential role of the parietal lobe in schizophrenia. *Epidemiol. Psychiatr. Soc.* 2010; 19(2):118–119. [PubMed: 20815295]
- Bishop JR, Ellingrod VL, et al. Association between the polymorphic GRM3 gene and negative symptom improvement during olanzapine treatment. *Schizophr. Res.* 2005; 77(2–3):253–260. [PubMed: 15913960]
- Boutte, D.; Liu, J. Sparse canonical correlation analysis applied to fMRI and genetic data fusion.; *IEEE International Conference on Bioinformatics and Biomedicine*; 2010.
- Buonanno A. The neuregulin signaling pathway and schizophrenia: from genes to synapses and neural circuits. *Brain Res. Bull.* 2010; 83(3–4):122–131. [PubMed: 20688137]
- Callicott JH, Straub RE, et al. Variation in DISC1 affects hippocampal structure and function and increases risk for schizophrenia. *Proc. Natl. Acad. Sci. USA.* 2005; 102(24):8627–8632. [PubMed: 15939883]
- Chen X, Liu H. An efficient optimization algorithm for structured sparse CCA, with applications to eQTL Mapping. *Stat. Biosci.* 2012; 4:3–26.
- Chen J, Bushman FD, et al. Structure-constrained sparse canonical correlation analysis with an application to microbiome data analysis. *Biostatistics.* 2013; 14(2):244–258. [PubMed: 23074263]
- Clinton SM, Meador-Woodruff JH. Thalamic dysfunction in schizophrenia: neurochemical, neuropathological, and in vivo imaging abnormalities. *Schizophr. Res.* 2004; 69(2–3):237–253. [PubMed: 15469196]
- Egan MF, Straub RE, et al. Variation in GRM3 affects cognition, prefrontal glutamate, and risk for schizophrenia. *Proc. Natl. Acad. Sci. USA.* 2004; 101(34):12604–12609. [PubMed: 15310849]
- Fransson P, Marrelec G. The precuneus/posterior cingulate cortex plays a pivotal role in the default mode network: evidence from a partial correlation network analysis. *Neuroimage.* 2008; 42(3): 1178–1184. [PubMed: 18598773]

