

Aberystwyth University

Measures for unsupervised fuzzy-rough feature selection

MacParthaláin, Neil Seosamh; Jensen, Richard

Published in:

International Journal of Hybrid Intelligent Systems

DOI:

[10.3233/HIS-2010-0118](https://doi.org/10.3233/HIS-2010-0118)

Publication date:

2010

Citation for published version (APA):

MacParthaláin, N. S., & Jensen, R. (2010). Measures for unsupervised fuzzy-rough feature selection. *International Journal of Hybrid Intelligent Systems*, 7(4), 249-259. <https://doi.org/10.3233/HIS-2010-0118>

General rights

Copyright and moral rights for the publications made accessible in the Aberystwyth Research Portal (the Institutional Repository) are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the Aberystwyth Research Portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the Aberystwyth Research Portal

Take down policy

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

tel: +44 1970 62 2400
email: is@aber.ac.uk

Measures for Unsupervised Fuzzy-Rough Feature Selection

Neil Mac Parthaláin and Richard Jensen

Department of Computer Science,

Aberystwyth University, Wales, UK

{ncm, rkj}@aber.ac.uk

Abstract

For supervised learning, feature selection algorithms attempt to maximise a given function of predictive accuracy. This function usually considers the ability of feature vectors to reflect decision class labels. It is therefore intuitive to retain only those features that are related to or lead to these decision classes. However, in unsupervised learning, decision class labels are not provided, which poses questions such as; which features should be retained? and, why not use all of the information? The problem is that not all features are important. Some of the features may be redundant, and others may be irrelevant and noisy. In this paper, some new fuzzy-rough set-based approaches to unsupervised feature selection are proposed. These approaches require no thresholding or domain information, can operate on real-valued data, and result in a significant reduction in dimen-

sionality whilst retaining the semantics of the data.

1. Introduction

Large dimensionality presents a problem for handling data due to the fact that the complexity of many commonly used operations are highly dependent (e.g. exponentially) on the level of dimensionality. The problems associated with such large dimensionality however mean that any attempt to use machine learning or data-mining tools to extract knowledge, results in very poor performance. Feature selection (FS) [7] is a process which attempts to select features which are information-rich whilst retaining the original meaning of the features following reduction. Indeed, FS has been applied to problems which have very large dimensionality ($>10\,000$) [2]. Most learning algorithms are unable to consider problems of such size, whilst those that

can will usually perform poorly [7].

Rough set theory (RST) [18] is an approach that can be used for dimensionality reduction, whilst simultaneously preserving the semantics or meaning of the features. Also, as RST operates only on the data and does not require any thresholding information, it is completely data-driven. RST however has one main disadvantage: its inability to deal with real-valued data. In order to tackle this problem, methods of discretising the data were employed prior to the application of RST. The use of such methods can result in information loss however, and a number of extensions to RST have emerged [10], [22], [27] which have attempted to address this inability to operate on real-valued domains. One such approach is fuzzy-rough sets (FRS) which have the ability to operate effectively on real-valued (and crisp) data, thus minimising any information loss [14]. This is achieved by extending traditional RST, to the fuzzy-rough case.

Conventional supervised FS methods evaluate various feature subsets using an evaluation function or metric to select only those features which are related to, or lead to, the decision classes of the data under consideration. However, for many data mining applications, decision class labels are often unknown or incomplete, thus indicating the significance of unsupervised feature selection. In a broad sense, two different types of approach to unsupervised FS have been adopted: Those which maximise clustering performance using an index function [8], [17], and those which

consider features for selection on the basis of dependency or relevance. The central idea behind the latter, is that any single feature which carries little or no further information than that subsumed by the remaining features is redundant and can therefore be eliminated [6], [11], [16]. The approach described in this paper is related to these techniques since it involves the removal of features which are considered to be redundant.

The work presented here is based on fuzzy-rough sets and allows the consideration of real-valued data. It employs the fuzzy-rough discernibility measure to examine the level of discernibility between a single feature and subsets of other features. Where a single feature can be discerned completely by a subset of features, that single feature is considered to be redundant and can be removed from the feature set. FS is conducted through the removal of features until no further inter-dependency can be found. The resulting subset of original features can then be used to define the original complete feature set.

The remainder of the paper is structured as follows. Section 2 introduces the theoretical background to RST and FRS and their application to FS. Section 3 presents the new unsupervised fuzzy-rough feature selection metrics. The proposed approach is compared with an advanced supervised FS technique [15] which is also based on FRS, and results are presented in Section 4. The paper is then concluded in Section 5.

2. Supervised rough approaches

There has been great interest in developing methodologies which are capable of dealing with imprecision and uncertainty. The success of rough set theory in this respect, is due in part to the fact that it operates only on the data and does not require any external information. As RST handles only one type of imperfection found in data, it is complementary to other concepts, such as fuzzy set theory. These two fields may be considered analogous in the sense that both can tolerate inconsistency and uncertainty - the difference being the type of uncertainty and their approach to it; fuzzy sets are concerned with vagueness, rough sets are concerned with indiscernibility.

2.1. Rough set feature selection

Let $I = (\mathbb{U}, \mathbb{A})$ be an information system, where \mathbb{U} is a non-empty set of finite objects (the universe of discourse) and \mathbb{A} is a non-empty finite set of attributes such that $a : \mathbb{U} \rightarrow V_a$ for every $a \in \mathbb{A}$. V_a is the set of values that attribute a may take. With any $P \subseteq \mathbb{A}$ there is an associated equivalence relation $IND(P)$:

$$IND(P) = \{(x, y) \in \mathbb{U}^2 \mid \forall a \in P, a(x) = a(y)\} \quad (1)$$

The partition of \mathbb{U} , generated by $IND(P)$ is denoted $\mathbb{U}/IND(P)$ (or \mathbb{U}/P for simplicity) and can be calculated as

follows:

$$\mathbb{U}/IND(P) = \otimes \{\mathbb{U}/IND(\{a\}) \mid a \in P\}, \quad (2)$$

where \otimes is specifically defined as follows for sets A and B :

$$A \otimes B = \{X \cap Y \mid X \in A, Y \in B, X \cap Y \neq \emptyset\} \quad (3)$$

If $(x, y) \in IND(P)$, then x and y are indiscernible by attributes from P . The equivalence classes of the P -indiscernibility relation are denoted $[x]_P$.

