# Label-Informed Non-negative Matrix Factorization with Manifold Regularization for Discriminative Subnetwork Detection

Takanori Watanabe[1(✉)], Birkan Tunc[1], Drew Parker[1],
Junghoon Kim[2], and Ragini Verma[1]

[1] Section of Biomedical Image Analysis, University of Pennsylvania,
Philadelphia, PA, USA
watanabe@uphs.upenn.edu
[2] The City College of New York, New York, NY, USA

**Abstract.** In this paper, we present a novel method for obtaining a low dimensional representation of a complex brain network that: (1) can be interpreted in a neurobiologically meaningful way, (2) emphasizes group differences by accounting for label information, and (3) captures the variation in disease subtypes/severity by respecting the intrinsic manifold structure underlying the data. Our method is a supervised variant of non-negative matrix factorization (NMF), and achieves dimensionality reduction by extracting an orthogonal set of subnetworks that are interpretable, reconstructive of the original data, and also discriminative at the group level. In addition, the method includes a manifold regularizer that encourages the low dimensional representations to be smooth with respect to the intrinsic geometry of the data, allowing subjects with similar disease-severity to share similar network representations. While the method is generalizable to other types of non-negative network data, in this work we have used structural connectomes (SCs) derived from diffusion data to identify the cortical/subcortical connections that have been disrupted in abnormal neurological state. Experiments on a traumatic brain injury (TBI) dataset demonstrate that our method can identify subnetworks that can reliably classify TBI from controls and also reveal insightful connectivity patterns that may be indicative of a biomarker.

## 1 Introduction

Substantial evidence suggests that many major psychiatric and neurological disorders are associated with aberrations in the network structure of the brain [5, 7]. With the availability of modern neuroimaging modalities such as diffusion tensor (DTI) and functional (fMRI) imaging, there is currently an exciting potential for researchers to identify connectivity-based biomarkers of disease states. Since brain networks are known to exhibit complex interactions, multivariate pattern analysis (MVPA) methods are particularly suitable here, as they aim to identify the site of the pathology by examining the data as a whole, accounting for the correlations among the network features.

In this work, we are interested in applying MVPA methods on diffusion-based structural connectomes (SCs) to identify the patterns of structural dysconnectivity induced by a brain disorder. However, due to the high dimensionality of SCs, standard MVPA methods such as the support vector machine (SVM) become prone to over-fitting and thus tend to generalize poorly to test data. Even when generalizability is achieved, SVM lacks clinical interpretability since it returns a dense, high dimensional weight vector. One way to address this is by adding an L1-regularizer to the SVM objective for feature selection [6], but this approach is known to perform poorly when the features are highly correlated. Thus dimensionality reduction becomes critical for improving classification performance and interpretability. Some well-established dimensionality reduction methods in neuroimaging include the principal and independent component analysis (PCA and ICA). However, these approaches do not preserve the non-negativity of the SCs, thus return global representations of brain network that are highly overlapping and lack interpretability since negative structural connection is biologically ill-defined.

Non-negative matrix factorization (NMF) [9] is a relatively recent method that addresses this problem by incorporating non-negativity as a constraint. This constraint leads to a more localized "parts-based" representation where the data is decomposed into purely additive combinations of non-negative basis components. For our work, the bases can be interpreted as data-driven subnetworks, and the corresponding coefficients provide a low-dimensional representation of the SC that can be used in a classifier.

However, despite its success, NMF possesses several limitations. First, NMF does not guarantee the basis components to be local and parts-based, i.e., the subnetworks may be global representations that are overlapping and redundant. Moreover, standard NMF and many of its variants are unsupervised, thus they ignore discriminative structures that may signify important group differences. Finally, NMF assumes that the data are sampled from a Euclidean space, thus does not account for the intrinsic manifold structure underlying the data. While this last issue was addressed in a recent work by Ghanbari et al. [7] under a graph-embedding framework, their method is also unsupervised and thus ignores label information. On the other hand, although supervised subnetwork detection frameworks have been introduced in some recent works [2, 8], these methods do not account for the manifold structure underlying the data.

