# Bridging Imaging, Genetics, and Diagnosis in a Coupled Low-Dimensional Framework

Sayan Ghosal[1]([✉]), Qiang Chen[2], Aaron L. Goldman[2], William Ulrich[2], Karen F. Berman[3], Daniel R. Weinberger[2,4], Venkata S. Mattay[2,4], and Archana Venkataraman[1]

[1] Department of Electrical and Computer Engineering, Johns Hopkins University, Baltimore, USA
sghosal3@jhu.edu
[2] Lieber Institute for Brain Development, Baltimore, USA
[3] Clinical and Translational Neuroscience Branch, NIMH, NIH, Bethesda, USA
[4] Department of Neurology and Radiology, Johns Hopkins School of Medicine, Baltimore, USA

**Abstract.** We propose a joint dictionary learning framework that couples imaging and genetics data in a low dimensional subspace as guided by clinical diagnosis. We use a graph regularization penalty to simultaneously capture inter-regional brain interactions and identify the representative set anatomical basis vectors that span the low dimensional space. We further employ group sparsity to find the representative set of genetic basis vectors that span the same latent space. Finally, the latent projection is used to classify patients versus controls. We have evaluated our model on two task fMRI paradigms and single nucleotide polymorphism (SNP) data from schizophrenic patients and matched neurotypical controls. We employ a ten fold cross validation technique to show the predictive power of our model. We compare our model with canonical correlation analysis of imaging and genetics data and random forest classification. Our approach shows better prediction accuracy on both task datasets. Moreover, the implicated brain regions and genetic variants underlie the well documented deficits in schizophrenia.

## 1 Introduction

Neuropsychiatric disorders, such as autism and schizophrenia are hereditary, which suggests a genetic underpinning. These disorders are characterized by behavioral and cognitive deficits linked to atypical neural functioning. Identifying the brain mechanisms through which the genomes confer risk is essential to find targeted biomarkers for these disorders. Functional MRI (fMRI) and Single Neucleotipe Polymorphisms (SNPs) are the most common modalities used to study brain activity and genetic variation, respectively. However, integrating them is hard due to large dimensionality and the complex interactions between them.

Prior work in imaging-genetics can be grouped into three general categories. The first are multivariate regression methods [8] that uses SNPs as feature vectors and the imaging phenotype as the response variables in penalized least square setting. Some of these methods also induce structured sparsity both at the SNP level and gene levels to find unique interactions between the imaging and the genetic components. However, they do not consider inter-regional brain interactions or the impact of diagnosis. The second category uses canonical correlation analysis (CCA) to maximize the correlation between linear projections of the imaging and genetics data [5]. Once again these unsupervised methods do not incorporate the clinical factor, so the implicated features may not be related to the disease. Finally, the recent approach of [1] develops a probabilistic framework that incorporates imaging, genetics and diagnosis. Specifically, they consider the imaging data as an intermediate phenotype between genetics and disease. However, this method cannot identify genetic variants associated with the disease that do not also express themselves in the imaging data.

In contrast to prior work, we propose coupled dictionary learning framework to bridge the three data domains. Our model assumes that the imaging and genetics data share a joint latent space. The shared projection coefficients are used as a low-dimensional feature vector to predict diagnosis. We couple the imaging, genetics and diagnosis terms in a regularized optimization framework. We use alternating minimization to estimate the unknown dictionary atoms, projection coefficients, and regression weights. The coupling between these variables overcomes the drawbacks of previous methods and yields better diagnosis prediction in a nested cross validation setting. We validate our framework on a population study of schizophrenia and compare our model with standard machine learning baselines. Our framework achieves the best classification accuracy while finding the interpretable and clinically relevant biomarkers.

## 2    Joint Modeling of Imaging, Genetics, and Diagnosis

Figure 1 presents a overview of our joint modelling approach. Let $M$ be the total number of subjects in our study. Our input data for each subject $m$ consists of fMRI activation maps $\mathbf{f_m}$, SNP variants $\mathbf{g_m}$, and binary disease diagnosis $y_m$. We model the fMRI and SNPs in a parallel dictionary learning framework, where the matrices $\mathbf{A}$ and $\mathbf{B}$ contain the associated dictionary elements. The projection onto these dictionaries is controlled by the latent vector $\mathbf{z}_m$ for both imaging and genetics. Likewise, classification is performed using the projection vector $\mathbf{z}_m$. This joint optimization method allows us to learn the related basis features of both imaging and genetics that are associated with the disease.

