

Unravelling communities of ALS patients using network mining

André V. Carreiro

Sara C. Madeira

Alexandre P. Francisco

INESC-ID and Instituto Superior Técnico, Technical University of Lisbon, Portugal
{acarreiro,smadeira,aplf}@kdbio.inesc-id.pt

ABSTRACT

Amyotrophic Lateral Sclerosis is a devastating neurodegenerative disease characterized by a usually fast progression of muscular denervation, generally leading to death in a few years from onset. In this context, any significant improvement of the patient's life expectancy and quality is of major relevance. Several studies have been made to address problems such as ALS diagnosis, and more recently, prognosis. However, these analysis have been mostly restricted to classical statistical approaches used to find the most associated features to a given outcome of interest. In this work we explore an innovative approach to the analysis of clinical data characterized by multivariate time series. We use a distance measure between patients as a reflection of their relationship, to build a network of patients, that in turn can be studied from a modularity point of view, in order to search for communities (groups of similar patients). Preliminary results show that it is possible to extract relevant information from such groups, each presenting a particular behavior for some of the features (patient characteristics) under analysis.

Categories and Subject Descriptors

G.2.2 [Graph Theory]: Network problems; H.2.8 [Database applications]: Data Mining; H.3.3 [Information Search and Retrieval]: Clustering; J.3 [Computer Applications]: Life and Medical Sciences

General Terms

Algorithms, Experimentation

Keywords

Graph mining, ALS, patient network, patient similarity.

1. INTRODUCTION

Included in the category of neurodegenerative diseases, Amyotrophic Lateral Sclerosis (ALS) is characterized by a

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

KDD-DMH'13 August 11, 2013, Chicago, Illinois, USA.
Copyright 2013 ACM 978-1-4503-2174-7/13/08 ...\$15.00.

progressive muscular deterioration, leading to a (usually) fast progressing weakness due to the denervation of critical skeletal and respiratory muscles, which can ultimately result in death [1]. Nonetheless, while devastating from a motor impairment point of view, ALS is somewhat (or even completely) unconnected to any cognitive decline [2]. Taking these aspects into account, and given that ALS currently ranks third in the incidence of neurodegenerative diseases, maintaining, or even improving, the patients' quality of life is a problem of crucial importance.

1.1 Related Work on Data Mining in ALS

In the context of ALS, the related work is mostly associated to a population based approach, focusing on common features significantly associated to reduced survival. In fact, we can divide the ALS studies in two problems. The first is related to the patients' diagnosis, investigating the heterogeneity in ALS subtypes [3], or the relevance of certain clinical features in the diagnosis, such as the paraspinal muscle EMG and motor-unit potentials (MUP) [4].

The second problem concerns the prognostic prediction, which, in turn, can be divided in two different analysis. The most explored is the study of ALS survival, and the most significantly associated features, including respiratory measures [5], but also the site of onset [6] and the ALS Functional Rating Scale (ALSFRS) score [6]. The other, least explored, type of studies is related to the prediction of auxiliary respiration requirement, either with tracheostomy, or non-invasive ventilation (NIV).

A patient-driven model for ALS prognosis prediction of respiratory failure has been recently proposed [7]. We note, however, that the strategies adopted for the most part of these studies rely on statistical tests, Kaplan-Meier survival tables, and multivariable Cox proportional hazard regression models, which are typical of population based studies.

1.2 Related Work on Patient Similarity

The topic of patient similarity has been increasingly explored in the last years, mainly motivated by the need of patient cohort identification for comparative clinical trials or decision support applications [8], but also due to the arising of the so called patient social networks, where patients try to find people with similar experiences and conditions. Nevertheless, patient similarity, or distance, poses several different challenges, where the subjective notion of similarity rises as one of the most critical. In fact, each physician may have a different perspective about how similar two patients are, as they assign different weights to different features. Some studies have proposed a way of learning these

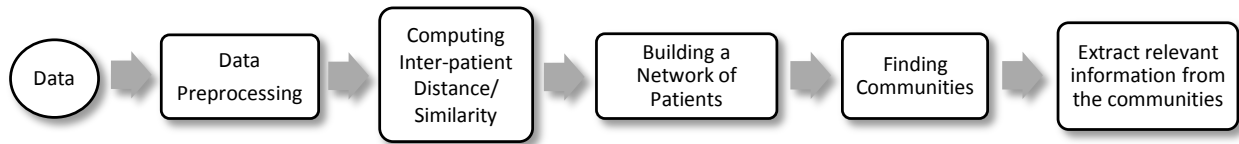


Figure 1: Workflow used in this work for community finding and interpretation using a network of ALS patients.

