

Supplemental Material to “Identifying Alzheimer’s Disease-Related Brain Regions from Multi-Modality Neuroimaging Data using Sparse Composite Linear Discrimination Analysis”

I – How does the formulation (5) serve the purpose of the composite parameterization

It is interesting to see how the penalty function, $\sqrt{\sum_{l=1}^q \sum_{m=1}^M |\theta_{k,l}^{(m)}|}$, will serve the purpose of the composite parameterization, i.e., common information and specific information can be estimated separately and simultaneously. To illustrate this point, let’s look at LASSO first. It has been pointed out in [40] that, the variable selection ability of LASSO depends on the singularity of the parameters when the parameters are equal to zero. For example, consider the LASSO problem [27], whose penalized likelihood function is $LS(\boldsymbol{\beta}) + \lambda \sum_i |\beta_i|$, and the ridge regression, whose penalized likelihood function is $LS(\boldsymbol{\beta}) + \lambda \sum_i \beta_i^2$. Here, $LS(\boldsymbol{\beta})$ represents the least square error function. It is well known that LASSO is capable of performing variable selection, i.e., forcing many parameters in the $\boldsymbol{\beta}$ estimate to be exactly zero, while ridge regression is not. The reason is that the penalized likelihood function of LASSO is singular when $\beta_i = 0$ [40]. Similarly, the penalty function used in SCLDA is singular when all the variables under the square root are zero (corresponding to $\delta_k = 0$), which facilitates the variable selection on the common information level. Furthermore, this penalty function is also singular when any variable under the square root is zero (corresponding to any $\gamma_{k,l}^{(m)} = 0$), which further facilitates the variable selection on the specific information level. It is worth mentioning that L2/L1 regularization employed in multitask feature selection methods, which will be $\sqrt{\sum_{l=1}^q \sum_{m=1}^M (\theta_{k,l}^{(m)})^2}$ in our case, is only singular when $\sum_{l=1}^q \sum_{m=1}^M (\theta_{k,l}^{(m)})^2 = 0$. This tends to produce an “all-in-all-out” solution, i.e., when one variable is selected, all the other variables in the square root will be selected, since there is no more singularity. Thus, the L2/L1 regularization is not capable of performing the variable selection on the specific information level, which limits its practical use in some applications.

II – The details on the DC programming for solving (5)

The optimization problem (5) is a non-convex optimization problem that is difficult to solve. We address this problem by using an iterative two-stage procedure known as Difference of Convex functions (DC) programming [39]. The basic idea behind the DC programming is that, at the first stage of every iteration, a surrogate convex objective function is proposed to bound the non-convex objective function at the current solution; then, at the second stage, a new solution is obtained by maximizing this surrogate convex objective function. This process iterates until a certain convergence rule is met. It is worth mentioning that the DC programming shares the same spirit as the Expectation-Maximization (EM) algorithm or Minorization-Maximization (MM) algorithms that have been widely used in statistics and machine learning.

We adopt the DC algorithm to solve (5) following [39], in which the DC programming is to solve a least-square problem with non-convex penalty terms. In our problem, the essential task is to find a decomposition of $\sum_k \sqrt{\sum_{l=1}^q \sum_{m=1}^M |\theta_{k,l}^{(m)}|}$ as a sum of two convex functions. As suggested in [39], the decomposition, $\sqrt{\theta_{k,l}^{(m)}} = |\theta_{k,l}^{(m)}| - \left(|\theta_{k,l}^{(m)}| - \sqrt{\theta_{k,l}^{(m)}} \right)$, can be used for the non-convex penalty term $\sqrt{\theta_{k,l}^{(m)}}$. Use

$$\sqrt{\sum_{l=1}^q \sum_{m=1}^M |\theta_{k,l}^{(m)}|} = \sum_{l=1}^q \sum_{m=1}^M |\theta_{k,l}^{(m)}| - \left(\sum_{l=1}^q \sum_{m=1}^M |\theta_{k,l}^{(m)}| - \sqrt{\sum_{l=1}^q \sum_{m=1}^M |\theta_{k,l}^{(m)}|} \right),$$