**Generative Model for Imaging:** Mathematically, let $N$ denotes the total number of ROIs in the brain. The imaging data has dimensionality $\mathbf{f}_m \in \mathbb{R}^{N \times 1}$. We assume that this data can be represented by a sparse set of anatomical basis vectors that lie in a lower dimensional latent space:

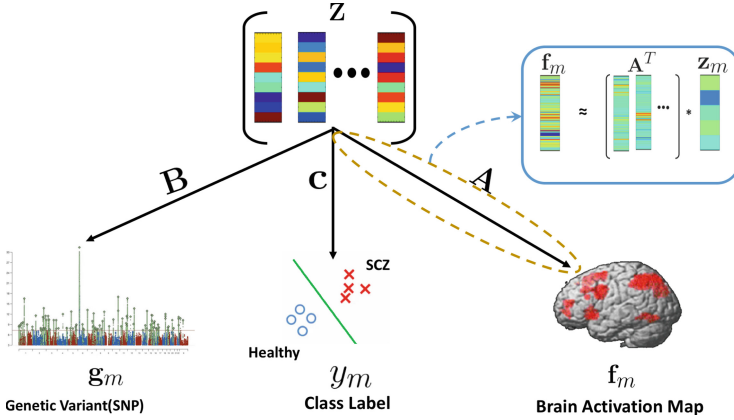$$\mathbf{f}_m \approx \mathbf{A}^T \mathbf{z}_m \quad \text{s.t. } \mathbf{A}\mathbf{A}^T = \mathbf{I}$$

**Fig. 1.** The generative process linking imaging ($\mathbf{f}_m$), genetics ($\mathbf{g}_m$), and diagnosis ($y_m$). The imaging model is shown as a linear projection. The genetics model parallels this, whereas the classification is based on a logistic regression.

where rows of $\mathbf{A} \in \mathbb{R}^{d \times N}$ correspond to the basis vectors (common across the group) and $\mathbf{z}_m \in \mathbb{R}^{d \times 1}$ is the subject specific projection to the latent space. The orthogonality constraint ensures that the basis vectors capture different facets of the data. We also introduce a graph based regularizer on $\mathbf{A}$ so that highly correlated brain regions play a similar role in projection:

$$\frac{\lambda_1}{2} Tr(\mathbf{A}\mathbf{L}\mathbf{A}^T) = \frac{\lambda_1}{2} \sum_{(i,j)} \mathbf{W}_{ij} ||\mathbf{a}_i - \mathbf{a}_j||_2^2$$

where $\mathbf{a}_i$ denotes the $i^{th}$ column of $\mathbf{A}$ and $\mathbf{W}_{ij}$ is the Pearson correlation between the activation map of region $i$ and region $j$ across $M$ patients. The standard graph Laplacian, $\mathbf{L}$, is computed from the sample correlation matrix, $\mathbf{W}$.

**Generative Model for Genetics:** Let $G$ denote the number of genetic variants measured in each subject. This data is captured by the vector $\mathbf{g_m} \in \mathbb{R}^{G \times 1}$. We assume that the genetic data express itself through a sparse set of basis vectors, $\mathbf{g}_m \approx \mathbf{B}^T \mathbf{z}_m$; furthermore, the projection is tied to that of the imaging data. Here $\mathbf{B} \in \mathbb{R}^{d \times G}$ is the genetic basis matrix. We employ a group sparsity penalty over $\mathbf{B}$, in the form of $\ell_{2,1}$ norm. This penalty selects a sparse set of genetic variants through the $\ell_1$ penalty across rows. At the same time, the $\ell_2$ penalty across columns preserve the genetic representation across basis vectors.