To overcome these limitations, in this paper we introduce a novel supervised NMF framework for identifying an orthogonal set of subnetworks that is interpretable and emphasizes group differences in structural connectivity. The method also respects the intrinsic geometric structure in the data through manifold regularization [7, 10], which encourages subnetwork representations to be smooth with respect to the data manifold. To solve the proposed objective function, we introduce an optimization algorithm based on the alternating direction method (ADM), which has recently been demonstrated to solve NMF with superior performance over other state-of-the-art algorithms [12]. The proposed framework was evaluated on a TBI dataset, and the results demonstrate the interpretability and the discriminative capacity of the subnetworks.

## 2   Method

**Projective NMF.** Let $X = [x_1, \cdots, x_n]$ and $y = [y_1, \cdots, y_n]^T$ denote a set of training samples consisting of SCs $x_i \in \mathbb{R}^p_+$, $i = 1, \cdots, n$, and $y_i \in \{\pm 1\}$ indicates the label of subject $i$. An SC is a vector representation of the brain network obtained via tractography, where each vector elements represents the strength of structural connection between distinct pair of brain regions (see Sect. 3 for details). Given a target dimension $r \ll p$, NMF learns a decomposition of the form $X \approx WH$ by minimizing the Frobenius norm error $\|X - WH\|_F^2$, where $W = [w_1, \cdots, w_r] \in \mathbb{R}^{p \times r}_+$ is the basis matrix and $H = [h_1, \cdots, h_n] \in \mathbb{R}^{r \times n}_+$ is the coefficient matrix. In the context of our work, the columns of $W$ are connectivity bases that represent subnetworks.

Following [10], we assume that $H$ is obtained from a linear projection of $X$, i.e., $H = PX$, where $P \in \mathbb{R}^{r \times p}_+$ is a nonnegative projection matrix that embeds the data onto the intrinsic subspace. Under this assumption, the objective function for NMF becomes

$$\min_{W, P \geq 0} \frac{1}{2} \|X - W(PX)\|_F^2. \tag{1}$$

A key advantage of this projective NMF is that once an optimal projection $P^*$ is learned from solving (1), the trained model can be readily generalized to unseen data. That is, given a new test data $x^*$, we can immediately obtain its low dimensional representation by $h^* = P^* x^*$. This is extremely important for running cross-validation (CV).

**Orthogonal NMF with Manifold Regularization and Label Information.** Despite the merits of the projective NMF, it has three key deficiencies. Firstly, it is often reported that NMF does not necessarily return meaningful parts-based decompositions for some datasets. Secondly, although many real-world data are found to lie in a low dimensional manifold, NMF assumes that the data are sampled from a Euclidean space, neglecting the intrinsic geometric structure in the data. Thirdly, traditional NMF models are unsupervised and thus ignore the discriminative information from the different label groups.

In light of these limitations, we propose to include the following terms in our model:

1. *Orthogonality constraint*: $F_1(W) = I_\Omega(W)$, where $\Omega := \{W \in \mathbb{R}^{p \times r} | W^T W = I_r\}$ and $I_C(\cdot)$ is the indicator function of a set $C : I_C(W) = 0$ if $W \in C$ and $I_C(W) = \infty$ elsewise.

2. *Manifold regularization*: $F_2(P) = \sum_{i=1}^{n} \sum_{j=1}^{n} \|Px_i - Px_j\| S_{ij}$.

3. *Classification error*: $F_3(P, \beta, b) = \|y - (PX)^T \beta - b1_n\|_2^2$, where $\beta \in \mathbb{R}^r$ and $b \in \mathbb{R}$ defines a hyperplane in the intrinsic subspace, and $1_n \in \mathbb{R}^n$ is a vector of all ones.