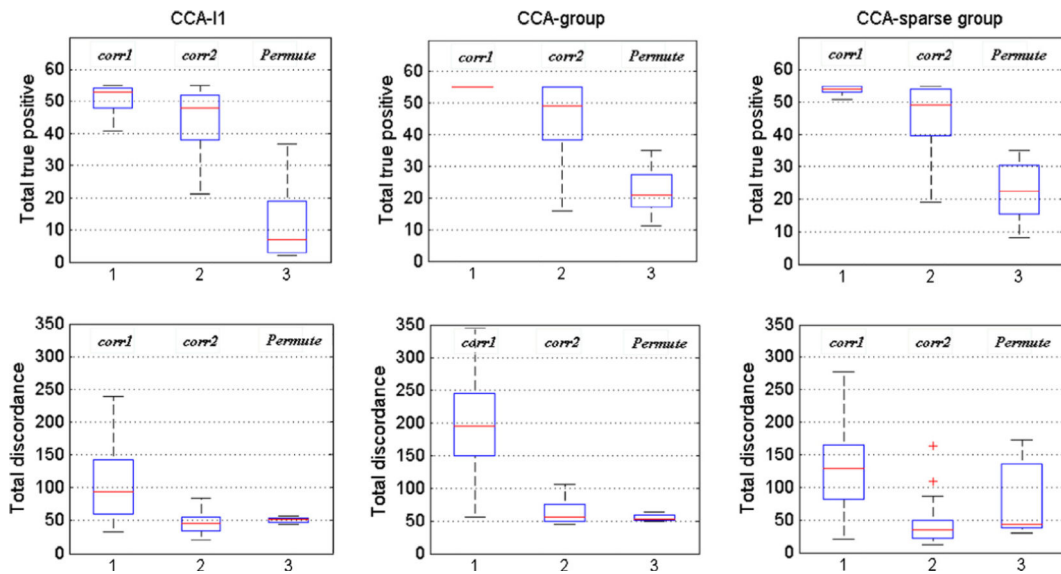


- Friedman, J.; Hastie, T., et al. A Note on the Group Lasso and a Sparse Group Lasso. 2010. <<http://www.arxiv.org/pdf/1001.0736>>
- Friedman J, Hastie T, et al. Regularization paths for generalized linear models via coordinate descent. *J. Stat. Softw.* 2010b; 33(1):1–22. [PubMed: 20808728]
- Hamid JS, Hu P, et al. Data integration in genetics and genomics: methods and challenges. *Hum. Genom., Proteom.* 2009
- Honey GD, Pomarol-Clotet E, et al. Functional dysconnectivity in schizophrenia associated with attentional modulation of motor function. *Brain.* 2005; 128(Pt 11):2597–2611. [PubMed: 16183659]
- Hotelling H. Relations between two sets of variates. *Biometrika.* 1936; 28:321–377.
- Jansma JM, Ramsey NF, et al. Working memory capacity in schizophrenia: a parametric fMRI study. *Schizophr. Res.* 2004; 68(2–3):159–171. [PubMed: 15099600]
- Karlsson R, Graae L, et al. MAGII copy number variation in bipolar affective disorder and schizophrenia. *Biol. Psychiatr.* 2012; 71(10):922–930.
- Kiehl KA, Liddle PF. An event-related functional magnetic resonance imaging study of an auditory oddball task in schizophrenia. *Schizophr. Res.* 2001; 48(2–3):159–171. [PubMed: 11295369]
- Kiehl KA, Stevens MC, et al. Abnormal hemodynamics in schizophrenia during an auditory oddball task. *Biol. Psychiatr.* 2005; 57(9):1029–1040.
- Kim DI, Mathalon DH, et al. Auditory oddball deficits in schizophrenia: an independent component analysis of the fMRI multisite function BIRN study. *Schizophr. Bull.* 2009; 35(1):67–81. [PubMed: 19074498]
- Kumari V, Gray JA, et al. Procedural learning in schizophrenia: a functional magnetic resonance imaging investigation. *Schizophr. Res.* 2002; 57(1):97–107. [PubMed: 12165380]
- Le Cao KA, Rossouw D, et al. A sparse PLS for variable selection when integrating omics data. *Stat. Appl. Genet. Mol. Biol.* 2008; 7 Article 35.
- Le Cao KA, Martin PGP, et al. Sparse canonical methods for biological data integration: application to a cross-platform study. *BMC Bioinform.* 2009; 10:34.
- Li X, Branch CA, et al. fMRI study of language activation in schizophrenia, schizoaffective disorder and in individuals genetically at high risk. *Schizophr. Res.* 2007; 96(1–3):14–24. [PubMed: 17719745]
- Lidow MS. Calcium signaling dysfunction in schizophrenia: a unifying approach. *Brain Res. Brain Res. Rev.* 2003; 43(1):70–84. [PubMed: 14499463]
- Lin D, Zhang J, et al. Group sparse canonical correlation analysis for genomic data integration. *BMC Bioinform.* 2013; 14:245.
- Liu J, Huang J, et al. Incorporating group correlations in genome-wide association studies using smoothed group Lasso. *Biostatistics.* 2013; 14(2):205–219. [PubMed: 22988281]
- Meda SA, Bhattarai M, et al. An fMRI study of working memory in first-degree unaffected relatives of schizophrenia patients. *Schizophr. Res.* 2008; 104(1–3):85–95. [PubMed: 18678469]
- Ng B, Abugharbieh R. Generalized sparse regularization with application to fMRI brain decoding. *Inform. Process. Med. Imag.* 2011; 22:612–623.
- Ohtsuki T, Sakurai K, et al. Mutation analysis of the NMDAR2B (GRIN2B) gene in schizophrenia. *Mol. Psychiatr.* 2001; 6(2):211–216.
- Onitsuka T, Shenton ME, et al. Middle and inferior temporal gyrus gray matter volume abnormalities in chronic schizophrenia: an MRI study. *Am. J. Psychiatr.* 2004; 161(9):1603–1611. [PubMed: 15337650]
- Parkhomenko E, Trichler D, et al. Sparse canonical correlation analysis with application to genomic data integration. *Stat. Appl. Genet. Mol. Biol.* 2009; 8(1)
- Pinault D. Dysfunctional thalamus-related networks in schizophrenia. *Schizophr. Bull.* 2011; 37(2): 238–243. [PubMed: 21307040]
- Porteous DJ, Thomson P, et al. The genetics and biology of DISC1 – an emerging role in psychosis and cognition. *Biol. Psychiatr.* 2006; 60(2):123–131.

