



Published in final edited form as:

*Neuroinformatics*. 2014 April ; 12(2): 229–244. doi:10.1007/s12021-013-9204-3.

## A review of feature reduction techniques in neuroimaging

Benson Mwangi<sup>1</sup>, Tian Siva Tian<sup>2</sup>, and Jair C. Soares<sup>1</sup>

<sup>1</sup>UT Center of Excellence on Mood Disorders, Department of Psychiatry and Behavioral Sciences, UT Houston Medical School, Houston, TX, 77054 USA

<sup>2</sup>Department of psychology, University of Houston, Houston, TX, 77024 USA

### Abstract

Machine learning techniques are increasingly being used in making relevant predictions and inferences on *individual* subjects neuroimaging scan data. Previous studies have mostly focused on categorical discrimination of patients and matched healthy controls and more recently, on prediction of *individual* continuous variables such as clinical scores or age. However, these studies are greatly hampered by the large number of predictor variables (voxels) and low observations (subjects) also known as the *curse-of-dimensionality* or *small-n-large-p* problem. As a result, feature reduction techniques such as feature subset selection and dimensionality reduction are used to remove redundant predictor variables and experimental noise, a process which mitigates the *curse-of-dimensionality* and *small-n-large-p* effects. Feature reduction is an essential step before training a machine learning model to avoid *overfitting* and therefore improving model prediction *accuracy* and generalization ability. In this review, we discuss feature reduction techniques used with machine learning in neuroimaging studies.

### Keywords

Feature reduction; Feature selection; Dimensionality reduction; Machine learning; Multivariate; Predictive Modeling; Neuroimaging

## 1.0 Introduction

Conventional *group-level* neuroimaging data analysis techniques (e.g. voxel-based morphometry) have been used in neuroimaging research for decades. However, these techniques have only been able to identify average *between-group* differences and unable to make predictions on *individual* subjects. The inability to make predictions from *individual* subjects neuroimaging scan data may have greatly hampered the ability to translate neuroimaging research results into clinical practice (Brammer, 2009; Linden, 2012).

---

Correspondence to: Benson Mwangi, PhD, Department of Psychiatry and Behavioral Sciences, University of Texas Health Science Center at Houston, 1941 East Road, Houston, TX 77054, United States, Tel: +17139412616, benson.irungu@uth.tmc.edu.

#### Conflict of Interest

J.C.S has participated in research funded by Forest, Merck, BMS and GSK. He has been a speaker for Pfizer and Abbot.

#### Information Sharing Statement

Public repositories and software routines explicitly mentioned in this review were available in the public domain at the time of publication of this review.

Recently though, machine learning (ML) techniques also variously known as multivariate pattern analysis (MVPA) or pattern recognition (PR) techniques, have been used to decode or predict *individual* subjects' brain states using neuroimaging scan data. These studies have increased significantly and equally become successful in decoding brain states which raises the possibility of being deployed for potential clinical use and particularly in making *personalized* clinical decisions (Linden, 2012; Mwangi et al. 2012a).

Machine learning applications in neuroimaging can be divided into two broad categories namely *classification* and *regression*. In classification, neuroimaging data with corresponding categorical labels (e.g. Healthy Controls vs Patients or Treatment responders vs Non-responders) are used to develop a predictive classifier. The resulting classifier is used to make predictions on subjects' data not present during the training stage. Previous classification studies include; Alzheimer's disease (AD) (Kloppel et al. 2008; Magnin et al. 2009; Zhang et al. 2011b), Major Depressive Disorder (MDD) (Mwangi et al. 2012a; Zeng et al. 2012), Autism Spectrum Disorder (ASD) (Ecker et al. 2010; Ingalhalikar et al. 2011), Schizophrenia (Koutsouleris et al. 2009), Mild Cognitive Impairment (MCI) (Haller et al. 2010b) and Attention Deficient Hyperactivity Disorder (ADHD) (Lim et al. 2013; Sato et al. 2012; Zhu et al. 2005). Conversely, in regression neuroimaging data with corresponding continuous targets (e.g. clinical scores or age) are used to 'train' a regression model. Similar to above, the model is used to make predictions on novel *individual* subjects' data not present during training. These techniques have recently been used to predict *individual* subjects' age (Brown et al. 2012; Dosenbach et al. 2010; Franke et al. 2012; Franke et al. 2010; Mwangi et al. 2013) and clinical scores in MDD, AD and chronic pain (Marquand et al. 2010a; Mwangi et al. 2012b; Stonnington et al. 2010; Wang et al. 2010).

Notably, previous studies in both *classification* and *regression* have used neuroimaging scan data from multiple modalities such as T<sub>1</sub>-weighted magnetic resonance imaging (MRI), functional MRI (fMRI), diffusion tensor imaging (DTI), positron emission tomography (PET) and single-photon emission computed tomography (SPECT). In this review, we generalize our discussions to all imaging modalities and assume that all scans have been subjected to standardized pre-processing routines (e.g. spatial normalisation, segmentation and or smoothing). Neuroimaging data pre-processing frameworks in different modalities are discussed elsewhere (Ashburner, 2007; Ashburner, 2009; Johansen-Berg and Behrens, 2009; Penny et al. 2007). In this review, pre-processed neuroimaging scans' voxels are referred to as 'features' or 'predictor variables'. In addition, we generalize our feature reduction discussion to all machine learning techniques such as support vector machines (SVMs) (Mwangi et al. 2012a; Orru et al. 2012; Vapnik, 1999), relevance vector machines (RVMs) (Mwangi et al. 2012a; Tipping, 2001) and Gaussian process classifiers (GPCs) (Marquand et al. 2010b; Rasmussen and Williams, 2006). Discussions of these methods with respect to neuroimaging data are also given elsewhere (Johnston et al. 2012; Misaki et al. 2010; Orru et al. 2012).

### 1.1 Rationale for feature reduction

In many neuroimaging studies, the sample size (number of subjects or observations) is often less than 1000. In comparison, pre-processed brain scans may contain (>100,000) non-zero

voxels. As a result, the numbers of features (voxels) greatly outnumber the number of observations (sample size). This is a common problem in machine learning literature known as the *curse-of-dimensionality* (Bellman, 1961) or *small-n-large-p* (Fort and Lambert-Lacroix, 2005). Consequently, without pre-selecting the ‘most relevant’ features and effectively discarding redundant features plus noise, a predictive machine learning model has a marked risk of ‘overfitting’ (Guyon and Elisseeff, 2003; Hua et al. 2009). Overfitting implies that the model training process yields a machine learning model with poor generalization ability which is interpreted as inability to make accurate predictions on novel subjects’ data (Guyon and Elisseeff, 2003; Hua et al. 2009; Kohavi and John, 1997).

In view of the above, there is a general consensus from previous neuroimaging machine learning studies that feature reduction is a fundamental step before applying a predictive model (e.g. SVM or RVM) to neuroimaging data (Bray et al. 2009; Cheng et al. 2012; De Martino et al. 2008; Duchesnay et al. 2007; Duchesnay et al. 2004; Duff et al. 2011; Franke et al. 2010; Liu et al. 2012; Lohmann et al. 2007; Mitchell et al. 2004; Mwangi et al. 2012a; Norman et al. 2006; Pereira et al. 2009; Rizk-Jackson et al. 2011; Schrouff et al. 2013; Valente et al. 2011; Van De Ville and Lee, 2012). Other than mitigating the *curse-of-dimensionality* effect as above, the feature reduction process may also facilitate a deeper understanding of the scientific question of interest. For example, in a highly accurate predictive classifier able to predict patient treatment responders against non-responders, brain regions identified during feature reduction may offer an insight into different neural substrates of non-responders. In this review we divide feature reduction techniques into two broad categories namely; *supervised* and *unsupervised* techniques respectively.

## 2.0 Supervised feature reduction techniques

Supervised feature reduction techniques use highly dimensional neuroimaging data and the required outcome labels (e.g. +1 treatment responders, -1 treatment non-responders) to select relevant features and discard redundant features plus noise. However, as in the wider machine learning literature, we further subdivide these techniques into three categories namely; ‘filter’, ‘wrapper’ and ‘embedded’ methods, respectively (Guyon and Elisseeff, 2003; Saeys et al. 2007). A clear distinction between these categories exists. First, filter techniques such as t-tests, Anova and Pearson correlation coefficient use simple statistical measures (e.g. mean, variance, correlation coefficients) to rank features according to their relevance in detecting *group-level* differences. Second, wrapper techniques use an *objective* function from a classification or regression machine learning model to rank features according to their relevance to the model. Lastly, embedded methods select relevant features as ‘part’ of the machine learning process by enforcing certain ‘penalties’ on a machine learning model thus yielding a small subset of relevant features.

In passing, we note that relevant features are usually selected using training data only to avoid *double-dipping* which entails using both training and testing data partitions to select relevant features (Kriegeskorte et al. 2010; Kriegeskorte et al. 2009). In this review, we assume both (training and testing datasets) have been adequately separated using a *cross-validation* procedure such as hold-out or leave-one-out which are explored elsewhere (Johnston et al. 2012; Pereira et al. 2009; Strother et al. 2002; Theodoridis and

Koutroumbas, 2009). Figure 1 illustrates a feature reduction process without double dipping. Conversely, Figure 2 illustrates a feature reduction procedure with double dipping - which should be avoided.

## 2.1 Filter techniques

**2.1.1 Pearson correlation coefficient**—The Pearson correlation coefficient (PCC) ranks features by calculating linear correlations between individual features and class labels in classification or continuous targets in regression (Guyon and Elisseeff, 2003). Here, we assume a two group classification problem with predictor variables  $x$  and diagnostic labels  $y_i$ , (e.g. Patients-1 vs Controls -2). The Pearson correlation coefficient between predictor variables and diagnostic labels is calculated as:

$$P_i = \frac{\sum_i (x_i - \bar{x})(y_i - \bar{y}_i)}{\sqrt{\sum_i (x_i - \bar{x})^2 \sum_i (y_i - \bar{y}_i)^2}} \quad (1)$$

Where  $x_i$  stands for the feature value of the  $i^{th}$  sample and  $\bar{x}$  is the mean of these feature values.  $y_i$  are diagnostic labels and  $\bar{y}_i$  is the mean of all  $y_i$  in the sample (Fan et al. 2007; Guyon and Elisseeff, 2003). The relevance of a feature in separating both classes is evaluated by considering the absolute value of the correlation coefficient  $P_i$  with higher values indicating the feature's greater relevance in discriminating between classes. Lastly, relevant features above a *user-defined* threshold of correlation coefficients are selected for subsequent machine learning analyses. However, to select the optimal *user-defined* threshold, cross-validation procedures (e.g. k-fold or leave-one-out) are used to iteratively evaluate a range of threshold values and the threshold with a low generalization error selected.

Numerous studies have applied PCC filters to select relevant features. For example, in MCI classification (Wee et al. 2011), cocaine exposed individuals classification (Fan et al. 2007), AD classification (Dai et al. 2012; Davatzikos et al. 2008; Grana et al. 2011) and in gender (Men vs Women) classification (Fan et al. 2006).

However, although PCC filters are applicable to both multi-group (>2) classification and regression tasks, they are only able to detect linear dependencies between features and corresponding targets (Guyon and Elisseeff, 2003). This is a major setback for PCC filters especially in tasks involving selecting relevant features from high-dimensional data with multivariate relationships (Wang et al. 2010).

**2.1.2 T-tests and analysis of variance (ANOVA)**—Statistical hypothesis testing techniques such as t-tests and ANOVA have extensively been used to detect *average group-level* differences in neuroimaging studies. Recently, though, these techniques have been used for supervised feature selection in neuroimaging machine learning studies with success. Here, we assume a two class classification task of categorising (treatment responders +1 vs. treatment Non-responders -1). The assumption that predictor variables are from pre-processed neuroimaging scans still holds. As a result, a two-sample t-test evaluates the null

hypothesis that the two-group (Responders, Non-responders) sample means are equal given observed variance. Let  $\bar{x}_p$  and  $\bar{x}_c$  be voxel means for patient and control groups respectively and  $s_p, s_c$  denote corresponding standard deviations. The t-score of a feature  $v$  is calculated as (Sheskin, 2004);

$$t = \frac{|\bar{x}_p - \bar{x}_c|}{\sqrt{\frac{(V_p - 1)s_p^2 + (V_c - 1)s_c^2}{N_p + N_c - 2} \cdot \left[ \frac{1}{N_p} + \frac{1}{N_c} \right]}} \quad (2)$$

Where  $N_p$  and  $N_c$  denote number of subjects in each group. Once this calculation is performed, a threshold of significance (e.g. p-value) which represents the probability of obtaining a statistic greater in magnitude than  $t$  under the null hypothesis is defined (Ashburner and Friston, 2000; Penny et al. 2007; Theodoridis and Koutroumbas, 2009). Subsequently, an optimal *user-defined* threshold of significance (p-value) representing *relevant* features is selected through a cross-validation process and relevant features used for subsequent machine learning analyses (Mwangi et al. 2012a). Several classification studies have recently used t-tests to select relevant features for machine learning in neuroimaging. For example, AD (Chaves et al. 2009; Chu et al. 2012; Hinrichs et al. 2009; Zhang et al. 2011a), ASD (Duchesnay et al. 2011), MDD (Mwangi et al. 2012a), gender (Duchesnay et al. 2007), schizophrenia (Kovalev et al. 2003), MCI (Wee et al. 2012) and in other non-clinical studies (Balci et al. 2008; Haynes and Rees, 2005). Applications of t-tests in feature selection are computationally fast and easy to implement and scale well to high dimensional data meaning they are able to select a small subset of *relevant* features from the original high-dimensional feature set (Hua et al. 2009; Mwangi et al. 2012a; Saeys et al. 2007).