Let $X \subseteq \mathbb{U}$. X can be approximated using only the information contained within P by constructing the P -lower and P -upper approximations of X :

$$\underline{P}X = \{x \in \mathbb{U} \mid [x]_P \subseteq X\} \quad (4)$$

$$\overline{P}X = \{x \in \mathbb{U} \mid [x]_P \cap X \neq \emptyset\} \quad (5)$$

The tuple $\langle \underline{P}X, \overline{P}X \rangle$ is called a rough set.

Let P and Q be sets of attributes inducing equivalence relations over \mathbb{U} , then the positive region can be defined as:

$$POS_P(Q) = \bigcup_{x \in \mathbb{U}/Q} \underline{P}X \quad (6)$$

The positive region contains all objects of \mathbb{U} that can be classified to classes of \mathbb{U}/Q using the information in attributes P . Based on this definition, dependencies between attributes can be determined. For $P, Q \subseteq \mathbb{A}$, it is said that Q

depends on P in a degree k ($0 \leq k \leq 1$), denoted $P \Rightarrow_k Q$, if

$$k = \gamma_P(Q) = \frac{|POS_P(Q)|}{|\mathbb{U}|} \quad (7)$$

A *reduct* R_{min} is defined as a minimal subset R of the initial attribute set \mathbb{C} such that for a given set of attributes D , $\gamma_R(\mathbb{D}) = \gamma_{\mathbb{C}}(\mathbb{D})$. From the literature, R is a minimal subset if $\gamma_{R-\{a\}}(\mathbb{D}) \neq \gamma_R(\mathbb{D})$ for all $a \in R$. The search for reducts is achieved by comparing equivalence relations generated by sets of attributes [14]. Attributes are removed so that the reduced set provides the same predictive capability of the decision attribute as the original. Other rough set FS approaches involve the use of discernibility matrices [21], [23], etc. to search for reducts (in contrast to the comparison of equivalence classes mentioned previously). These approaches can be computationally complex however, although there are some ways in which this can be alleviated to a certain degree [12], [23].

2.2. Fuzzy-rough feature selection

Fuzzy-rough sets (FRS) [10] encapsulate the related but distinct concepts of vagueness (for fuzzy sets) and indiscernibility (for rough sets), both of which occur as a result of uncertainty in knowledge.

Definitions for the fuzzy lower and upper approximations can be found in [20], where a \mathcal{T} -transitive fuzzy simi-

ilarity relation is used to approximate a fuzzy concept X :

$$\mu_{R_P X}(x) = \inf_{y \in \mathbb{U}} \mathcal{I}(\mu_{R_P}(x, y), \mu_X(y)) \quad (8)$$

$$\mu_{\overline{R_P X}}(x) = \sup_{y \in \mathbb{U}} \mathcal{T}(\mu_{R_P}(x, y), \mu_X(y)) \quad (9)$$

Here, \mathcal{I} is a fuzzy implicator and \mathcal{T} a t-norm. A fuzzy implicator is any $[0, 1]^2 \rightarrow [0, 1]$ -mapping \mathcal{I} satisfying $\mathcal{I}(0, 0) = 1, \mathcal{I}(1, x) = x$ for all x in $[0, 1]$. R_P is the fuzzy similarity relation induced by the subset of features P :

$$\mu_{R_P}(x, y) = \mathcal{T}_{a \in P} \{ \mu_{R_a}(x, y) \} \quad (10)$$

$\mu_{R_a}(x, y)$ is the degree to which objects x and y are similar for feature a , and may be defined in many ways, for example:

$$\mu_{R_a}(x, y) = 1 - \frac{|a(x) - a(y)|}{|a_{max} - a_{min}|} \quad (11)$$

$$\mu_{R_a}(x, y) = \max\left(\min\left(\frac{a(y) - (a(x) - \sigma_a)}{\sigma_a}, \frac{((a(x) + \sigma_a) - a(y))}{\sigma_a}, 0\right), 0\right) \quad (12)$$

where σ_a^2 is the variance of feature a . As these relations do not necessarily display \mathcal{T} -transitivity, the fuzzy transitive closure can be computed for each attribute. The choice of relation is largely determined by the intended application. For feature selection, a relation such as (12) may be appropriate as this permits only small differences between attribute values of differing objects. For classification tasks,

a more gradual and inclusive relation such as (11) should be used.

In a similar way to the original crisp rough set approach, the fuzzy positive region [15] can be defined as :

$$\mu_{POS_P(\mathbb{D})}(x) = \sup_{X \in \mathbb{U}/\mathbb{D}} \mu_{R_P X}(x) \quad (13)$$

An important issue in data analysis is discovering dependencies between attributes. The fuzzy-rough degree of dependency of \mathbb{D} on the attribute subset P can be defined in the following way:

$$\gamma'_P(\mathbb{D}) = \frac{\sum_{x \in \mathbb{U}} \mu_{POS_P(\mathbb{D})}(x)}{|\mathbb{U}|} \quad (14)$$

A fuzzy-rough reduct R can be defined as a minimal subset of features that preserves the dependency degree of the entire dataset, i.e. $\gamma'_R(\mathbb{D}) = \gamma'_C(\mathbb{D})$. Based on this, a fuzzy-rough greedy hill-climbing algorithm can be constructed that uses equation (14) to gauge subset quality. In [15], it has been shown that the dependency function is monotonic and that fuzzy discernibility matrices may also be used to discover reducts.

3. Unsupervised fuzzy-rough feature selection

In the previous section, it was demonstrated how RST and FRS can be applied to the problem of supervised feature selection. One of the most important aspects relating to fea-

ture set reduction is the fuzzy-rough dependency measure, and it is this measure that is also employed for the new unsupervised fuzzy-rough FS (UFRFS) method described in this section. A short worked example is also provided here to illustrate the approach. This section introduces the new unsupervised subset evaluation measures based on fuzzy-rough set theory, and the corresponding reduction algorithm. It is worth noting that a method which is termed ‘unsupervised rough set feature selection’ is described in [19]. However, on closer examination it is revealed that no use is made of the upper and lower approximation rough set concepts (which as demonstrated previously, are central to RST), but merely a discernibility relation. Furthermore, this approach only has the ability to consider crisp or discrete valued data which is one of the main limitations of crisp rough set-based approaches. For the processing of real-valued data, a discretisation step must first be carried out which may result in information loss. This motivates the use of fuzzy-rough sets for feature selection as described in the following section.

