

Feature clustering and ranking for selecting stable features from high dimensional remotely sensed data

Dugal Harris ^a & Adriaan van Niekerk ^b

^a *Department of Geography and Environmental Studies, Stellenbosch University, Stellenbosch 7602, South Africa, (email: dugalh@gmail.com, tel: +27 82 843 9679, postal address: PO Box 180, Newlands, Cape Town, 7725)*

^b *Centre for Geographical Analysis, Stellenbosch University, Stellenbosch 7602, South Africa, (email: avn@sun.ac.za, tel: +27 21 808 3101, postal address: Private Bag X1, Matieland, Stellenbosch, 7602)*

Corresponding author: Dugal Harris

Feature clustering and ranking for selecting stable features from high dimensional remotely sensed data

High dimensional remote sensing data sets typically contain redundancy amongst the features. Traditional approaches to feature selection are prone to instability and selection of sub-optimal features in these circumstances. They can also be computationally expensive, especially when dealing with very large remote sensing datasets. This article presents an efficient, deterministic feature ranking method that is robust to redundancy. ~~Average-linkage-hierarchical-clustering~~Affinity propagation is used to group correlated features into clusters. A relevance criterion is evaluated for each feature. Clusters are then ranked based on the median of the relevance values of their constituent features. The most relevant ~~Individual-individual~~ features can then be selected automatically from the best clusters. Other criteria, such as computation time or measurement cost, can optionally be considered interactively when making this selection. The proposed feature selection method is compared to ~~traditional-competing~~ filter approach methods on a number of remote sensing data sets containing feature redundancy. Mutual information and naive Bayes relevance criteria were evaluated in conjunction with the feature selection methods. Using the proposed method it was shown that the stability of selected features improved under different data samplings, while similar or better classification accuracies were achieved compared to ~~traditional~~ competing methods.

Keywords: feature selection; clustering; high dimensionality; redundancy; stability

1 Introduction

In image classification, the amount of training data required to adequately represent class distributions in feature space increases exponentially as the number of features (variables) is increased – a phenomena known as the ‘curse of dimensionality’ (Bishop 2003). For finite training samples, increasing the features beyond a certain point results in overtraining and a decrease in the classifier accuracy. This so-called ‘peaking phenomenon’ (Jain, Duin, and Mao 2000) requires the size of the feature set to be reduced to a salient minimum in order to achieve an accurate classification. Support vector machine (SVM) (Burges 1998) and

random forest (Breiman 2001) classifiers have become popular in remote sensing, partly because of their lack of sensitivity to the peaking phenomenon (Guyon et al. 2002), but several studies have demonstrated the benefits of feature reduction when these classifiers were applied to high-dimensional data (Guyon et al. 2002; Strobl et al. 2008; Tolosi and Lengauer 2011). Reducing the number of features is also beneficial from the perspective of measurement costs and feature computation time. This is particularly relevant in large scale remote sensing studies involving Very High Resolution (VHR) imagery, as vast quantities of data require processing. ~~A further motivation for reducing features to an informative minimum is the ‘ugly duckling theorem’, which implies that the more redundant features contained in a data set, the less separable classes become (Jain, Duin, and Mao 2000).~~

Two basic approaches to feature set reduction exist: feature selection and feature extraction. In feature extraction the feature set is mapped into a new feature space of reduced dimensionality (Webb 2002). ~~Various criteria, such as separability and variance, are used to define the dimensions of the new ‘optimal’ space. Principal Components Analysis (PCA) (Webb 2002) is an example of a popular feature extraction method. It uses a linear transform composed of the largest m eigenvectors of the covariance matrix to project the input features into a reduced space. Other methods such as projection pursuit and independent component analysis (ICA) also incorporate linear projections, but do not rely on covariance, and are thus better suited to non-Gaussian distributed data than PCA (Jain, Duin, and Mao 2000). Commonly used non-linear methods include kernel PCA (Schölkopf, Smola, and Müller 1998) and multidimensional scaling (MDS) (Webb 2002).~~ A major disadvantage of the feature extraction approach is that it requires measurements and computations to produce the full feature set, which can be prohibitively costly. Feature extraction also hinders

interpretability as it alters the original representation of the features. A feature extraction approach was thus not followed in this study.

Feature selection involves the selection of a subset of features from the original set according to some criterion of subset performance. The number of possible subsets increases combinatorially with the size of the feature set and it is seldom practical to evaluate all possible subsets (Jain, Duin, and Mao 2000). A variety of search ~~methods-schemes~~ exists for reducing the portion of feature space searched. ~~Of these, only the branch and bound method is globally optimal, the rest achieve reduced computation at the price of optimality. The complexity of the branch and bound method increases exponentially with the size of the feature set and is as such still computationally impractical for large feature sets (Jain, Duin, and Mao 2000).~~ The fastest and most straightforward search ~~method-scheme~~ is simply to rank features based on their individual performance and select the best N . However, feature ranking approaches are problematic for data sets containing feature redundancy. In these situations, correlated features are ranked similarly, resulting in sub-optimal and redundant feature sets.

More advanced search ~~methods-schemes~~ use greedy sequential approaches, such as forward selection and backward elimination. Compared to the feature ranking approach, greedy search ~~methods-procedures~~ are more likely to find the globally optimal feature set as they explore more of the search space and are less inclined to select multiple redundant features (Webb 2002). The forward selection (FS) approach starts with an empty feature set and proceeds in a number of steps where one feature is added to the selected set at each step. The feature whose selection most improves an accuracy criterion is the one that is selected for that step. The selection process proceeds for a set number of steps or until a stopping

criterion is reached (Bishop 2003). The backward elimination (BE) method starts with the full set of features and proceeds in a number of steps where one feature is eliminated from the selected set at each step. The feature whose removal produces the best accuracy according to some criterion is the one eliminated for that step. Again, the BE selection process proceeds for a set number of steps or until a stopping criterion is reached (Bishop 2003). – BE is computationally more costly than FS as it begins evaluation on the full feature set. For the same reason, it also requires adequate data to represent the full feature set.