However, t-tests are also weighed down by a number of shortcomings. First, they are *univariate*, meaning they do not take into account interactions between multiple features and spatial patterns (non-multivariate). Second, t-tests are only able to detect two-group differences, although this shortcoming is compensated by the equivalent analysis of variance (ANOVA) technique.

The ANOVA technique is used to select *relevant* features in multiple (>2) groups by generalizing a t-test to more than two groups (Cohen, 1998). However, similar to t-tests, ANOVAs are *univariate* but have also been used in selecting relevant features in neuroimaging classification tasks. For example, Costafreda and colleagues (Costafreda et al. 2011) recently used ANOVA for feature selection in classifying patients with bipolar disorder (BD), schizophrenia and healthy controls. Other neuroimaging classification studies using ANOVA for feature reduction include; ASD (Coutanche et al. 2011), MDD (Costafreda et al. 2009), schizophrenia (Yoon et al. 2008) and fMRI visual activation classification task (Cox and Savoy, 2003). In passing, we note that ANOVAs offer the same benefit as t-tests, but an optimal threshold of relevant features is selected using a cross-validation process with the training data. Most notably, we note that the multivariate analysis of covariance (MANCOVA) is a multivariate extension of ANOVA and has recently been used in a number of neuroimaging feature reduction tasks (Friston et al. 1996, Calhoun et al. 2011, Allen et al. 2011). Lastly, Wilcoxon test (De Martino et al. 2008), signal-to-noise ratio (Ingalhalikar et al. 2011), Gini-contrast (Langs et al. 2011) and relief

(Haller et al. 2010a) are other filter techniques, albeit with few applications in neuroimaging machine learning studies so far.

However, whilst these supervised filter techniques have extensively been used in mitigating the *small-n-large-p* or *curse-of-dimensionality* problems in neuroimaging, they collectively suffer from a significant setback by not considering interactions between multiple features or spatial patterns (not multivariate).

## 2.2 Multivariate wrapper techniques

Multivariate wrapper techniques use an *objective-function* from a multivariate machine learning model (e.g. classification or regression) to rank features according to their relevance or importance to the model (Guyon and Elisseeff, 2003; Kohavi and John, 1997). Multivariate wrapper approaches are further subdivided into two sub-categories namely; forward selection and backward elimination (Kohavi and John, 1997). In forward selection, the search for *relevant* features begins with an empty set and features are iteratively added in ‘small’ pre-defined steps until an *optimal* number of features are found. On the contrary, in backward elimination, the search begins with all features in the training set which are iteratively removed in ‘small’ pre-defined steps until an *optimal* number of features are found. Recursive feature elimination (RFE) (Guyon et al. 2001), which is popularly used in neuroimaging studies is an example of a backward elimination technique.

**2.2.1 Recursive feature elimination**—We assume a two-group classification task with a set of features  $x_i$ , and corresponding labels  $y_i$ . Equally, we assume training data is now sub-divided into ‘training’ and ‘evaluation’ subsets respectively. A machine learning algorithm (e.g. linear support vector machine or linear relevance vector machine) is ‘trained’ resulting into observation weights. As a result, feature or voxel relevance weights  $\alpha_i$  are calculated as:

$$W = \sum_{x_i \in cf} \alpha_i y_i x_i \quad (3)$$

Where  $cf$  stands for observations or subjects with non-zero weights (e.g. support vectors in SVM or analogous relevance vectors in RVM). Subsequently, absolute values of the weights  $W$  are ranked in order of importance (low-less relevant, high-most relevant) and a *user-defined* percentage (e.g. 2%) of the lowest ranking features removed. In the next step, a new model is trained minus the non-relevant features and the new model’s generalization error (or accuracy) on the evaluation subset is assessed. This process is repeated until a termination criterion is reached or until the feature set is empty. Lastly, features leading to the best generalization ability (high accuracy) are selected for training the final machine learning model and the rest of the features discarded. Figure 3 illustrates a recursive feature elimination process.

However, RFE requires definition of two parameters. 1) Backward-eliminations termination or stopping criteria. The first possibility is to remove low ranking features iteratively until the feature set is empty and the iteration resulting to the best generalization ability or high accuracy is selected. The second possibility is to terminate the procedure when the model

performance of the current iteration is not significantly better than the previous step, as explored by De Martino and colleagues (De Martino et al. 2008). 2) Another *user-defined* parameter is the percentage of features removed at every backward-elimination step. Previous neuroimaging studies using RFE have variably used this parameter, for example 8% (De Martino et al. 2008), 10% (Craddock et al. 2009) and 2% (Mwangi et al. 2013). It is not currently well understood how the choice of this parameter affects the overall model performance and this remains an open research question. However, removing a very small percentage of features at every iteration is computationally expensive, whilst removing a higher percentage may result into inclusion of non-relevant features (Craddock et al. 2009; De Martino et al. 2008).

RFE offers several benefits by first, considering multivariate interactions between entire spatial patterns in the data. Second, the technique uses a predictive model to remove non-relevant or redundant features which may result in a better generalization ability (Guyon and Elisseeff, 2003). However potential setbacks should be noted. First, this technique is computationally intensive as it performs a complete heuristic search of the feature input space (Saeys et al. 2007). Second, as cross-validation methods are used to avoid a biased feature selection process, this may result in different features being selected in every cross-validation iteration (Craddock et al. 2009). This problem has recently been addressed by Dosenbach and colleagues (Dosenbach et al. 2010) by recommending reporting of a *consensus* discrimination map which aggregates features selected in all cross-validation iterations.

Previous feature selection applications using RFE in neuroimaging classification tasks include; ASD (Calderoni et al. 2012; Duchesnay et al. 2011; Ecker et al. 2010; Ingahlalkar et al. 2011), MDD (Craddock et al., 2009), schizophrenia (Castro et al. 2011b), psychosis (Gothelf et al. 2011), object recognition in a fMRI task (Hanson and Halchenko, 2008), fragile X syndrome (Hoeft et al. 2011), ADHD (Marquand et al. 2011), MCI (Nho et al. 2010), fMRI spatial patterns (De Martino et al. 2008), mood disorders (Mourao-Miranda et al. 2012) and AD (Davatzikos et al. 2008; Mesrob et al. 2008). An interesting variant of RFE, which involves backward-elimination of voxel clusters rather than individual features, has recently been explored (Deshpande et al. 2010). In passing, we note that although the majority of RFE applications in neuroimaging are largely in predictive classification, recently RFE has been used for *regression* tasks (Fan et al. 2010; He et al. 2008; Mwangi et al. 2013).

Remarkably, although previous RFE applications have largely utilized multivariate wrappers in *input* spaces (e.g. linear SVM), continuing efforts in solving the ‘pre-image problem’ (Kwok and Tsang, 2004) have recently allowed extraction of *feature/voxel* weights in *non-linear* feature spaces (e.g. non-linear SVM) (Kjems et al. 2002; LaConte et al. 2005; Mwangi et al. 2013; Rasmussen et al. 2011).

**2.2.2 Searchlight**—The searchlight technique was introduced by Kriegeskorte and colleagues (Kriegeskorte et al. 2006) for multivariate feature reduction in neuroimaging data. This technique selects relevant features as follows. First, 3-dimensional (3D) spherical volumes of pre-defined radius or ‘searchlight’ (e.g. 4mm) are centered at every voxel and

populated through neuroimage volumes in the training data. Second, a machine learning classifier (e.g. SVM) is trained at every searchlight volume using only voxels within the spherical searchlight volume and classifier accuracies from searchlight volumes centered at every voxel recorded. Through permutation methods (Kriegeskorte et al. 2006; Pereira et al. 2009), the searchlight accuracies at every voxel are converted into p-values, which are subsequently thresholded to remove non-relevant voxels. However, the searchlight technique requires a user to determine the spherical regions-of-interest *radius* (e.g. 4mm), and according to Kriegeskorte and colleagues (Kriegeskorte et al. 2006), there should be a balance between the size of neuroimaging scans and the spherical searchlight volume.

Previous applications of feature reduction using searchlight in neuroimaging machine learning tasks include: episodic memory decoding (Chadwick et al. 2010), scene representation decoding (Bonnici et al. 2011) and attention shifting cortical activation decoding (Greenberg et al. 2010). We note that majority of these studies have largely been task-based fMRI decoding solutions.

Notably, the searchlight method assumes that any discriminative information lies within the searchlight radius and according to (Formisano et al. 2008), this technique may be unreliable in detecting ‘distant pattern’ differences (e.g. bilateral cerebral regions). However, we note that this is not the case with other multivariate wrapper methods such as RFE, which are able to detect multivariate interactions across ‘distant patterns’ as above.

### 2.3 Embedded feature reduction techniques

In this section, we discuss three of the most popular embedded methods namely; least absolute shrinkage and selection operator (LASSO) (Tibshirani, 1996; Tibshirani, 2011), the Elastic Net (Zou and Hastie, 2005) and the partial least square (PLS) method (MacIntosh and Lobaugh 2004). The Elastic Net and LASSO techniques combine both machine learning and feature reduction steps by enlisting a regularization framework (e.g.  $L_1$  and  $L_2$  norm regularization) resulting to a reduced subset of relevant features (Zou and Hastie, 2005). On the contrary, PLS selects relevant features by establishing or analyzing associations between the independent and dependent variables (e.g. brain activity or structure and behavior) (Krishnan et al. 2011).

**2.3.1 Least absolute shrinkage and selection operator (LASSO)**—We assume a two-group classification task with a set of features  $x_{ij}$ , and corresponding target labels  $y_i$ , where  $i = 1, 2, \dots, N$  represents observations (subjects) and  $j = 1, 2, \dots, P$  represents the number of features (predictor variables). Additionally, we assume that the predictor variables are normalized by subtracting sample mean and dividing by the standard deviation. As a result, the LASSO computes model coefficients  $\hat{\beta}$  by minimizing the following function (Tibshirani, 1996; Tibshirani, 2011).

$$\sum_{i=1}^N \left( y_i - \sum_j x_{ij} \beta_j \right)^2 + \lambda \sum_{j=1}^P |\beta_j| \quad (4)$$



$\lambda$  is a *user-defined* parameter which controls the balance between the model having few non-zero coefficients  $\hat{\beta}$  (sparsity) and high prediction accuracy or generalization ability. Interestingly, as  $\lambda$  approaches 1, the model becomes increasingly sparse meaning few ‘relevant’ features whilst as  $\lambda$  approaches 0, the model is less sparse meaning more ‘relevant’ features (Bunea et al. 2011). To arrive at an optimal value for, cross-validation procedures (e.g. k-fold or leave-one-out) are used to test a range of  $\lambda$  parameters with pre-defined steps and the parameter resulting to high model accuracies is selected. The LASSO function is solved using an optimization procedure such as the coordinate descent algorithm as explored elsewhere (Friedman et al. 2010; Tseng and Yun, 2009). A solution to this optimization problem is provided by Friedman and colleagues (Friedman et al. 2010) in R and Matlab (The Mathworks, Inc) routines.

The benefits of this method in a feature reduction process are two fold. First, the LASSO yields a small set of model coefficients  $\hat{\beta}$  with majority of coefficients set to zero and corresponding features discarded from the subsequent machine learning process. Second, the LASSO is able to cope with situations where there are a large number of predictor variables (voxels) and fewer observations (subjects) as is the case in majority of neuroimaging studies (Bunea et al. 2011). Previous applications of feature selection using LASSO in neuroimaging machine learning tasks include: AD classification (Casanova et al. 2011; Kohannim et al. 2012b; Rao et al. 2011; Vounou et al. 2011; Yan et al. 2012), prediction of video stimulus scores in fMRI (Carroll et al. 2009), ASD classification (Duchesnay et al. 2011), prediction of brain characteristics using genetic data (Kohannim et al. 2012a; Kohannim et al. 2012b), prediction of pain stimuli in fMRI (Rish et al. 2010) and Gender classification (Casanova et al. 2012).