The approaches described here use three different measures to perform unsupervised FS: fuzzy-rough dependency, fuzzy-rough boundary region and fuzzy-rough discernibility. Supervised adaptations of these measures can also be employed for supervised fuzzy-rough FS.

3.1 Dependency measure

The discovery of dependencies between attributes, is in general, an important issue in data analysis. Intuitively, a set of attributes Q depends totally on a set of attributes P , denoted $P \rightarrow Q$, if all attribute values from Q are uniquely determined by values of attributes from P .

The central idea behind the present work is that, as with supervised fuzzy-rough FS [15], the fuzzy dependency measure can also be used to discover the inter-dependency of features. This can be achieved by substituting the decision feature(s) \mathbb{D} of the supervised approach for any given feature or group of features Q such that

$$\gamma'_P(Q) = \frac{\sum_{x \in \mathbb{U}} \mu_{POS_{R_P}(Q)}(x)}{|\mathbb{U}|} \quad (15)$$

where $P \cap Q = \emptyset$ and,

$$\mu_{POS_{R_P}(Q)}(x) = \sup_{z \in \mathbb{U}} \mu_{\underline{R}_P R_{Qz}}(x) \quad (16)$$

Here, R_{Qz} indicates the fuzzy tolerance class (or fuzzy equivalence class) for object z . The lower approximation becomes:

$$\mu_{\underline{R}_P R_{Qz}}(x) = \inf_{y \in \mathbb{U}} \text{mathcal{I}}(\mu_{R_P}(x, y), \mu_{R_Q}(y, z)) \quad (17)$$

3.2. Boundary region measure

Most approaches to crisp rough set FS and all approaches to fuzzy-rough FS use only the lower approximation for the evaluation of feature subsets. The lower approximation contains information regarding the extent of certainty of object membership to a given concept. However, the upper approximation contains information regarding the degree of uncertainty of objects and hence this information can be used to discriminate between subsets. For example, two subsets may result in the same lower approximation but one subset may produce a smaller upper approximation. This subset will be more useful as there is less uncertainty concerning objects within the boundary region (the difference between upper and lower approximations). The fuzzy-rough boundary region for a fuzzy tolerance class R_{Qz} X may thus be defined:

$$\mu_{BND_P(R_{Qz})}(x) = \mu_{\overline{R}_P R_{Qz}}(x) - \mu_{\underline{R}_P R_{Qz}}(x) \quad (18)$$

with the upper approximation defined as:

$$\mu_{\overline{R}_P R_{Qz}}(x) = \sup_{y \in \mathbb{U}} \mathcal{T}(\mu_{R_P}(x, y), \mu_{R_Q}(y, z)) \quad (19)$$

As the search for an optimal subset progresses, the object memberships to the boundary region diminish until a minimum is achieved. From this, the total certainty degree

given a feature subset P is defined as:

$$\lambda_P(Q) = 1 - \frac{\sum_{z \in \mathbb{U}} \sum_{x \in \mathbb{U}} \mu_{BND_{R_P}(R_Q z)}(x)}{|\mathbb{U}|^2} \quad (20)$$

It is this measure, λ , that can be used to guide an unsupervised subset selection process.

3.3. Discernibility measure

There are two main branches of research in crisp rough set-based FS: those based on the dependency degree and those based on discernibility matrices and functions. Therefore, it is natural to extend concepts in the latter branch to the fuzzy-rough domain [5].

The fuzzy tolerance relations that represent objects' approximate equality can be used to extend the classical discernibility function. For each combination of features P , a value is obtained indicating how well these attributes maintain the discernibility, relative to another subset of features Q , between all objects.

$$f(P, Q) = \mathcal{T}(\underbrace{c_{ij}(P, Q)}_{1 \leq i < j \leq |\mathbb{U}|}) \quad (21)$$

with

$$c_{ij}(P, Q) = \mathcal{I}(\mathcal{T}(\underbrace{\mu_{R_a}(x_i, x_j)}_{a \in P}), \mu_{R_Q}(x_i, x_j)) \quad (22)$$

Alternatively, rather than taking a minimum operation in

Eq. (21), one can also consider the average over all object pairs, i.e.,

$$g(P, Q) = \frac{2 \cdot \sum_{1 \leq i < j \leq |\mathbb{U}|} c_{ij}(P, Q)}{|\mathbb{U}|(|\mathbb{U}| - 1)} \quad (23)$$

This measure is less rigid than equation (21), which produces the value 0 as soon as one of the c_{ij} equals 0.

3.4. Finding reductions

For the supervised approach, search is conducted within $\mathcal{P}(\mathbb{C})$, the set of all possible subsets of the conditional feature set. However, for the unsupervised approach search is performed within $\mathcal{P}(\mathbb{C}) \times \mathcal{P}(\mathbb{C})$, as to search for reductions any subset can be compared with any other subset. This is a vastly more complex space in which to search. For the purposes of this paper, a linear backward search is employed that achieves reasonable reductions in a short space of time.

The algorithm (figure 1) starts by considering all of the features contained in the dataset. The removal of each feature is then examined iteratively, and the corresponding measure is calculated. If the measure is unaffected then the feature can be removed. This process continues until all features have been examined. If no interdependency exists, the algorithm will return the full set of features. The complexity for the search in the worst case is $O(n)$, where n is the number of original features.

Reduction is achieved for the three measures by replac-

UFRQUICKREDUCT(\mathbb{C})
 \mathbb{C} , the set of all features.

- (1) $R \leftarrow \mathbb{C}$
- (2) **foreach** $x \in \mathbb{C}$
- (3) $R \leftarrow R - \{x\}$
- (4) **if** $M(R, \{x\}) < 1$
- (5) $R \leftarrow R \cup \{x\}$
- (6) **return** R

Figure 1. The UFRQUICKREDUCT Algorithm

ing $M(T, \{x\})$ with either $\gamma'_T(\{x\})$, $\lambda_T(\{x\})$ or $g(T, \{x\})$.