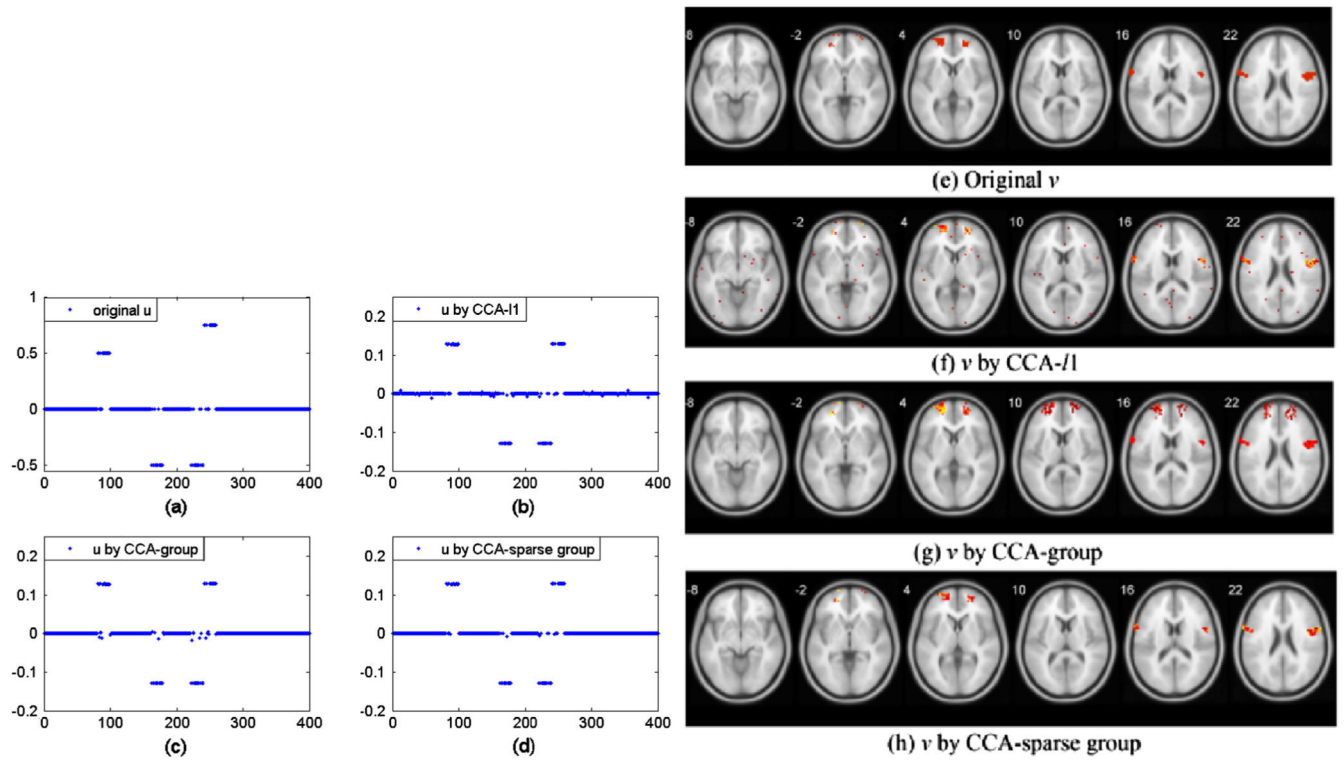
- Posada D, Buckley TR. Model selection and model averaging in phylogenetics: Advantages of akaike information criterion and Bayesian approaches over likelihood ratio tests. *Syst. Biol.* 2004; 53(5): 793–808. [PubMed: 15545256]
- Qin S, Zhao X, et al. An association study of the N-methyl-D-aspartate receptor NR1 subunit gene (GRIN1) and NR2B subunit gene (GRIN2B) in schizophrenia with universal DNA microarray. *Eur. J. Hum. Genet.* 2005; 13(7):807–814. [PubMed: 15841096]
- Shenton ME, Dickey CC, et al. A review of MRI findings in schizophrenia. *Schizophr. Res.* 2001; 49(1–2):1–52. [PubMed: 11343862]
- Simon N, Tibshirani R. Standardization and the group lasso penalty. *Stat. Sin.* 2012; 22:983–1001.
- Sui J, Pearlson G, et al. Discriminating schizophrenia and bipolar disorder by fusing fMRI and DTI in a multimodal CCA+ joint ICA model. *Neuroimage.* 2011; 57(3):839–855. [PubMed: 21640835]
- Sutrala SR, Goossens D, et al. Gene copy number variation in schizophrenia. *Schizophr. Res.* 2007; 96(1–3):93–99. [PubMed: 17826036]
- Szyck GR, Munte TF, et al. Audiovisual integration of speech is disturbed in schizophrenia: an fMRI study. *Schizophr. Res.* 2009; 110(1–3):111–118. [PubMed: 19303257]
- Torrey EF. Schizophrenia and the inferior parietal lobule. *Schizophr. Res.* 2007; 97(1–3):215–225. [PubMed: 17851044]
- Tyekucheva S, Marchionni L, et al. Integrating diverse genomic data using gene sets. *Genome Biol.* 2011; 12(10):R105. [PubMed: 22018358]
- Tzourio-Mazoyer N, Landeau B, et al. Automated anatomical labeling of activations in SPM using a macroscopic anatomical parcellation of the MNI MRI single-subject brain. *Neuroimage.* 2002; 15(1):273–289. [PubMed: 11771995]
- Waaijenborg S, Zwinderman AH. Correlating multiple SNPs and multiple disease phenotypes: penalized non-linear canonical correlation analysis. *Bioinformatics.* 2009; 25(21):2764–2771. [PubMed: 19689958]
- Waaijenborg S, Hamer PCVDW, et al. Quantifying the association between gene expressions and DNA-Markers by penalized canonical correlation analysis. *Stat. Appl. Genet. Mol. Biol.* 2008; 7(1)
- Wiley HS. Integrating multiple types of data for signaling research: challenges and opportunities. *Sci. Signal.* 2011; 4(160):pe9. [PubMed: 21325205]
- Wilson GM, Flibotte S, et al. DNA copy-number analysis in bipolar disorder and schizophrenia reveals aberrations in genes involved in glutamate signaling. *Hum. Mol. Genet.* 2006; 15(5):743–749. [PubMed: 16434481]
- Witten DM, Tibshirani RJ. Extensions of sparse canonical correlation analysis with applications to genomic data. *Stat. Appl. Genet. Mol. Biol.* 2009; 8(1)
- Witten DM, Tibshirani R, et al. A penalized matrix decomposition, with applications to sparse principal components and canonical correlation analysis. *Biostatistics.* 2009; 10(3):515–534. [PubMed: 19377034]
- Wu TT, Chen YF, et al. Genome-wide association analysis by lasso penalized logistic regression. *Bioinformatics.* 2009; 25(6):714–721. [PubMed: 19176549]