The LASSO technique has been successful in mitigating *small-n-large-p* or *curse-of-dimensionality* problems in neuroimaging albeit with several setbacks. First, if predictor variables in a group are highly correlated, the LASSO selects only one variable from the group and ignores the rest (Bunea et al. 2011; Zou and Hastie, 2005). Second, the number of selected relevant features may not exceed the number of observations or samples before the model begins to saturate (Zou and Hastie, 2005). To overcome these two limitations, Zou and Hastie (Zou and Hastie, 2005) introduced the Elastic Net technique, which is relatively similar to LASSO, albeit with a few modifications, as explored below.

**2.3.2 Elastic Net**—The Elastic Net is formulated in a similar manner as the LASSO, but with an additional quadratic term (Zou and Hastie, 2005). We assume a two-class classification problem with similar observations and predictor variables and a preceding normalization step, as above. As a result, the Elastic Net computes model coefficients  $\hat{\beta}$  by minimizing the following *objective* function (Bunea et al. 2011; Ogotu et al. 2012; Zou and Hastie, 2005).

$$\sum_{i=1}^N \left( y_i - \sum_j x_{ij} \beta_j \right)^2 + \lambda_1 \sum_{j=1}^P |\beta_j| + \lambda_2 \sum_{j=1}^P \beta_j^2 \quad (5)$$

Contrary to LASSO, the Elastic Net requires definition of two *user-defined* model regularization parameters ( $\lambda_1$  and  $\lambda_2$ ), which control the degree of penalization. The  $L_1$  penalty  $\sum_{j=1}^P \beta_j^2$  promotes sparsity in the solution, resulting in few features with non-zero weights, whilst  $L_2$  penalty encourages stability in the solution and acts as a bound on the number of features selected (Bunea et al. 2011; Kohannim et al. 2012a; Ogutu et al. 2012; Zou and Hastie, 2005). These parameters are often selected using an *objective* parameter grid-search process which evaluates a ‘range’ of parameters in two-dimensions (grid-search) and parameters giving the best performance selected. However, the grid-search process can be computationally intensive (Bunea et al. 2011). Similar to LASSO, the Elastic Net function is solved using an optimization procedure such as the coordinate descent algorithm also explored elsewhere (Friedman et al. 2010). An implementation of the Elastic Net solution is freely provided in both R and Matlab (The Mathworks, Inc) by Friedman and colleagues (Friedman et al. 2010).

Previous applications of feature reduction using Elastic Net in neuroimaging machine learning tasks include: AD classification (Rao et al. 2011; Shen et al. 2011; Wan et al. 2011) and treatment response predictions in ADHD (Marquand et al. 2012).

**2.3.3 Partial Least Squares**—The partial least square (PLS) feature reduction method is divided into two major categories namely; partial least squares correlation (PLSC; McIntosh et al. 1996, Krishnan et al. 2011) and partial least squares regression (PLSR; Wold et al 2001). Here, we discuss PLSC which is by far the most popular variant of the partial least squares method used in neuroimaging feature reduction tasks (Krishnan et al. 2011).

We recall our two-group classification task with a set of features  $x_{ij}$  stored in a matrix  $X$  and corresponding target labels  $y_i$ , stored in a matrix  $Y$  where  $i$  1,2,..,  $N$  represents observations (subjects) and  $j$  1,2,..,  $P$  represents the number of features (predictor variables). Additionally, we assume that the predictor variables are normalized by subtracting sample mean and dividing by the standard deviation (z-scores). PLSC begins by computing the cross product of the predictor variables and target vectors as below;

$$M=Y^T X \quad (6)$$

The resulting matrix  $M$  is decomposed using singular value decomposition (SVD) (Krishnan et al. 2011, McIntosh and Misis 2013). The SVD matrix decomposition process is explored elsewhere (Bishop 2006, Krishnan et al. 2011). The SVD decomposition yields the following.

$$M=USV^T \quad (7)$$

The above decomposition results into a set of singular vectors ( $U$ -left,  $V$ -right) whilst  $S$  is a diagonal matrix representing the ‘singular values’ (McIntosh and Misis 2013). The left singular vectors  $U$  contains the weights or coefficients identifying the variables or voxels in matrix  $X$  making the highest contribution in explaining the relationship between the predictor variables (e.g. brain scans) and the targets (e.g. diagnostic labels or behavioural scores) (McIntosh and Misis 2013). Lastly, a set of latent variables for both the predictor

variable matrix  $X$  and target labels  $Y$  are reconstructed by computing the dot product of the singular vectors and the original data as below and also explored elsewhere (Krishnan et al. 2011).

$$L_x = XV \quad (8)$$

Where  $L_x$  is a reduced set of latent variables representing the original predictor variables in  $X$  (Krishnan et al. 2011). On the contrary,  $L_y$  represents the latent variables for the target variables and computed as follows (Krishnan et al. 2011).

$$L_y = YU \quad (9)$$

At this point, the original high-dimensional predictor variables (e.g. brain scans) are now represented by a set of low dimensional latent variables. Lastly, bootstrap and permutation tests are used to identify the optimal number of latent variables as explored elsewhere (Krishnan et al. 2011, Ziegler et al. 2013, McIntosh and Misic. 2013). An implementation of the PLSC method for neuroimaging data analysis is freely provided in Matlab (The Mathworks, Inc) by Shen and colleagues at ([research.baycrest.org/pls/source/](http://research.baycrest.org/pls/source/)). A rigorous introduction to the PLS method and accompanying pseudocodes are also provided elsewhere (Krishnan et al. 2011). Notably, as PLSC correlates predictor variables (e.g. brain scans) and target variables (e.g. diagnosis or behavioral design), prediction tasks are handled by its variant - partial least squares regression (PLSR) (Wold et al. 2001). The PLSR is explored elsewhere (Wold et al. 2001, Krishnan et al. 2011).

Previous applications of feature reduction using the partial least square method in neuroimaging machine learning tasks include: Age classification (Young vs Old) (Chen et al. 2009) and prediction of cognitive behavioral scores (Ziegler et al. 2013, Menzies et al. 2007, Nester et al. 2002). Notably, PLS has also been used in multimodal feature reduction tasks (Sui et al. 2012, Martinez et al. 2004).

### 3.0 Unsupervised feature reduction techniques

Unsupervised feature reduction techniques also known as dimensionality reduction or feature extraction techniques *construct* relevant features through *linear* or *non-linear* combinations of the *original* predictor variables (features) (Lee and Verleysen, 2007). In this review, we restrict our discussion to two of the most popular unsupervised dimensionality reduction techniques in neuroimaging namely, principal component analysis and independent component analysis.

#### 3.1 Principal component analysis

Principal component analysis (PCA) constructs relevant features by linearly transforming correlated variables (e.g. raw voxels in a brain scan) into a smaller number of uncorrelated variables also known as principal components (Jolliffe, 2002). The resulting principal components are essentially linear combinations of the original data capturing most of the variance in the data (Jolliffe, 2002; Mourao-Miranda et al. 2005; Mourao-Miranda et al. 2006; Zhu et al. 2005). Construction of principal components from high-dimensional data begins by first, normalizing the original features by subtracting sample mean and dividing

by the standard deviation. Second, Eigen decomposition of the covariance matrix from the standardized data is performed resulting in eigenvalues and eigenvectors of the covariance matrix. Third, eigenvalues are sorted in a decreasing order effectively representing decreasing variance in the data (Jolliffe, 2002). Lastly, principal components are constructed by multiplying the originally *normalized* data with the ‘leading’ eigenvectors whose exact number is a *user-defined* parameter. The ‘leading’ eigenvectors explain most of the variance in the data. As a result, highly-dimensional neuroimaging scans with potentially many correlated voxels are now represented by relatively few uncorrelated principal components which are later used for machine learning analyses. A detailed and rigorous mathematical derivation of PCA is given elsewhere (Jolliffe, 2002).

PCA has been successful in extracting relevant features in neuroimaging classification studies. For example in Schizophrenia (Caprihan et al. 2008; Radulescu and Mujica-Parodi, 2009; Yoon et al. 2007), AD (Lopez et al. 2011), Psychosis (Koutsouleris et al. 2009), MDD (Fu et al. 2008), ADHD (Zhu et al. 2008) and face recognition in fMRI (Mourao-Miranda et al. 2005). Most notably principal components have also been used in neuroimaging machine learning regression studies. For example, AD clinical scores prediction (Wang et al. 2010) and age predictions (Franke et al. 2010).

PCA offers two major benefits to dimensionality reduction in neuroimaging machine learning studies. First, the technique is easy to implement and computationally efficient. Second, the technique is *un-supervised*- meaning it does not require corresponding categorical or continuous labels or targets to extract relevant features. However, PCA also suffers from several setbacks. First, the user is required to define the number of principal components although Hansen and colleagues (Hansen et al. 1999) propose using the *generalization error* from a cross-validation process to select an *optimal* number of principal components. Second, as principal components are linear combinations of the original features, they may not be easily interpretable (Bunea et al. 2011). Third, whilst constructing relevant components, PCA may ignore the required outcome (e.g. discrimination of disease vs healthy) (Bunea et al. 2011). Lastly, as principal components are constructed through a linear transformation, this process may not adequately detect more complex non-linear feature interactions (Bishop, 1995), such as those that may occur in neuroimaging scan data. This setback has led to the formulation of kernel-pca- which in brief is the application of PCA in a feature space created through a kernel function (Scholkopf and Smola, 2002). Several studies have recently explored the application of kernel-pca to dimensionality reduction problems in neuroimaging (Rasmussen et al. 2012; Sidhu et al. 2012; Thirion and Fugeras, 2003; Wang et al. 2011).

### 3.2 Independent component analysis

Independent component analysis (ICA) is a multivariate data-driven technique which belongs to the broader category of *blind-source separation* methods used to separate data into underlying *independent* information components (Stone, 2004). ICA separates a set of ‘mixed signals’ (e.g. raw data from an fMRI scan) into a set of *independent* and relevant features (e.g. paradigm-related signals in fMRI). To achieve this goal, ICA assumes the

source signals are statistically independent from an unknown but linear mixing process (Calhoun et al. 2009).

To elucidate this, first we consider a  $i \times j$  matrix  $X$  where ( $i$  = number of fMRI scans in a study and  $j$  = number of voxels from the pre-processed scans). Second, we consider a  $n \times j$  matrix  $S$  with rows representing ( $n$ ) number of ‘expected’ independent components. Additionally, we represent  $A$  as a  $i \times n$  ‘mixing’ matrix with columns containing associated time-courses of the  $n$  components. Effectively, this becomes a ‘blind-source separation’ problem (Stone, 2004) which can be represented by a linear model  $X = AS$  (Calhoun and Adali, 2006). In addition, we consider the function  $Y = WX$ . Consequently, the goal of an ICA algorithm is to estimate the  $j \times i$  ‘unmixing’ matrix  $W$  such that  $Y$  becomes a good approximation of the true signal sources  $S$ . ICA algorithms use high-order statistical methods to solve for independent components (ICs) and a gentle introduction to these methods is given elsewhere (Calhoun and Adali, 2006; Stone, 2004).

There are two broad variants of ICA applications in fMRI. The first variant depends on the dimension of priority (e.g. spatial dimension-spatial ICA, temporal dimension-temporal ICA). Majority of ICA dimensionality reduction studies in fMRI have mostly extracted *relevant* independent components from the spatial dimension (Calhoun et al. 2009) and a detailed evaluation of both spatial and temporal options is given elsewhere (Calhoun and Adali, 2006). The second category of ICA further sub-divides the technique into either *individual-subject* ICA or *group-level* ICA (Calhoun et al. 2001; Calhoun and Adali, 2006). In *individual-subject* ICA, each subject’s data are entered into a separate ICA analysis while in *group-level* ICA one set of components for the groups are estimated and back-reconstructed from an aggregate mixing matrix to obtain *individual-subject* independent components and discriminative maps. The group-level ICA back-reconstruction step is explored elsewhere (Erhardt et al. 2011). A Matlab (The Mathworks, Inc) implementation of group ICA for fMRI is made freely available elsewhere (Calhoun, 2011).