If a greater reduction in features is required (at the expense of accuracy), line (4) in the algorithm can be replaced by:

$$\mathbf{if} \ M(R, \{x\}) < \alpha$$

with $\alpha \in (0, 1]$.

3.5 Worked Example

To illustrate the ideas which have been described in the preceding sections, a small dataset shown in Table 1 is employed. As recommended in [4], the Łukasiewicz t-norm ($\max(x + y - 1, 0)$) and the Łukasiewicz fuzzy implicator ($\min(1 - x + y, 1)$) are adopted to implement the fuzzy connectives. Other interpretations may also be used.

Object	a	b	c	d
1	0.0	0.1	0.1	0.5
2	0.2	0.1	0.6	0.9
3	0.6	0.7	0.3	0.9
4	0.3	0.4	0.8	0.6
5	0.2	0.7	0.9	0.2

Using the fuzzy similarity measure defined in (12), the

resulting relations for each feature in the dataset are shown (for brevity) in Table 2.

Initially, the lower approximations of the concepts of a given feature must be computed for each of the other features in the dataset. This is then used to calculate the dependency degree. For the example dataset, consider the dependency of the feature b on the feature a :

$$\mu_{R_{\{a\}}R_{ub}}(x) = \inf_{y \in \mathbb{U}} \mathcal{I}(\mu_{R_{\{a\}}}(x, y), \mu_{R_{ub}}(y)) \quad (24)$$

Thus, for a particular instance where object $x = 2$, and $u = 2$, this is (as highlighted in Table 2):

$$\begin{aligned} \mu_{R_{\{a\}}R_{2b}}(2) &= \inf_{y \in \mathbb{U}} \mathcal{I}(\mu_{R_a}(2, y), \mu_{R_{2b}}(y)) = \\ &= \inf\{\mathcal{I}(0.83, 1), \mathcal{I}(1, 1)\mathcal{I}(0.33, 0)\mathcal{I}(0.83, 0.5)\mathcal{I}(1, 0)\} = 0 \end{aligned}$$

and for the remaining objects regarding a (i.e. $u \in \{1, 3, 4, 5\}$) this is:

$$\begin{aligned} \mu_{R_{\{a\}}R_{1b}}(2) &= \\ &= \inf\{\mathcal{I}(1, 1), \mathcal{I}(0.83, 1)\mathcal{I}(0, 0)\mathcal{I}(0.5, 0.5)\mathcal{I}(0.83, 0)\} = 0.0 \end{aligned}$$

$$\begin{aligned} \mu_{R_{\{a\}}R_{3b}}(2) &= \\ &= \inf\{\mathcal{I}(0, 1), \mathcal{I}(0.33, 1)\mathcal{I}(1, 0)\mathcal{I}(0.5, 0.5)\mathcal{I}(0.33, 0)\} = 0.0 \end{aligned}$$

$$\begin{aligned} \mu_{R_{\{a\}}R_{4b}}(2) &= \\ &= \inf\{\mathcal{I}(0.5, 1)\mathcal{I}(0.83, 1)\mathcal{I}(0.5, 0)\mathcal{I}(1, 0.5)\mathcal{I}(0.83, 0)\} = 0.17 \end{aligned}$$

$R_a(x, y)$					$R_b(x, y)$					$R_c(x, y)$					$R_d(x, y)$				
1.0	0.83	0.0	0.50	0.83	1.0	1.0	0.0	0.50	0.0	1.0	0.375	0.75	0.125	0.0	1.0	0.429	0.429	0.857	0.572
0.83	1.0	0.33	0.83	1.0	1.0	1.0	0.0	0.50	0.0	0.375	1.0	0.625	0.75	0.625	0.429	1.0	1.0	0.572	0.0
0.0	0.33	1.0	0.50	0.33	0.0	0.0	1.0	0.50	1.0	0.75	0.625	1.0	0.375	0.25	0.429	1.0	1.0	0.572	0.0
0.50	0.83	0.50	1.0	0.83	0.50	0.50	0.50	1.0	0.50	0.125	0.75	0.375	1.0	0.375	0.857	0.572	0.572	1.0	0.429
0.83	1.0	0.33	0.83	1.0	0.0	0.0	1.0	0.50	1.0	0.0	0.625	0.25	0.375	1.0	0.572	0.0	0.0	0.429	1.0

Table 2. Fuzzy similarity relations

$$\mu_{R_{\{a\}}R_{5b}}(2) =$$

$$\inf\{\mathcal{I}(0.83, 1), \mathcal{I}(1, 1)\mathcal{I}(0.33, 0)\mathcal{I}(0.83, 0.5)\mathcal{I}(1, 0)\} = 0.0$$

This process is repeated for every object regarding b in order to calculate the remaining lower approximations for each object. These can then be used to calculate the positive regions:

$$\mu_{POS_{R_{\{a\}}(\{b\})}}(1) = 0.5$$

$$\mu_{POS_{R_{\{a\}}(\{b\})}}(2) = 0.5$$

$$\mu_{POS_{R_{\{a\}}(\{b\})}}(3) = 0.67$$

$$\mu_{POS_{R_{\{a\}}(\{b\})}}(4) = 0.67$$

$$\mu_{POS_{R_{\{a\}}(\{b\})}}(5) = 0.67$$

Therefore the resulting dependency degree is:

$$\gamma'_{\{a\}}(\{b\}) = \frac{\sum_{x \in \mathbb{U}} \mu_{POS_{R_P}}(x)}{|\mathbb{U}|} = \frac{3.01}{6} = 0.602$$

In the interests of brevity only the computation of the dependency of feature b upon feature a is illustrated here. However, in the actual implementation of the UFRFS algorithm, the first step is to consider the dependency of $\{a\}$ on

the subset $\{b, c, d\}$. For the example dataset this leads to the following result:

$$\gamma'_{\{b,c,d\}}(\{a\}) = 1.0 \quad (T = \{b, c, d\})$$

$$\gamma'_{\{c,d\}}(\{b\}) = 0.9569 \quad (T = \{c, d\})$$

$$\gamma'_{\{b,d\}}(\{c\}) = 1.0 \quad (T = \{d\})$$

$$\gamma'_{\{b\}}(\{d\}) = 0.2 \quad (T = \emptyset)$$

Note that each time $\gamma' = 1$, the feature in question is eliminated resulting in the final subset $\{b, d\}$, after all features have been examined.