**Joint Objective with Diagnosis Prediction:** We use the patient-specific projection coefficients $\{\mathbf{z}_m\}_{m=1}^M$ to predict the class labels. Mathematically, $y_m \approx \sigma(\mathbf{z}_m^T \mathbf{c})$ where $\sigma(\cdot)$ is the standard sigmoid function and $\mathbf{c} \in \mathbb{R}^{d \times 1}$ is the regression vector. We combine the dictionary learning and logistic regression terms in a joint objective. For convenience, we concatenate the imaging

and genetics data into the matrices $\mathbf{F} = [\mathbf{f}_1, \ldots, \mathbf{f}_M]$ and $\mathbf{G} = [\mathbf{g}_1, \ldots, \mathbf{g}_M]$, respectively. Likewise, we define the latent projection matrix $\mathbf{Z} = [\mathbf{z}_1, \ldots, \mathbf{z}_M]$.

$$
\begin{aligned}
\mathcal{J}(\mathbf{A}, \mathbf{B}, \mathbf{Z}, \mathbf{c}) = & \ ||\mathbf{F} - \mathbf{A}^T\mathbf{Z}||_F^2 + ||\mathbf{G} - \mathbf{B}^T\mathbf{Z}||_F^2 \\
& - \sum_{m=1}^{M} (y_m \log(\sigma(\mathbf{z}_m^T\mathbf{c})) + (1 - y_m) \log(1 - \sigma(\mathbf{z}_m^T\mathbf{c}))) \\
& + \frac{\lambda_1}{2} Tr(\mathbf{A}\mathbf{L}\mathbf{A}^T) + \lambda_2 ||\mathbf{B}||_{2,1} + \frac{\lambda_3}{2} ||\mathbf{Z}||_F^2 + \frac{\lambda_4}{2} ||\mathbf{c}||_2^2 \ \ \text{s.t. } \mathbf{A}\mathbf{A}^T = \mathbf{I}
\end{aligned}
\tag{1}
$$

where $\{\lambda_1, \lambda_2, \lambda_3, \lambda_4\}$ are the associated regularization parameters. The penalties, $\frac{\lambda_3}{2} ||\mathbf{Z}||_F^2$, and $\frac{\lambda_4}{2} ||\mathbf{c}||_2^2$ are introduced to make the optimization well-posed.

## 2.1 Optimization Strategy

We use an alternating minimization strategy to optimize the unknown variables $\{\mathbf{A}, \mathbf{B}, \mathbf{Z}, \mathbf{c}\}$ in Eq. (1) from the data $\{\mathbf{f}_m, \mathbf{g}_m, y_m\}$. This procedure iteratively updates each variable while holding the remaining variables constant.

**Update for A via Interior Point Solver:** The objective $\mathcal{J}(\cdot)$ is a convex function in $\mathbf{A}$ with a nonconvex equality constraint $\mathbf{A}\mathbf{A}^T = \mathbf{I}$. We use an interior point solver to incrementally optimize over $\mathbf{A}$. Specifically, each iteration solves the following modified problem:

$$
\mathbf{A}_{t+1} = \min_{\mathbf{A}} \ ||\mathbf{F} - \mathbf{A}^T\mathbf{Z}||_F^2 + \frac{\lambda_1}{2} Tr(\mathbf{A}\mathbf{L}\mathbf{A}^T) \quad \text{subject to: } \mathbf{A}\mathbf{A}^T - \mathbf{I} = 0
\tag{2}
$$

The interior point solver uses two methods to solve Eq. (2). It first tries to take a Newton step by solving the augmented Lagrangian problem. If this step fails the algorithm approximates the augmented Lagragian in the least squares sense to solve for the lagrange multipliers and takes a conjugate gradient step to approximately solve Eq. (2) using a trust region method.

**Update B, Z, and c Using Trust Region Method:** The objective $\mathcal{J}(\cdot)$ is convex in each of the variables $\{\mathbf{B}, \mathbf{Z}, \mathbf{c}\}$ while keeping the other variables constant. So, we solve the minimization problem for each variable in an iterative fashion using an unconstrained trust region solver. At each iteration the trust region method estimates the step size and direction $\mathbf{s}_k$ to optimize the function $f(\mathbf{x})$ through the following quadratic program:

$$
\mathbf{s}_k = \text{argmin}_{\mathbf{s}} \ f(\mathbf{x}_k) + \mathbf{g}_k^T\mathbf{s} + \frac{1}{2}\mathbf{s}^T\mathbf{H}_k\mathbf{s} \ \ \text{subject to: } ||\mathbf{s}|| < \delta
\tag{3}
$$

where $\mathbf{g}_k$ and $\mathbf{H}_k$ are the gradient and hessian of $f(\mathbf{x})$ at $\mathbf{x}_k$. The update $\mathbf{x}_k \rightarrow \mathbf{x} + \mathbf{s}_k$ is taken such that $f(\mathbf{x}_k + \mathbf{s}_k) < f(\mathbf{x}_k)$. This method is guaranteed to converge to a local minimum in polynomial time.