Feature selection methods can be grouped into filter, wrapper and embedded approaches. In the filter approach, generic measures of separability or importance are used to evaluate feature subsets, while in the wrapper approach, the accuracy of a specific classifier trained on the feature subset is used as the selection criterion (Duin and Tax 2005). An embedded approach is one where feature selection is incorporated into the classifier training procedure, such as with random forests (Breiman 2001). Filter approaches have the advantage over wrapper and embedded approaches of making feature selection independent of the classifier, thus allowing for greater flexibility in the choice of classifier (Brown et al. 2012). In general, filter approaches are also computationally more efficient than wrapper approaches. This is an important consideration for large and high dimensional data sets such as those often encountered in remote sensing. This study thus focuses on filter approaches.

High dimensional feature spaces typically contain feature redundancy (Cukur et al. 2015; Tolosi and Lengauer 2011; Yu and Liu 2004). Although feature correlation and redundancy are related, they are not strictly the same thing (Brown et al. 2012; Guyon and Elisseeff 2003). Features can help improve separability when the within class correlation is stronger than the between class correlation. We use the term ‘redundancy’ to refer to

correlation of features between classes. The raw bands of multispectral imagery often have significant spectral overlap and consequently are correlated with one another. This spectral overlap will exacerbate the redundancy amongst features derived from these raw bands (Cukur et al. 2015). Hyperspectral imagery is also well-known for containing redundancy amongst the bands (Yuan, Zhu, and Wang 2015).

A number of authors have noted difficulties in selecting features from high dimensional data sets. Kononenko et al. (1997), Guyon et al. (2002), Yu & Liu (2004) and Yousef et al. (2007) noted that feature redundancy can have a negative impact on the optimality of feature selection. Feature redundancy not only leads to sub-optimal feature selection, but also makes selected features unstable and sensitive to small changes in the data used for selection (Tolosi and Lengauer 2011; Guyon and Elisseeff 2003; Li, Harner, and Adjero 2011; Kalousis, Prados, and Hilario 2007).

~~Redundancy can be effectively dealt with using a feature extraction approach, such as PCA, but this requires computation of the full feature set. This is not practical in computationally demanding applications such as analysing very high resolution (VHR) imagery over large areas.~~

~~In gene expression studies, such as DNA microarray studies, all features have similar measurement costs and computation times (Inza et al. 2004), but this is not the case for remote sensing problems, where some features might carry significantly larger computation time or measurement cost burdens than others (Blaschke 2010).~~

The increasing availability of high resolution imagery ~~and associated feature extraction and classification techniques are rapidly increasing (Chi et al. 2016) resulting in large quantities of high spatial and spectral resolution data requiring interpretation. As the size of remote sensing data increases, necessitates~~ computationally efficient feature selection techniques robust to high dimensional redundant spaces ~~will become increasingly important.~~ In this ~~paper article~~ we propose a computationally-efficient filter approach feature selection method for addressing the problems of sub-optimality and instability associated with high dimensional, redundant feature spaces of remotely sensed data. We adopt the filter approach due its ~~relative speed and~~ separation of feature selection and classification tasks. The method ~~follows the category 1 approach above and employs hierarchical clustering affinity propagation to identify clusters of redundant features and redundancy is reduced by selecting a single representative feature from the most relevant clusters. This has the advantage of producing deterministic results, not~~The method ~~requires~~ no prior knowledge of the number of clusters ~~and allowing user selection of the final partitioning.~~ Correlation is used to measure feature similarity, which allows a broader encapsulation of feature redundancy than distance measures such as Euclidean distance (X. Chen et al. 2017) ~~, used by X. Chen et al. (2017).~~ Assumptions of linear dependence between features and class labels made in structured sparsity regularisation approaches (Gui et al. 2016), are also avoided by selecting features with a relevance heuristic, based on the use of naïve Bayes or mutual information criteria.

~~Where a group of correlated features are similarly relevant, it is beneficial to select the feature with the lowest computation time or measurement cost. Consideration of these costs can potentially result in a substantial reduction in classification computation time for this form of remote sensing problem.~~

The proposed method can be fully automated, or it can be used interactively to ~~is~~ unique in that it allows for consideration of computation time and measurement cost ~~in~~ selecting features from correlated clusters of similarly relevant features. ~~We compare it~~ The performance of the proposed method is compared to popular feature selection approaches, on a number of remote sensing data sets. The various

~~While many~~ feature selection methods are evaluated ~~editions only consider based on~~ classification accuracy (Yu and Liu 2004; Strobl et al. 2008; X. Chen et al. 2017), ~~we present~~ measures of computation time, classification accuracy and stability of selected features under different data samplings.

2 Methods

2.1 Related Work

A number of feature selection approaches have been developed to address the issues of stability and sub-optimality encountered in high dimensional and redundant data. These methods consider the trade-off between feature relevance (i.e. how much information the feature contains about the class labels) and redundancy.

A means of selecting good features from redundant spaces was devised by Yousef et al. (2007), who. ~~They~~ used a k -means algorithm to produce a fixed number of clusters of correlated features. The accuracy of a SVM classifier is found for all the features of each cluster and. ~~The~~ lowest performing clusters are eliminated, the remaining features combined, and the process is repeated until a desired number of clusters is reached. ~~A~~ related feature selection method that finds and removes redundancy by clustering features into similar groups was presented by Mitra et al. (2002). They used a novel clustering

algorithm to group correlated features based on a similarity measure they call ‘maximal information compression index’, which is the smallest eigenvalue of the feature covariance. Sahu & Mishra (2011) also used *k*-means clustering to group redundant features. The best feature, according to an importance measure, is then selected from each cluster. Cukur et al. (2015) proposed a similar method, where redundant features are clustered and top ranked features selected from each cluster using an importance measure called ‘minimum redundancy maximum relevance’ (mRMR).

A two-step procedure called the ‘Fast Correlation Based Filter’ (FCBF), was developed by Yu & Liu (2004). It first creates a reduced set of relevant features, and then removes redundant features from this set using a search procedure based on Markov blanket filtering. A non-linear correlation measure, called symmetrical uncertainty, is used to measure both feature relevance and redundancy. Relevance is measured by how well features are correlated with class labels and redundancy is measured by how well features are correlated with each other.