as a decomposition for $\sqrt{\sum_{l=1}^q \sum_{m=1}^M |\theta_{k,l}^{(m)}|}$. Based on the theory developed in [39], this decomposition results in the following objective function

$$\begin{aligned} \tilde{\Theta}^{(t+1)} &= \operatorname{argmin}_{\Theta} l_3(\Theta | \{\mathbf{Z}^{(1)}, \mathbf{Z}^{(2)}, \dots, \mathbf{Z}^{(M)}\}) = \\ &\operatorname{argmin}_{\Theta} \left\{ -l_0(\Theta | \{\mathbf{Z}^{(1)}, \mathbf{Z}^{(2)}, \dots, \mathbf{Z}^{(M)}\}) + \sum_k \lambda_k^{(t+1)} \sum_{l=1}^q \sum_{m=1}^M |\theta_{k,l}^{(m)}| \right\}, \end{aligned} \quad (6)$$

as the surrogate convex objective function, where

$$\lambda_k^{(t+1)} = \lambda - h'(z) \Big|_{z = \sum_{l=1}^q \sum_{m=1}^M |\theta_{k,l}^{(m)(t)}|}, \quad h(z) = \lambda(z - \sqrt{z}),$$

$\tilde{\Theta}^{(t+1)}$ is the solution at iteration $t + 1$ and $\theta_{k,l}^{(m)(t)}$ is a corresponding element. It is shown that this decomposition produces a surrogate convex objective function with L1- penalty, which can be solved by many existing efficient algorithms developed for LASSO-type problems. A complete procedure for the DC programming is depicted in Figure 1.

Initialize: Let $t = 0$;
Repeat
 Calculate $\lambda_k^{(t)} = \lambda - h'(z) \Big|_{z = \sum_{l=1}^q \sum_{m=1}^M |\theta_{k,l}^{(m)(t)}|}$;
 for $1 \leq k, l \leq p, 1 \leq m \leq M$;
 Solve (6) and get $\tilde{\Theta}^{(t)}$;
 Let $t = t + 1$;
Until converge

Figure 1: The DC programming for solving (5)

The optimization problem (6) is a standard L1-regularization type problem, whose objective function is a sum of a smooth likelihood/least square error function and an L1- penalty on the parameters. This problem can be solved by many efficient numeric algorithms in the literature [25, 26]. In our case, the two-metric method is employed [25].

It has been pointed out that the DC programming for solving many non-convex regularization problems can be closely linked to the adaptive Lasso formulation [39]. To illustrate this point in our case, we note that

$$h'(z) \Big|_{z = \sum_{l=1}^q \sum_{m=1}^M |\theta_{k,l}^{(m)}|} = \lambda - \frac{\lambda}{2 \left(\sqrt{\sum_{l=1}^q \sum_{m=1}^M |\theta_{k,l}^{(m)}|} + \epsilon \right)}$$

is in our DC programming, where ϵ is a user-specified number for numerical stability consideration. This produces a new regularization parameter for iteration $t + 1$,

$$\lambda_k^{(t+1)} = \frac{\lambda}{2 \left(\sqrt{\sum_{l=1}^q \sum_{m=1}^M |\theta_{k,l}^{(m)(t)}|} + \epsilon \right)},$$

which is inversely proportional to the magnitude of $\sum_{l=1}^q \sum_{m=1}^M |\theta_{k,l}^{(m)(t)}|$. This implies that, at each iteration, the DC programming essentially reweights the regularization parameters of each

L1-regularization type problem (6). Specifically, at iteration $t + 1$, the new regularization parameters associated with the zero $\theta_{k,l}^{(m)}$'s identified at the previous iteration will increase drastically, while the new regularization parameters associated with the non-zero $\theta_{k,l}^{(m)}$'s identified at the previous iteration will decrease proportionally to the sum of the absolute magnitudes of these $\theta_{k,l}^{(m)}$'s, which belong to the same variable. In this manner, the shrinkage effect imposed on the non-zero $\theta_{k,l}^{(m)}$'s by the L1-regularization is effectively alleviated. This explains why the proposed SCLDA has the capability of preserving weak-effect features.