weights automatically [9], and others suggested that the expert knowledge can be integrated, and used to learn a new unified similarity measure [8]. Nonetheless, in the network context, research has been focused on the visualization, or on the aforementioned social network, where patient similarity is used to build a recommendation system for the user [9]. Moreover, these studies rely on supervised learning approaches, and it is not clear how temporal data can be addressed.

In this context, this work proposes a novel unsupervised learning strategy using a distance measure capable of dealing with multivariate time series in order to build a network of patients, which can then be analyzed from a modular point of view. The found modules, or patient communities, can be studied according to their particular characteristics, possibly reflecting ALS subtypes, which might help to better understand the disease. Moreover, such modules can then be used in a supervised learning fashion to train expert models for discriminating subgroups of patients.

2. METHODS

This section presents and discusses the methodology used in this work. The workflow is in Figure 1.

2.1 Data and Preprocessing

In this work we use the data made available at the DREAM-Phil Bowen ALS Prediction Prize4Life Challenge (or ALS Prediction Prize) (<https://www.innocentive.com/ar/challenge/9933047>), where the goal was to predict the rate of disease progression as measured by the slope of the functional scale ALSFRS between the 3 and 12 months after the beginning of the trial. However, in this paper we go beyond the regression problem, moving towards the search for groups of similar patients, or communities. We note that this dataset was later extended, and now counts over 8500 clinical ALS patient records, available in the Pro-ACT database (<https://nctu.partners.org/ProACT>). This particular dataset, consists of information for 918 ALS patients, and comprises features such as demographics, discriminated ALSFRS features, medical and family history, respiratory measurements and other general lab data. Given that these data come from merged trials, the first step of preprocessing consisted in cleaning the outlier values (e.g. unrealistic age or weight), followed by unit standardization (e.g. height in *cm* or *in*). This is a major challenge in clinical data and reinforces the need for a standardized input and storage model, as it is being slowly implemented with electronic health records (EHR). Another challenge was the presence of missing values. In fact, many of the available features present a missing value percentage close to or over 90%. We chose to proceed with the features with less than 50% of missing values. For the experiments that required it, we per-

formed missing values estimation using a probabilistic PCA method implemented in Matlab [10]. One important aspect of this dataset is that it can be divided into two distinct subsets: static data, including demographics, medical and family history as well as the disease onset data; and temporal data, with all the measurements varying between appointments, as the ALSFRS, respiratory features and other lab data. The number of evaluations (time points) spans from 7 to 23, with an average of 13 evaluations per patient.

2.2 Distance/Similarity Measures

In order to build a network of patients, we need a measure which reflects the relationship between two patients. However, since we have both static and temporal partitions, in the latter we are dealing with multivariate time series (MTS). Thus, and according to the considered partition, we use different distance metrics to summarize how closely related two different patients are.

2.2.1 Static data

In the static partition, each patient is represented by a row containing the values of the static features including type and time of onset, first symptoms, demographics, etc. As such, we compute the normalized difference for numerical features and perform a string comparison for categorical. We note, however, that when a feature has a set of simultaneous categories, we compute the minimum proportion of shared strings (“Weakness, Speech” vs. “Weakness” results in a distance score of 1/2). The static distance is then the average value computed for all static features.

2.2.2 Eros Distance for Temporal Data

For temporal data, characterized by MTS with different number of time points, we use the Extended Frobenius Norm (Eros) as a distance measure [11]. Essentially, this measure is based on the application of PCA to the MTS data, thus generating principal components which in turn can be used to compute a distance based on the angle formed by the corresponding principal components for both patients. We note that for this measure, the missing values are estimated as aforementioned. This measure was shown to outperform the more traditional distance measures for MTS (for more details refer to [11]). In this work, the Eros distance was normalized so that it would be contained in the interval $[0,1]$.

2.2.3 Combined Distance

We consider that a global distance measure can be a simple average of the distance computed for both partitions, although a weighted version can be used to assign a higher relevance to one type of data. Assigning different weights to the individual features, either automatically or based on expert knowledge, is also possible.

