# INFINITE FEATURE SELECTION ON SHORE-BASED BIOMARKERS REVEALS CONNECTIVITY MODULATION AFTER STROKE

*S. Obertino[1], G.Roffo[1], C. Granziera[2], G. Menegaz[1]*

[1] Dept. of Computer Science, University of Verona, Italy
[2] Martinos Center for Biomedical Imaging, Massachusetts General Hospital and Harvard Medical School, Chalestown, MA, United States

## ABSTRACT

*Connectomics* is gaining increasing interest in the scientific and clinical communities. It consists in deriving models of structural or functional brain connections based on some local measures. Here we focus on structural connectivity as detected by diffusion MRI. Connectivity matrices are derived from microstructural indices obtained by the 3D-SHORE. Typically, graphs are derived from connectivity matrices and used for inferring node properties that allow identifying those nodes that play a prominent role in the network. This information can then be used to detect network modulations induced by diseases. In this paper we take a complementary approach and focus on link as opposed to node properties. We hypothesize that network modulation can be better described by measuring the connectivity alteration directly in the form of modulation of the properties of white matter fiber bundles constituting the network communication backbone. The goal of this paper is to detect the paths that are most altered by the pathology by exploiting a feature selection paradigm. Temporal changes on connection weights are treated as features and those playing a leading role in a patient versus healthy controls classification task are detected by the *Infinite Feature Selection* (Inf-FS) method. Results show that connection paths with high discriminative power can be identified that are shared by the considered microstructural descriptors allowing a classification accuracy ranging between 83% and 89%.

*Index Terms*— Feature Selection, Stroke, 3D-SHORE

## 1. INTRODUCTION

Recent advances in diffusion MRI led to the definition of numerical parameters describing microstructural properties. The so called Ensemble Average Propagator (EAP) derived models approximate the diffusion signal by series expansion over basis functions providing an analytical solution from which indices can be derived in closed form. Among the most promising ones is the 3D-SHORE [1] model form which, under some ideal conditions, geometrical descriptors of the diffusion compartments can be inferred [2, 3].

Combined with quantitative tractography, such indices allow modeling structural connectivity through the construction of the connectivity matrix from which a representation of the network in the form of a graph is derived. Typically, cortical/subcortical regions represent the nodes and the matrix elements represent the weight of the link between pairs of regions. Then, graph theory is exploited for deriving node properties and thus identifying those playing a leading role in the network. This information can be exploited for detecting network modulations due to some pathological conditions. In this paper we take a complementary approach focusing on *connections* instead than on nodes, with the aim of identifying the paths that are more prominent in the network modulation. The reason behind this choice is the observation that pathologies compromise the communication among regions by disrupting the white matter fibers that constitute the backbone for communication. In consequence, we hypothesize that network modulation can be better described by measuring the connectivity alteration directly in the form of modulations of the properties of white matter fiber bundles. The objective is to check this hypothesis by detecting the paths that are most altered by the pathology in a feature selection paradigm. Temporal changes on connection weights are treated as features and those playing a leading role in a patients versus healthy controls classification task are detected by the *Infinite Feature Selection* (Inf-FS) method. The novelty of this approach is twofold. First, the focus is on links instead than on nodes, as previously stated. Second, the set of connections that are altered by the pathology results from a feature selection task. Previous works exploiting microstructural indices for assessing neuronal plasticity after stroke have been focusing on a predefined set of manually selected cortical and subcortical regions as the end-points of the considered links [4, 5, 6]. The proposed approach relaxes the constraint on the choice of the regions and connections such that the relevant links naturally emerge by feature selection. This has the potential of providing new insights on the changes induced by pathologies both locally and globally.

## 2. METHODS

Four microstructural indices are derived from the 3D-SHORE model. Connectivity matrices are derived from each index by quantitative tractography. Each entry of the connectivity matrix represents the absolute percent temporal variation of the mean value of each index along the connection linking the corresponding cortical/subcortical regions and plays the role of feature in the classification task. Feature selection by Inf-FS is performed for detecting the set of connections playing a dominant role in group discrimination.

### 2.1. Dataset

A total of 18 subjects (9 patients and 9 age and gender matched controls) were imaged following Diffusion Spectrum Imaging (DSI) scans [TR/TE = 6600/138 msec, FoV = $212 \times 212$ mm$^2$, 34 slices, $2.2 \times 2.2 \times 3$ mm$^3$ resolution, 258 diffusion directions, $b$-value = 8000 s/mm$^2$, $\sim 25$ min scan time] within one week (*tp1*) and one month ($\pm$ one week, *tp2*) after stroke for patients and one month apart for controls. Pre-processing was performed as in [4]. All subjects provided written informed consent and the Lausanne University Hospital review board approved the study protocol.