**Fig. 1.** A schematic illustration of combining both fMRI image and SNP data by the group sparse CCA model to identify correlated genes and ROIs.

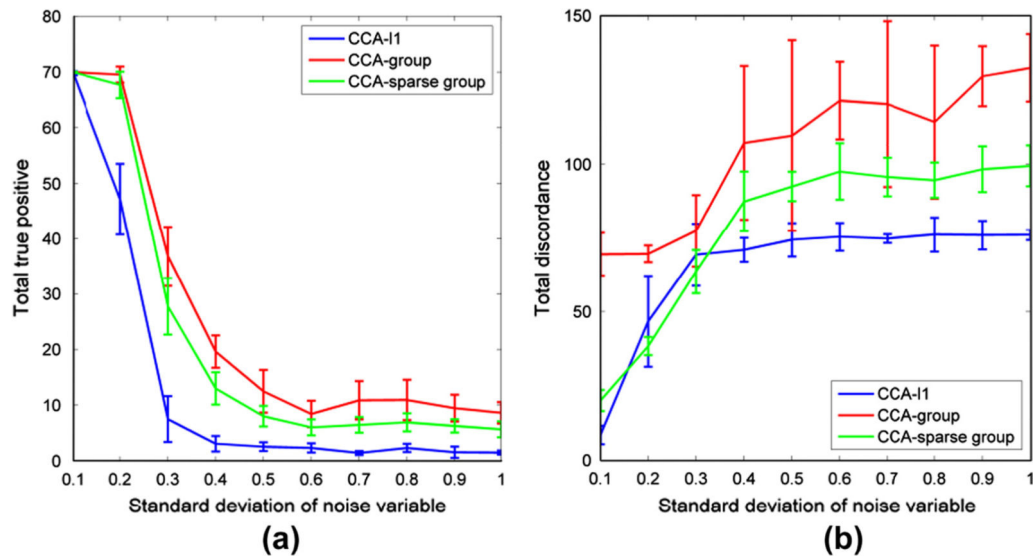


**Fig. 2.** The total true positive number (TPP) and total discordance (TD) of estimating correlations using CCA-l1, CCA-group and CCA-sparse group methods using three different parameter selection criteria: corr1 (Eq. (10)) proposed by Parkhomenko et al. (2009), corr2 (Eq. (11)) proposed by Waaijenborg et al. (2008) and permutation based criterion (Eq. (12)) proposed by Witten and Tibshirani (2009).