2.3 Building a Network of Patients

After computing the distance matrix for all patients, we can build the network (graph) using the distances. However, in network analysis, it is usual to use similarities instead of distances, and thus we get these in two distinct ways. The first is computing the similarities as $S_{ij} = 1 - D_{ij}$, where S is the similarities matrix and D the distances matrix. The other is based on the binary adjacencies matrix A , defined as

$$A_{ij} = 1 \iff D_{ij} \leq \tau, \quad (1)$$

where τ is a given threshold. Both matrices can be used as input to Gephi [12], an open-source software for network analysis. In the network context, the graph nodes represent the patients, whereas the edges represent their connection. When the similarities matrix is used, each edge has an associated weight representing the similarity between both patients. Eventually, some edges with lower similarities can be filtered out. In the case where we use the adjacencies matrix, each edge states that the two patients it connects are similar (their distance is below the threshold τ). With this tool, we can quickly analyze and visualize the built networks, apply several filters to delete outlier nodes or edges, and compute metrics such as the network density, modularity, average path length and many others.

2.4 Finding Communities in a Network of Patients

The goal of this work is to investigate if we can find communities within networks of patients, each presenting particular characteristics which might bring new insights to the disease. Hence, the most relevant metric is the modularity, since it is related to how well the network can be divided into modules. Essentially, modularity can be seen as the difference between the number of edges within identified communities in given network and the random expectation [13]. A higher value means that the network presents a more modular structure and, hence, vertices in each community are more similar. The algorithm used in Gephi is time-efficient [14], and assigns a modularity class (or community) label to each patient, which we can then visualize with different colors.

2.5 Feature Selection and Retrieving Meaningful Information

Gephi allows us to export the similarities/adjacencies matrix together with the assigned community label. Nevertheless, we still have to extract relevant information from the communities found. Which features (patients characteristics) are more important to these modules? Can this information be used to help the clinicians in the disease diagnosis and/or prognosis? In order to answer these questions, we started by performing feature selection to assess which features were more correlated with the discovered communities. We chose to use the minimum redundancy maximum relevance (mRMR) method [15], with a parallelized implementation in R: mRMRe package [16]. The basic idea is to select features that are highly correlated with the target class (the community) and mostly uncorrelated between themselves. Besides the selected features, this implementation can also return the causality values for each feature (as defined in [16]), where more negative scores indicate a stronger putative causality of the feature to the target. We

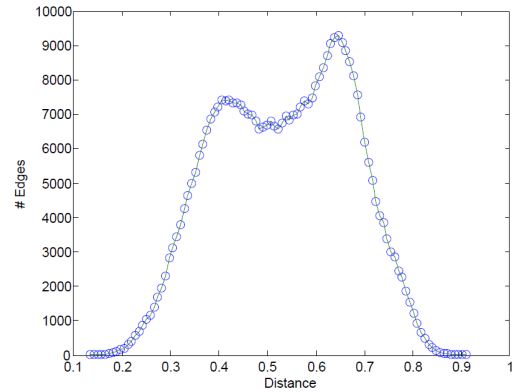


Figure 2: Distribution of number of edges vs. distance (average of static and temporal - Eros - distances).

started by applying this method to the dataset with all the communities, and then proceeded with a pairwise analysis, searching for features that better explain the differences between each pair of modules.

3. RESULTS

As shown in the previous section, we used a distinct distance measure for static and temporal data, and the final distance was the average for the two values, although weights can be used, whether to reinforce some type of data or, more specifically, to give more importance to certain features. Resorting to Gephi, we analyzed the two types of networks, either based on the similarities or adjacencies matrices. In fact, when filtering the edges with a given weight threshold to match the threshold used for the creation of the adjacencies matrix, the results of modularity are very similar, as well as the found communities. As such, we decided to present here the results obtained when adjacencies matrices are used.

Figure 2 shows the distribution of the computed distances by the number of links between patients (or edges). This distribution seems to follow a mixture of two Gaussian distributions, with peaks around the distance of 0.4 and 0.65, approximately.