III – Simulation procedure

For a given combination of values of these parameters, simulation data can be generated in the following way: First, we generate a $1 \times p$ vector, $\boldsymbol{\beta}$, with l elements being randomly selected to be 1 and the other elements being zero. Second, to generate a feature vector for dataset i , $\boldsymbol{\beta}_i$, we randomly change $(100 - s)\%$ of the non-zero elements of $\boldsymbol{\beta}$ to be zero, and, at the same time, we randomly change the same number of zero elements to be 1. The nonzero elements of $\boldsymbol{\beta}_i$ correspond to the features of dataset i . In this manner, the larger the $s\%$, the more overlapping of the features across the data sources. We further randomly pick up half of the nonzero elements of each $\boldsymbol{\beta}_i$ and modify them in such a way: if the element is 1, then its new value is sampled from the uniform distribution $U[0,0.5]$; if the element is -1, then its new value is sampled from the uniform distribution $U[-0.5,0]$. This is to mimic the true situation in the real application in section 5, where a large portion of weak-effect features are present. After that, we use the resulting $\boldsymbol{\beta}_i$ and $-\boldsymbol{\beta}_i$ as the mean vectors of two classes of dataset i (for simplicity, we only investigate 2-class problems for the simulation). And, as we assume heterogeneous covariance matrices for different classes, two covariance matrices are independently generated from a Wishart distribution with degrees of freedom to be n and scale matrix to be an Identity matrix. Each covariance matrix is further diagonalized to make the diagonal elements being 0.2 (the corresponding signal-to-noise ratio is about 5). With these mean vectors and covariance matrices, we generate the dataset i by sampling one sub-dataset for each class from its corresponding multivariate Gaussian distribution, and then combining these two sub-datasets as dataset i .

IV – Compare SCLDA, SLDA and MSLDA under different sample sizes

Furthermore, we conduct an experiment to compare SCLDA, SLDA and MSLDA under different sample sizes. As shown in Figure 3, SCLDA performs significantly better when sample sizes are small, which confirms that SCLDA is more statistically efficient in recovering the features.

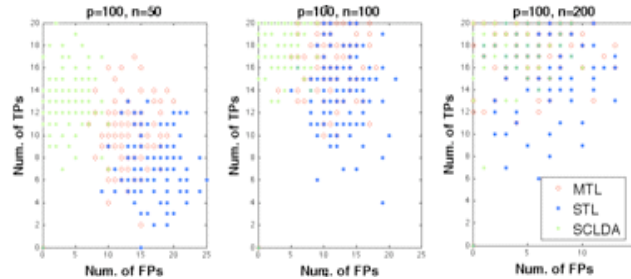


Figure 2: Average numbers of TPs vs FPs for proposed SCLDA (green symbols “+”), SLDA (blue symbols “*”) and MSLDA (red symbols “o”) with task relatedness $s\% = 90\%$, $m = 2$

V– Selection of λ and q

We first discuss the optimal selection of λ and q for the case when there is only one data source. This optimal selection can be obtained by maximizing Akaike's Information Criteria (AIC) [1]:

$$AIC = s(\lambda, q) - \frac{d}{n_{avg}},$$

where $s(\lambda, q)$ is the cross validation classification accuracy associated with λ and q , d is the number of nonzero parameters in $\tilde{\Theta}$, and n_{avg} is the average number of observations per class. For the general case where there are M data sources, the optimal selection of λ and q can be obtained by maximizing the average AIC, which is:

$$AIC_{avg} = \sum_{i=1}^M AIC_i / M.$$

where AIC_i is the AIC value for the i^{th} data source. Both the simulation studies in section 4 and real application in section 5 reveal that this criterion can produce accurate and meaningful model selection results.

[1] Chiang, L., Russell, E.R. 2001. *Fault Detection and Diagnosis in Industrial Systems*. Springer, London.