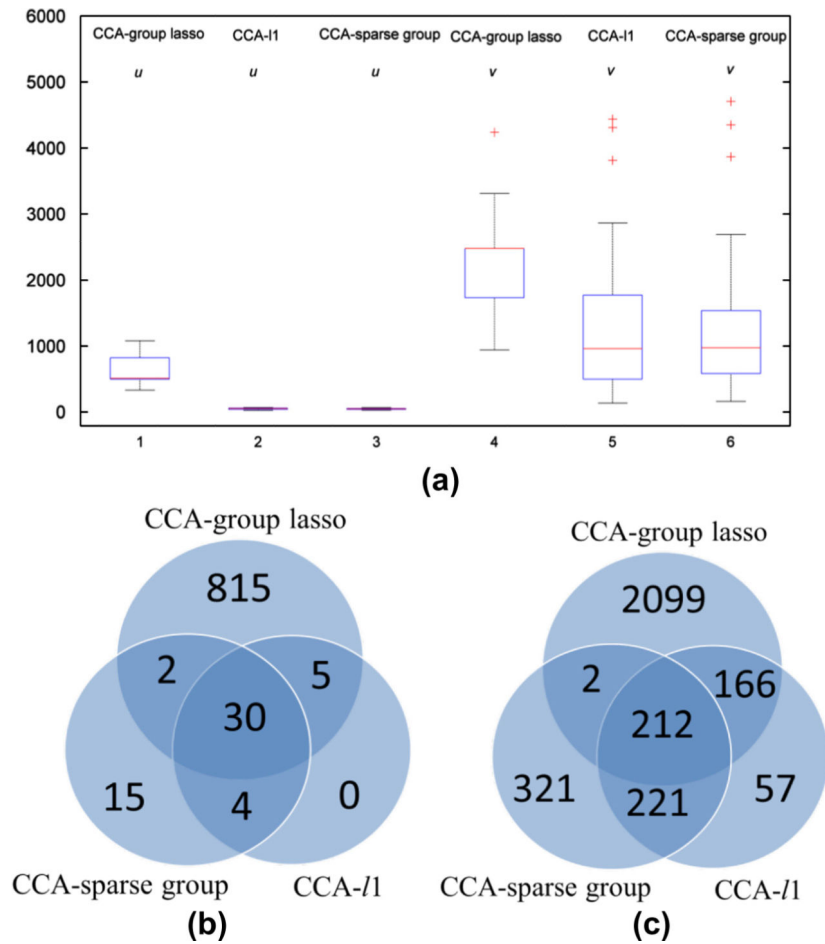


**Fig. 3.**

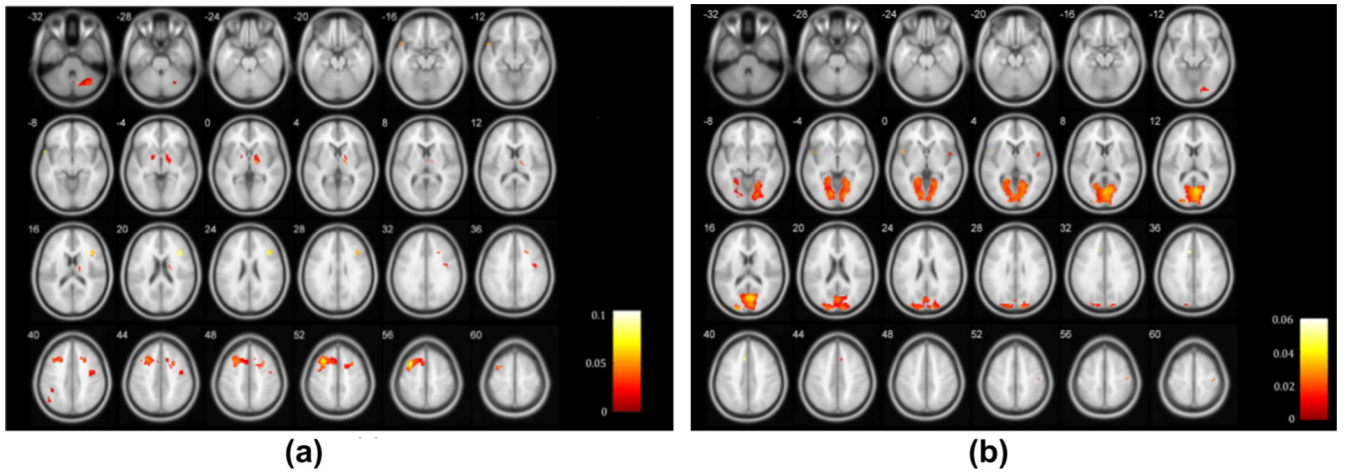
A comparison of the performance of group sparse CCA and three sCCA methods. (a) True  $u$ ; (b–d)  $u$  recovered by CCA-l1, CCA-group and CCA-sparse group respectively. (e) True  $v$ ; (f–h)  $v$  recovered by three sCCA methods.