Wu et al. (2013) compared a number of filter approach feature selection methods ~~based on their application to~~ for reducing redundancy in three hyperspectral data sets. They used a number of performance measures for comparison, including classifier accuracy, feature stability, and their own criterion called the maximal minimal associated index quotient (MMAIQ). MMAIQ uses Cramer’s *V*-test values to trade feature relevance against redundancy and is applied in a forward selection type routine. While the authors concluded that MMAIQ provided the best overall performance, it did not provide good stability for high dimensional data.

A number of feature importance measures (including the FCBF) were incorporated by Brown et al. (2012) into a common theoretical framework. These measures all consider both relevance and redundancy in some way. A comprehensive empirical study was used to compare the performance (in terms of stability and classifier accuracy) of these measures. The study tested the criteria in a FS ~~search scheme~~^{approach}, under varying conditions, including redundancy in high dimensional feature spaces. They concluded that joint mutual information (JMI) (Yang and Moody 1999) provides the best feature selection performance overall.

In recent years, a number of feature selection approaches based on structured sparsity regularisation have been developed (Wang et al. 2010; X. Chen and Gu 2015; X. Chen et al. 2017; Gui et al. 2016). Structured sparsity regularisation modifies the traditional sparsity regularisation approach by incorporating prior knowledge of the group structure of features to improve performance (Gui et al. 2016). The supervised multiview feature selection (SMFS) method of X. Chen et al. (2017) uses a structured sparsity approach that groups features by similarity and uses this similarity structure to address the trade-off between feature relevance and redundancy. In SMFS, features are clustered into homogenous groups or “views” using affinity propagation (AP) with a squared Euclidean distance similarity measure. A sparse set of features is selected from these views by a joint $\ell_{1,2}$ -norm minimisation of an objective function comprised of loss function and regularisation terms. The formation of the loss function assumes a linear dependence between features and class labels. Feature view structure is incorporated into the $\ell_{1,2}$ -norms so as to encourage the sparsity of selected features within views, while retaining the information of multiple heterogeneous views. The objective function is minimised with quadratic programming, which is computationally expensive compared to greedy search type feature selection methods such as FS and JMI.

Feature weights produced by the optimisation procedure can be considered an importance measure that trades relevance against redundancy.

~~With the exception of FCBF, the above feature selection procedures can be grouped into two categories:~~

~~Approaches that use some form of clustering of similar features to identify and isolate redundancy, followed by a measure of importance to select features with low redundancy and high relevancy.~~

~~Approaches that use a single measure of feature importance that incorporates the trade-off between feature relevance and redundancy, after which a FS or simple ranking procedure is used to select the best features.~~

2.12.2 Formulation Feature Clustering and Ranking

The proposed method is described as follows:

Within the context of the related research overviewed in the previous section, our proposed feature selection method consists of the following three steps:

- (1) Perform ~~average-linkage hierarchical~~affinity propagation clustering (~~Szekely and Rizzo 2005~~)(Frey and Dueck 2007) of the feature set using the absolute value of the correlation coefficient as the ~~dis~~similarity metric.
- (2) ~~Select a dissimilarity threshold at which to extract a natural number of clusters containing high correlation (in this study it is done by visual inspection of the dendrogram, but selection can conceivably also be automated).~~

- ~~(3)~~(2) Rank each cluster's importance by finding the value of a relevance criterion for each individual feature and then finding the median of the feature relevance values in the cluster.
- ~~(4)~~(3) Select a single feature from each of the N clusters with best importance scores.

Affinity propagation is a clustering technique that identifies cluster representatives ('exemplars'), and their corresponding clusters, by an iterative scheme of message passing between data points (Frey and Dueck 2007). A matrix of pair-wise similarities and a 'preference' parameter are required as inputs. The preference parameter affects the number of identified clusters and may be chosen automatically based on the values of the similarities. The proposed feature selection method sets the preference parameter to the median of the similarities, which results in a moderate number of clusters (Frey and Dueck 2007). Unlike clustering techniques such as k -means, affinity propagation does not require prior knowledge of the number of clusters.

Two kinds of messages, 'availability' and 'responsibility', are passed between data points at each iteration. The values of these messages express the current affinity one point has for choosing another as its exemplar. The responsibility $r(i, k)$ reflects the accumulated evidence that feature k is the exemplar for feature i , taking into consideration other possible exemplars for feature i . The availability $a(i, k)$ reflects the accumulated evidence for how appropriate it would be for feature i to choose feature k as its exemplar, taking into consideration support from other features for choosing k as their exemplar.

To initialise, the availabilities are set to zero, $a(i, k) = 0$. In our method, the similarity $s(i, k)$ between feature i and k , is set to the absolute value of the correlation

coefficient, and the self-similarities, $s(k, k)$, are set to the preference value. At each iteration, the responsibilities are updated using the rule:

$$r(i, k) \leftarrow s(i, k) - \max_{k' \text{ s.t. } k' \neq k} \{a(i, k') + s(i, k')\} \quad (1)$$

The availabilities are correspondingly updated using the rule:

$$a(i, k) \leftarrow \begin{cases} \min \left\{ 0, r(k, k) + \sum_{i' \text{ s.t. } i' \notin \{i, k\}} \max\{0, r(i', k)\} \right\}, & i \neq k \\ \sum_{i' \text{ s.t. } i' \neq k} \max\{0, r(i', k)\}, & i = k \end{cases} \quad (2)$$

The exemplar for feature i is identified by the value of k that maximises $a(i, k) + r(i, k)$.

The iterations continue until the clusters (and their corresponding exemplars) remain stable for ten consecutive updates.