Benefits of employing ICA in dimensionality reduction tasks should be noted. First, unlike *univariate* methods, ICA does not require an investigator to specify a regressor of interest which may require prior knowledge and assumptions about the experiment (e.g. paradigm or model regressor in fMRI) (Calhoun and Adali, 2006). Second, ICA has proved to be successful in disentangling otherwise mixed brain signals (e.g. separating physiological, motion and scanner related components) (Calhoun and Adali, 2006). However, potential setbacks of ICA should also be noted. First, ICA algorithms are computationally intensive (Correa et al. 2007). Second, according to Birn and colleagues (Birn et al. 2008), current ICA algorithms may not adequately separate respiration and default mode network signals in fMRI.

Notable neuroimaging machine learning studies using ICA in dimensionality reduction include (Castro et al. 2011a; Chai et al. 2010; De Martino et al. 2007; Douglas et al. 2011; Duff et al. 2011; Ince et al. 2008; Sato et al. 2012; Tagliazucchi et al. 2012; Toussaint et al. 2012; Yang et al. 2010). A review of ICA applications in multi-modal feature reduction tasks is given elsewhere (Sui et al. 2012). Lastly we note that a significant methodological difference between PCA and ICA exists. In PCA a set of possibly correlated variables are

converted into a set of uncorrelated features, whilst in ICA original variables are transformed into a set of nearly statistically independent features (Calhoun and Adali, 2006). However, the main commonality between ICA and PCA and possibly a distinctive characteristic that separates them from filter, wrapper and embedded feature reduction techniques is that the former are *unsupervised* while the latter are *supervised*.

### 3.3 Coordinate-based meta-analysis (CBMA) techniques

In the previous sections, we have mostly focused on ‘data-driven’ feature reduction techniques. Here though, we discuss techniques which may rely on existing ‘domain knowledge’ for feature reduction. Meta-analysis techniques have previously been used to model, analyze and report brain activation or group-level differences across neuroimaging studies (Eickhoff et al. 2009). Popular meta-analysis techniques include; activation likelihood estimation (ALE) (Laird et al. 2005b), kernel-density estimation (Scott, 1992) and multi-level kernel density estimation (Wager et al. 2007). A more detailed discussion of these techniques in relation to the wider meta-analysis reporting in neuroimaging is given elsewhere (Eickhoff et al. 2009; Laird et al. 2005a; Laird et al. 2005b).

Recently though, CBMA techniques have been used in modeling distributions of reported *foci* from fMRI activation studies and resulting regions-of-interest used as input features for machine learning analyses. For instance, Yarkoni and colleagues (Yarkoni et al. 2011) recently applied a CBMA technique to select relevant features in classifying working memory, emotion and pain using fMRI. The same group has recently developed a CBMA feature reduction framework known as *Neurosynth* (Mitchell, 2011; Yarkoni et al. 2011). *Neurosynth* is a tool for automated synthesis of fMRI data as reported in published studies (Yarkoni et al. 2011). Durkat and colleagues (Dukart et al. 2012) have recently applied a CBMA technique in classifying AD subjects in multicenter studies with high generalization ability. Other studies exploring CBMA techniques in feature reduction have also been reported. For example, Dosenbach and colleagues (Dosenbach et al. 2010) used regions-of-interest (ROIs) derived from a fMRI meta-analysis as relevant features for predicting individual subject’s age. Doyle and colleagues (Doyle et al. 2013) used existing domain knowledge to select the thalamus, anterior cingulate cortex and occipital cortex as relevant regions to predict the effect of a drug (ketamine) on brain activity. Recently, Chu and colleagues (Chu et al. 2012), reported that features from ROIs selected using *a priori* domain knowledge (e.g. hippocampal degeneration in Alzheimer’s disease) resulted to better generalization ability as compared to features selected from a data-driven approach such as t-test or RFE.

Notably, as CBMA techniques pool activation *foci* coordinates from numerous neuroimaging studies, this approach may improve a-posteriori certainty, increase statistical power and therefore make neuroimaging studies less susceptible to type II errors (Wager et al. 2007). The latter is a common problem in neuroimaging studies with relatively small sample sizes. However, CBMA techniques may also suffer from information loss as they represent data in published studies with a high degree of sparseness (Salimi-Khorshidi et al. 2009).

## 4.0 Summary and Discussion

There is general consensus within the neuroimaging machine-learning community that feature reduction is an important process before training a machine learning model. The main benefits of this process are two fold. First, to remove any redundant features (voxels) plus noise a process which may improve prediction accuracy or generalization ability as well as support interpretability of study results. The latter may in turn help in generating *post-hoc* inferences.

In this review, we have noticeably distinguished feature reduction methods into two major categories namely; *supervised and unsupervised* techniques respectively. We have further subdivided supervised techniques into three subcategories namely; filter, wrapper and embedded techniques. In summary, filter techniques use statistical feature ranking criterions (e.g. labels vs feature correlations in PCC) to discard redundant features. However, these techniques have two common setbacks. First, they are not multivariate and as such they do not take into account interactions between multiple features and spatial patterns. Second, these techniques need a user to define an optimal relevant feature threshold value (e.g. absolute correlation coefficient in PCC or p-value in t-test), although previous studies have prevailed over this impediment by optimizing this parameter through cross-validation procedures. On the other hand, wrapper techniques are multivariate although computationally intensive. Notably, De Martino and colleagues recently attempted a *hybrid* combination of both filter and wrapper techniques (t-test and RFE) which had superior prediction accuracy as compared to both techniques alone. Lastly, we note that the performance of *embedded* feature reduction methods (LASSO and Elastic Net) strongly depends on the choice of penalization parameters which should be chosen through a cross-validation process (Bunea et al. 2011; Shi et al. 2007).

We note that a significant difference between supervised and unsupervised techniques exists. Unsupervised techniques *construct* features independent of the outcome of interest whilst supervised techniques choose relevant features based on their ability to detect group-level differences.

Pertinent issues that in general may determine the performance of many feature reduction techniques at a high-level should be noted. First, the need to set a threshold of optimal number of features. Second, In the event one is using a cross-validation process to train and test the model (e.g. leave-one-out), relevant features may differ from fold-to-fold. Notably, although feature reduction methods ultimately help in removing redundant data and noise, equally important and relevant features may be inadvertently removed during feature selection (Guyon and Elisseeff, 2003).

In conclusion, feature reduction techniques are frequently being used in neuroimaging machine learning studies to mitigate the curse-of-dimensionality or small-n-large-p problems and maximize prediction accuracies. Lastly, whilst there are a number of studies empirically comparing different feature reduction approaches (Craddock et al. 2009; De Martino et al. 2008; Ryalı et al. 2010) in neuroimaging studies, none of these studies recommend any technique as the best in all neuroimaging machine learning tasks.

## Acknowledgments

This research was funded by NIMH R01085667 and Pat Rutherford, Jr. Endowed Chair in Psychiatry (UT Medical School) grants to J.C.S.

## References

- Ashburner J. A fast diffeomorphic image registration algorithm. *Neuroimage*. 2007; 38:95–113. [PubMed: 17761438]
- Ashburner J. Computational anatomy with the SPM software. *Magnetic Resonance Imaging*. 2009; 27:1163–1174. [PubMed: 19249168]
- Ashburner J, Friston K. Voxel-Based Morphometry-The methods. *Neuroimage*. 2000; 11:805–821. [PubMed: 10860804]
- Allen EA, Erhardt EB, Damaraju E, Gruner W, Segall JM, Silva RF, Havlicek M, Rachakonda S, Fries J, Kalyanam R, Michael AM, Caprihan A, Turner JA, Eichele T, Adelsheim S, Bryan AD, Bustillo J, Clark VP, Ewing SWF, Filbey F, Ford CC, Hutchison K, Jung RE, Kiehl KA, Koditwakku P, Komesu YM, Mayer AR, Pearlson GD, Phillips JR, Sadek JR, Michael S, Teuscher U, Thoma RJ, Calhoun VD. A baseline for the multivariate comparison of resting-state networks. *Frontiers in systems neuroscience*. 2011; 5:2. [PubMed: 21442040]
- Balci S, Sabuncu M, Yoo J, Gosh S, Gabrieli W, Gabrieli J, Golland P. Prediction of Successful Memory Encoding from fMRI Data. *Med Image Comput Comput Assist Interv*. 2008; 11:97–104. [PubMed: 20401334]
- Bellman, R. Princeton University. 1961. Adaptive control process: A guided tour.
- Birn RM, Murphy K, Bandettini PA. The effect of respiration variations on independent component analysis results of resting state functional connectivity. *Human Brain Mapping*. 2008; 29:740–750. [PubMed: 18438886]
- Bishop, C. *Neural Networks for Pattern Recognition*. Oxford University Press; New York: 1995.
- Bishop, C. *Pattern recognition and machine learning*. Springer; New York: 2006.
- Bonnici HM, Kumaran D, Chadwick MJ, Weiskopf N, Hassabis D, Maguire EA. Decoding representations of scenes in the medial temporal lobes. *Hippocampus*. 2011; 22:1143–1153. [PubMed: 21656874]
- Brammer M. The role of neuroimaging in diagnosis and personalized medicine--current position and likely future directions. *Dialogues In Clinical Neuroscience*. 2009; 11:389–396. [PubMed: 20135896]
- Bray S, Chang C, Hoefft F. Applications of multivariate pattern classification analyses in developmental neuroimaging of healthy and clinical populations. *Frontiers In Human Neuroscience*. 2009; 3:32. [PubMed: 19893761]
- Brown TT, Kuperman JM, Chung Y, Erhart M, McCabe C, Hagler DJ, Venkatraman VK, Akshoomoff N, Amaral DG, Bloss CS. Neuroanatomical assessment of biological maturity. *Current Biology*. 2012
- Bunea F, She Y, Ombao H, Gongvatana A, Devlin K, Cohen R. Penalized least squares regression methods and applications to neuroimaging. *NeuroImage*. 2011; 55:1519–1527. [PubMed: 21167288]
- Calderoni S, Retico A, Biagi L, Tancredi R, Muratori F, Tosetti M. Female children with autism spectrum disorder: An insight from mass-univariate and pattern classification analyses. *Neuroimage*. 2012; 59:1013–1022. [PubMed: 21896334]
- Calhoun, VD. Group ICA Of fMRI Toolbox(GIFT). 2011. <http://mialab.mrn.org/software/gift/http://mialab.mrn.org/software/gift/>
- Calhoun VD, Adali T, Pearlson GD, Pekar JJ. A method for making group inferences from functional MRI data using independent component analysis. *Human Brain Mapping*. 2001; 14:140–151. [PubMed: 11559959]
- Calhoun VD, Adali T. I. Unmixing fMRI with independent component analysis. *IEEE Engineering In Medicine And Biology Magazine: The Quarterly Magazine Of The Engineering In Medicine & Biology Society*. 2006; 25:79–90.