### 2.2. 3D-SHORE model

The SHORE model decomposes the signal $E(\mathbf{q})$ as a linear combination of basis functions that are the solutions of the 3D harmonic oscillator. The orthonormal formulation of the 3D-SHORE model is expressed as

$$\mathbf{E}(q\mathbf{u}) = \sum_{l=0,even}^{N_{max}} \sum_{n=l}^{(N_{max}+l)/2} \sum_{m=-l}^{l} c_{nlm}\Phi_{nlm}(q\mathbf{u}) \quad (1)$$

where $N_{max}$ is the maximal order of the functions in the truncated series and $\Phi_{nlm}(\mathbf{q})$ is the orthonormal 3D-SHORE basis, defined as

$$\Phi_{nlm}(q\mathbf{u}) = \left[\frac{2(n-l)!}{\zeta^{3/2}\Gamma(n+3/2)}\right]^{1/2} \left(\frac{q^2}{\zeta}\right)^{l/2} \quad (2)$$
$$\times \quad \exp\left(\frac{-q^2}{2\zeta}\right) \times L_{n-l}^{l+1/2}\left(\frac{q^2}{\zeta}\right) Y_l^m(\mathbf{u})$$

where $\Gamma$ is the Gamma function and $\zeta$ is a scaling parameter dependent on the diffusion time $\tau$ and the diffusivity $D$. All these models provide close approximations of the diffusion signal and allow deriving important EAP features such as the Orientation Distribution Function (ODF) in a reliable manner. From EAP models, it is also possible to derived microstructure related indices, namely the Return To the Origin Probability (RTOP), the Return To the Axis Probability (RTAP), and the Return To the Plane Probability (RTPP). RTOP, RTAP and RTPP represent the zero net displacement probabilities in

the three, two and mono-dimensional cases, respectively. Under some ideal conditions regarding the acquisition protocol they provide an estimation of the mean pore geometry being respectively proportional to the reciprocal of the mean volume, cross-sectional area and length of the pore [1]. In this work, only RTAP was used

$$RTAP = \int_{\mathbb{R}^2} E(\mathbf{q}_{\perp})d^2\mathbf{q}_{\perp} = \int_{\mathbb{R}} P(\vec{r}_{\parallel})dr \quad (3)$$

Additionally, Generalized Fractional Anysotropy (GFA) and Propagator Anysotropy (PA), expressing the distance between the EAP and its isotropic component, were calculated. These indices reflect the degree of restriction of the water molecules in the voxel, which is directly linked to the underlying pore shape.
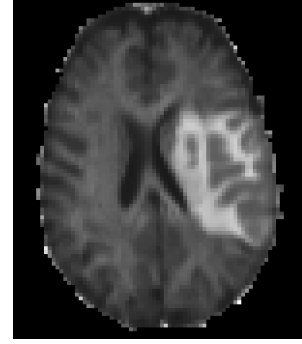


**Fig. 1**: Image of RTAP in presence of a stroke lesion.

### 2.3. Feature extraction by tract-based analysis

The ODFs were reconstructed and fiber-tracking was performed via a streamline algorithm (www.cmtk.org). Brain segmentation was performed using the Desikan-Killany atlas provided by Freesurfer (www.surfer.nmr.mgh.harvard.edu) resulting in as a set of 32 cortical and 7 subcortical regions. The mean value of the microstructural indices along the fibers were calculated and the absolute percentage changes across time points were used for generating one connectivity matrix of size $39 \times 39$ per subject per index as

$$\Delta_{t_{p12,i}} = \frac{|F_{i,tp1} - F_{i,tp2}|}{F_{i,tp1}} \quad (4)$$

where $F_i$ denotes the index in use.

### 2.4. Feature Selection

Let $G_n = (V, E)$ be an undirected graph where $V$ is the set of vertices corresponding to ROIs, and $E$ codifies weighted edges among connections. We represent the graph $G_n$ by the adjacency matrix $D_n$, where each element $d_{ij}^n$, $1 \leq i, j \leq N$, $N = 39$ is the corresponding entry of the connectivity matrix

of subject $n$, $n = 1, \ldots, 18$. In order to measure how well each connection separates the two classes of patients ($P$) and controls ($C$), we define a discriminant matrix $M$ by using a simple heuristic for measuring class separation. The heuristic is based on the separation of the class means. In the two-class problem at hand there are $N_c = 9$ adjacency matrices of class $C$ and $N_p = 9$ adjacency matrices of class $P$ for each microstructural index RTAP, R, GFA and PA. For each entry, that is for each feature, the mean and variance are estimated across subjects to generate the matrix $M$ whose entries are

$$M_{i,j} = \frac{\mu_{i,j}^C - \mu_{i,j}^P}{(\sigma_{i,j}^C)^2 + (\sigma_{i,j}^P)^2} \quad (5)$$