~~Hierarchical clustering provides a simple way of grouping features and does not require prior knowledge of the number of clusters (Webb 2002). The method begins with each individual feature in its own cluster and proceeds in a number of steps. At each step, the pair of clusters that minimise a cluster dissimilarity criterion are merged into a single cluster. The hierarchy of clusters can be graphically represented with a dendrogram (MathWorks 2016), as illustrated with an example in Figure 1. The inverted U-shaped lines show which clusters are combined into new clusters at each step. The height of the horizontal line indicates the magnitude of the dissimilarity between clusters.~~

~~The average linkage criterion was used to measure cluster dissimilarity. This is the average of the pairwise distances between the objects in the two clusters. The pairwise object distances were calculated as $|1 - \rho|$, where ρ is the Pearson correlation coefficient between~~

two objects. Cluster stability and strength of correlation within each cluster are the key factors to consider when choosing the number of clusters and can be visually interpreted from the dendrogram. In Figure 1, the dotted line shows an example dissimilarity threshold at which to extract clusters from the hierarchy. At this threshold, the clusters are highly correlated (i.e. the dissimilarity is small) and the cluster contents are stable (i.e. the next level in the hierarchy only occurs at a substantially larger dissimilarity).

[Figure 1. Example dendrogram showing chosen threshold at which to extract clusters]

We investigated the performance of two different feature relevance measures: the accuracy of a naive Bayes classifier and the mutual information (MI) between the feature and the class labels. The naive Bayes classifier, using a histogram to model class densities, was chosen primarily because it makes no assumptions about the form of the class distributions and can thus provide a generic measure of separability. It is simple, fast and recognised as being accurate for a variety of problems (Hand and Yu 2001). The ‘naive’ assumption of feature independence is of no consequence when testing individual features.

MI is a measure of the dependence between two random variables (Brown et al. 2012). Given two random variables X and Y , with probability distributions $p(x)$ and $p(y)$ and joint probability distribution $p(x, y)$, the MI between X and Y is defined as:

$$I(X; Y) = \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p(x, y) \log \frac{p(x, y)}{p(x)p(y)} \quad (34)$$

The MI between a feature and the class labels gives a useful indication of that feature's relevance or importance (Brown et al. 2012). The probability distributions in Equation ~~(3)~~(4) are not known and are estimated using histograms.

The cluster importance measure for the k^{th} cluster is expressed as

$$C_k = \text{median}_{X_j \in G_k} R(X_j, Y) \quad (4)$$

where G_k is the set of features in cluster k and $R(X_j, Y)$ is the MI or naïve Bayes feature relevance measure for feature X_j and class labels Y . Once the clusters of similarly relevant features have been ranked according to their importance measures, single features may be selected from the best clusters using an automatic procedure or by the user, taking measurement cost and computation time into account. The process of selecting a representative feature from each of the top ranked clusters reduces redundancy while retaining relevance. For automatic feature selection, the feature with the maximum relevance measure is selected from each of the N best clusters. When criteria of measurement cost or computation time require consideration, the user should hand select features minimising these values from the N best clusters.

The number of clusters to select, N , can be specified by the user based on the size of the training set or by using a grid search with the final classifier accuracy as performance measure. To avoid biased accuracy estimates, all classifier accuracy evaluations, for cluster ranking or selection of N , are done on unseen test data using a five-fold cross validation (Bishop 2003).

2.22.3 Data Sets

Five remote sensing and one synthetic data set ([Table 1](#)) were used for comparing the proposed method against popular existing feature selection methods. ~~The data sets are detailed in Table 1.~~ The ‘difficulty’ in the last column [of Table 1](#) is calculated as $N/(mc)$, as in Brown et al. (2012), where N is the number of objects, m the number of features and c the number of classes. Smaller values indicate that the data is less representative of the underlying class distributions, which results in more challenging feature selection and classification tasks. The Spekboom set consists of 46 spectral and textural features derived from four band multispectral, 0.5 m spatial resolution aerial imagery. The classes represent three types of vegetation found in the Little Karoo, a semi-arid region in South Africa. It was created as part of a vegetation mapping project being conducted by the authors.

The two class synthetic data set was generated to have redundancy amongst the features. The first five features for class j were generated from a normal distribution, $N(\mu_j, 1)$, with ~~mean~~ mean μ_j and standard deviation of one ($j = 1..2$). The mean, μ_j , of each distribution, was generated from the standard normal distribution, $N(0,1)$. The same number of objects were generated for each class. To introduce redundancy, an additional five features were generated by adding normally distributed noise, $N(0,0.25)$, to the original five features. A further five redundant features were similarly generated, but by adding normally distributed noise, $N(0,0.5)$, to the original features. Finally, two spurious features, sampled from a standard normal distribution, $N(0,1)$, were added to the data set.

The Statlog Landsat and Urban Land Cover data sets were obtained from the UCI Machine Learning Repository (Lichman 2013). The Statlog Landsat features are generated from six band multispectral pixel values in three by three neighbourhoods. The data set

consists of six land cover classes. The features of the Urban Land Cover data set are comprised of multi-scale spectral, size, shape and textural measures, derived from high resolution aerial imagery (Johnson and Xie 2013).

Kennedy Space Centre (KSC) and Botswana are public hyperspectral data sets with vegetation and land cover classes (GIC 2014). The Botswana data were acquired by the Hyperion sensor on board the NASA EO-1 satellite and consist of 145 bands in the 400-2500 nm portion of the spectrum, at a 30 m pixel resolution. The KSC data were acquired by the NASA AVIRIS (Airborne Visible/Infrared Imaging Spectrometer) sensor and consist of 176 bands in the 400-2500 nm range, acquired at a spatial resolution of 18 m.

[Table 1. Data sets]

2.32.4 Experimental design~~valuation~~

The proposed Feature Clustering and Ranking (FCR) method was compared to a number of other established and competing feature selection methods. We adopted a similar, although reduced, evaluation approach to that of Brown et al. (2012) and Wu et al. (2013). The compared methods included the standard selection-search approaches of ranking, FS and BE. These standard approaches and FCR, were each evaluated with two different feature relevance criteria: MI and the naïve Bayes classification accuracy. The MI relevance criterion for FCR and ranking approaches finds the MI between individual features and the class labels. To integrate the MI relevance criterion into FS and BE, it is necessary to compute the MI of a set of multiple candidate features with the class labels. In this situation, the candidate features are first merged into a joint variable and then the MI of the class labels

with this joint variable is computed (Brown et al. 2012). We used histograms with ten bins along each dimension to approximate probability densities for both the MI and naive Bayes criteria (Webb 2002). ~~This approximation is made~~ to avoid difficulties and inefficiencies associated with estimating probability densities for continuous variables (Brown et al. 2012).