- Calhoun VD, Liu J, Adali T. I. A review of group ICA for fMRI data and ICA for joint inference of imaging, genetic, and ERP data. *Neuroimage*. 2009; 45:S163–S172. [PubMed: 19059344]
- Calhoun VD, Sui J, Kiehl K, Turner J, Allen E, Pearlson G. Exploring the psychosis functional connectome: aberrant intrinsic networks in schizophrenia and bipolar disorder. *Frontiers in psychiatry*. 2011; 2:75. [PubMed: 22291663]
- Caprihan A, Pearlson GD, Calhoun VD. Application of principal component analysis to distinguish patients with schizophrenia from healthy controls based on fractional anisotropy measurements. *Neuroimage*. 2008; 42:675–682. [PubMed: 18571937]
- Carroll MK, Cecchi GA, Rish I, Garg R, Rao AR. Prediction and interpretation of distributed neural activity with sparse models. *NeuroImage*. 2009; 44:112–122. [PubMed: 18793733]
- Casanova R, Whitlow C, Wagner B, Espeland M, Maldjian J. Combining Graph and Machine Learning Methods to Analyze Differences in Functional Connectivity Across Sex. *The open neuroimaging journal*. 2012; 6:1. [PubMed: 22312418]
- Casanova R, Whitlow CT, Wagner B, Williamson J, Shumaker SA, Maldjian JA, Espeland MA. High dimensional classification of structural MRI Alzheimer's disease data based on large scale regularization. *Frontiers in Neuroinformatics*. 2011; 5
- Castro E, Martinez-Ramon M, Pearlson G, Sui J, Calhoun VD. Characterization of groups using composite kernels and multi-source fMRI analysis data: Application to schizophrenia. *NeuroImage*. 2011a; 58:526–536. [PubMed: 21723948]
- Castro E, Martinez-Raman M, Pearlson G, Sui J, Calhoun VD. Characterization of groups using composite kernels and multi-source fMRI analysis data: application to schizophrenia. *Neuroimage*. 2011b; 58:526–536. [PubMed: 21723948]
- Chadwick MJ, Hassabis D, Weiskopf N, Maguire EA. Decoding Individual Episodic Memory Traces in the Human Hippocampus. *Current Biology*. 2010; 20:544–547. [PubMed: 20226665]
- Chai J-W, Chi-Chang Chen C, Chiang C-M, Ho Y-J, Chen H-M, Ouyang Y-C, Yang C-W, Lee S-K, Chang C-I. Quantitative analysis in clinical applications of brain MRI using independent component analysis coupled with support vector machine. *Journal Of Magnetic Resonance Imaging: JMIR*. 2010; 32:24–34. [PubMed: 20578007]
- Chaves R, Ramirez J, Garriz JM, Lopez M, Salas-Gonzalez D, Alvarez I, Segovia F. SVM-based computer-aided diagnosis of the Alzheimer's disease using t-test NMSE feature selection with feature correlation weighting. *Neuroscience Letters*. 2009; 461:293–297. [PubMed: 19549559]
- Cheng W, Ji X, Zhang J, Feng J. Individual classification of ADHD patients by integrating multiscale neuroimaging markers and advanced pattern recognition techniques. *Frontiers In Systems Neuroscience*. 2012; 6
- Chen K, Reiman EM, Huan Z, Caselli RJ, Bandy D, Ayutyanont N, Alexander GE. Linking functional and structural brain images with multivariate network analyses: A novel application of the partial least square method. *Neuroimage*. 2009; 47:602–610. [PubMed: 19393744]
- Chu C, Hsu A-L, Chou K-H, Bandettini P, Lin C. Does feature selection improve classification accuracy? Impact of sample size and feature selection on classification using anatomical magnetic resonance images. *NeuroImage*. 2012; 60:59–70. [PubMed: 22166797]
- Cohen, J. *Statistical power analysis for the behavioural sciences*. 2nd edition. Lawrence Erlbaum Associates; New Jersey: 1998.
- Correa N, Adalı T, Calhoun VD. Performance of blind source separation algorithms for fMRI analysis using a group ICA method. *Magnetic resonance imaging*. 2007; 25:684–694. [PubMed: 17540281]
- Costafreda S, FU C, Picchioni M, Touloupoulou T, McDonald C, Kravariti E, Walsge M, Prata D, Murray R, McGuire P. Pattern of neural responses to verbal fluency shows diagnostic specificity for schizophrenia and bipolar disorder. *BMC Psychiatry*. 2011; 11:18. [PubMed: 21276242]
- Costafreda SG, Chu C, Ashburner J, Fu CHY. Prognostic and diagnostic potential of the structural neuroanatomy of depression. *Plos One*. 2009; 4:e6353. [PubMed: 19633718]
- Coutanche MN, Thompson-Schill SL, Schultz RT. Multi-voxel pattern analysis of fMRI data predicts clinical symptom severity. *NeuroImage*. 2011; 57:113–123. [PubMed: 21513803]

- Cox DD, Savoy RL. Functional magnetic resonance imaging (fMRI) “brain reading”•: detecting and classifying distributed patterns of fMRI activity in human visual cortex. *NeuroImage*. 2003; 19:261–270. [PubMed: 12814577]
- Craddock RC, Holtzheimer PE 3rd, Hu XP, Mayberg HS. Disease state prediction from resting state functional connectivity. *Magnetic Resonance In Medicine: Official Journal Of The Society Of Magnetic Resonance In Medicine/Society Of Magnetic Resonance In Medicine*. 2009; 62:1619–1628.
- Dai Z, Yan C, Wang Z, Wang J, Xia M, Li K, He Y. Discriminative analysis of early Alzheimer’s disease using multi-modal imaging and multi-level characterization with multi-classifier (M3). *NeuroImage*. 2012; 59:2187–2195. [PubMed: 22008370]
- Davatzikos C, Fan Y, Wu X, Shen D, Resnick SM. Detection of prodromal Alzheimer’s disease via pattern classification of magnetic resonance imaging. *Neurobiology Of Aging*. 2008; 29:514–523. [PubMed: 17174012]
- De Martino F, Gentile F, Esposito F, Balsi M, Di Salle F, Goebel R, Formisano E. Classification of fMRI independent components using IC-fingerprints and support vector machine classifiers. *Neuroimage*. 2007; 34:177–194. [PubMed: 17070708]
- De Martino F, Valente G, Staeren N. I. Ashburner J, Goebel R, Formisano E. Combining multivariate voxel selection and support vector machines for mapping and classification of fMRI spatial patterns. *Neuroimage*. 2008; 43:44–58. [PubMed: 18672070]
- Deshpande G, Li Z, Santhanam P, Coles CD, Lynch ME, Hamann S, Hu X. Recursive cluster elimination based support vector machine for disease state prediction using resting state functional and effective brain connectivity. *Plos One*. 2010; 5:e14277. [PubMed: 21151556]
- Dosenbach NUF, Nardos B, Cohen AL, Fair DA, Power JD, Church JA, Nelson SM, Wig GS, Vogel AC, Lessov-Schlaggar CN, et al. Prediction of individual brain maturity using fMRI. *Science*. 2010; 329:1358–1361. [PubMed: 20829489]
- Douglas PK, Harris S, Yuille A, Cohen MS. Performance comparison of machine learning algorithms and number of independent components used in fMRI decoding of belief vs. disbelief. *Neuroimage*. 2011; 56:544–553. [PubMed: 21073969]
- Doyle OM, Ashburner J, Zelaya FO, Williams SCR, Mehta MA, Marquand AF. Multivariate decoding of brain images using ordinal regression. *Neuroimage*. 2013; 81:347–357. [PubMed: 23684876]
- Duchesnay E, Cachia A, Boddaert N, Chabane N, Mangin J-F, Martinot J-L, Brunelle F, Zilbovicius M. Feature selection and classification of imbalanced datasets: Application to PET images of children with autistic spectrum disorders. *NeuroImage*. 2011; 57:1003–1014. [PubMed: 21600290]
- Duchesnay E, Cachia A, Roche A, Riviere D, Cointepas Y, Papadopoulos-Orfanos D, Zilbovicius M, Martinot J-L, Regis J, Mangin J-F. Classification Based on Cortical Folding Patterns. *Medical Imaging, IEEE Transactions*. 2007; 26:553–565.
- Duchesnay, E.; Roche, A.; Riviere, D.; Papadopoulos, D.; Cointepas, Y.; Mangin, J-F. Population classification based on structural morphometry of cortical sulci; Paper presented at: *Biomedical Imaging: Nano to Macro*; 2004; 2004
- Duff EP, Trachtenberg AJ, Mackay CE, Howard MA, Wilson F, Smith SM, Woolrich MW. Task-driven ICA feature generation for accurate and interpretable prediction using fMRI. *NeuroImage*. 2011
- Dukart J, Mueller K, Barthel H, Villringer A, Sabri O, Schroeter ML. Meta-analysis based SVM classification enables accurate detection of Alzheimer’s disease across different clinical centers using FDG-PET and MRI. *Psychiatry Research: Neuroimaging*. 2012
- Ecker C, Rocha-Rego V, Johnston P, Mourao-Miranda J, Marquand A, Daly EM, Brammer MJ, Murphy C, Murphy DG. Investigating the predictive value of whole-brain structural MR scans in autism: a pattern classification approach. *Neuroimage*. 2010; 49:44–56. [PubMed: 19683584]
- Eickhoff S, Laird A, Grefkes C, Wang L, Zilles K, Fox P. Coordinate-based activation likelihood estimation meta-analysis of neuroimaging data: a random-effects approach based on empirical estimates of spatial uncertainty. *Human Brain Mapping*. 2009; 30:2907–2926. [PubMed: 19172646]

- Erhardt EB, Rachakonda S, Bedrick EJ, Allen EA, Adali T, Calhoun VD. Comparison of multi-subject ICA methods for analysis of fMRI data. *Human brain mapping*. 2011; 32:2075–2095. [PubMed: 21162045]
- Fan, Y.; Kaufer, D.; Shen, D. Joint estimation of multiple clinical variables of neurological diseases from imaging patterns. Paper presented at: *Biomedical Imaging: From Nano to Macro*; IEEE International Symposium on (IEEE); 2010; 2010.
- Fan Y, Rao H, Hurt H, Giannetta J, Korczykowski M, Shera D, Avants B, Gee J. Multivariate examination of brain abnormality using both structural and functional MRI. *Neuroimage*. 2007; 36:1189–1199. [PubMed: 17512218]
- Fan Y, Shen D, Gur RC, Gur RE, Davatzikos C. COMPARE: Classification of Morphological Patterns Using Adaptive Regional Elements. *Medical Imaging, IEEE Transactions*. 2006; 26:93–105.
- Formisano E, De Martino F, Valente G. Multivariate analysis of fMRI time series: classification and regression of brain responses using machine learning. *Magnetic Resonance Imaging*. 2008; 26:921–934. [PubMed: 18508219]
- Fort G, Lambert-Lacroix S. Classification using partial least squares with penalized logistic regression. *Bioinformatics*. 2005; 21:1104–1111. [PubMed: 15531609]
- Franke K, Luders E, May A, Wilke M, Gaser C. Brain maturation: Predicting individual BrainAGE in children and adolescents using structural MRI. *NeuroImage*. 2012; 63:1305–1312. [PubMed: 22902922]
- Franke K, Ziegler G, Kloppel S, Gaser C. Estimating the age of healthy subjects from T1-weighted MRI scans using kernel methods: exploring the influence of various parameters. *Neuroimage*. 2010; 50:883–892. [PubMed: 20070949]
- Friedman J, Hastie T, Tibshirani R. Regularization paths for generalized linear models via coordinate descent. *Journal of statistical software*. 2010; 33:1. [PubMed: 20808728]
- Friston KJ, Poline JB, Holmes AP, Frith CD, Frakowiack RS. A multivariate analysis of PET activation studies. *Human Brain Mapping*. 1996; 4:140–151. [PubMed: 20408193]
- Fu CHY, Mourao-Miranda J, Costafreda SG, Khanna A, Marquand AF, Williams SCR, Brammer MJ. Pattern classification of sad facial processing: toward the development of neurobiological markers in depression. *Biological Psychiatry*. 2008; 63:656–662. [PubMed: 17949689]
- Gothelf D, Hoefl F, Ueno T, Sugiura L, Lee AD, Thompson P, Reiss AL. Developmental changes in multivariate neuroanatomical patterns that predict risk for psychosis in 22q11.2 deletion syndrome. *Journal Of Psychiatric Research*. 2011; 45:322–331. [PubMed: 20817203]
- Grana M, Termenon M, Savio A, Gonzalez-Pinto A, Echeveste J, Perez JM, Besga A. Computer Aided Diagnosis system for Alzheimer Disease using brain Diffusion Tensor Imaging features selected by Pearson's correlation. *Neuroscience Letters*. 2011; 502:225–229. [PubMed: 21839143]
- Greenberg AS, Esterman M, Wilson D, Serences JT, Yantis S. Control of spatial and feature-based attention in frontoparietal cortex. *The Journal of Neuroscience*. 2010; 30:14330–14339. [PubMed: 20980588]
- Guyon I, Elisseeff A. An Introduction to Variable and Feature Selection. *Journal of Machine Learning*. 2003; 7/8:1157–1182.
- Guyon I, Weston J, Barnhill S, Vapnik V. Gene Selection for Cancer Classification using Support Vector Machines. *Machine Learning*. 2001; 46:389–422.
- Haller S, Bartsch A, Nguyen D, Rodriguez C, Emch J, Gold G, Lovblad KO, Giannakopoulos P. Cerebral microhemorrhage and iron deposition in mild cognitive impairment: susceptibility-weighted MR imaging assessment. *Radiology*. 2010a; 257:764–773. [PubMed: 20923870]
- Haller S, Nguyen D, Rodriguez C, Emch J, Gold G, Bartsch A, Lovblad KO, Giannakopoulos P. Individual Prediction of Cognitive Decline in Mild Cognitive Impairment Using Support Vector Machine-Based Analysis of Diffusion Tensor Imaging Data. *Journal of Alzheimer's disease*. 2010b; 22:315–327.
- Hansen LK, Larsen J, Nielsen FÅ, Strother SC, Rostrup E, Savoy R, Lange N, Sidtis J, Svarer C, Paulson OB. Generalizable patterns in neuroimaging: How many principal components? *NeuroImage*. 1999; 9:534–544. [PubMed: 10329293]