**Fig. 4.**

A comparison of three methods for different correlation level influenced by noise. (a) The value of total true positive obtained by three methods when the standard deviation of noise increases from 0.1 to 1, showing that the highest correlation of true variables between two data sets is within the range from 0.958 to 0.18. (b) The total discordance by three methods when the standard deviation of noise changes from 0.1 to 1.

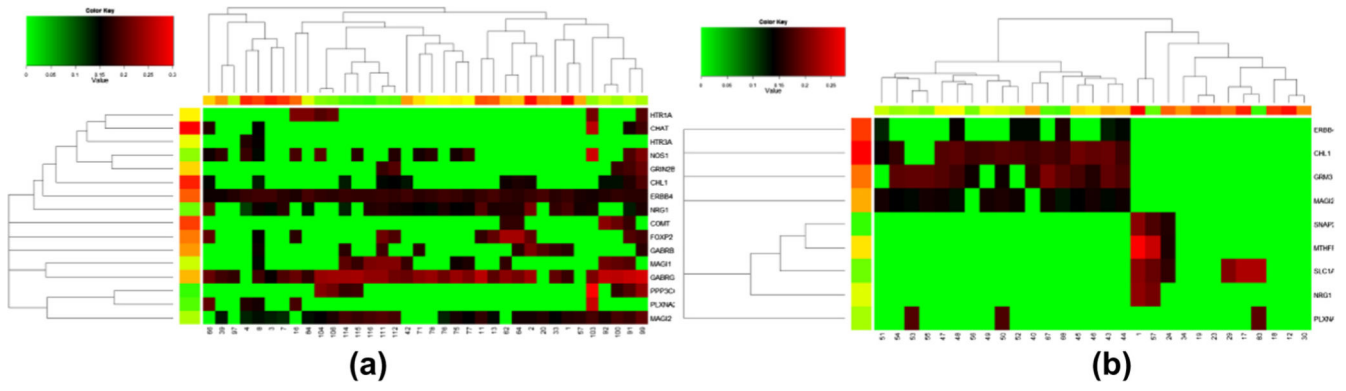


**Fig. 5.** The comparison of the number of features in loading vectors  $u$  and  $v$  selected by three methods. (a) The number of SNPs and voxels selected by three methods with 50 sub-samplings. (b–c) Show the top ranked SNPs (b) and voxels (c) in the first pair of canonical variates with cut-off threshold = 0.3 selection probability.