3.1 Community Finding

We tested several values for the threshold τ , while assessing the modularity metric, as well the size of the found communities and of the giant connected component (largest connected group of patients). In fact, with lower values of τ , we end up with many disconnected patients. These net-

Table 1: Modularity, Giant Component (GC) size and Edges statistics vs. distance threshold τ

τ	Mod.	GC	# Edges	Weight ($\mu \pm \sigma$)	σ^2/μ
0.25	0.67	383	3791	0.23 ± 0.019	0.0016
0.30	0.52	647	15362	0.27 ± 0.028	0.0029
0.35	0.35	773	41632	0.31 ± 0.036	0.0042
0.40	0.23	862	83074	0.34 ± 0.045	0.0060
1.00	0.05	918	420903	0.53 ± 0.136	0.0349

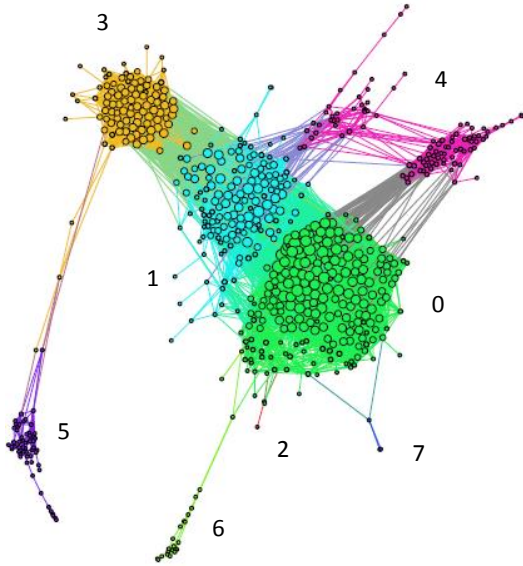


Figure 3: Layout of the patients network when using an adjacency matrix with distance threshold τ of 0.3, and using only the giant component. Different colors represent different patient communities.

work nodes can be filtered out, and a usual way to achieve this is to keep only the giant component, although we tested also with smaller disconnected groups of patients. In these situations, however, there are never disconnected groups of significant size (over 10 patients), and so we proceeded with only the giant component.

Table 1 shows the modularity value and giant component size for each distance threshold τ (or similarity threshold $1 - \tau$). From these results, we chose to proceed with a τ distance threshold of 0.3, since it presented the best ratio $GCsize/(1 - Modularity)$, as we want to preserve as many patients as possible, while achieving a good modularity. As explained in Section 2.4, the modularity metric was computed and communities were assigned using Gephi. One of the strengths of this network analysis tool is that we can color the nodes with different colors corresponding to the different modules, and use an automatic layout to better visualize them. Such result is shown in Figure 3, for an adjacency matrix with distance threshold of 0.3 and keeping only the giant component, which resulted in 8 different communities (although two of them are very small, with only 2 patients).

Regarding this subnetwork, when we compare its index of dispersion ($\sigma^2/\mu = 0.0029$) to the one of the original network (0.0349) we can see that the latter is approximately 12 times larger, although this value already suggests that the original network is under-dispersed. Edges in the filtered network are only within the 2.5% of edges among more similar nodes, which provides a good degree of confidence on our ability to avoid spurious connections.

3.2 Feature Selection

As aforementioned in Section 2.5, after the communities are discovered, we still cannot understand what particular characteristics each group presents. To achieve this goal,

we use the classic implementation of the mRMR feature selection method in the mRMRe package [16]. The first analysis was performed with all the communities, trying to find which features best explain the different groups. Table 2 shows the causality values, according to [16], computed for the top 10 features (more negative values). Most of the features presenting a more putative causality to the communities are related to the functional scale for the disease (ALSFRS). However, we can see that the Treatment Group Delta ranks first with the Swallowing portion of the functional scale ALSFRS. This is interesting, since this feature provides information not only about the time when the patient started undergoing treatment or placebo, but also if it participated in such a trial, since a missing value means that the patient was not assigned to a Treatment or Placebo group. Other features include the Respiratory rate and Gamma-Glutamyltransferase. We now know which features are more relevant to the formation of the discovered communities. Nonetheless, we decided to proceed with a pairwise analysis, so we can better identify which features are best correlated with each pair of patient communities. We note that only 6 of the 8 discovered communities were investigated. Two of them (2 and 7) were too small to be considered (only 2 or 3 patients). Due to space constraints, it is not possible to show all the features and causalities for each pair of modules, and since a more negative causality value is associated with a stronger putative causality to the community assignment, Table 3 shows only the 1st and 10th causality values, in order to provide a general idea of which pairs of communities are more easily distinguished.