- Hanson SJ, Halchenko YO. Brain reading using full brain support vector machines for object recognition: there is no “face” identification area. *Neural Computation*. 2008; 20:486–503. [PubMed: 18047411]
- Haynes J-D, Rees G. Predicting the orientation of invisible stimuli from activity in human primary visual cortex. *Nature Neuroscience*. 2005; 8:686–691.
- He W, Wang Z, Jiang H. Model optimizing and feature selecting for support vector regression in time series forecasting. *Neurocomputing Machine Learning for Signal Processing (MLSP 2006)/Life System Modelling, Simulation, and Bio-inspired Computing (LSMS 2007)*. 2008; 72:600–611.
- Hinrichs C, Singh V, Mukherjee L, Xu G, Chung MK, Johnson SC. Spatially augmented LPboosting for AD classification with evaluations on the ADNI dataset. *NeuroImage*. 2009; 48:138–149. [PubMed: 19481161]
- Hoeft F, Walter E, Lightbody AA, Hazlett HC, Chang C, Piven J, Reiss AL. Neuroanatomical differences in toddler boys with fragile × syndrome and idiopathic autism. *Archives Of General Psychiatry*. 2011; 68:295–305. [PubMed: 21041609]
- Hua, Tembe W, Dougherty E. Performance of feature-selection methods in the classification of high-dimension data. *Pattern Recognition*. 2009; 42:409–424.
- Ince, NF.; Goksu, F.; Pellizzer, G.; Tewfik, A.; Stephane, M. Selection of spectro-temporal patterns in multichannel MEG with support vector machines for schizophrenia classification; Paper presented at: Engineering in Medicine and Biology Society; 2008; 2008
- Ingalhalikar M, Parker D, Bloy L, Roberts TPL, Verma R. Diffusion based abnormality markers of pathology: toward learned diagnostic prediction of ASD. *Neuroimage*. 2011; 57:918–927. [PubMed: 21609768]
- Johansen-Berg, H.; Behrens, T. Diffusion MRI - from quantitative measurements to in-vivo neuroanatomy. Academic Press; London, UK: 2009.
- Johnston, B.; Mwangi, B.; Matthews, K.; Coghill, D.; Steele, J. European Child & Adolescent Psychiatry (Springer Berlin/Heidelberg). 2012. Predictive classification of individual magnetic resonance imaging scans from children and adolescents; p. 1-12.
- Joliffe, I. Principle component Analysis. Springer-Verlag; New York: 2002.
- Kjems U, Hansen LK, Anderson J, Frutiger S, Muley S, Sittis J, Rottenberg D, Strother S. The quantitative evaluation of functional neuroimaging experiments: Mutual information learning curves. *NeuroImage*. 2002; 15:772–786. [PubMed: 11906219]
- Kloppel S, Stonnington C, Chu C, Draganski B, Scahill R, Rohrer J, Fox N, Jack C Jr, Ashburner J, Frackowiak R. Automatic classification of MR scans in Alzheimer’s disease. *Brain: A Journal of Neurology*. 2008; 131:681–689. [PubMed: 18202106]
- Kohannim, O.; Hibar, D.; Jahanshad, N.; Stein, J.; Hua, X.; Toga, A.; Jack, C.; Weinen, M.; Thompson, P. Predicting temporal lobe volume on MRI from genotypes using L1-L2 regularized regression; 2012a; Paper presented at: Biomedical Imaging (ISBI), 2012 9th IEEE International Symposium on (IEEE);
- Kohannim O, Hibar DP, Stein JL, Jahanshad N, Hua X, Rajagopalan P, Toga AW, Jack CR Jr, Weiner MW, de Zubicaray GI. Discovery and replication of gene influences on brain structure using LASSO regression. *Frontiers in Neuroscience*. 2012b; 6
- Kohavi R, John G. Wrappers for Feature Subset Selection. *Artificial Intelligence*. 1997; 97:1–2.
- Koutsouleris N, Meisenzahl EM, Davatzikos C, Bottlender R, Frodl T, Scheuerecker J, Schmitt G, Zetzsche T, Decker P, Reiser M, et al. Use of neuroanatomical pattern classification to identify subjects in at-risk mental states of psychosis and predict disease transition. *Archives Of General Psychiatry*. 2009; 66:700–712. [PubMed: 19581561]
- Kovalev VA, Petrou M, Suckling J. Detection of structural differences between the brains of schizophrenic patients and controls. *Psychiatry Research: Neuroimaging*. 2003; 124:177–189.
- Kriegeskorte N, Goebel R, Bandettini P. Information-based functional brain mapping. *Proceedings of the National Academy of Sciences of the United States of America*. 2006; 103:3863–3868. [PubMed: 16537458]
- Kriegeskorte N, Lindquist MA, Nichols TE, Poldrack RA, Vul E. Everything you never wanted to know about circular analysis, but were afraid to ask. *Journal of Cerebral Blood Flow & Metabolism*. 2010; 30:1551–1557. [PubMed: 20571517]

- Kriegeskorte N, Simmons WK, Bellgowan PSF, Baker CI. Circular analysis in systems neuroscience: the dangers of double dipping. *Nature Neuroscience*. 2009; 12:535–540.
- Krishnan A, Williams LJ, McIntosh AR, Abdi H. Partial Least Squares (PLS) methods for neuroimaging: A tutorial and review. *NeuroImage*. 2011; 56:455–475. [PubMed: 20656037]
- Kwok JTY, Tsang IWH. The pre-image problem in kernel methods. *Neural Networks, IEEE Transactions*. 2004; 15:1517–1525.
- LaConte S, Strother S, Cherkassky V, Anderson J, Hu X. Support vector machines for temporal classification of block design fMRI data. *NeuroImage*. 2005; 26:317. [PubMed: 15907293]
- Laird A, Lancaster J, Fox P. BrainMap: the social evolution of a functional neuroimaging database. *Neuroinformatics*. 2005a; 3:65–78. [PubMed: 15897617]
- Laird AR, McMillan KM, Lancaster JL, Kochunov P, Turkeltaub PE, Pardo JV, Fox PT. A comparison of label-based review and ALE meta-analysis in the Stroop task. *Human Brain Mapping*. 2005b; 25:6–21. [PubMed: 15846823]
- Langs G, Menze BH, Lashkari D, Golland P. Detecting stable distributed patterns of brain activation using Gini contrast. *NeuroImage*. 2011; 56:497–507. [PubMed: 20709176]
- Lee, J.; Verleysen, M. *Nonlinear Dimensionality Reduction*. Springer Publishing Co.; New York, USA: 2007.
- Lim L, Marquand A, Cubillo AA, Smith AB, Chantiluke K, Simmons A, Mehta M, Rubia K. Disorder-Specific Predictive Classification of Adolescents with Attention Deficit Hyperactivity Disorder (ADHD) Relative to Autism Using Structural Magnetic Resonance Imaging. *PLOS ONE*. 2013; 8:e63660. [PubMed: 23696841]
- Linden DEJ. The challenges and promise of neuroimaging in psychiatry. *Neuron*. 2012; 73:8–22. [PubMed: 22243743]
- Liu M, Zhang D, Shen D. Ensemble sparse classification of Alzheimer's disease. *NeuroImage*. 2012; 60:1106–1116. [PubMed: 22270352]
- Lohmann G, Volz KG, Ullsperger M. Using non-negative matrix factorization for single-trial analysis of fMRI data. *NeuroImage*. 2007; 37:1148–1160. [PubMed: 17662621]
- Lopez M, Ramirez J, Garriz J, Álvarez I, Salas-Gonzalez D, Segovia F, Chaves R, Padilla P, Gomez-Rao M. Principal component analysis-based techniques and supervised classification schemes for the early detection of Alzheimer's disease. *Neurocomputing*. 2011; 74:1260–1271.
- MacIntosh AR, Bookstein F, Haxby J, Grady C. Spatial pattern analysis of functional brain images using partial least squares. *NeuroImage*. 1996; 3:143–157. [PubMed: 9345485]
- McIntosh AR, Lobaugh NJ. Partial least squares analysis of neuroimaging data: applications and advances. *NeuroImage*. 2004; 23:250–263.
- McIntosh AR, Misić B. Multivariate statistical analyses for neuroimaging data. *Annu. Rev. Psychol.* 2013; 64:499–525. [PubMed: 22804773]
- Magnin B, Mesrob L, Kinkingnahun S, Palagrini-Issac M, Colliot O, Sarazin M, Dubois B, Leharicy S, Benali H. Support vector machine-based classification of Alzheimer's disease from whole-brain anatomical MRI. *Neuroradiology*. 2009; 51:73–83. [PubMed: 18846369]
- Marquand A, Howard M, Brammer M, Chu C, Coen S, Mourao-Miranda J. Quantitative prediction of subjective pain intensity from whole-brain fMRI data using Gaussian processes. *NeuroImage*. 2010a; 49:2178–2189. [PubMed: 19879364]
- Marquand A, Howard M, Brammer M, Chu C, Coen S, Mourao-Miranda J. Quantitative prediction of subjective pain intensity from whole-brain fMRI data using Gaussian processes. *NeuroImage*. 2010b; 49:2178–2189. [PubMed: 19879364]
- Marquand AF, De Simoni S, O'Daly OG, Williams SCR, Mourao-Miranda J, Mehta MA. Pattern classification of working memory networks reveals differential effects of methylphenidate, atomoxetine, and placebo in healthy volunteers. *Neuropsychopharmacology*. 2011; 36:1237–1247. [PubMed: 21346736]
- Marquand AF, O'Daly OG, De Simoni S, Alsop DC, Maguire RP, Williams SCR, Zelaya FO, Mehta MA. Dissociable effects of methylphenidate, atomoxetine and placebo on regional cerebral blood flow in healthy volunteers at rest: A multi-class pattern recognition approach. *NeuroImage*. 2012; 60:1015–1024. [PubMed: 22266414]