where

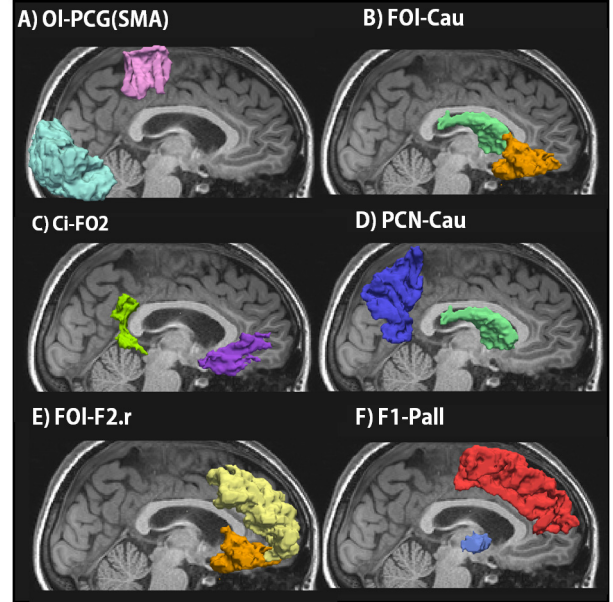$$\mu_{i,j}^k = \frac{1}{N_k} \sum_{n \in N_k} d_{i,j}^n, \qquad k \in \{C, P\}.$$

In the same way, we calculate the standard deviation vectors $\sigma_{i,j}^k$ for each feature $d_{i,j}^k$ of class $k$.

Our approach proposes to rank the features by importance regarding the patients versus controls classification task. To this end, we use the matrix $M$ as input of the infinite feature selection (Inf-FS) [7] algorithm, where the percent absolute changes of the microstructural values along the connections are seen as features. By construction, the Inf-FS method allows to use convergence properties of the power series of matrices, and evaluate the relevance of a feature with respect to all the other ones taken together. In the Inf-FS formulation, each path of a certain length $l$ over the graph is seen as a possible selection of features. Letting these paths tend to an infinite number permits the investigation of the importance of each feature. As a result, this method assigns a score of "importance" to each feature by taking into account all the possible feature subsets, therefore the higher the final score, the most important the feature. In this work a simplified version where only the Fisher distance of the features across classes was used. The final rank was then used in our experimental section, where we proved that the selected connections turn out to be effective from the classification point of view. In order to obtain some measure of relevance of the subset of features (connections), a classification approach was followed. Performance was defined in terms of accuracy, precision and recall. Moreover, the ROC curve was obtained as well as the corresponding the area under the curve (AUC). Training and testing pools were created using a cross-validation leave-1-out method, while a SVM was used for classification.

## 3. RESULTS & DISCUSSION

Figure 3 illustrates the performance of the classifier as a function of the number of features that are retained after the Inf-FS based selection in terms of accuracy, precision, recall and area under the curve (AUC). Good performance was obtained using a relatively low number of features, suggesting that few

key connections could be the key for discriminating patients from controls. Among the set of the first 20 features, six were common to the four indices. These correspond to connections between the pairs of regions illustrated in Figure 2.
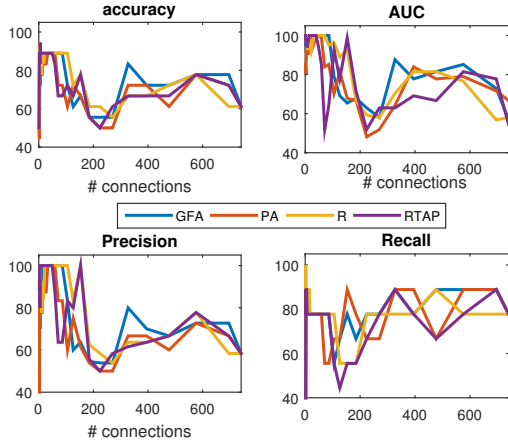


**Fig. 2**: A) lateral occipital (Ol) - para central (PCG) (as supplementary motor area); B) lateral orbito frontal (Fol) - Caudate (Cau); C) isthmus cingulate (Ci) - medial orbito frontal (F02); D) pre cuneus (PCN) - Caudate (Cau); E) lateral orbito frontal (Fol) - rostral middle frontal (F2.r); F) superior frontal (F1) - Pallidum (Pall)

Reducing the feature set to this ensemble the classification performance is slightly degraded especially for GFA and RTAP. However, the still good performance could be an indication of the relevance of such connections in the considered task, pointing to a network modulation involving areas in different cortical and subcortical regions.