VI – Proof of Theorem 1

Proof: We first demonstrate that any local optima of $l_1(\Gamma, \Psi | \{\mathbf{Z}^{(1)}, \mathbf{Z}^{(2)}, \dots, \mathbf{Z}^{(M)}\})$, $\hat{\Gamma}$ and $\hat{\Psi}$, corresponds to a local optima of

$$\begin{aligned} l'_1(\Gamma, \Psi | \{\mathbf{Z}^{(1)}, \mathbf{Z}^{(2)}, \dots, \mathbf{Z}^{(M)}\}) &= \\ \sum_{m=1}^M \left\{ \frac{N^{(m)}}{2} \log \left| \boldsymbol{\theta}_{p-q}^{(m)T} \mathbf{T}^{(m)} \boldsymbol{\theta}_{p-q}^{(m)} \right| + \sum_{j=1}^J \frac{N_j}{2} \log \boldsymbol{\theta}_q^{(m)} \mathbf{W}_j^{(m)} \boldsymbol{\theta}_q^{(m)} - N^{(m)} \log |\boldsymbol{\theta}^{(m)}| \right\} &+ \sum_k \delta_k + \\ \eta \sum_k \sum_{l=1}^q \sum_{m=1}^M \gamma_{k,l}^{(m)}, &\text{ which is denoted } \tilde{\Gamma}, \tilde{\Psi}. \text{ To show this, we insert } \tilde{\Gamma}, \tilde{\Psi} \text{ into } l'_1, \text{ which is} \\ l'_1(\tilde{\Gamma}, \tilde{\Psi} | \{\mathbf{Z}^{(1)}, \mathbf{Z}^{(2)}, \dots, \mathbf{Z}^{(M)}\}) &= \\ = \sum_{m=1}^M \left\{ \frac{N^{(m)}}{2} \log \left| \tilde{\boldsymbol{\theta}}_{p-q}^{(m)T} \mathbf{T}^{(m)} \tilde{\boldsymbol{\theta}}_{p-q}^{(m)} \right| + \sum_{j=1}^J \frac{N_j}{2} \log \tilde{\boldsymbol{\theta}}_q^{(m)} \mathbf{W}_j^{(m)} \tilde{\boldsymbol{\theta}}_q^{(m)} - N^{(m)} \log |\tilde{\boldsymbol{\theta}}^{(m)}| \right\} &+ \sum_k \delta_k + \\ \eta \sum_k \sum_{l=1}^q \sum_{m=1}^M \tilde{\gamma}_{k,l}^{(m)}, &\text{ which can be further written as} \\ l'_1(\tilde{\Gamma}, \tilde{\Psi} | \{\mathbf{Z}^{(1)}, \mathbf{Z}^{(2)}, \dots, \mathbf{Z}^{(M)}\}) &= \\ \sum_{m=1}^M \left\{ \frac{N^{(m)}}{2} \log \left| \tilde{\boldsymbol{\theta}}_{p-q}^{(m)T} \mathbf{T}^{(m)} \tilde{\boldsymbol{\theta}}_{p-q}^{(m)} \right| + \sum_{j=1}^J \frac{N_j}{2} \log \tilde{\boldsymbol{\theta}}_q^{(m)} \mathbf{W}_j^{(m)} \tilde{\boldsymbol{\theta}}_q^{(m)} - N^{(m)} \log |\tilde{\boldsymbol{\theta}}^{(m)}| \right\} &+ \lambda_1 \sum_k \frac{\tilde{\delta}_k}{\lambda_1} + \\ \lambda_2 \sum_k \sum_{l=1}^q \sum_{m=1}^M \lambda_1 \tilde{\gamma}_{k,l}^{(m)}, & \\ = l_1 \left(\lambda_1 \tilde{\Gamma}, \frac{\tilde{\Psi}}{\lambda_1} | \{\mathbf{Z}^{(1)}, \mathbf{Z}^{(2)}, \dots, \mathbf{Z}^{(M)}\} \right) &\geq l_1(\hat{\Gamma}, \hat{\Psi} | \{\mathbf{Z}^{(1)}, \mathbf{Z}^{(2)}, \dots, \mathbf{Z}^{(M)}\}). \text{ The last inequality holds} \\ \text{since } \hat{\Gamma} \text{ and } \hat{\Psi} \text{ is a local optima of } l_1. \text{ Similarly, we insert } \hat{\Gamma}, \hat{\Psi} \text{ into } l_1 \text{ and we can prove} & \\ l_1(\hat{\Gamma}, \hat{\Psi} | \{\mathbf{Z}^{(1)}, \mathbf{Z}^{(2)}, \dots, \mathbf{Z}^{(M)}\}) = l'_1 \left(\frac{\hat{\Gamma}}{\lambda_1}, \lambda_1 \hat{\Psi} | \{\mathbf{Z}^{(1)}, \dots, \mathbf{Z}^{(M)}\} \right) &\geq l'_1(\tilde{\Gamma}, \tilde{\Psi} | \{\mathbf{Z}^{(1)}, \mathbf{Z}^{(2)}, \dots, \mathbf{Z}^{(M)}\}), \end{aligned}$$