In addition to the MI and naive Bayes criteria, two other criteria, namely JMI and maximum relevance minimum redundancy (mRMR), were included in our study to represent ‘state of the art’ performance.

Brown et al. (2012) compared the performance of several feature selection criteria in redundant high dimensional spaces, and found the JMI criterion gave the best overall performance in terms of classification accuracy and stability. ~~Based on these results, FS with the JMI criterion was included in our study to represent the ‘state of the art’ performance.~~

The JMI measure for feature X_k is

$$J_{\text{JMI}}(X_k) = \sum_{X_j \in S} I(X_k X_j; Y) \quad (52)$$

where Y are the class labels and S is the set of previously selected features. JMI considers the MI between the class labels and the joint variables $X_k X_j$, which are the pairwise combinations of the candidate feature with each feature already selected. It measures how well the candidate feature complements selected features in describing the class labels.

The popular mRMR criterion, introduced by (Peng, Long, and Ding (2005), expresses the trade-off between feature relevance and redundancy using mutual information measures. The mRMR measure for candidate feature X_k is

$$J_{mRMR}(X_k) = I(X_k; Y) - \frac{1}{|S|} \sum_{X_j \in S} I(X_k; X_j) \quad (6)$$

The first term expresses relevance as the dependence between the candidate feature and class labels, while the second term approximates the redundancy as the mean pair-wise dependencies between the candidate and previously selected features. The JMI and mRMR criteria are used in a FS search scheme. The evaluated methods ~~and their criteria~~ are detailed in Table 2. In this evaluation, the term ‘method’ is used to refer to a combination of search scheme, such as FS, and criterion, such as JMI.

[Table 2. Methods as combinations of search schemes and criteria~~Method Search scheme and criteria combination~~]

To quantify the stability of the selected features, we used the consistency index developed by Kuncheva (2007). If A and B are subsets of the full feature set X , with $|A| = |B| = k$, $r = |A \cap B|$ and $0 < k < |X| = n$, the consistency index is

$$C(A, B) = \frac{rn - k^2}{k(n - k)} \quad (73)$$

It's value lies in the range $[-1, +1]$, where positive values indicate similar sets, zero indicates a random relation and negative values indicate an anti-correlation between the feature sets (Kuncheva 2007). To evaluate stability for a particular method, we select features from bootstrap samples of the data. The consistency index is found for each pairwise combination of selected features over ten bootstraps of the data and averaged to give a measure of overall stability.

A k -nearest-neighbour (k -NN) classifier (with $k = 3$) was used to evaluate the accuracy of the features selected by each method. k -NN is a generic classifier that makes no assumptions about the data and requires no tuning. While other classifiers may be more accurate in particular situations, k -NN allows a relative comparison of the feature selection methods, independent of the influence of classifier tuning for specific data. For each of the feature sets found from the bootstrap samples, the k -NN accuracy was found as the average per-class accuracy from a ten-fold cross validation. For each method and data set combination, an overall accuracy was computed as the average of the bootstrap accuracies.

The number of features to select for each data set was fixed across methods. This parameter was selected by using the accuracy of a k -NN classifier ($k = 3$), trained on the first N features selected by FS-NaiveBC, as the criterion in a grid search. A low value of N that achieved good accuracy was selected for each data set. The number of features selected for each data set are detailed in Table 3.

[Table 3. Feature selection parameters]

The FCR methods (FCR-MI and FCR-NaiveBC) required some specific treatment to integrate them into the evaluation. ~~The dissimilarity threshold at which to extract the feature clusters from the hierarchy was determined by visual inspection of the dendrogram for each dataset. Thresholds at which the clusters were both stable and strongly correlated were favoured (see section 2.1 for a description of these concepts). The chosen threshold was used across all bootstraps of the data set.~~ After bootstrapping, clusters were assigned unique

indices, ensuring identical clusters had the same index. The consistency index was then found using the selected cluster indices rather than feature indices. This was done to simulate hand-selection of preferred features from the best clusters for each bootstrap, while, ~~in~~ practice, allowing the FCR algorithm to automatically choose the top ranked feature from each cluster (for the sake of simplicity and speed). In other words, the cluster index was used to represent the index of the preferred feature that could otherwise have been selected by hand from the cluster contents.

We followed a similar approach to that of Brown et al. (2012), ~~for~~ for computing a single ‘non-dominated’ ranking of the methods that considers stability and accuracy performance simultaneously. The concept of ‘Pareto optimality’, ~~is~~ was used to find a single optimal solution in terms of multiple criteria. In the context of our evaluation, the ‘Pareto front’ is the set of methods on which no other method can improve without degrading either the accuracy or stability. The methods in this set are called ‘non-dominated’ (Mishra and Harit 2010). Successive Pareto fronts can be formed iteratively by finding the current Pareto front of the set of methods that excludes members of the previous fronts. A method was thus given a non-dominated rank of N if it was a member of the N^{th} Pareto front. The average of the non-dominated ranks for each method over the six data sets was used to produce an overall ranking.

The bulk of the software implementation was done in MatlabTM, making use of the PRTools toolbox (TU Delft 2015). The MI , ~~and~~ JMI and mRMR criteria were computed using the FEAST (FEAture Selection Toolbox) C++ implementation (Brown et al. 2012).

3 Results and Discussion

The chosen FCR correlation dissimilarity thresholds and corresponding number of features selected for each data set are detailed in Table 3.