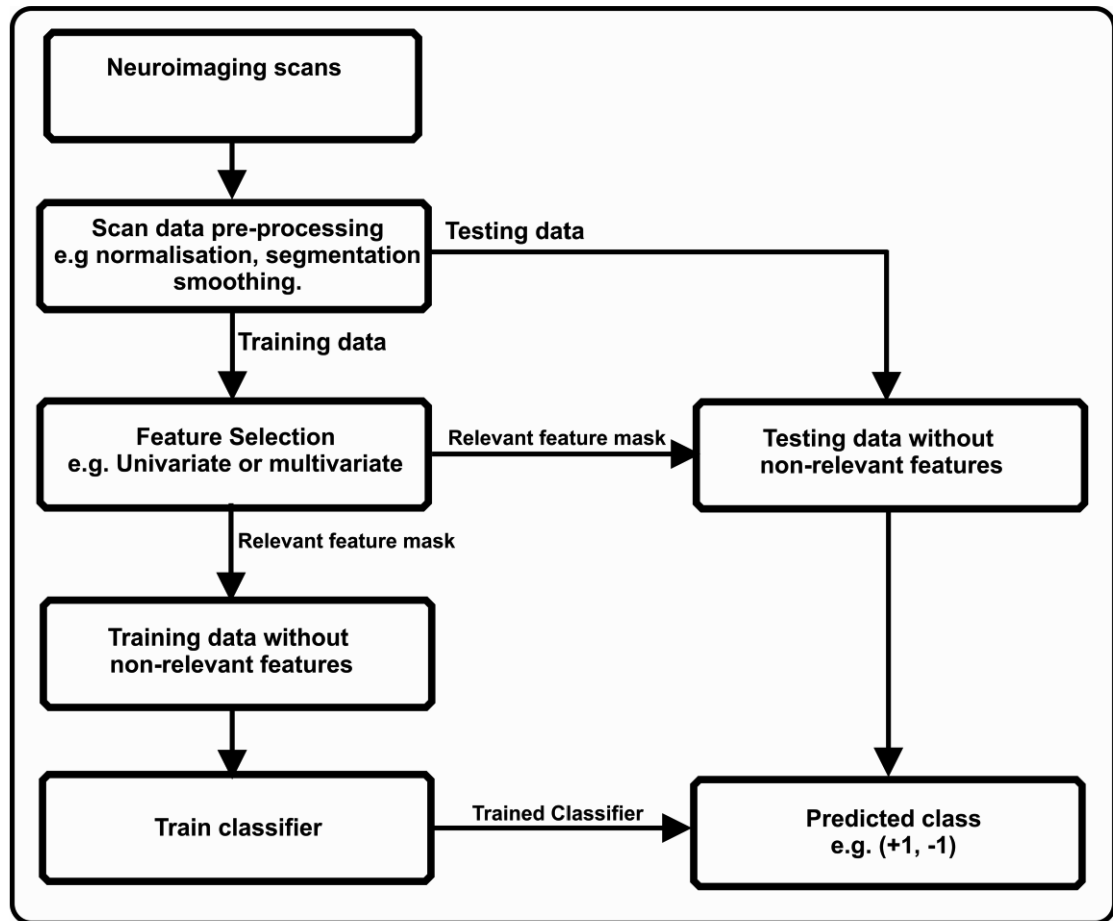
- Martinez-Montes E, Valdes-Sosa PA, Miwakeichi F, Goldman RI, Cohen MS. EEG/fMRI analysis by multiway partial least squares. *Neuroimage*. 2004; 22:1023–34. [PubMed: 15219575]
- Mesrob L, Magnin B, Colliot O, Sarazin M, Hahn-Barma V, Dubois B, Gallinari P, Leharicy S, Kinkinghuhn S, Benali H. Identification of Atrophy Patterns in Alzheimers Disease Based on SVM Feature Selection and Anatomical Parcellation. *Medical Imaging and Augmented Reality*. 2008:124–132.
- Menzies L, Achard S, Chamberlain SR, Fineberg N, Chen C-H, delCampo N, Sahakian BJ, Robbins TW, Bullmore E. Neurocognitive endophenotypes of obsessive-compulsive disorder. *Brain*. 2007; 130:3223–3236. [PubMed: 17855376]
- Misaki M, Kim Y, Bandettini PA, Kriegeskorte N. Comparison of multivariate classifiers and response normalizations for pattern-information fMRI. *NeuroImage*. 2010; 53:103–118. [PubMed: 20580933]
- Mitchell TM. From journal articles to computational models: a new automated tool. *Nature methods*. 2011; 8:627. [PubMed: 21799495]
- Mitchell TM, Hutchinson R, Niculescu RS, Pereira F, Wang X, Just M, Newman S. Learning to decode cognitive states from brain images. *Machine Learning*. 2004; 57:145–175.
- Mourao-Miranda J, Bokde ALW, Born C, Hampel H, Stetter M. Classifying brain states and determining the discriminating activation patterns: Support Vector Machine on functional MRI data. *NeuroImage*. 2005; 28:980–995. [PubMed: 16275139]
- Mourao-Miranda J, Oliveira L, Ladouceur CD, Marquand A, Brammer M, Birmaher B, Axelson D, Phillips ML. Pattern recognition and functional neuroimaging help to discriminate healthy adolescents at risk for mood disorders from low risk adolescents. *Plos One*. 2012; 7:e29482. [PubMed: 22355302]
- Mourao-Miranda J, Reynaud E, McGlone F, Calvert G, Brammer M. The impact of temporal compression and space selection on SVM analysis of single-subject and multi-subject fMRI data. *NeuroImage*. 2006; 33:1055–1065. [PubMed: 17010645]
- Mwangi B, Ebmeier K, Matthews K, Douglas Steele J. Multi-centre diagnostic classification of individual structural neuroimaging scans from patients with major depressive disorder. *Brain: A Journal of Neurology*. 2012a; 135:1508–1521. [PubMed: 22544901]
- Mwangi B, Hasan KM, Soares JC. Prediction of individual subject's age across the human lifespan using diffusion tensor imaging: A machine learning approach. *NeuroImage*. 2013
- Mwangi B, Matthews K, Steele J. Prediction of illness severity in patients with major depression using structural MR brain scans. *Journal of Magnetic Resonance Imaging*. 2012b; 35:64–71. [PubMed: 21959677]
- Nester PG, O'Donnell BF, Mccarley RW, Niznikiewicz M, Barnard J, Shen ZJ, Bookstein FL, Shenton ME. A new statistical method for testing hypotheses of neuropsychological/MRI relationships in schizophrenia: partial least squares analysis. *Schizophr. Res*. 2002; 53:57–66. [PubMed: 11728838]
- Nho, K.; Shen, L.; Kim, S.; Risacher, SL.; West, JD.; Foroud, T.; Jack, CR.; Weiner, MW.; Saykin, AJ. Automatic Prediction of Conversion from Mild Cognitive Impairment to Probable Alzheimer's Disease using Structural Magnetic Resonance Imaging; AMIA Annual Symposium Proceedings/AMIA Symposium AMIA Symposium 2010; 2010; p. 542-546.
- Norman KA, Polyn SM, Detre GJ, Haxby JV. Beyond mind-reading: multi-voxel pattern analysis of fMRI data. *Trends in cognitive sciences*. 2006; 10:424–430. [PubMed: 16899397]
- Ogutu, JO.; Schulz-Streeck, T.; Piepho, HP. Genomic selection using regularized linear regression models: ridge regression, lasso, elastic net and their extensions; Paper presented at: BMC proceedings (Springer); 2012;
- Orru G, Pettersson-Yeo W, Marquand AF, Sartori G, Mechelli A. Using Support Vector Machine to identify imaging biomarkers of neurological and psychiatric disease: A critical review. *Neuroscience And Biobehavioral Reviews*. 2012; 36:1140–1152. [PubMed: 22305994]
- Penny, W.; Friston, K.; Ashburner, J.; Nicols, T. Academic Press. 2007. Statistical parametric mapping: The Analysis of functional Brain images.
- Pereira F, Mitchell T, Botvinick M. Machine learning classifiers and fMRI: A tutorial overview. *NeuroImage Mathematics in Brain Imaging*. 2009; 45:S199–S209.

- Radulescu AR, Mujica-Parodi LR. A principal component network analysis of prefrontal-limbic functional magnetic resonance imaging time series in schizophrenia patients and healthy controls. *Psychiatry Research*. 2009; 174:184–194. [PubMed: 19880294]
- Rao, A.; Lee, Y.; Gass, A.; Monsch, A. Classification of Alzheimer's Disease from structural MRI using sparse logistic regression with optional spatial regularization; Paper presented at: Engineering in Medicine and Biology Society, EMBC, 2011 Annual International Conference of the IEEE (IEEE); 2011;
- Rasmussen, C.; Williams, C. MIT Press. 2006. *Gaussian Processes for Machine Learning*.
- Rasmussen PM, Abrahamsen TJ, Madsen KH, Hansen LK. Nonlinear denoising and analysis of neuroimages with kernel principal component analysis and pre-image estimation. *NeuroImage*. 2012; 60:1807–1818. [PubMed: 22305952]
- Rasmussen PM, Madsen KH, Lund TE, Hansen LK. Visualization of nonlinear kernel models in neuroimaging by sensitivity maps. *NeuroImage*. 2011; 55:1120–1131. [PubMed: 21168511]
- Rish I, Cecchi G, Baliki M, Apkarian A. Sparse regression models of pain perception. *Brain Informatics*. 2010:212–223.
- Rizk-Jackson A, Stoffers D, Sheldon S, Kuperman J, Dale A, Goldstein J, Corey-Bloom J, Poldrack RA, Aron AR. Evaluating imaging biomarkers for neurodegeneration in pre-symptomatic Huntington's disease using machine learning techniques. *NeuroImage*. 2011; 56:788–796. [PubMed: 20451620]
- Ryali S, Supekar K, Abrams DA, Menon V. Sparse logistic regression for whole brain classification of fMRI data. *NeuroImage*. 2010; 51:752. [PubMed: 20188193]
- Saeys Y, Inza I, Larranaga P. A review of feature selection techniques in bioinformatics. *Bioinformatics*. 2007
- Salimi-Khorshidi G, Smith SM, Keltner JR, Wager TD, Nichols TE. Meta-analysis of neuroimaging data: a comparison of image-based and coordinate-based pooling of studies. *Neuroimage*. 2009; 45:810–823. [PubMed: 19166944]
- Sato JR, Hoexter MQ, Fujita A, Rohde LA. Evaluation of pattern recognition and feature extraction methods in ADHD prediction. *Frontiers in Systems Neuroscience*. 2012:6. [PubMed: 22438838]
- Scholkopf, B.; Smola, A. *Learning with Kernels: Support Vector Machines, Regularization, Optimization, and Beyond*. MIT Press; Cambridge-MA: 2002.
- Schrouff J, Rosa MJ, Rondina J, Marquand A, Chu C, Ashburner J, Phillips C, Richiardi J, Mourão-Miranda J. PRoNTO: Pattern Recognition for Neuroimaging Toolbox. *Neuroinformatics*. 2013:1–19. [PubMed: 23224666]
- Scott, D. *Multivariate density estimation: Theory, practice and visualization*. Wiley; New York: 1992.
- Shen L, Kim S, Qi Y, Inlow M, Swaminathan S, Nho K, Wan J, Risacher S, Shaw L, Trojanowski J. Identifying neuroimaging and proteomic biomarkers for MCI and AD via the elastic net. *Multimodal Brain Image Analysis*. 2011:27–34.
- Sheskin, D. *Handbook of parametric and nonparametric statistical procedures Florida*. Chapman & Hall; 2004.
- Shi, W.; Lee, KE.; Wahba, G. Detecting disease-causing genes by LASSO-Patternsearch algorithm; Paper presented at: BMC proceedings, (BioMed Central Ltd); 2007;
- Sidhu GS, Asgarian N, Greiner R, Brown MRG. Kernel Principal Component Analysis for dimensionality reduction in fMRI-based diagnosis of ADHD. *Frontiers in Systems Neuroscience*. 2012; 6:74. [PubMed: 23162439]
- Stone, J. *Independent Component Analysis*. MIT Press; Cambridge, MA: 2004.
- Stonnington CM, Chu C, Kloppel S, Jack CR Jr, Ashburner J, Frackowiak RSJ. Predicting clinical scores from magnetic resonance scans in Alzheimer's disease. *Neuroimage*. 2010; 51:1405–1413. [PubMed: 20347044]
- Strother SC, Anderson J, Hansen LK, Kjems U, Kustra R, Sidtis J, Frutiger S, Muley S, LaConte S, Rottenberg D. The quantitative evaluation of functional neuroimaging experiments: The NPAIRS data analysis framework. *NeuroImage*. 2002; 15:747–771. [PubMed: 11906218]
- Sui J, Adali T, Yu Q, Chen J, Calhoun VD. A review of multivariate methods for multimodal fusion of brain imaging data. *Journal of Neuroscience Methods*. 2012; 204:68–81. [PubMed: 22108139]