The $F_1$ term constrains the basis matrix to reside within the set $\Omega$, which is the set of orthogonal matrices known as the Stiefel manifold [11]. Since $W$ is non-negative,

orthogonality implies that the bases representing the subnetworks are non-overlapping, which enhances interpretability and eliminates redundancy.

The $F_2$ term ensures smoothness of the low dimensional representation with respect to the manifold structure encoded in affinity matrix $S \in \mathbb{R}^{n \times n}$. Intuitively, this regularizer preserves the intrinsic geometric structure in the data by encouraging representations $Px_i$ and $Px_j$ to be close if $S_{i,j}$ is large, i.e., subjects $i$ and $j$ are similar under some notion. This regularizer can also be expressed in terms of the trace operator: $F_2(P) = \text{Tr}\left((PX)L(PX)^T\right)$, where $L \in \mathbb{R}^{n \times n}$ is the graph Laplacian defined by $L = D - S$, and $D$ is a diagonal matrix with $D_{i,i} = \sum_{j=1}^{n} S_{i,j} \forall i$. While the type of inter-subject relationship that can be encoded via the affinity matrix $S$ is general, in this work, we will take advantage of the clinical scores that are used to evaluate patients, and create a "disease-severity graph" to capture the disease-induced variation in the SCs. Specifically, we will assign higher value to $S_{i,j}$ if subjects $i$ and $j$ share similar severity scores.

Finally, the classification error term $F_3$ enhances the discriminatory power of NMF by encouraging the label groups in the low dimensional embedding $PX$ to be separated by a hyperplane $\beta$ (for clarity, the intercept term $b$ is dropped from our presentation hereon after). Thus, our proposed NMF model seeks to identify subnetwork bases that are not only reconstructive of data but also discriminative of label groups (note that the squared error is used here to allow the ADM algorithm to admit a closed form solution).

Integrating the above constraint terms into the projective NMF Eq. (1) gives us our final objective function ($\lambda_1, \lambda_2 \geq 0$ below are regularization parameters):

$$\min_{W,P \geq 0, \beta} \|X - W(PX)\|_F^2 + \lambda_1 \text{Tr}\left((PX)L(PX)^T\right) + \lambda_2 \|y - (PX)^T \beta\|_2^2 + I_\Omega(W). \quad (2)$$

**ADM Algorithm.** We now introduce an optimization algorithm based on the ADM algorithm [12] for solving the proposed cost function. Before applying ADM, we first convert objective function (2) into the following equivalent constrained form by introducing auxiliary variables $\{H, \tilde{H}, \tilde{W}_1, \tilde{W}_2, \tilde{P}\}$ (a technique called *variable splitting*):

$$\min_{\substack{W, P, H, \beta, \\ \tilde{P}, \tilde{H}, \tilde{W}_1, \tilde{W}_2}} \|X - WH\|_F^2 + \lambda_1 \text{Tr}\left(\tilde{H}L\tilde{H}^T\right) + \lambda_2 \|y - H^T \beta\|_2^2 + I_+\left(\tilde{W}_1\right) + I_\Omega\left(\tilde{W}_2\right) + I_+\left(\tilde{P}\right)$$

$$\text{such that } H = PX, W = \tilde{W}_1, W = \tilde{W}_2, P = \tilde{P}, H = \tilde{H},$$

where $I_+(\cdot)$ denotes the indicator function of the non-negative orthant. Although the auxiliary variables introduced from variable splitting may appear redundant, this strategy is commonly used in ADM frameworks (see [12] for example), as it allows the ADM subproblems to be solved in closed form. In the context of our work, the augmented Lagrangian (AL) function for the above constrained problem is given by:

$$\mathcal{L}_{\text{AL}}\left(W, P, \beta, \tilde{P}, H, \tilde{H}, \tilde{W}_1, \tilde{W}_2, \Lambda_{\tilde{W}_1}, \Lambda_{\tilde{W}_2}, \Lambda_{\tilde{P}}, \Lambda_H, \Lambda_{\tilde{H}}\right) = \|X - WH\|_F^2$$