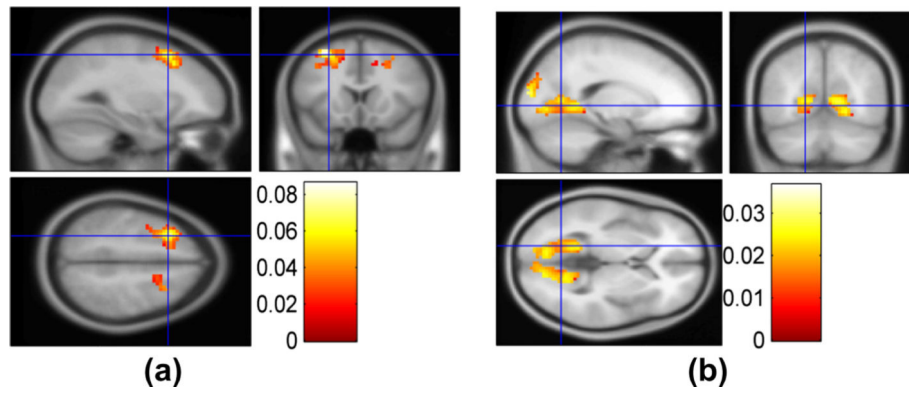


**Fig. 6.** Maps showing regions correlated with genetic factors in the first (a) and second (b) canonical variates.



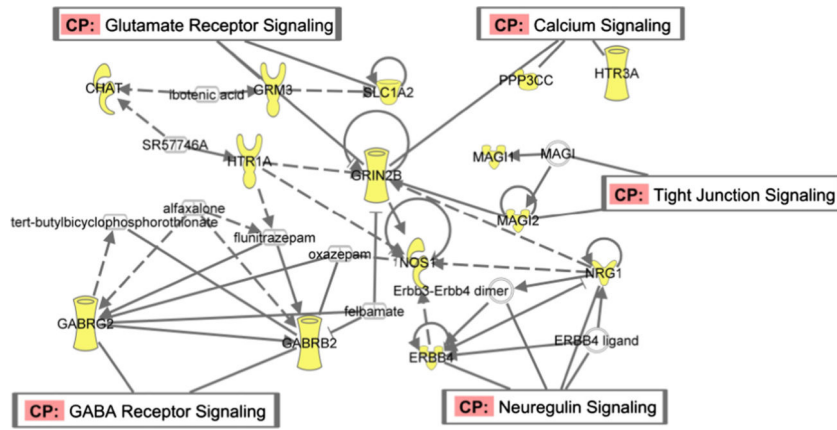


**Fig. 7.** Heatmap of the gene-ROI correlations derived from the first (a) and second (b) pairs of canonical variates.



**Fig. 8.**

(a) ROIs 3, 4, 7, 8 which are identified to be mainly correlated with genes ERBB4, NRG1, MAGI1 and GABRG2. (b) ROIs 40, 43, 44, 45, 46, 47, 48, 49, 50, 51, 52, 53, 54, 55, 56, 67, 68 which are found to correlate with Gene ERBB4, CHL1, GRM3 and MAGI2. Note: the index of ROI is given by the aal template (Tzourio-Mazoyer et al., 2002).



**Fig. 9.** Functional network built based on the genes identified within the two pairs of canonical variates.

**Table 1**

The iterative algorithm of group sparse CCA.

---

1. Initialize  $\mathbf{u}^0$  and  $\mathbf{v}^0$  by traditional CCA decomposition,  $\|\mathbf{u}^0\|_2 = 1$ ,  $\|\mathbf{v}^0\|_2 = 1$ .
2. Solve  $\mathbf{U}^j, \mathbf{V}^j$  using the following iterations until it convergence:
  - (a) Fix  $\mathbf{v} = \mathbf{v}^{j-1}$ ,  $\mathbf{u}^j \leftarrow \arg \min_{\mathbf{u}, d} \| \mathbf{K} - d\mathbf{u}\mathbf{v}^t \|_F^2 + \Psi(\mathbf{u}) \quad s.t. \quad \|\mathbf{u}\|_2 = 1$
  - (b) Fix  $\mathbf{u} = \mathbf{u}^j$ ,  $\mathbf{v}^j \leftarrow \arg \min_{\mathbf{v}, d} \| \mathbf{K} - d\mathbf{u}\mathbf{v}^t \|_F^2 + \Phi(\mathbf{v}) \quad s.t. \quad \|\mathbf{v}\|_2 = 1$
  - (c)  $d^j \leftarrow \text{tr}(\mathbf{K}\mathbf{v}^j(\mathbf{u}^j)^t)$  or  $\text{tr}(\mathbf{K}\mathbf{u}^j(\mathbf{v}^j)^t)$
3. Update the remaining matrix  $\mathbf{K} \leftarrow \mathbf{K} - \text{tr}(\mathbf{K}\mathbf{v}\mathbf{u}^t)\mathbf{u}\mathbf{v}^t$ ; go to Step (1) to obtain the next pair of loading vectors  $(\mathbf{u}, \mathbf{v})$ .

---

**Table 2**

Block coordinate decent for group sparse CCA.

---

**Input:** iteration step  $j$ ,  $\mathbf{u}^j$ ,  $\mathbf{v}^j$ ,  $\|\mathbf{u}^j\|_2 = 1$ ,  $\|\mathbf{v}^j\|_2 = 1$ ,  $\lambda_1$ ,  $\tau_1$ .