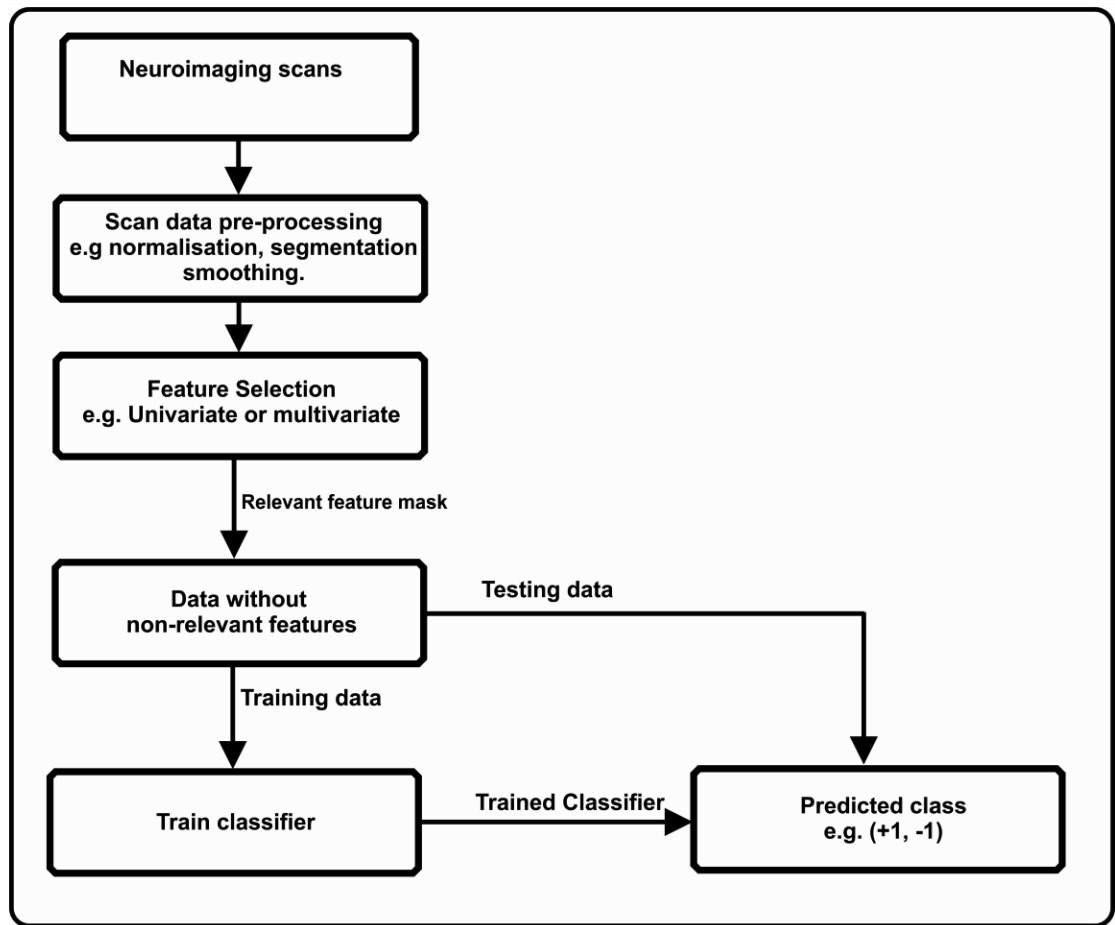
- Tagliazucchi E, von Wegner F, Morzelewski A, Borisov S, Jahnke K, Laufs H. Automatic sleep staging using fMRI functional connectivity data. *Neuroimage*. 2012
- Theodoridis, S.; Koutroumbas, K. *Pattern Recognition*. 4th edition. Elsevier; San Diego, California: 2009.
- Thirion B, Fugeras O. Dynamical components analysis of fMRI data through kernel PCA. *NeuroImage*. 2003; 20:34–49. [PubMed: 14527568]
- Tibshirani R. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society Series B (Methodological)*. 1996;267–288.
- Tibshirani R. Regression shrinkage and selection via the lasso: a retrospective. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*. 2011; 73:273–282.
- Tipping M. Sparse Bayesian Learning and the Relevance Vector Machine. *Journal of Machine Learning Research*. 2001; 1:211–244.
- Toussaint PJ, Perlberg V, Bellec P, Desarnaud S, Lacomblez L, Doyon J, Habert MO, Benali H. Resting state FDG-PET functional connectivity as an early biomarker of Alzheimer’s disease using conjoint univariate and independent component analyses. *NeuroImage*. 2012
- Tseng P, Yun S. Block-coordinate gradient descent method for linearly constrained nonsmooth separable optimization. *Journal of optimization theory and applications*. 2009; 140:513–535.
- Valente G, De Martino F, Esposito F, Goebel R, Formisano E. Predicting subject-driven actions and sensory experience in a virtual world with Relevance Vector Machine Regression of fMRI data. *NeuroImage*. 2011; 56:651–661. [PubMed: 20888922]
- Van De Ville D, Lee S-W. Brain decoding: Opportunities and challenges for pattern recognition. *Pattern Recognition Brain Decoding*. 2012; 45:2033–2034.
- Vapnik, V. *The nature of statistical learning theory*. 2nd edition. Springer-verlag; New York: 1999.
- Vounou M, Janousova E, Wolz R, Stein JL, Thompson PM, Rueckert D, Montana G. Sparse reduced-rank regression detects genetic associations with voxel-wise longitudinal phenotypes in Alzheimer’s disease. *Neuroimage*. 2011
- Wager TD, Lindquist M, Kaplan L. Meta-analysis of functional neuroimaging data: current and future directions. *Social Cognitive And Affective Neuroscience*. 2007; 2:150–158. [PubMed: 18985131]
- Wan J, Kim S, Inlow M, Nho K, Swaminathan S, Risacher S, Fang S, Weiner M, Beg M, Wang L. Hippocampal surface mapping of genetic risk factors in AD via sparse learning models. *Medical Image Computing and Computer-Assisted Intervention, MICCAI*. 2011;2011:376–383.
- Wang, X.; Jiao, Y.; Lu, Z. Discriminative analysis of resting-state brain functional connectivity patterns of Attention-Deficit Hyperactivity Disorder using Kernel Principal Component Analysis; Paper presented at: Fuzzy Systems and Knowledge Discovery (FSKD), 2011 Eighth International Conference on (IEEE); 2011;
- Wang Y, Fan Y, Bhatt P, Davatzikos C. High-dimensional pattern regression using machine learning: from medical images to continuous clinical variables. *Neuroimage*. 2010; 50:1519–1535. [PubMed: 20056158]
- Wee C-Y, Yap P-T, Li W, Denny K, Browndyke JN, Potter GG, Welsh-Bohmer KA, Wang L, Shen D. Enriched white matter connectivity networks for accurate identification of MCI patients. *Neuroimage*. 2011; 54:1812–1822. [PubMed: 20970508]
- Wee C-Y, Yap P-T, Zhang D, Denny K, Browndyke JN, Potter GG, Welsh-Bohmer KA, Wang L, Shen D. Identification of MCI individuals using structural and functional connectivity networks. *NeuroImage*. 2012; 59:2045–2056. [PubMed: 22019883]
- Wold S, Sjostrom M, Eriksson L. PLS-regression: a basic tool of chemometrics. *Chemom. Intell. Lab. Syst*. 2001; 58:109–130.
- Yan J, Risacher S, Kim S, Simon J, Li T, Wan J, Wang H, Huang H, Saykin A, Shen L. Multimodal Neuroimaging Predictors for Cognitive Performance Using Structured Sparse Learning. *Multimodal Brain Image Analysis*. 2012:1–17.
- Yang H, Liu J, Sui J, Pearlson G, Calhoun VD. A hybrid machine learning method for fusing fMRI and genetic data: combining both improves classification of schizophrenia. *Frontiers in human neuroscience*. 2010;4. [PubMed: 20198130]



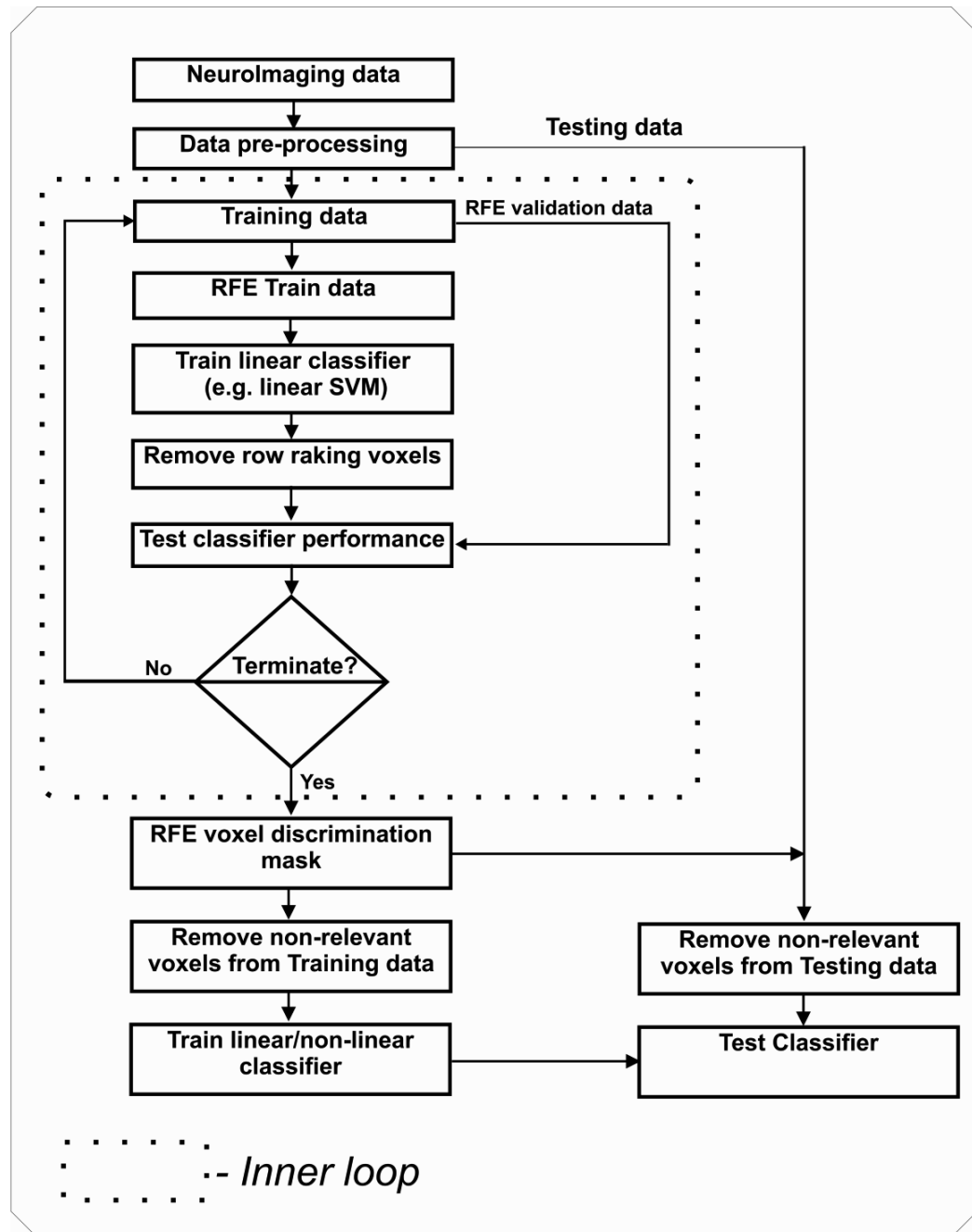
- Yarkoni T, Poldrack RA, Nichols TE, Van Essen DC, Wager TD. Large-scale automated synthesis of human functional neuroimaging data. *Nature methods*. 2011; 8:665–670. [PubMed: 21706013]
- Yoon JH, Tamir D, Minzenberg MJ, Ragland JD, Ursu S, Carter CS. Multivariate Pattern Analysis of Functional Magnetic Resonance Imaging Data Reveals Deficits in Distributed Representations in Schizophrenia. *Biological Psychiatry Schizophrenia: Neural Circuitry and Molecular Mechanisms*. 2008; 64:1035–1041.
- Yoon U, Lee JM, Im K, Shin YW, Cho BH, Kim IY, Kwon JS, Kim SI. Pattern classification using principal components of cortical thickness and its discriminative pattern in schizophrenia. *Neuroimage*. 2007; 34:1405–1415. [PubMed: 17188902]
- Zeng L-L, Shen H, Liu L, Wang L, Li B, Fang P, Zhou Z, Li Y, Hu D. Identifying major depression using whole-brain functional connectivity: a multivariate pattern analysis 10.1093/brain/aws059. *Brain*. 2012; 135:1498–1507. [PubMed: 22418737]
- Zhang D, Wang Y, Zhou L, Yuan H, Shen D. Multimodal classification of Alzheimer’s disease and mild cognitive impairment. *NeuroImage*. 2011a; 55:856–867. [PubMed: 21236349]
- Zhang D, Wang Y, Zhou L, Yuan H, Shen D. Multimodal classification of Alzheimer’s disease and mild cognitive impairment. *Neuroimage*. 2011b; 55:856–867. [PubMed: 21236349]
- Zhu CZ, Zang YF, Cao QJ, Yan CG, He Y, Jiang TZ, Sui MQ, Wang YF. Fisher discriminative analysis of resting-state brain function for attention-deficit/hyperactivity disorder. *Neuroimage*. 2008; 40:110–120. [PubMed: 18191584]
- Zhu, CZ.; Zang, YF.; Liang, M.; Tian, LX.; He, Y.; Li, XB.; Sui, MQ.; Wang, YF.; Jiang, TZ. Discriminative analysis of brain function at resting-state for attention-deficit/hyperactivity disorder; *Medical Image Computing And Computer-Assisted Intervention: MICCAI International Conference On Medical Image Computing And Computer-Assisted Intervention*; 2005; p. 468-475.
- Zigler G, Dahnke R, Winkler AD, Gaser C. Partial least squares correlation of multivariate cognitive abilities and local brain structure in children and adolescents. *Neuroimage*. 2013; 82:284–294. [PubMed: 23727321]
- Zou H, Hastie T. Regularization and variable selection via the elastic net. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*. 2005; 67:301–320.



**Figure 1.** Feature reduction without *double-dipping*. Training and testing datasets are separated before the feature reduction process.



**Figure 2.** Feature reduction with *double-dipping*. Features are selected from the same set of training and testing data.



**Figure 3.**  
Flow diagram illustrating the recursive feature elimination (RFE) process.

**Table 1**

A tabular summary of both supervised and unsupervised feature reduction techniques.

<b>Supervised</b>	<b>Unsupervised</b>
a) Filter techniques <ul style="list-style-type: none"> <li>- T-test, Anova, pearson correlation coefficient</li> </ul>	a) Data driven <ul style="list-style-type: none"> <li>- Principal component analysis</li> <li>- Independent component analysis</li> </ul>
b) Wrapper techniques <ul style="list-style-type: none"> <li>- Recursive feature elimination</li> <li>- Searchlight</li> </ul>	
c) Embedded techniques <ul style="list-style-type: none"> <li>- Least absolute shrinkage and selection operator</li> <li>- Elastic Net</li> <li>- Partial least squares method</li> </ul>	b) Domain Knowledge driven <ul style="list-style-type: none"> <li>- coordinate-based meta-analysis</li> </ul>