$$+ \lambda_1 \text{Tr}\left(\tilde{H}L\tilde{H}^T\right) + \lambda_2\|y - H^T\beta\|_2^2 + I_+\left(\tilde{W}_1\right) + I_\Omega\left(\tilde{W}_2\right) + I_+\left(\tilde{P}\right)$$

$$+ \left\langle \Lambda_{\tilde{W}_1}, W - \tilde{W}_1\right\rangle + \left\langle \Lambda_{\tilde{W}_2}, W - \tilde{W}_2\right\rangle + \left\langle \Lambda_P, P - \tilde{P}\right\rangle + \left\langle \Lambda_H, H - PX\right\rangle + \left\langle \Lambda_{\tilde{H}}, H - \tilde{H}\right\rangle$$

$$+ \frac{\rho}{2}\left\{\|W - \tilde{W}_1\|_F^2 + \|W - \tilde{W}_2\|_F^2 + \|P - \tilde{P}\|_F^2 + \|H - PX\|_F^2 + \|H - \tilde{H}\|_F^2\right\},$$

where $\{W, P, \beta, \tilde{W}_1, \tilde{W}_2, \tilde{P}, H, \tilde{H}\}$ and $\{\Lambda_{\tilde{W}_1}, \Lambda_{\tilde{W}_2}, \Lambda_{\tilde{P}}, \Lambda_H, \Lambda_{\tilde{H}}\}$ are primal and dual variables, $\rho > 0$ is the AL penalty parameter, and $\cdot, \cdot$ denotes the trace inner product. The ADM algorhm is derived by alternately minimizing $\mathcal{L}_{\text{AL}}$ with respect to each primal variable while holding others fixed, followed by a gradient ascent step on dual variables. The overall ADM algorithm can be summarized as follows:

```
Repeat until convergence after variable initialization:
```

| Primal updates (1) | Primal updates (2) | Dual updates |
|---|---|---|
| $P \leftarrow \arg\min_P \mathcal{L}_{AL}$ | $\tilde{P} \leftarrow \arg\min_{\tilde{P}} \mathcal{L}_{AL}$ | $\Lambda_{\tilde{P}} \leftarrow \Lambda_{\tilde{P}} + \rho(P - \tilde{P})$ |
| $W \leftarrow \arg\min_W \mathcal{L}_{AL}$ | $\tilde{W}_1 \leftarrow \arg\min_{\tilde{W}_1} \mathcal{L}_{AL}$ | $\Lambda_{\tilde{W}_1} \leftarrow \Lambda_{\tilde{W}_1} + \rho(W - \tilde{W}_1)$ |
| $H \leftarrow \arg\min_H \mathcal{L}_{AL}$ | $\tilde{W}_2 \leftarrow \arg\min_{\tilde{W}_2} \mathcal{L}_{AL}$ | $\Lambda_{\tilde{W}_2} \leftarrow \Lambda_{\tilde{W}_2} + \rho(W - \tilde{W}_2)$ |
| $\beta \leftarrow \arg\min_\beta \mathcal{L}_{AL}$ | $\tilde{H} \leftarrow \arg\min_{\tilde{H}} \mathcal{L}_{AL}$ | $\Lambda_H \leftarrow \Lambda_H + \rho(H - PX)$ |
| | | $\Lambda_{\tilde{H}} \leftarrow \Lambda_{\tilde{H}} + \rho(H - \tilde{H})$ |