In our setting $f(\cdot)$ involves the terms of the objective function, $\mathcal{J}(\cdot)$ that involves the variable in consideration. A typical example is minimization of $\mathbf{B}$ where, $f(\mathbf{B}) = ||\mathbf{G} - \mathbf{B}^T\mathbf{Z}||_F^2 + \lambda_2 ||\mathbf{B}||_{2,1}$. $f(\mathbf{Z})$ and $f(\mathbf{c})$ are similarly specified.

**Prediction on Unseen Data:** We use a tenfold cross validation setup to evaluate the predictive power of our framework. In this case, we optimize the variables $\{\mathbf{A}^*, \mathbf{B}^*, \mathbf{Z}^*, \mathbf{c}^*\}$ based just on the training data. For testing, we use just the imaging and genetics data $\{\mathbf{f}_{test}, \mathbf{g}_{test}\}$ of a new subject to obtain the subject-specific projection $\mathbf{z}_{test}$ while setting the cross-entropy term to zero (since the diagnosis $y_{test}$ in unknown). We then use the same logistic regression $y_{test} = \sigma(\mathbf{z}_{test}^T \mathbf{c}^*)$ to predict the class label. Unlike prior work [1], our setup performs feature selection in a nested fashion, since the bases matrices $\mathbf{A}$ and $\mathbf{B}$, which in turn govern the latent projection, are estimated only from the training data.

**Baseline Comparisons:** We compare out model with three baseline methods, as described below:

– **Random Forest (RF) Classification:** We construct a RF classifier for diagnosis based on the concatenated imaging and genetics data, $[\mathbf{f}_m^T, \mathbf{g}_m^T]^T$.
– **CCA + RF Classification:** Canonical correlation analysis (CCA) identifies bi-multivariate associations between imaging and genetics data. This approach is similar in spirit to our coupled latent projection, but it does not include the diagnosis to guide the association. The input to the RF classifier is the aligned imaging & genetics data after performing CCA. Once again, we concatenate the modalities into a single feature vector.
– **Imaging Only Variant of Eq. 1:** Finally we consider a variant of our own model with just the imaging terms and ignore the terms that involve genetic information. This baseline will help us to quantify the performance gain for adding genetic data. We again evaluate this model in a nested fashion, where we optimize the variables $\{\mathbf{A}^*, \mathbf{Z}^*, \mathbf{c}^*\}$ over training set and use them to estimate subject specific projection and diagnosis on the testing data.

We use a grid search to fix parameters for each method. Based on this experiment, we fix the genetic regularizer $\lambda_2$, the projection regularizer $\lambda_3$, and the regression regularizer $\lambda_4$ to $\lambda_2 = 7.5$, $\lambda_3 = 0.6$, $\lambda_4 = 0.04$ for both our model variants. The imaging regularizer $\lambda_1$ and latent dimension $d$ are sensitive to the complexity of the fMRI paradigm used in the study. We use $\{d = 9, \lambda_1 = 17.5\}$ for the Nback task, and $\{d = 7, \lambda_1 = 12.5\}$ for SDMT task. The baseline RFs produce stable results for 9000 trees, which is what we use in the analysis.