4. Experimentation

In this section results for the new unsupervised FS method are presented. The approaches are compared with some advanced supervised methods. It must be remembered that that the approaches proposed in this paper are unsupervised and it is assumed that there is no knowledge about the label or class to which each data object belongs. Furthermore it would be very difficult, if not impossible for any unsupervised approach to consistently equal the performance of a supervised method or discover the absent classification functions which are represented by the class labels. The comparison is included here to show that despite missing or

incomplete labels, UFRFS can still reduce dimensionality and discover useful subsets of features.

The experimental setup is shown in Fig. 2. It consists of three main steps: feature selection, dataset reduction (using selected subsets), and classifier learning. Note the class label removal step when employing UFRFS, so that FS is only performed on the unlabelled data.

Following feature selection, the datasets are reduced according to the discovered reducts. These reduced datasets are then classified using the relevant classifier learning method (as described below) and evaluated with 10-fold cross validation.

Four learning mechanisms have been employed to create classifiers for the purpose of evaluating the resulting subsets from the feature selection phase: JRip [3], J48 [26], PART [25] and FRNN [13]. JRip learns propositional rules by repeatedly growing rules and pruning them. During the growth phase, features are added greedily to fit training samples. Once the ruleset is generated, a further optimisation is performed where rules are evaluated and deleted, based on their performance on randomised data. J48 creates decision trees by choosing the most informative features via an entropy measure, and recursively partitions the data into subtables based on their values. Each node in the tree represents a feature with branches from a node representing the alternative values this feature can take according to the current subtable. Partitioning stops when all data items in the

subtable have the same classification. PART generates rules by means of repeatedly creating partial decision trees from data. The algorithm adopts a divide-and-conquer strategy such that it removes instances covered by the current ruleset during processing. Essentially, a rule is created by building a pruned tree for the current set of instances; the leaf with the highest coverage is promoted to a rule. FRNN uses the membership values of test objects to the fuzzy upper and lower approximation concepts of each of the training data decision classes in order to predict test object classes. Test objects are assigned to the training class to which has the highest lower approximation membership value.

The supervised feature selection methods employed are: correlation-based (CFS) [11], consistency-based [26], fuzzy-rough lower approximation-based (FRFS)[15], boundary region-based (B-FRFS)[15], discernibility-based (D-FRFS)[15]. The search method employed for all of these approaches is a forward greedy step-wise approach. The unsupervised methods are all based on UFRFS but employ the new measures as described previously: fuzzy-rough lower approximation-based (UFRFS), unsupervised boundary region-based (B-UFRFS) and unsupervised discernibility-based (D-UFRFS). The classification accuracies for the unreduced data are also included for comparison.

All of the data used in this experimental investigation is labelled. However, before applying the unsupervised meth-

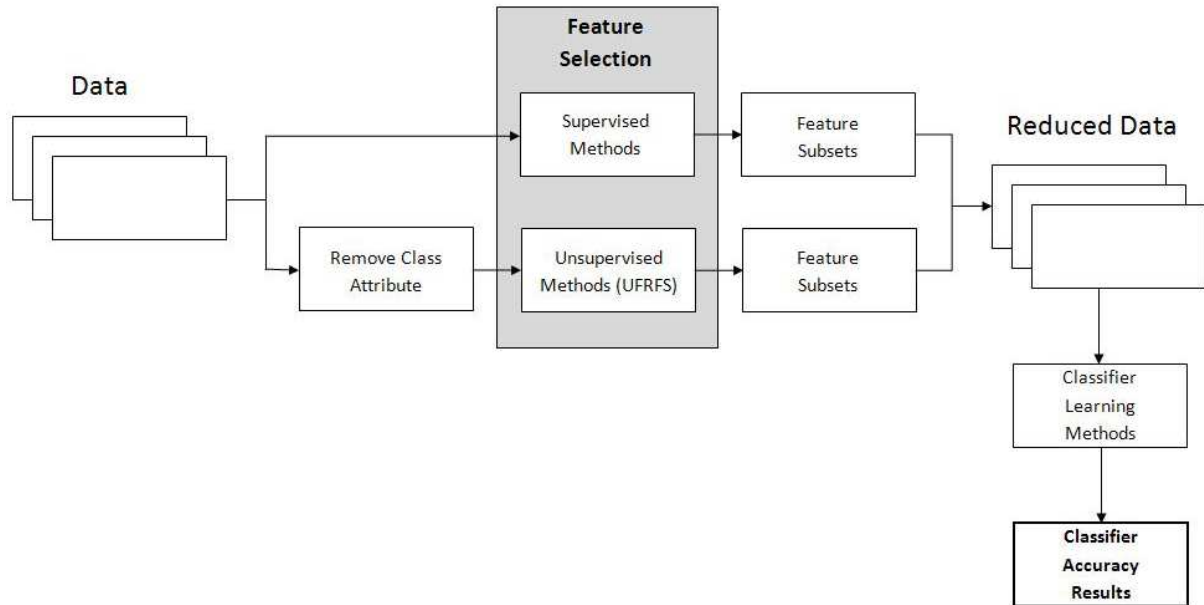


Figure 2. Experimental Evaluation Setup

ods, the decision feature is removed from the data, and the approaches operate on the unlabelled data only. When learning classifiers, or applying supervised FRFS, the complete dataset is used.

The results presented in Table 3 and Table 4 show the subset sizes discovered by the both the supervised and new unsupervised methods. Two different approaches were adopted when generating these results: in the first a subset is selected from all of the data before performing classification, the second involves selecting a subset of features for each fold of the 10 fold cross validation, and getting the average for all ten folds. It can be seen that the proposed methods manage reduction in all cases and return substantial levels of dimensionality reduction for some datasets. These results compare well with the supervised approach

and show that the unsupervised approaches may even find smaller subsets in some cases. In particular, UFRFS manages to outperform the CFS and Consistency-based supervised FS methods for 6 of the 9 datasets for both the average and absolute subset sizes. UFRFS also manages to find a smaller subset for the *Glass* dataset than the supervised FRFS method. It is important to note once again that UFRFS does not consider the class labels and these subsets are obtained only by examining the redundancy of features.