The primal updates above can all be carried out efficiently in closed form:

| | |
|---|---|
| $P \leftarrow \left(HX^T + \tilde{P} + [\Lambda_H X^T - \Lambda_P]/\rho\right)\left(XX^T + I_p\right)^{-1}$ | $\tilde{P} \leftarrow \max\left(0, P + \Lambda_{\tilde{P}}/\rho\right)$ |
| $W \leftarrow \left(XH^T + \rho[\tilde{W}_1 + \tilde{W}_2] - \Lambda_{\tilde{W}_1} - \Lambda_{\tilde{W}_2}\right)\left(HH^T + 2\rho I_r\right)^{-1}$ | $\tilde{W}_1 \leftarrow \max\left(0, W_1 + \Lambda_{\tilde{W}_1}/\rho\right)$ |
| $H \leftarrow \left(W^T W + 2\rho I_r + \lambda_2 \beta\beta^T\right)^{-1}\left(W^T X + \rho PX - \Lambda_H + \lambda_2 \beta y^T\right)$ | $\tilde{H} \leftarrow \left(\rho H + \Lambda_{\tilde{H}}\right)\left(\lambda_1 L + \rho I_n\right)^{-1}$ |
| $\beta \leftarrow \left(HH^T\right)^{-1} y$ | $\tilde{W}_2 \leftarrow \text{Proj}_\Omega\left(W + \Lambda_{\tilde{W}_2}/\rho\right)$ |

Note $\text{Proj}_\Omega(\cdot)$ for the $\tilde{W}_2$ update denotes the Euclidean projection of a matrix onto the Stiefel manifold. Letting $A \in \mathbb{R}^{p \times r}$ $(r \leq p)$ denote a rank-$r$ matrix, this is given by:

$$\text{Proj}_\Omega(A) = \underset{Q \in \Omega}{\arg\min} \|A - Q\|_F^2 = U\begin{bmatrix} I_r \\ 0 \end{bmatrix} V^H \tag{3}$$

Here $U\Sigma V^H$ represents the SVD of $A$ and $0 \in \mathbb{R}^{(p-r) \times r}$ is a matrix of all zeros; solution (3) is unique as long as $A$ is full column rank (see Proposition 7 in [11]).

## 3   Experiments and Conclusions

**Dataset.** We apply our method to a TBI dataset consisting of 34 TBI patients and 32 age-matched controls. While the control subjects were scanned only once, the TBI patients were scanned and evaluated at three different time points: 3, 6, and 12 months post-injury. Of the 34 TBI patients, 18 had all 3 time points, 9 had 2 and 7 had only one timepoint. The functional outcome of patients was evaluated using the Glasgow Outcome Scale Extended (GOSE) and Disability Rating Scale (DRS), which are commonly used in TBI. GOSE ranges from 1 = dead to 8 = good recovery, whereas DRS ranges from 0 = normal to 29 = extremely vegetated. In total, the dataset comprises 111 total scans, with 32 labeled *control* and 79 labeled *TBI*. All scans are accompanied with 11 clinical scores that are intended to assess the cognitive functioning of the subject.

**Creating the SCs.** DTI data was acquired for each subject (Siemens 3T TrioTim, 8 channel head coil, single shot spin echo sequence, TR/TE = 6500/84 ms, b = 1000 s/mm$^2$, 30 gradient directions). 86 ROIs from the Desikan atlas were extracted to represent the nodes of the structural network. Probabilistic tractography [3] was performed from each of these regions with 100 streamline fibers sampled per voxel, resulting in an 86 × 86 matrix of weighted connectivity values, where each element represents the conditional probability of a pathway between regions, normalized by the active surface area of the seed ROI. Finally, the 86 × 86 connectivity matrix of each subject was vectorized to its $p = 3655$ lower triangular elements, resulting in $\boldsymbol{x} \in \mathbb{R}^{\boldsymbol{p}}_+$ representing the SC.