Therefore, we have

$$l_1(\hat{\Gamma}, \hat{\Psi} | \{\mathbf{Z}^{(1)}, \mathbf{Z}^{(2)}, \dots, \mathbf{Z}^{(M)}\}) = l'_1 \left(\frac{\hat{\Gamma}}{\lambda_1}, \lambda_1 \hat{\Psi} | \{\mathbf{Z}^{(1)}, \dots, \mathbf{Z}^{(M)}\} \right) = l'_1(\tilde{\Gamma}, \tilde{\Psi} | \{\mathbf{Z}^{(1)}, \mathbf{Z}^{(2)}, \dots, \mathbf{Z}^{(M)}\}).$$

This has an implication that $\left\{ \frac{\hat{\Gamma}}{\lambda_1}, \lambda_1 \hat{\Psi} \right\}$ is also a local optimizer of $l'_1(\Gamma, \Psi | \{\mathbf{Z}^{(1)}, \mathbf{Z}^{(2)}, \dots, \mathbf{Z}^{(M)}\})$

and $\tilde{\delta}_k \tilde{\gamma}_{k,l}^{(m)} = \hat{\delta}_k \hat{\gamma}_{k,l}^{(m)}$, $\eta = \lambda_1 \lambda_2$.

Now we demonstrate that any local optima of $l'_1(\mathbf{\Gamma}, \mathbf{\Psi}|\{\mathbf{Z}^{(1)}, \mathbf{Z}^{(2)}, \dots, \mathbf{Z}^{(M)}\})$ corresponds to a local optima of $l_2(\mathbf{\Theta}|\{\mathbf{Z}^{(1)}, \mathbf{Z}^{(2)}, \dots, \mathbf{Z}^{(M)}\})$, $\hat{\mathbf{\Theta}}$. Using the same idea above, we can obtain that

$$\begin{aligned} & l'_1(\tilde{\mathbf{\Gamma}}, \tilde{\mathbf{\Psi}}|\{\mathbf{Z}^{(1)}, \mathbf{Z}^{(2)}, \dots, \mathbf{Z}^{(M)}\}), \\ & \geq \sum_{m=1}^M \left\{ \frac{N^{(m)}}{2} \log \left| \tilde{\mathbf{\Theta}}_{p-q}^{(m)T} \mathbf{T}^{(m)} \tilde{\mathbf{\Theta}}_{p-q}^{(m)} \right| + \sum_{j=1}^J \frac{N_j}{2} \log \tilde{\mathbf{\Theta}}_q^{(m)} \mathbf{W}_j^{(m)} \tilde{\mathbf{\Theta}}_q^{(m)} - N^{(m)} \log |\tilde{\mathbf{\Theta}}^{(m)}| \right\} + \\ & 2 \sum_k \sqrt{\eta \tilde{\delta}_k \sum_{l=1}^q \sum_{m=1}^M |\tilde{\gamma}_{k,l}^{(m)}|}, \\ & = l_2(\tilde{\mathbf{\Theta}}|\{\mathbf{Z}^{(1)}, \mathbf{Z}^{(2)}, \dots, \mathbf{Z}^{(M)}\}) \geq l_2(\hat{\mathbf{\Theta}}|\{\mathbf{Z}^{(1)}, \mathbf{Z}^{(2)}, \dots, \mathbf{Z}^{(M)}\}). \end{aligned}$$