The resulting classification accuracies for the classifiers can be seen in Tables 5 – 8. These demonstrate that the unsupervised methods retain useful features, without considering the decision feature. This is borne out by comparison to the classification accuracy of the unreduced data, showing that the greatest decrease amongst all of the reduced data is

Dataset	Features	Objects	CFS	Consis.	FRFS	B-FRFS	D-FRFS	UFRFS	B-UFRFS	D-UFRFS
Cleveland	13	297	7	9	8	8	8	11	11	11
Glass	9	214	7	9	8	8	8	7	7	7
Heart	13	270	7	10	7	7	7	11	11	11
Ionosphere	34	230	11	7	8	8	7	9	9	9
Olitos	25	120	16	11	5	5	5	6	6	6
Water 2	38	390	9	14	6	6	6	7	8	7
Water 3	38	390	11	11	6	6	6	7	7	7
Wine	13	178	11	5	5	5	5	7	7	6

Table 3. Subset sizes for all approaches using the complete dataset

Dataset	Features	Objects	CFS	Consis.	FRFS	B-FRFS	D-FRFS	UFRFS	B-UFRFS	D-UFRFS
Cleveland	13	297	6.7	10.4	7.7	7.7	7.7	10.5	10.5	10.5
Glass	9	214	6.3	7.6	9	8.2	8.2	7.1	7.1	7.1
Heart	13	270	7.4	11	7.1	7.1	7.1	10.2	10.2	10.2
Ionosphere	34	230	15	8.6	5	5	5	6.2	6	6.2
Olitos	25	120	10.8	10.1	7.1	7.1	6.9	9	9	9
Water 2	38	390	9.1	15	6	6	6	7	7.4	7
Water 3	38	390	10.6	11.6	6	6	5.9	7.1	7.4	7.1
Wine	13	178	10.7	7.3	5	4.9	4.8	6	6.2	6

Table 4. Average Subset Sizes for 10-fold Cross Validation

only in the order of 10% overall. There are also cases where the use of unsupervised-reduced data outperforms the unreduced data and that of the supervised-reduced data. UFRFS demonstrates generally good results when compared to the unreduced data. In particular the *water2* and *Heart* datasets show that UFRFS actually increased classification accuracy when compared to the unreduced data. The *ionosphere* dataset shows a decrease in accuracy for UFRFS using JRip but also shows good comparative results for the other 3 classifier learners. The *olitos* dataset also shows a decrease for UFRFS for PART and FRNN but shows comparable results to those of the supervised methods for JRip and J48. This demonstrates the power of the unsupervised methods as they perform drastic dimensionality reduction that generally maintains the classification accuracy whilst ignoring

the class information. This would imply that using the proposed approaches, the quality of reduction should be high for datasets with missing or incomplete class labels.

The FS process helps to remove measurement noise as a positive by-product of the actual selection. A pertinent question therefore is whether other subsets of dimensionality 6 (e.g. for the *olitos* dataset) would perform similarly as those identified by UFRFS. In order avoid a biased answer to this question, and without resorting to exhaustive computation 25 sets of randomly chosen subsets of size 6 are used to build 25 classifiers. The classification results that are achieved are shown in Fig. 3 along with the error rate of the classifier that uses the UFRFS selected subset.

The average error of the classifiers that each employ five randomly selected features is 40.90%, far higher than

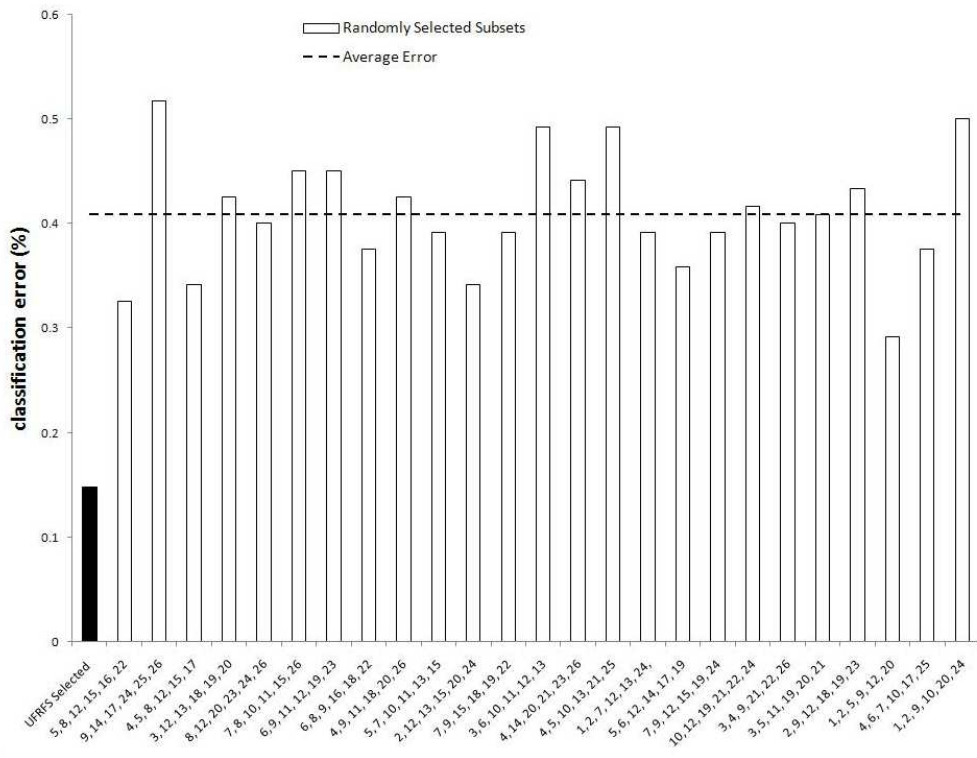


Figure 3. UFRFS vs. randomly selected subsets

Dataset	Unred	CFS	Consis.	FRFS	B-FRFS	D-FRFS	UFRFS	B-UFRFS	D-UFRFS
Cleveland	52.18	56.93	55.51	54.52	54.52	54.52	54.9	53.57	54.9
Glass	71.39	68.23	65.42	71.39	66.9	66.9	64.94	64.94	64.94
Heart	77.41	78.89	78.88	77.04	77.04	77.04	80	80.74	80
Ionosphere	70.83	71.67	89.56	65.83	69.17	63.33	60	67.5	60
Olitos	86.52	90.87	67.50	86.52	86.52	85.22	86.09	84.78	86.09
Water 2	82.82	82.31	84.35	84.1	84.1	81.28	84.87	82.82	84.87
Water 3	81.79	82.56	83.84	78.97	80.26	81.28	80	78.97	80
Wine	95.00	92.12	90.13	88.27	91.57	90.46	74.67	79.9	74.67