**Implementation Details.** We applied our method to SCs computed from the TBI dataset to compute the subnetwork bases and their corresponding NMF coefficients; here we let y = + 1 indicate TBI and y = - 1 indicate control. The disease-severity graph was created using the functional outcome indices of GOSE/DRS as follows. First, we constructed a symmetrized *k*-nearest-neighbor (*k*-NN) graph with $k = 5$, where the distance between scans $i$ and $j$ was measured as $d_{i,j} = (\text{GOSE}_i\text{-}\text{GOSE}_j)^2 + (\text{DRS}_i\text{ - }\text{DRS}_j)^2$. Then a binary affinity graph was created by setting $S_{i,j}$ to 1 if and only if scans $i$ and $j$ were connected by the *k*-NN graph and did not represent the same subject (to avoid connecting same TBI patients who underwent multiple scans); controls were left un-connected.

   We identified $r = 5$ subnetwork bases using this affinity graph, and the regularization parameters were set at $\lambda_1 = \lambda_2 = 0.25$, as the model became stable around this value (degradation in classification performance was observed when parameters were set at $\lambda_1 = \lambda_2 = 0$, i.e., a setup equivalent to traditional NMF). To initialize the ADM variables, we use the strategy introduced in [4] to deterministically initialize $\boldsymbol{W}$ and $\boldsymbol{H}$ and set all other variables to zero for replicability. The AL parameter value was set to $\rho = 1000$ based on empirical test runs, and the ADM algorithm was terminated when the relative change in the objective function value (Eq. 2) at successive iterations fell below $10^{-4}$ and the following primal residual condition was met:

$$\max\left(\frac{\left\|\boldsymbol{W} - \tilde{\boldsymbol{W}}_1\right\|_F}{\|\boldsymbol{W}\|_F}, \frac{\left\|\boldsymbol{W} - \tilde{\boldsymbol{W}}_2\right\|_F}{\|\boldsymbol{W}\|_F}, \frac{\|\boldsymbol{H} - \boldsymbol{PX}\|_F}{\|\boldsymbol{H}\|_F}, \frac{\left\|\boldsymbol{H} - \tilde{\boldsymbol{H}}\right\|_F}{\|\boldsymbol{H}\|_F}, \frac{\left\|\boldsymbol{P} - \tilde{\boldsymbol{P}}\right\|_F}{\|\boldsymbol{P}\|_F}\right) < 10^{-4}.$$

To remove features that are likely non-biological, we applied feature selection using the aforementioned 11 clinical scores. Precisely, we first correlated individual SC features with each clinical score to obtain 11 separate p-value rankings (rank = 1 the smallest), and summed these rankings to obtain a rank-sum value for each feature. We then selected 1000 features having the smallest rank-sum that were then standardized via linear scaling to the range [0,1]. This feature selection and standardization procedures were conducted within the CV-folds to avoid biasing the classification performance.

We compared the performance for the following classifiers (implemented using Liblinear [6]). The first three methods are applied to the 1000 features selected using the above procedure: (1) L1-loss L2-regularized SVM (SVM), (2) L2-loss, L1 regularized SVM (SVM + L1), (3) L1-regularized Logistic regression (LogReg + L1), and (4) L1-loss L2-regularized SVM applied to the projected NMF coefficients with our method. A weighted loss function was used for all classifiers, where the weights assigned to each label class is inversely proportional to the class frequency. Since subjects have multiple timepoints, the classification accuracy was assessed using a Leave-One-Subject-Out CV (LOSO-CV) procedure, where all scans from a test subject are iteratively left out during training. Finally, the hyperparameter $C$, which is common to all classifiers, were tuned via an internal LOSO-CV over the range $C \in \{2^{-10}, 2^{-9}, \cdots, 2^{10}\}$.