Table 3 suggests that some pairs are more separable than others. For instance, community 3 seems to be the more separable community, given that the causality values of pairs containing these patients are always much more negative. Moreover, modules 5 and 6 appear to be somewhat distinct from each other. On the other hand, community 4 seems to be less separable than the others, although it may be more distinct from communities 3 and 5. Nevertheless, after this analysis we still cannot answer the more specific question: What particular characteristics do these communities present?

3.3 Studying the Communities

To analyze what characterizes the different communities and distinguishes them from each other, we studied the distribution of the evaluations for each of the features selected in the previous step. We restricted the shown results to the first four features, which can be seen in Figure 4.

The distribution for the Treatment Group Delta feature is the most peculiar, since the modules 0 and 6 only presented null values for this feature, meaning that these patients did not participate in the trial where individuals were divided into Treatment and Placebo groups. For the other four communities, we can observe that almost all the patients in module 3 were assigned into Treatment or Placebo at day 0, which is very different from modules 1 and 4, where patients were only assigned beyond the 100 days from the beginning of the trial. In what concerns the feature Swallowing of the functional scale, we can see an almost clear separation between two sets of modules: 1, 3 and 4 present a lower functional score than modules 0, 5 and 6. In a similar analysis, for features representing the Salivation and Speech scores, modules 3 and 4 present a lower value, while the pa-

Table 2: Causalities for the 10 top scoring features. More negative causality values mean a stronger putative causality of the feature to the target (community) [16].

Treatment Group Delta	3.Swallowing	2.Salivation	1.Speech	10.Respiratory	8.Walking	4.Handwriting	Respiratory Rate	Gamma-Glutamyl-transferase	9.Climbing stairs
-0.1408	-0.1408	-0.1237	-0.0829	-0.0427	-0.0285	-0.0190	-0.0158	-0.0130	-0.0121

Table 3: Causalities for the 1st and 10th top scoring features. More negative causality values mean a stronger putative causality of the feature to the target (community) [16].

	0 vs 1	0 vs 3	0 vs 4	0 vs 5	0 vs 6	1 vs 3	1 vs 4	1 vs 5	1 vs 6	3 vs 4	3 vs 5	3 vs 6	4 vs 5	4 vs 6	5 vs 6
1 st	-0.047	-1.960	-0.058	-0.197	-0.127	-1.820	-0.052	-0.307	-0.303	-0.435	-1.731	-1.456	-0.204	-0.085	-0.810
10 th	-0.025	-0.372	-0.008	-0.084	-0.031	-0.404	-0.022	-0.010	-0.080	-0.193	-0.218	-0.487	-0.073	-0.054	-0.138