[Table 3. Feature selection parameters]

The results of the stability and accuracy evaluations for each method and data set combination are shown in Figure 1~~Figure 2~~ and Figure 2~~Figure 3~~ respectively. The methods appear along the x axis in order of their mean stability in Figure 1~~Figure 2~~, and mean accuracy in Figure 2~~Figure 3~~, over the six data sets. ~~FS-MI was the most stable overall, but had one of the poorest accuracies. Similarly, FS-NaiveBC is the most accurate overall, but is the least stable. While neither FCR-NaiveBC nor FCR-MI achieve the best overall accuracy or stability, they are amongst the top three methods for both performance measures.~~ The wide range of stabilities confirms the sensitivity of some methods to variations in the data. The method accuracies span a smaller range than the method stabilities. Nonetheless, there ~~are~~ were substantial differences in accuracy between the best and worst methods. Compared to the other data sets, the stability of the Spekboom, Synthetic and Landsat data ~~is~~ was noticeably superior. As reflected in the ‘difficulty’ values in Table 1, these data sets are more representative of the underlying distributions and are thus less sensitive to disturbances.

FCR-NaiveBC and FCR-MI occupy the top two positions for both performance measures. The ranking methods, Rank-MI and Rank-NaiveBC both ~~have~~ had poor accuracy performance. This ~~is~~ was expected as these methods do not consider feature complementarity and only measure relevance of features in isolation. The relatively poor accuracy and

stability of FS-JMI was surprising in the context of the results of Brown et al. (2012), where it produced the best overall performance. ~~Perhaps FS-JMI is more competitive when applied to higher dimensional data sets, containing hundreds or thousands of features, such as those that were used in Brown et al. (2012).~~ The FS-JMI results nevertheless provide a benchmark that helps confirm the usefulness of the FCR method for the type of data investigated in our study.

As with classifier design, there is a ‘curse of dimensionality’ problem with computing the MI of joint variables. As the number of features increases, the number of objects needed to adequately represent the feature distribution increases exponentially (Brown et al. 2012). For this reason, the MI criterion is not well suited for evaluating the BE ~~method~~search scheme, which requires computation of the relevance criterion for the full feature set. This likely explains the poor performance of BE-MI in terms of both accuracy and stability. ~~Note that p~~Part of the motivation for the JMI and mRMR formulations is to circumvent this kind of representivity issue by using a low dimensional approximation to MI.

[Figure 12. Method stability per data set (methods along the x axis are ordered by their mean stability over the data sets)]

[Figure 23. Method accuracy per data set (methods along the x axis are ordered by their mean accuracy over the data sets)]

The method execution times, summed over the six data sets, are provided ~~for~~ reference in Table 4. The execution time of FCR ~~compete~~ds well with the other methods.

although mRMR was the fastest overall. The ~~NaiveBC-naïve Bayes~~ criterion is slower to compute than the MI criterion as it uses a five-fold cross-validation, ~~implemented in MatlabTM,~~ to evaluate the classification accuracy, ~~while,~~ MI is computed once-off ~~using the efficient FEAST-C++ implementation.~~ Methods using the ~~NaiveBC-naïve Bayes~~ criterion are consequently slower than their MI counterparts. FS-JMI is faster than the related FS-MI method, as the criterion only requires MI computations between pairwise combinations of features and the class labels, while the MI criterion is evaluated on the combination all selected features. BE is known to be less efficient than FS (Guyon and Elisseeff 2003), and was the slowest of the tested ~~methods~~ search schemes.

[Table 4. Method cumulative execution time over all data]

Table 5 presents the non-dominant ranking of the methods, in terms of both accuracy and stability. The best ranked method overall was FCR-MI, followed by ~~with~~-FCR-NaiveBC, ~~FS-MI and FS NaiveBC occupying the second rank position.~~ ~~While the FS-MI and FS NaiveBC produced the best performance for stability and accuracy respectively,~~ ~~FCR-MI and FCR NaiveBC achieved a better compromise between these two measures.~~ The Rank-NaiveBC, Rank-MI, BE-NaiveBC and BE-MI methods ~~are~~ were ranked lowest due to the known limitations of these methods. The FS-MRMR method was competitive on all measures, and was the third ranked method overall.

If the clustering step were omitted, FCR-MI and FCR-NaiveBC would simplify to Rank-MI and Rank-NaiveBC respectively. FCR-MI and FCR-NaiveBC showed a substantial improvement in performance compared to Rank-MI and Rank-NaiveBC, which lends support to the effectiveness of the clustering step. Considering the combination of the MI and naive

Bayes criteria with each ~~method~~ search scheme in isolation, there ~~is~~ was a general trend for MI to produce better stability and naive Bayes to produce better accuracy. While FCR ~~workeds~~ well with either criterion, the results favoured ed the use of MI as it is faster and produce ds a better non-dominant ranking than naive Bayes. On the whole, the evaluations ~~study show~~ demonstrate thats the proposed FCR method ~~to be~~ is effective at selecting accurate and stable features from high dimensional remote sensing data containing redundancy.

[Table 5. Non-dominated ranking of methods by accuracy and stability]

4 Conclusions

Small changes in data sets containing redundancy can result in substantial changes in selected features. Feature redundancy is also known to cause selection of sub-optimal features. This study presented and evaluated a ~~A new~~ method for selecting stable and informative features from redundant data by ranking correlated clusters of features, ~~was presented~~. The method uses ~~Using affinity propagation hierarchical clustering to identify,~~ a ~~natural moderate~~ number of clusters of correlated and similarly relevant features ~~can be selected by observing the stability of correlation relationships in the data using a dendrogram~~. It then ranks the ~~c~~ Clusters ~~are then ranked~~ using an importance measure, calculated as the median of a relevance criterion evaluated on each individual feature in the cluster. By selecting individual features from the best clusters, a set of informative features is found while simultaneously removing redundancy from the data. These features may be selected automatically based on their relevance, or interactively, taking criteria such as computation time and measurement cost into account. The ~~ability option~~ to include factors other than relevance and redundancy in determining selected features ~~hand pick features from the best~~

~~clusters~~ distinguishes FCR from related feature selection methods; although this option is currently limited to a manual procedure. Future work will investigate ways of combining relevance with other criteria to allow feature selection without user input. This ability is beneficial as it allows other factors, such as speed of computation and physical interpretability, to be considered when determining an effective feature set.