Table 5. Classification accuracies: JRip (%)

Dataset	Unred	CFS	Consis.	FRFS	B-FRFS	D-FRFS	UFRFS	B-UFRFS	D-UFRFS
Cleveland	51.87	56.92	56.57	49.84	49.84	49.84	52.91	50.21	52.91
Glass	67.29	69.98	64.47	67.29	65.87	65.87	65.91	65.91	65.91
Heart	76.67	80.74	78.88	77.04	77.04	77.04	78.89	79.26	78.89
Ionosphere	67.5	57.5	89.56	62.5	61.67	70.83	59.17	64.17	59.17
Olitos	87.83	88.7	68.33	86.96	86.96	86.52	85.22	83.91	85.22
Water 2	83.08	84.1	83.58	84.36	84.36	83.59	83.59	82.05	83.59
Water 3	83.08	81.54	81.02	79.49	80.26	80.77	81.54	80.51	81.54
Wine	94.41	94.41	97.10	94.97	96.08	94.41	79.74	81.99	79.74

Table 6. Classification accuracies: J48 (%)

that attained by the classifier which utilises the UFRFS selected subset of the same dimensionality. This implies that those randomly selected entail important information loss in the course of feature selection; this is not the case for the UFRFS unsupervised selection-based approach.

In addition to the empirical evaluation, a statistical analysis was carried out in order to check the significance of the obtained results. A paired t-test was performed on the unreduced data and the UFRFS data to ensure that the generated subsets were not obtained through chance. The results for this are shown in Table 9. For each dataset 10 times 10-fold cross validation was used to generate the test data, then a paired t-test was employed to analyse the results. Again, the UFRFS methods perform well with only the results for the *wine* dataset showing results that are sta-

tistically worse than the unreduced data. However, it should be remembered that with the absence of labels much discriminative information is lost that can otherwise be easily utilised by supervised methods.

5. Conclusion

This paper has presented novel techniques for unsupervised feature selection, based on the fuzzy-rough dependency measure. These approaches are data-driven, and no user-defined thresholds or domain-related information is required, although a choice must be made regarding fuzzy similarity relations and connectives. Note that these choices must also be made for the existing supervised FS approaches that employ the same underlying mathemati-

Dataset	Unred	CFS	Consis.	FRFS	B-FRFS	D-FRFS	UFRFS	B-UFRFS	D-UFRFS
Cleveland	51.85	57.91	55.21	53.19	53.19	53.19	53.87	53.87	53.87
Glass	67.28	68.69	71.96	67.76	70.56	70.56	66.36	66.36	66.35
Heart	76.66	77.03	74.04	76.30	76.30	76.30	81.49	81.49	81.49
Ionosphere	87.82	90.00	88.69	85.23	86.09	76.29	86.08	89.56	86.08
Olitos	67.50	71.67	65.00	64.17	67.50	64.16	60.00	63.33	57.50
Water 2	83.33	83.07	85.64	85.60	84.61	85.64	81.53	85.64	80.00
Water 3	77.43	82.05	82.56	79.74	79.74	79.74	76.67	78.46	78.71
Wine	93.82	93.82	97.1	94.38	94.38	90.46	94.38	86.51	74.71

Table 7. Classification accuracies: PART (%)

Dataset	Unred	CFS	Consis.	FRFS	B-FRFS	D-FRFS	UFRFS	B-UFRFS	D-UFRFS
Cleveland	52.52	53.53	55.55	51.17	51.17	51.17	52.18	52.18	53.87
Glass	73.83	76.63	75.23	73.83	73.36	73.36	71.97	71.96	71.96
Heart	75.55	74.81	76.65	77.40	77.40	75.56	77.03	77.03	77.03
Ionosphere	89.56	90	90.00	88.26	88.26	89.56	86.52	89.13	86.52
Olitos	79.16	81.67	79.16	70.23	62.50	63.33	63.33	50.83	63.33
Water 2	83.58	97.19	85.38	83.58	82.56	82.82	77.17	75.64	77.17
Water 3	80.25	85.38	83.07	80.25	79.23	78.71	76.16	75.64	76.15
Wine	97.19	84.35	97.10	97.19	97.19	91.57	82.02	86.51	82.02

Table 8. Classification accuracies: FRNN (%)

cal theory. The results show that the approach can reduce dataset dimensionality considerably whilst retaining useful features when class labels are unknown or missing.

At present the unsupervised search algorithm utilises a simple but nevertheless effective backwards elimination method for search. The problem with such search techniques is that they often return a result which is a local optimum as the search can proceed down non-optimal paths. In addition to this the adoption of such approaches can be computationally complex. The investigation of other more efficient search techniques such as ant colony optimisation (ACO) [14], particle swarm optimisation (PSO) [24] and propositional satisfiability (SAT) [9], may help in alleviating this problem and thus further improving the efficiency of the approaches. Also, a more complete comparison of

UFRFS and other unsupervised FS techniques for clustering performance, would form the basis for a series of topics for future investigation.

As mentioned previously the fuzzy similarity relations and connectives must be chosen for UFRFS. As only one choice of fuzzy connective (Łukasiewicz), and also a single fuzzy similarity measure (as defined in 11) are explored in this paper, the evaluation of other options in this regard would form the basis for a further more comprehensive investigation. A further interesting topic is the use of a method which would combine both the unsupervised, and supervised measures. The supervised measure determines relevance and the unsupervised measure determines redundancy, hence a method that combines these should be particularly powerful for subset evaluation.