## 3   Experimental Results

We evaluate our model on two fMRI datasets and a SNP dataset of schizophrenia patients and normal controls. The first fMRI paradigm is a working memory task (NBack), comprised of 2-back working memory trial blocks alternating with 0-back trial blocks. During 0-back trials participants were instructed to press a button corresponding to a number displayed on the screen and during the 2-back working memory trials the participants were instructed to press the button corresponding to the number they saw two stimuli previously. This dataset includes

**Table 1.** Performance of each of the methods during nested cross validation. We abbreviated Sensitivity to Sens, Specificity to Spec and Accuracy to Acc.

| Method | NBack task | | | SDMT task | | |
|---|---|---|---|---|---|---|
| | Sens | Spec | Acc | Sens | Spec | Acc |
| RF | 0.58 | 0.56 | 0.57 | **0.68** | 0.56 | 0.62 |
| CCA + RF | 0.41 | 0.47 | 0.44 | 0.55 | 0.69 | 0.62 |
| Our method (Imaging only) | 0.66 | 0.52 | 0.60 | 0.63 | 0.63 | 0.63 |
| Our method (Imaging+Genetics) | **0.71** | **0.68** | **0.70** | 0.60 | **0.76** | **0.68** |

53 patients and 53 controls, matched on age, IQ and education. The second fMRI paradigm is a simple declarative memory task (SDMT), which involved incidental encoding of complex aversive visual scenes. This dataset includes 46 patients and 47 controls, matched on age, IQ, and education. All fMRI data was acquired on 3-T General Electric Sigma scanner (EPI, TR/TE = 2000/28 ms; flip angle = 90; field of view = 24 cm, res: 3.75 mm in x and y dimensions and 6 mm in the z dimension for NBack and 5 mm for SDMT;). FMRI preprocessing include slice timing correction, realignment, spatial normalization to an MNI template, smoothing and motion regression. SPM12 is used to generate activation and contrast maps for each paradigm. We use the Brainnetome atlas [6] to define 246 cortical and subcortical regions. The input to our model is the contrast map over these 246 ROIs. In parallel, genotyping was done using variate Illumina Bead Chips including 510K/610K/660K/2.5M. Quality control and imputation were performed using PLINK and IMPUTE2 respectively. The resulting 102K linkage disequilibrium independent SNPs are used to calculate the polygenic risk score of schizophrenia via a log-odds ratio of the imputation probability for the reference allele [3]. By selecting $P < 10^{-4}$, we obtain 1252 linkage disequilibrium independent SNPs. We remove the effect of first principal component from the SNP training data and use the same estimated projection matrix to remove the effect of first principal component from the test data.

Table 1 quantifies the classification performance of each method. Notice that both dictionary learning frameworks (with and without genetics) outperform the standard machine learning baselines. Additionally, the genetic information improves the overall performance of our model. This result suggest that our coupled dictionary learning framework is able to identify meaningful features from both the imaging and genetics data that distinctly separates patients and controls.

Figure 2 shows the most significant set of brain region as obtained from the template basis vectors of **A**. We observe that in the Nback dataset the set of regions include the dorsolateral prefrontal cortex that is well known to underlie executive function including working memory. This region has been implicated in the executive functioning deficits of Schizophrenia [2]. The SDMT task implicates hippocampal and parahippocampal regions also thought to be disrupted in
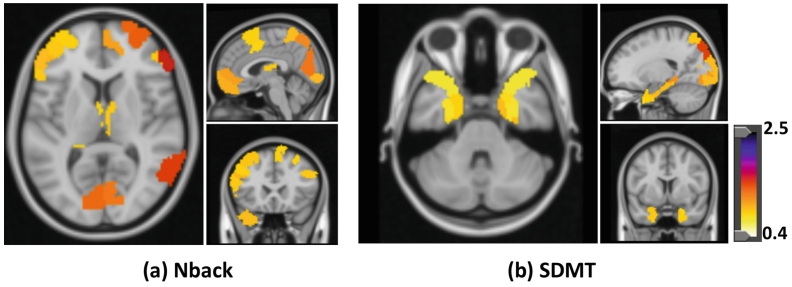
**(a) Nback**                     **(b) SDMT**

**Fig. 2.** The representative set of regions captured by the matrix, **A**. The color bar shows the level of contribution of each region for classification. (Color figure online)
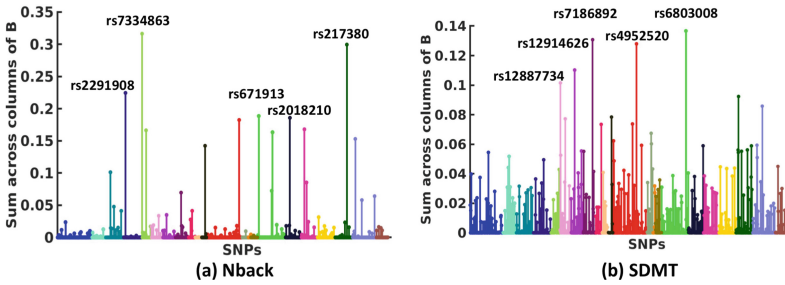