~~The need for user specification of dissimilarity threshold can be seen as a weakness of the FCR method. This is a subjective choice and different thresholds can lead to different sets of selected features. A possible way to automate this choice would be to extract clusters from all levels in the hierarchy, select a set of features from the best clusters at each level, and then use a performance measure such as the accuracy of a k-NN classifier to choose the best set of features overall. The need for visual inspection of the dendrogram to make the choice of dissimilarity threshold limits the dimensionality of data that the FCR method can practically be applied to. For data sets of hundreds or thousands of features, the dendrogram may be too cluttered to make a sensible choice and a feature selection algorithm other than FCR may be more appropriate. It is worth noting that for problems where feature stability and user specification of preferred variables are not required, FS-NaiveBC may be a more sensible choice of feature selection method. It achieved the best accuracy results and does not require any user intervention as with FCR.~~

~~Despite these limitations, FCR performed well on a diverse range of high-dimensional redundant data sets.~~ The effectiveness of the proposed FCR method was evaluated by comparing its accuracy, stability and execution time to a set of popular feature selection methods. ~~The feature selection methods were each tested in combination with t~~Two criteria were tested for measuring feature relevance: the MI between the candidate feature(s) and the

class labels, and the accuracy of a naive Bayes classifier trained on the candidate feature(s). Unlike structured sparsity regularisation approaches, these relevance criteria do not assume a linear dependence between features and class labels. ~~FS-NaiveBC provided the best accuracy performance but the worst stability performance. In a similar vein, FS-MI provided the best stability performance but the second worst accuracy performance. The FCR method~~ performed well overall, with both naive Bayes and MI criteria. ~~Although FCR did not quite achieve the best performance in accuracy or stability alone, it and~~ was the highest ranked method when considering the accuracy and stability measures in combination. Another benefit of FCR is its relative speed compared to greedy search FS and BE type methods. Ever increasing quantities of high spatial and spectral resolution remote sensing data are being produced and require interpretation (Chi et al. 2016). In this context, instability and sub-optimality associated with feature selection from high dimensional redundant data will become increasingly ~~critical issues~~important. Computationally efficient techniques, such as FCR, are required to address these challenges.

5 Acknowledgements

This work was supported by funding from the Department of Environmental Affairs (DEA) via the Working for Natural Resources Programme. DEA was otherwise not involved in this research. The financial assistance of the National Research Foundation (NRF) towards this research is hereby acknowledged. Opinions expressed and conclusions arrived at, are those of the authors and are not necessarily to be attributed to the NRF. We would like to thank www.linguafix.net for language editing.

6 References

- Bishop, Christopher M. 2003. *Neural Networks for Pattern Recognition*. New York: Oxford University Press. doi:10.1002/0470854774.
- Breiman, Leo. 2001. "Random Forests." *Machine Learning* 45 (1): 5–32.

doi:10.1023/A:1010933404324.

Brown, Gavin, Adam Pocock, Ming-Jie Zhao, and Mikel Lujan. 2012. "Conditional Likelihood Maximisation: A Unifying Framework for Mutual Information Feature Selection." *Journal of Machine Learning Research* 13: 27–66.

doi:10.1016/j.patcog.2015.11.007.

Burges, Christopher J.C. 1998. "A Tutorial on Support Vector Machines for Pattern Recognition." *Data Mining and Knowledge Discovery* 2 (2): 121–167.

doi:10.1023/A:1009715923555.

Chen, Xi, and Yanfeng Gu. 2015. "Class-Specific Feature Selection With Local Geometric Structure and Discriminative Information Based on Sparse Similar Samples." *IEEE Geoscience and Remote Sensing Letters* 12 (7): 1392–1396.

doi:10.1109/LGRS.2015.2402205.

Chen, Xi, Gongjian Zhou, Yushi Chen, Guofan Shao, and Yanfeng Gu. 2017. "Supervised Multiview Feature Selection Exploring Homogeneity and Heterogeneity with L1,2 - Norm and Automatic View Generation." *IEEE Transactions on Geoscience and Remote Sensing* 55 (4): 2074–2088. doi:10.1109/TGRS.2016.2636329.

Chi, Mingmin, Antonio Plaza, Jon Atli Benediktsson, Zhongyi Sun, Jinsheng Shen, and Yangyong Zhu. 2016. "Big Data for Remote Sensing: Challenges and Opportunities." *Proceedings of the IEEE* 104 (11): 2207–2219. doi:10.1109/JPROC.2016.2598228.

Cukur, Huseyin, Hamidullah Binol, Faruk Sukru Uslu, Yusuf Kalayci, and Abdullah Bal. 2015. "Cross Correlation Based Clustering for Feature Selection in Hyperspectral Imagery." In *2015 9th International Conference on Electrical and Electronics Engineering (ELECO)*, 232–236. Bursa: IEEE. doi:10.1109/ELECO.2015.7394552.

Duin, R P W, and David M. J. Tax. 2005. "Statistical Pattern Recognition." In *Handbook of Pattern Recognition and Computer Vision, 3rd Ed.*, edited by CH Chen and PSP Wang, 1–21. Singapore: World Scientific. doi:10.1142/9789812775320_0001.

Frey, Brendan J, and Delbert Dueck. 2007. "Clustering by Passing Messages between Data Points." *Science* 315 (5814): 972–976. doi:10.1126/science.1136800.

GIC. 2014. "Hyperspectral Remote Sensing Scenes."

http://www.ehu.eus/ccwintco/index.php?title=Hyperspectral_Remote_Sensing_Scenes.

Gui, Jie, Zhenan Sun, Shuiwang Ji, Dacheng Tao, and Tieniu Tan. 2016. "Feature Selection Based on Structured Sparsity: A Comprehensive Study." *IEEE Transactions on Neural Networks and Learning Systems* 28 (7): 1–18. doi:10.1109/TNNLS.2016.2551724.