Dataset	Unred	UFRFS	B-UFRFS	D-UFRFS
Cleveland	53.59	55.74	54.63	55.74
Glass	68.08	69.26	69.26	69.26
Heart	78.15	78.85	79.26	78.85
Ionosphere	89.56	84.61	86.04	84.61
Olitos	65.75	59.08	63.50	59.08
Water 2	83.18	81.95	79.28	81.95
Water 3	81.59	80.92	80.92	80.92
Wine	93.37	82.02*	80.31*	79.74*

* denotes a result that is statistically worse than the unreduced data

Table 9. Paired t-test for UFRFS generated results - J48 (% correct)

Acknowledgement

The authors would like to acknowledge the financial support for this research through *The Research institute of Visual Computing - Wales (RIVIC)*.

References

- [1] C.L. Blake, C.J. Merz, *UCI Repository of Machine Learning Databases*. Irvine, University of California, 1998. <http://archive.ics.uci.edu/ml/>
- [2] A. Chouchoulas and Q. Shen, "Rough set-aided keyword reduction for text categorisation", *Applied Artificial Intelligence*, vol. 15, no. 9, pp. 843–873, 2001.
- [3] W.W. Cohen, "Fast effective rule induction," In *Proceedings of the 12th International Conference on Machine Learning*, pp. 115–123, 1995.
- [4] M. De Cock, C. Cornelis, and E.E. Kerre, "Fuzzy Rough Sets: The Forgotten Step", *IEEE Transactions on Fuzzy Systems*, vol.15, no.1, pp.121-130, 2007.
- [5] C. Cornelis, G. Hurtado Martín, R. Jensen, D. Ślęzak, "Feature Selection with Fuzzy Decision Reducts," *3rd Int. Conf. on Rough Sets and Knowledge Technology (RSKT'08)*, pp. 284–291, 2008.
- [6] S.K. Das, "Feature Selection with a Linear Dependence Measure", *IEEE Transactions on Computers*, pp. 1106–1109. 1971.
- [7] M. Dash and H. Liu, "Feature Selection for Classification," *Intelligent Data Analysis*, vol. 1, no. 3, 1997.
- [8] M. Dash and H. Liu, "Unsupervised Feature Selection," *Proceedings of the Pacific and Asia Conference on Knowledge Discovery and Data Mining*, pp. 110–121, 2000.
- [9] M. Davis, G. Logemann, and D. Loveland. "A machine program for theorem proving", *CACM*, vol. 5, pp. 394–397, 1962.
- [10] D. Dubois, H. Prade, "Rough fuzzy sets and fuzzy rough sets," *International Journal of General Systems*,

- vol. 17, 91–209, 1990.
- [11] M.A. Hall, “Correlation-based Feature Selection for Discrete and Numeric Class Machine Learning,” Proceedings of the 17th International Conference on Machine Learning, pp. 359–366, 2000.
- [12] X. Hu, N. Cercone, “Learning in Relational Databases: a Rough Set Approach”, Computational Intelligence, vol.2 pp.323–337, 1995.
- [13] R. Jensen and C. Cornelis. “A New Approach to Fuzzy-Rough Nearest Neighbour Classification”. Proceedings of the 6th International Conference on Rough Sets and Current Trends in Computing, pp. 310–319, 2008.
- [14] R. Jensen, Q. Shen, *Computational Intelligence and Feature Selection: Rough and Fuzzy Approaches*, IEEE Press/Wiley & Sons, 2008.
- [15] R. Jensen, Q. Shen, “New approaches to fuzzy-rough feature selection,” *IEEE Transactions on Fuzzy Systems*, in press, 2009.
- [16] P. Mitra, C.A. Murthy, and S.K. Pal, “Unsupervised Feature Selection Using Feature Similarity,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 24, no. 4, pp. 1–13, 2002.
- [17] S.K. Pal, R.K. De, and J. Basak, “Unsupervised Feature Evaluation: A Neuro-Fuzzy approach,” *IEEE Transactions in Neural Networks*, vol. 11, pp. 366–376, 2000.
- [18] Z. Pawlak, “Rough sets,” *International Journal of Computing and Information Sciences*, vol. 11, pp. 341–356, 1982.
- [19] F. Questier, I. Arnaut-Rollier, B. Walczak, and D. L. Massarta, “Application of rough set theory to feature selection for unsupervised clustering”, *Chemometrics and Intelligent Laboratory Systems*, vol. 63, no. 2, pp. 155–167, 2002
- [20] A.M. Radzikowska and E.E. Kerre, “A comparative study of fuzzy rough sets,” *Fuzzy Sets and Systems*, vol. 126, no. 2, pp. 137–155, 2002.
- [21] A. Skowron, and C. Rauszer, “The Discernibility Matrices and Functions in Information Systems”, in: R.Slowiski(ed.), *Intelligent Decision Support Handbook of Applications and Advances of the Rough Set Theory*, Kluwer Academic Publishers, Dordrecht, pp. 311–362, 1992.
- [22] A. Skowron, and J. Stepaniuk, “Tolerance Approximation Spaces”, *Fundamenta Informaticae*, vol. 27, pp. 245–253, 1996.
- [23] J. Wang, and J.Wang, “Reduction Algorithms Based on Discernibility Matrix: The Ordered Attributes

- Method” *Journal of Computer Science & Technology*, vol. 16, no. 6, pp. 489–504, 2001.
- [24] X. Wang, J. Yang, X. Teng, W. Xia and R. Jensen, “Feature Selection based on Rough Sets and Particle Swarm Optimization,” *Pattern Recognition Letters*, vol. 28, no. 4, pp. 459–471, 2007.
- [25] I.H. Witten and E. Frank, “Generating Accurate Rule Sets Without Global Optimization,” *Proceedings of the 15th International Conference on Machine Learning*, Morgan Kaufmann Publishers, San Francisco, 1998.
- [26] I.H. Witten and E. Frank, *Data Mining: Practical machine learning tools with Java implementations*. Morgan Kaufmann, San Francisco, 2000.
- [27] W. Ziarko, “Variable Precision Rough Set Model”, *Journal of Computer and System Sciences*, vol. 46, no. 1, pp. 39–59, 1993.