On the other hand, let $\hat{\mathbf{\Psi}} = \left\{ \hat{\delta}_k = \sqrt{\eta \sum_{l=1}^q \sum_{m=1}^M |\hat{\theta}_{k,l}^{(m)}|}, 1 \leq k \leq p \right\}$ and $\hat{\mathbf{\Gamma}} = \left\{ \hat{\gamma}_{k,l}^{(m)} = \frac{\hat{\theta}_{k,l}^{(m)}}{\hat{\delta}_k}, 1 \leq k \leq p, 1 \leq l \leq p, 1 \leq m \leq M \right\}$, then we can obtain

$$\begin{aligned} & l_2(\hat{\mathbf{\Theta}}|\{\mathbf{Z}^{(1)}, \mathbf{Z}^{(2)}, \dots, \mathbf{Z}^{(M)}\}), \\ & = \sum_{m=1}^M \left\{ \frac{N^{(m)}}{2} \log \left| \hat{\mathbf{\Theta}}_{p-q}^{(m)T} \mathbf{T}^{(m)} \hat{\mathbf{\Theta}}_{p-q}^{(m)} \right| + \sum_{j=1}^J \frac{N_j}{2} \log \hat{\mathbf{\Theta}}_q^{(m)} \mathbf{W}_j^{(m)} \hat{\mathbf{\Theta}}_q^{(m)} - N^{(m)} \log |\hat{\mathbf{\Theta}}^{(m)}| \right\} + \\ & \sum_k \sqrt{\eta \sum_{l=1}^q \sum_{m=1}^M |\hat{\theta}_{k,l}^{(m)}|} + \sqrt{\eta} \sum_k \sum_{l=1}^q \sum_{m=1}^M \frac{|\hat{\theta}_{k,l}^{(m)}|}{\sqrt{\sum_{l=1}^q \sum_{m=1}^M |\hat{\theta}_{k,l}^{(m)}|}}, \\ & = \sum_{m=1}^M \left\{ \frac{N^{(m)}}{2} \log \left| \hat{\mathbf{\Theta}}_{p-q}^{(m)T} \mathbf{T}^{(m)} \hat{\mathbf{\Theta}}_{p-q}^{(m)} \right| + \sum_{j=1}^J \frac{N_j}{2} \log \hat{\mathbf{\Theta}}_q^{(m)} \mathbf{W}_j^{(m)} \hat{\mathbf{\Theta}}_q^{(m)} - N^{(m)} \log |\hat{\mathbf{\Theta}}^{(m)}| \right\} + \sum_k \delta_k + \\ & \lambda \sum_k \sum_{l=1}^q \sum_{m=1}^M \hat{\gamma}_{k,l}^{(m)}, \\ & = l'_1(\hat{\mathbf{\Theta}}|\{\mathbf{Z}^{(1)}, \mathbf{Z}^{(2)}, \dots, \mathbf{Z}^{(M)}\}) \geq l'_1(\tilde{\mathbf{\Theta}}|\{\mathbf{Z}^{(1)}, \mathbf{Z}^{(2)}, \dots, \mathbf{Z}^{(M)}\}). \end{aligned}$$

Thus, we have

$$l_2(\hat{\mathbf{\Gamma}}, \hat{\mathbf{\Psi}}|\{\mathbf{Z}^{(1)}, \mathbf{Z}^{(2)}, \dots, \mathbf{Z}^{(M)}\}) = l'_1(\hat{\mathbf{\Gamma}}, \hat{\mathbf{\Psi}}|\{\mathbf{Z}^{(1)}, \mathbf{Z}^{(2)}, \dots, \mathbf{Z}^{(M)}\}) = l'_1(\tilde{\mathbf{\Gamma}}, \tilde{\mathbf{\Psi}}|\{\mathbf{Z}^{(1)}, \mathbf{Z}^{(2)}, \dots, \mathbf{Z}^{(M)}\}),$$

which implies that $\{\hat{\mathbf{\Gamma}}, \hat{\mathbf{\Psi}}\}$ is also a local optimizer of $l_2(\mathbf{\Gamma}, \mathbf{\Psi}|\{\mathbf{Z}^{(1)}, \mathbf{Z}^{(2)}, \dots, \mathbf{Z}^{(M)}\})$, with

$$\lambda = 2\sqrt{\eta} = 2\sqrt{\lambda_1 \lambda_2} \text{ and } \tilde{\delta}_k \tilde{\gamma}_{k,l}^{(m)} = \hat{\theta}_{k,l}^{(m)}.$$