Guyon, Isabelle, and Andre Elisseeff. 2003. "An Introduction to Variable and Feature

- Selection.” *Journal of Machine Learning Research* 3: 1157–1182.
doi:10.1016/j.aca.2011.07.027.
- Guyon, Isabelle, Jason Weston, Stephen Barnhill, and Vladimir Vapnik. 2002. “Gene Selection for Cancer Classification Using Support Vector Machines.” *Machine Learning* 46 (1–3): 389–422. doi:10.1023/A:1012487302797.
- Hand, David J, and Kerning Yu. 2001. “Idiot’s Bayes - Not So Stupid After All?” *International Statistical Review* 69 (3): 385–398.
- Jain, Anil K, Robert P W Duin, and Jianchang Mao. 2000. “Statistical Pattern Recognition: A Review.” *IEEE Transactions on Pattern Analysis and Machine Intelligence* 22 (1): 4–37.
- Johnson, Brian, and Zhixiao Xie. 2013. “Classifying a High Resolution Image of an Urban Area Using Super-Object Information.” *ISPRS Journal of Photogrammetry and Remote Sensing* 83 (September): 40–49. doi:10.1016/j.isprsjprs.2013.05.008.
- Kalousis, Alexandros, Julien Prados, and Melanie Hilario. 2007. “Stability of Feature Selection Algorithms: A Study on High-Dimensional Spaces.” *Knowledge and Information Systems* 12 (1): 95–116. doi:10.1007/s10115-006-0040-8.
- Kononenko, Igor, E Šimec, and M Robnik-Šikonja. 1997. “Overcoming the Myopia of Inductive Learning Algorithms with RELIEFF.” *Applied Intelligence* 7 (1): 39–55. doi:10.1023/A:1008280620621.
- Kuncheva, L I. 2007. “A Stability Index for Feature Selection.” In *International Multi-Conference: Artificial Intelligence and Applications*, 390–395. Innsbruck, Austria: IASTED.
- Li, Shengqiao, E James Harner, and Donald a Adjero. 2011. “Random KNN Feature Selection - a Fast and Stable Alternative to Random Forests.” *BMC Bioinformatics* 12 (1). BioMed Central Ltd: 450. doi:10.1186/1471-2105-12-450.
- Lichman, M. 2013. “UCI Machine Learning Repository.” <http://archive.ics.uci.edu/ml>.
- Mishra, K K, and Sandeep Harit. 2010. “A Fast Algorithm for Finding the Non Dominated Set in Multi Objective Optimization.” *Multi-Objective Optimization Using Evolutionary Algorithms* 1 (25): 35–39.
- Mitra, Pabitra, C A Murthy, and Sankar K Pal. 2002. “Unsupervised Feature Selection Using Feature Similarity.” *IEEE Transactions on Pattern Analysis and Machine Intelligence PAMI* 24 (3): 301–312. doi:10.1109/34.990133.
- Peng, Hanchuan, Fuhui Long, and Chris Ding. 2005. “Feature Selection Based on Mutual Information: Criteria of Max-Dependency, Max-Relevance, and Min-Redundancy.”

- IEEE Trans. on Pattern Analysis and Machine Intelligence* 27 (8): 1226–1238.
doi:10.1109/TPAMI.2005.159.
- Sahu, Barnali, and Debahuti Mishra. 2011. “A Novel Approach for Selecting Informative Genes from Gene Expression Data Using Signal-to-Noise Ratio and t-Statistics.” In *2011 2nd International Conference on Computer and Communication Technology (ICCCCT-2011)*, 5–10. Allahabad, India: IEEE. doi:10.1109/ICCCCT.2011.6075207.
- Strobl, Carolin, Anne-Laure Boulesteix, Thomas Kneib, Thomas Augustin, and Achim Zeileis. 2008. “Conditional Variable Importance for Random Forests.” *BMC Bioinformatics* 9 (January): 307. doi:10.1186/1471-2105-9-307.
- Tolosi, Laura, and Thomas Lengauer. 2011. “Classification with Correlated Features: Unreliability of Feature Ranking and Solutions.” *Bioinformatics* 27 (14): 1986–1994. doi:10.1093/bioinformatics/btr300.
- TU Delft. 2015. “PRTools.” <http://prtools.org/prtools/>.
- Wang, Hua, Feiping Nie, Heng Huang, Shannon Risacher, Andrew J Saykin, Li Shen, and ADNI. 2010. “Efficient and Robust Feature Selection via Joint ℓ_2 , 1-Norms Minimization.” *Advances in Neural Information Processing Systems* 23 (March): 1813–1821. doi:10.1016/j.neuroimage.2010.10.081.
- Webb, Andrew R. 2002. *Statistical Pattern Recognition*. Chichester, UK: John Wiley & Sons, Ltd. doi:10.1002/0470854774.
- Wu, Bo, Chongcheng Chen, Tahar Mohand Kechadi, and Liya Sun. 2013. “A Comparative Evaluation of Filter-Based Feature Selection Methods for Hyper-Spectral Band Selection.” *International Journal of Remote Sensing* 34 (22): 7974–7990. doi:10.1080/01431161.2013.827815.
- Yang, Howard Hua, and John Moody. 1999. “Data Visualization and Feature Selection: New Algorithms for Nongaussian Data.” *Advances in Neural Information Processing Systems* 12 (Mi): 687–693.
- Yousef, Malik, Segun Jung, Louise C Showe, and Michael K Showe. 2007. “Recursive Cluster Elimination (RCE) for Classification and Feature Selection from Gene Expression Data.” *BMC Bioinformatics* 8 (144). doi:10.1186/1471-2105-8-144.
- Yu, Lei, and Huan Liu. 2004. “Efficient Feature Selection via Analysis of Relevance and Redundancy.” *Journal of Machine Learning Research* 5 (2004): 1205–1224.
- Yuan, Yuan, Guokang Zhu, and Qi Wang. 2015. “Hyperspectral Band Selection by Multitask Sparsity Pursuit.” *IEEE Transactions on Geoscience and Remote Sensing* 53 (2): 631–644. doi:10.1109/TGRS.2014.2326655.