### 3.1   Experimental Results and Conclusions

**Classification Results.** Table 1 reports the classification results from LOSO-CV for different methods, showing overall accuracy, specificity (type I error), sensitivity (type II error), and balanced score rate (BSR), which is the mean of specificity and sensitivity. The results show that the classification performance obtained using the proposed subnetwork features demonstrates a noticeable improvement over using the SC features in its original form, achieving accuracy of 82.0 % and a BSR of 81.8 %. The SVM achieves the next best performance, but the model is hard to interpret since all 1000 edge features contribute to the classifier. Finally, despite using a weighted loss function, we see the sparsity-promoting L1-regularized classifiers suffer from low sensitivity, which

**Table 1.** Classification results from "leave-one-subject-out" cross-validation.

| Classifier | Accuracy | Sensitivity | Specificity | BSR |
|---|---|---|---|---|
| SVM | 76.6 % | 77.2 % | 75.0 % | 76.1 % |
| SVM + L1 | 69.4 % | 73.4 % | 59.4 % | 66.4 % |
| LogReg + L1 | 67.6 % | 70.9 % | 59.4 % | 65.1 % |
| Proposed NMF + SVM | **82.0 %** | **82.3 %** | **81.3 %** | **81.8 %** |

is likely caused by data label imbalance, as well as the correlated structures among the features (a case where L1-regularizations tend to suffer).

**Effect of Manifold Regularization.** We next assessed whether the manifold regularizer with the disease-severity graph has successfully preserved the inter-patient relationship in terms of GOSE/DRS functional outcome indices. To do this, we computed Spearman's rank correlation between the subnetwork bases coefficients and GOSE/DRS indices from the 79 TBI scans. The results reported in Table 2 reveal that for all basis coefficients, consistently positive and negative correlations (statistically significant) are obtained for GOSE and DRS, respectively. This result indicates that subjects with similar level of disease-severity share similar representations in the embedding space, demonstrating the impact of manifold regularization.

**Table 2.** Spearman's correlation coefficients and corresponding p values between the $r = 5$ subnetwork basis coefficients and DRS/GOSE severity scores among TBI patients.

| Basis label | Basis coefs' correlation with DRS | Basis coefs' correlation with GOSE |
|---|---|---|
| 1 | −0.538 (p = 3.24e-7) | 0.596 (p = 6.75e-9) |
| 2 | −0.464 (p = 1.65e-5) | 0.584 (p = 1.63e-8) |
| 3 | −0.387 (p = 4.19e-4) | 0.408 (p = 1.88e-4) |
| 4 | −0.516 (p = 1.12e-6) | 0.607 (p = 3.00e-9) |
| 5 | −0.517 (p = 1.08e-6) | 0.605 (p = 3.54e-9) |

**Subnetwork Visualization.** Given the high predictive capacity of subnetwork coefficients, we next examine their corresponding subnetwork bases $W = [w_1, \cdots, w_5]$ to assess the pathological impact TBI may have induced on structural connectivity. For visualization and interpretation, we retrained the proposed NMF model using the entire dataset, and learned an SVM hyperplane $\beta \in \mathbb{R}^5$ in the corresponding embedding space. The resulting subnetworks are rendered in 3-D brain space in Fig. 1 (figures generated using Python module Nilearn [1]); the color of the edges represent the sign of the hyperplane coefficients in $\beta$, with red indicating contribution towards TBI (positive) and blue indicating contribution towards control (negative). From the figure, we can see
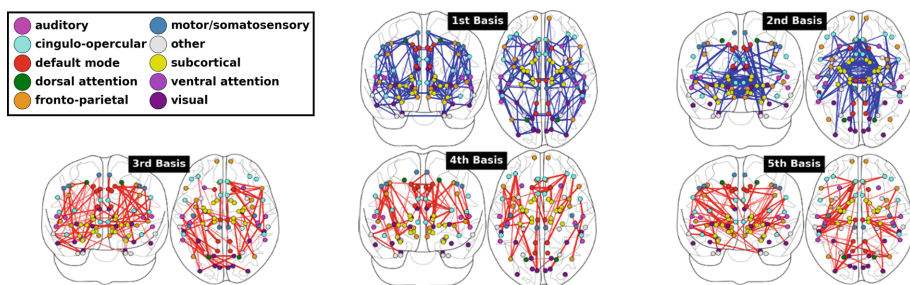