**Output:**  $\mathbf{u}^{j+1}$

Solve  $\mathbf{u}^j$  using block coordinate decent until it convergence:

1. For each group  $k = 0$  to  $L$
  2.  $\text{Soft}_k(\mathbf{Kv}) = S(\|\mathbf{Kv}\|_k, \tau_1)$ , where  $S(\cdot)$  is the soft-thresholding function.
  3. If  $\|\text{Soft}_k(\mathbf{Kv})\|_2 \leq \lambda_1$  then  $u_k^{j+1} = 0$ .
  4. Else
  5.  $Sg_k(\mathbf{Kv}) = \frac{1}{2} \left[ \text{Soft}_k(\mathbf{Kv}) - \lambda_1 \frac{\text{Soft}_k(\mathbf{Kv})}{\|\text{Soft}_k(\mathbf{Kv})\|_2} \right]$ .
  6. Update  $u_k^{j+1} = \frac{Sg_k(\mathbf{Kv})}{\|\text{Soft}_k(\mathbf{Kv})\|_2}$ .
  7. End if
  8. End for
  9. Update  $\mathbf{u}^{j+1} = \frac{[Sg^1(\mathbf{Kv}), Sg^2(\mathbf{Kv}), \dots, Sg^k(\mathbf{Kv})]}{\|\mathbf{u}^{j+1}\|_2}$  to make  $\|\mathbf{u}^{j+1}\|_2 = 1$ .
  10. Repeat (1-9), until  $\|\mathbf{u}^{j+1} - \mathbf{u}^j\|_2 \leq \varepsilon$ , else  $\mathbf{u}^j = \mathbf{u}^{j+1}$ .
-

**Table 3**

The count of true positive (TP), false positive (FP) and discordance by three methods.

<b>Model</b>	<b>TP (voxels/SNPs)</b>	<b>FP (voxels/SNPs)</b>	<b>Discordance (voxels/SNP)</b>
CCA-1	185/60	912/36	941/36
CCA-group	214/60	4097/20	4097/20
CCA-sparse group	200/60	113/4	127/4

**Table 4**

Correlation coefficients based on the variants estimated by three models in the test data.

	<b>CCA-group lasso</b>	<b>CCA-<i>l1</i></b>	<b>CCA-sparse group</b>
<i>Test sample correlation (mean <math>\pm</math>SD)</i>			
1st pair	0.4497 $\pm$ (0.043)	0.4235 $\pm$ (0.057)	0.4527 $\pm$ (0.048)
2nd pair	0.4137 $\pm$ (0.034)	0.3959 $\pm$ (0.045)	0.4292 $\pm$ (0.033)

**Table 5**

The list of genes selected in the first and second pair of canonical variates.

Canonical variate	Gene ID	SNPs number
1st pair	CHAT CHL1 COMT ERBB4 FOXP2 GABRB2 GABRG2 GRIN2B HTR1A HTR3A MAGI1 MAGI2 NOS1 NRG1 PLXNA2 PPP3CC	51
2nd pair	CHL1 ERBB4 GRM3 MAGI2 MTHFR NRG1 PLXNA2 SLC1A2 SNAP29	19



**Table 6**

The correlated brain ROIs selected in two pairs of canonical variates.

<b>1st pair</b>			<b>2nd pair</b>		
<b>Brain region</b>	<b>Brodmann area</b>	<b>L/R volume (cm<sup>3</sup>)</b>	<b>Brain region</b>	<b>Brodmann area</b>	<b>L/R volume (cm<sup>3</sup>)</b>
Superior Frontal Gyrus	6, 8	3.1/0.1	Lingual Gyrus	17, 18, 19	5.4/6.5
Inferior Parietal Lobule	40	0.6/*	Posterior Cingulate	18, 23, 30, 31	2.3/2.9
Medial Frontal Gyrus	6, 8, 32	0.7/0.3	Parahippocampal Gyrus	18, 19, 30	1.2/0.4
Precentral Gyrus	6	0.3/1.2	Precuneus	19, 31	0.7/0.8
Middle Frontal Gyrus	6, 8	4.5/0.6	Superior Temporal Gyrus	22	*/0.1
Thalamus	*	0.1/0.4	Cingulate Gyrus	32	0.5/0.2
Cingulate Gyrus	6	0.3/1.2	Cuneus	17, 18, 19, 23, 30	8.1/9.3
Culmen	*	0.1/0.5	Middle Occipital Gyrus	18	0.8/0.1
Declive	*	0.2/0.6	Insula	13	0.1/0.2