**(a) Nback**                     **(b) SDMT**

**Fig. 3.** The contribution of each SNP to discriminate the subjects between patients and controls. We have annotated the top five SNPs. The colors indicate the chromosomes on which the SNPs are located. (Color figure online)

schizophrenia [7]. Hence our model is able to find well reported and interpretable aberrations between schizophrenia patients and controls.

Figure 3 illustrates the contribution of individual SNP to the latent vector. It is calculated as the sum of absolute values of the columns in **B**. We have annotated the five most highly implicated SNPs for each dataset as a reference. According to GWAS the overlapping genes of these SNPs are closely related to schizophrenia. Additionally, we ran a gene ontology enrichment analysis based on the genes associated with the top 150 SNPs. The results are shown in Table 2. We found a common biological process for both the datasets implicated in *nervous system development* [4]. These findings verifies the ability of our coupled framework to find clinically relevant biomarkers from both the imaging and genetic data.

**Table 2.** The table shows the enriched biological processes along with their level of significance obtained via GO enrichment analysis. The processes are arranged by the most specific subclass first, with its parent terms indented directly below it.

| Dataset | Biological processes | FDR |
|---------|---------------------|-----|
| Nback | Central nervous system development | 0.03 |
| | → Nervous system development | 0.0002 |
| | → System development | 0.001 |
| | Generation of neurons | 0.03 |
| | → Neurogenesis | 0.02 |
| | → Cell differentiation | 0.003 |
| SDMT | Forebrain neuron differentiation | 0.04 |
| | → Nervous system development | 0.002 |
| | → Generation of neurons | 0.004 |
| | → Central nervous system neuron differentiation | 0.04 |
| | Central nervous system neuron development | 0.02 |
| | Regulation of neurogenesis | 0.03 |

## 4  Conclusion

We have introduced an elegant joint matrix decomposition framework that identifies imaging and genetic biomarkers guided by the clinical diagnosis. Unlike other conventional analysis this framework can robustly and efficiently integrate diverse datatypes while maintaining good prediction accuracy. Moreover, the biomarkers may help us understand the biology underlying cognitive deficits in patients with schizophrenia in relation to genetic variants. This model can easily be adapted to other imaging and genetic modalities. In this work we only explored a linear relationship between imaging, genetics and diagnosis, however in future work we will also explore the non linear relationships across them.

## References

1. Batmanghelich, N.K., et al.: Probabilistic modeling of imaging, genetics and diagnosis. IEEE Trans. Med. Imaging **35**(7), 1765–1779 (2016)
2. Callicott, J.H., et al.: Abnormal fMRI response of the dorsolateral prefrontal cortex in cognitively intact siblings of patients with schizophrenia. Am. J. Psychiatry **160**(4), 709–719 (2003)
3. Chen, Q., et al.: Schizophrenia polygenic risk score predicts mnemonic hippocampal activity. Brain **141**(4), 1218–1228 (2018)
4. Dean, B.: Is schizophrenia the price of human central nervous system complexity? Aust. New Zealand J. Psychiatry **43**(1), 13–24 (2009)

5. Du, L., et al.: Pattern discovery in brain imaging genetics via SCCA modeling with a generic non-convex penalty. Sci. Rep. **7**(1), 14052 (2017)
6. Fan, L., et al.: The human brainnetome atlas: a new brain atlas based on connectional architecture. Cereb. Cortex **26**(8), 3508–3526 (2016)
7. Rasetti, R., et al.: Altered hippocampal-parahippocampal function during stimulus encoding. JAMA Psychiatry **71**(3), 236 (2014)
8. Wang, H., et al.: Identifying quantitative trait loci via group-sparse multitask regression and feature selection: an imaging genetics study of the ADNI cohort. Bioinformatics **28**(2), 229–237 (2012)