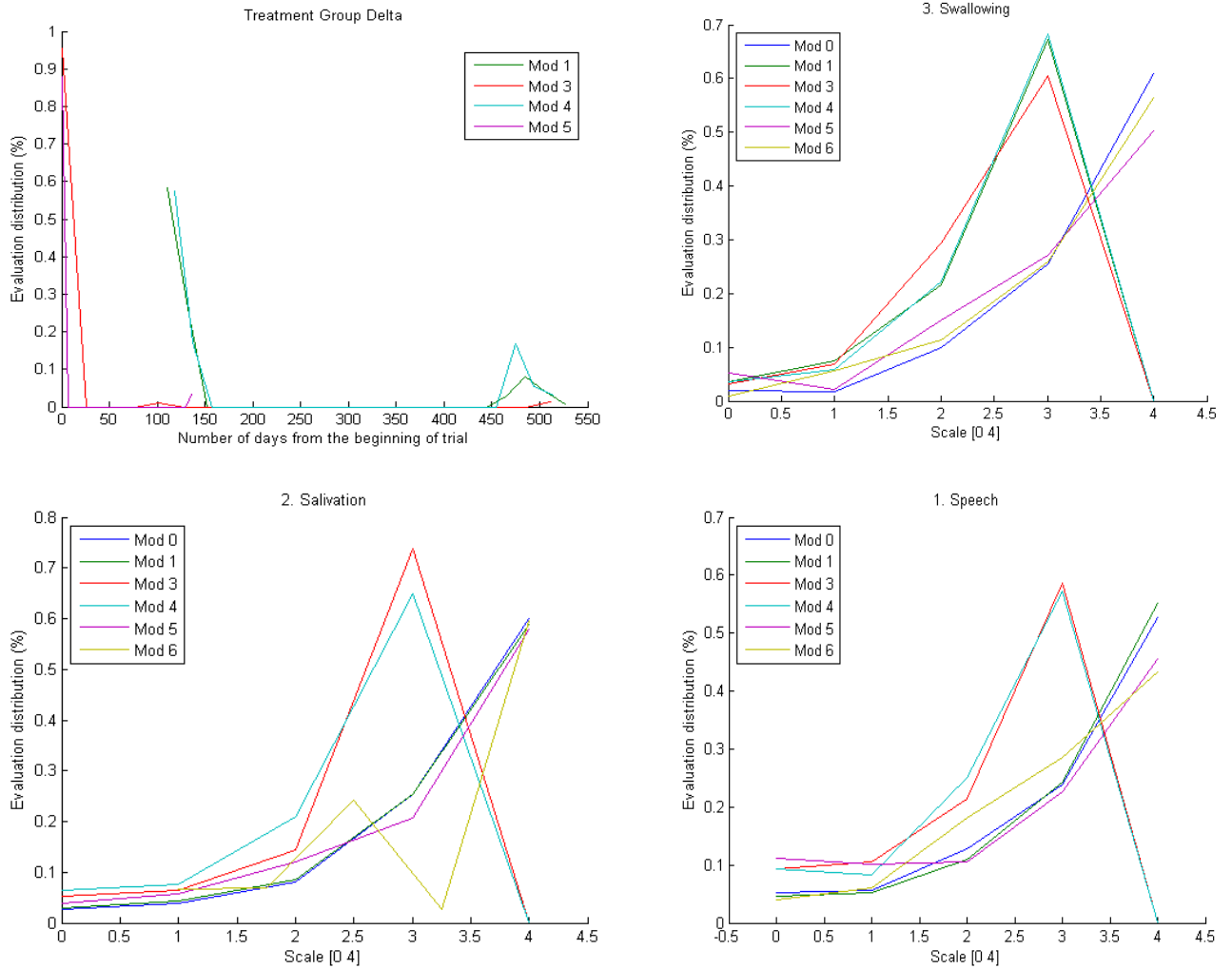


Figure 4: Normalized distribution of the patients' evaluations, according to the discovered modules (Mod), for the 4 top scoring features.

tients from module 1 are now in the same set as 0, 5 and 6, presenting a better salivation and speech conditions. Using only the four highest ranked selected features in terms of causality, we can already extract some characteristics that can distinguish between some of the modules. However, it

seems to be the combination of all these features which allows the modular structure in the network shown in Figure 3. In what concerns the pairwise analysis, we decided not to show new figures of distribution due to space constraints, and also because the features shown in Figure 4 are already

present when we compare the pairs of communities, with the characteristics already discussed.

4. CONCLUSIONS AND FUTURE WORK

In this work we propose a novel unsupervised learning approach to find communities in a network of patients, and then extract meaningful information from the discovered modules of similar patients. One of the crucial aspects of this study is the choice of the distance measure, as we deal with both static and temporal data. Nonetheless, Eros distance seems to be a suitable choice to handle temporal data, preserving particular characteristics which can be retrieved later, as shown in Figure 4. One of the future directions will be to further investigate what weights should be given to each type of data, and to each particular feature. This weighting could be performed automatically [9], or based on expert knowledge. Moreover, we would like to proceed with a feedback analysis, and recalculate the distance/similarity between patients using only the features we found relevant for the modular structure. Methodologies such as this look promising in terms of knowledge discovery with little or no prior knowledge, where the conclusions are achieved in a totally unsupervised fashion, and may help to gain new insights on different diseases. Furthermore, we intend to apply this strategy on different datasets, and if possible, to explore a supervised learning approach that could allow us to learn the distance metric from the data [8], and search for communities of different classes. Finally, such modules could be used to train expert models for classification problems regarding subgroups of patients, possibly discriminating the ones with different disease progression rates.