For the sake of comparison, Table 2 provides the performance of the classification algorithm when using the 23 connections involving the cortical and subcortical motor loops manually selected as in [4, 6]. As it can be observed, the discriminative power of those features is lower than the that obtained using the same number of features that are first ranked by the Inf-FS algorithm that are reported in Table 1. This could suggest that a more extended portion of the network is involved in the plasticity process and thus that a wider perspective should be taken for its assessment.

Finally, Table 4 shows the performance that is obtained by gathering the six features common to all indices together. Accuracy is not significantly affected while precision, recall and AUC reach the maximum value.

**Fig. 3**: Performance of classification after Inf-FS from absolute delta adjancecy matrix.

**Table 1**: Classification performance on the 23 first ranked features following Inf-FS.

| Index | Accuracy | AUC | Precision | Recall |
|---|---|---|---|---|
| GFA | 88.89 | 97.53 | 100 | 77.78 |
| PA | 83.33 | 92.59 | 87.50 | 77.78 |
| R | 88.89 | 97.56 | 87.50 | 77.78 |
| RTAP | 88.89 | 100 | 100 | 77.78 |

**Table 2**: Classification performance on the 23 manually selected features as in [4, 5, 6].

| Index | Accuracy | AUC | Precision | Recall |
|---|---|---|---|---|
| GFA | 66.67 | 56.79 | 61.53 | 88.89 |
| PA | 50 | 41.98 | 50 | 33.33 |
| R | 55.56 | 50.62 | 55.56 | 55.56 |
| RTAP | 50 | 54.32 | 50 | 22.22 |

**Table 3**: Classification performance on the 6 features within the first 20 ranked by Inf-FS for each index.

| Index | Accuracy | AUC | Precision | Recall |
|---|---|---|---|---|
| GFA | 83.33 | 95.06 | 80 | 88.89 |
| PA | 83.33 | 96.30 | 87.50 | 77.78 |
| R | 83.33 | 97.53 | 87.50 | 77.78 |
| RTAP | 77.78 | 95.06 | 77.78 | 77.78 |

**Table 4**: Classification performance on the 6 features within the first 20 ranked by Inf-FS for all the indices.

| Accuracy | AUC | Precision | Recall |
|---|---|---|---|
| 88.89 | 98.77 | 100 | 77.78 |

## 4. CONCLUSIONS

In this paper we proposed a novel approach for the characterization of the structural connectivity network after stroke focusing on link as opposed to node properties. The basic hypothesis is that network modulation can be better described by measuring the connectivity alterations directly in the form of modulation of the properties of white matter fiber bundles constituting the communication backbone. A feature selection paradigm was exploited where features were represented by percent absolute changes of mean values of a predefined set of microstructural indices across connections between pairs of regions. Results show that connection paths with high discriminative power can be identified that are shared by the considered microstructural descriptors allowing a classification accuracy ranging between 83% and 89% for the different indices.

## 5. REFERENCES

[1] E. Ozarslan, C. Koay, T. Shepherd, S. Blackband, and P. Basser, "Simple harmonic oscillator based reconstruction and estimation for three-dimensional q-space MRI," in *ISMRM*, 2009.

[2] E. Ozarslan, C. Koay, T. Shepherd, M. Komlosh, M. Irfanolu, C. Pierpaoli, and P. Basser, "Mean apparent propagator (map) mri: A novel diffusion imaging method for mapping tissue microstructure," *Neuroimage*, vol. 78, pp. 16–32, 2013.

[3] M. Zucchelli, E. Garyfallidis, M. Paquette, S. Merlet, G. Menegaz, and M. Deascoteaux, "Comparison between discrete and continuous propagator indices from cartesian q-space dsi sampling," in *ISMRM*, 2014.

[4] C. Granziera, A. Daducci, D.E. Meskaldji, A. Roche, P. Maeder, P. Michel, N. Hadjikhani, A.G. Sorensen, R.S. Frackowiak, J.P. Thiran, R. Meuli, and G. Krueger, "A new early and automated MRI-based predictor of motor improvement after stroke," *Neurology*, vol. 79, no. 1, pp. 39–46, 2012.

[5] Y. Lin, A. Daducci, D. Meskaldji, J. Thiran, P. Michel, R. Meuli, G. Krueger, G. Menegaz, and C. Granziera, "Quantitative Analysis of Myelin and Axonal Remodeling in the Uninjured Motor Network After Stroke," *Brain Connectivity*, vol. 5, no. 7, pp. 401–412, 2015.

[6] L. Brusini, S. Obertino, M. Zucchelli, I. Boscolo Galazzo, G. Krueger, C. Granziera, and G. Menegaz, "Assessment of Mean Apparent Propagator-based indices as biomarkers of axonal remodeling after stroke," in *MICCAI*, 2015.

[7] Giorgio Roffo, Simone Melzi, and Marco Cristani, "Infinite feature selection," in *ICCV*, 2015.