**Fig. 1.** The subnetwork bases obtained with $r = 5$. The edge color represents the sign of the corresponding hyperplane coefficient $\beta \in \mathbb{R}^r$ (blue = negative/control, red = positive/TBI).

that the network structure of the first basis exhibits strong bilateral symmetry with notable inter-hemispheric connections between the cerebellar, precuneus, and cingulate regions. Moreover, the second subnetwork basis resembles dense inter-hemispheric connections among the subcortical regions, with the sign indicating that these edges tend to be the weaker among TBI patients. On the other hand, subnetwork bases 3–5 represents connection towards TBI. Overall, the subnetworks exhibit a diffuse connectivity pattern that spans across the cortex, suggesting that damages from TBI results in a widespread disturbance in brain network. Interestingly, the connectivity patterns in the first two bases exhibit rich connectivity pattern within the subcortical and medial posterior regions, which are frequently reported to be vulnerable in TBI.

**Conclusions.** We have presented a supervised NMF framework for extracting a disjoint set of subnetworks that are interpretable and highlight group differences in structural connectivity. The method is also capable of preserving the manifold structure in the data encoded by an affinity graph, thereby respecting the intrinsic geometry of the data. Experiment on a TBI dataset shows that the subnetworks identified from our method can not only be used to reliably discriminate TBI from controls, but also exhibit tight correlation with TBI-outcome indices, indicating that subjects with similar level of TBI-severity share similar subnetwork representations due to manifold regularization.

# References

1. Abraham, A., Pedregosa, F., Eickenberg, M., Gervais, P., Mueller, A., et al.: Machine learning for neuroimaging with scikit-learn. Front. Neuroinformatics **8**(14) (2014)
2. Allahyar, A., Ridder, J.: FERAL: network-based classifier with application to breast cancer outcome prediction. Bioinformatics **31**(12), i311–i319 (2015)
3. Behrens, T., et al.: Non-invasive mapping of connections between human thalamus and cortex using diffusion imaging. Nat. Neurosci. **6**(7), 750–757 (2003)
4. Boutsidis, C., Gallopoulos, E.: SVD based initialization: a head start for nonnegative matrix factorization. Pattern Recognit. **41**, 1350–1362 (2008)
5. Cheplygina, V., Tax, D.M., Loog, M., Feragen, A.: Network-guided group feature selection for classification of autism spectrum disorder. In: Wu, G., Zhang, D., Zhou, L. (eds.) MLMI 2014. LNCS, vol. 8679, pp. 190–197. Springer, Heidelberg (2014)
6. Fan, R.E., Chang, K.W., Hsieh, C.J., Wang, X.R., Lin, C.J.: LIBLINEAR: a library for large linear classification. J. Mach. Learn. Res. **9**, 1871–1874 (2008)
7. Ghanbari, Y., Smith, A.R., Schultz, R.T., Verma, R.: Identifying group discriminative and age regressive sub-networks from DTI-based connectivity via a unified framework of non-negative matrix factorization and graph embedding. Med. Image Anal. **18**(8) (2014)
8. Kasenburg, N., et al.: Supervised hub-detection for brain connectivity. In: Proceedings of the SPIE, vol. 9784, Medical Imaging 2016: Image Processing, p. 978409 (2016)
9. Lee, D.D., Seung, H.S.: Learning the parts of objects by NMF. Nature **401**, 788–791 (1999)
10. Liu, X., et al., H.: Projective nonnegative graph embedding. IEEE Trans. Image Proc. (2010)
11. Manton, J.H.: Optimization algorithms exploiting unitary constraints. IEEE Trans. Signal Process. **50**(3), 635–650 (2002)
12. Xu, Y., Yin, W., Wen, Z., Zhang, Y.: An alternating direction algorithm for matrix completion with nonnegative factors. Front. Math. China **7**(2), 365–384 (2012)