5. ACKNOWLEDGMENTS

This work was supported by national funds through Fundação para a Ciência e a Tecnologia (FCT), under projects PEst-OE/EEI/LA0021/2013, Neuroclinomics (PTDC/EIA-EIA/111239/2009) and NetDyn (PTDC/EIA-CCO/118533/2010), and a doctoral grant SFRH/BD/82042/2011 to AVC.

6. REFERENCES

- [1] Cudkowicz, M., Qureshi, M. and Shefner, J., *Measures and markers in amyotrophic lateral sclerosis*, NeuroRx : the journal of the American Society for Experimental NeuroTherapeutics, 1:2 (2004), pp. 273–83.
- [2] Stein, J., Schettler, T., Rohrer, B., Valenti, M., *Environmental Threats to Healthy Aging With a Closer Look at Alzheimer’s & Parkinson’s Diseases*, Environmental Health, (2008), pp. 202.
- [3] Turner, M. R., Scaber, J., Goodfellow, J., Lord, M. E., Marsden, R. and Talbot, K., *The diagnostic pathway and prognosis in bulbar-onset amyotrophic lateral sclerosis*, Journal of the neurological sciences, 294:(1-2) (2010), pp. 81–85.
- [4] Carvalho, M., Pinto, S. and Swash, M., *Motor unit changes in thoracic paraspinal muscles in Amyotrophic Lateral Sclerosis*, Muscle & Nerve, 39:1 (2009), pp. 83–86.
- [5] Baumann, F., Henderson, R.D., Morrison, S.C., Brown, M., Hutchinson, N., Douglas, J., Robinson, P.J. and McCombe, P., *Use of respiratory function tests to predict survival in amyotrophic lateral sclerosis*, Amyotrophic lateral sclerosis : official publication of the World Federation of Neurology Research Group on Motor Neuron Diseases, 11:(1-2) (2010), pp. 194–202.
- [6] Kollwe, K., Mauss, U., Krampfl, K., Petri, S., Dengler, R. and Mohammadi, B., *ALSFRS-R score and its ratio: A useful predictor for ALS-progression*, Journal of the Neurological Sciences, 275 (2008), pp. 69–73.
- [7] Amaral, P. M. T., Pinto, S., de Carvalho, M., Tomás, P. and Madeira, S.C. *Predicting the need for non-invasive ventilation in patients with Amyotrophic Lateral Sclerosis*, ACM SIGKDD Workshop on Health Informatics (HI-KDD), (2012).
- [8] Sun, J., Wang, F., Hu, J., Edabollahi, S., *Supervised patient similarity measure of heterogeneous patient records*, ACM SIGKDD Explorations Newsletter, 14:1 (2012), pp. 16–24.
- [9] Klenk, S., Dippon, J., Fritz, P. and Heidemann, G., *Determining Patient Similarity in Medical Social Networks*, MedEx Workshop, WWW’10, (2010).
- [10] Verbeek, J.J., *Notes on probabilistic PCA with missing values*, Tech. Report, (2009).
- [11] Yang, K. and Shahabi, C., *A PCA-based Similarity Measure for Multivariate Time Series*, Proceedings of the 2nd ACM international workshop on Multimedia databases, (2004), pp. 65-74.
- [12] Bastian M., Heymann S., Jacomy M., *Gephi: an open source software for exploring and manipulating networks*, International AAAI Conference on Weblogs and Social Media, (2003).
- [13] Girvan, M. and Newman, M. E. J., *Community structure in social and biological networks*, Proceedings of the National Academy of Sciences, 99:12 (2002), pp. 782–791.
- [14] Blondel, V.D., Guillaume, J-L., Lambiotte, R. and Lefebvre, E., *Fast unfolding of communities in large networks*, Journal of Statistical Mechanics: Theory and Experiment , 10 (2009), P10008.
- [15] Peng, H., Fuhui L. and Ding, C., *Feature selection based on mutual information criteria of max-dependency, max-relevance, and min-redundancy*, Pattern Analysis and Machine Intelligence, IEEE Transactions, (2005), pp. 1226–38.
- [16] De Jay, N., Papillon-Cavanagh, S., Olsen, C., Bontempi, G. and Haibe-Kains, B., *mRMRe: an R package for parallelized mRMR ensemble feature selection*, (2012).