## JID: ESWA

# ARTICLE IN PRESS

Expert Systems With Applications xxx (2015) xxx-xxx

[m5G;July 24, 2015;17:13]



Contents lists available at ScienceDirect

**Expert Systems With Applications** 



journal homepage: www.elsevier.com/locate/eswa

# Feature selection using Joint Mutual Information Maximisation

#### Q1

School of Engineering, Cardiff University, Cardiff CF24 3AA, UK

Mohamed Bennasar, Yulia Hicks, Rossitza Setchi\*

#### ARTICLE INFO

Keywords: Feature selection Mutual information Joint mutual information Conditional mutual information Subset feature selection Classification Dimensionality reduction Feature selection stability

## ABSTRACT

Feature selection is used in many application areas relevant to expert and intelligent systems, such as data mining and machine learning, image processing, anomaly detection, bioinformatics and natural language processing. Feature selection based on information theory is a popular approach due its computational efficiency, scalability in terms of the dataset dimensionality, and independence from the classifier. Common drawbacks of this approach are the lack of information about the interaction between the features and the classifier, and the selection of redundant and irrelevant features. The latter is due to the limitations of the employed goal functions leading to overestimation of the feature significance.

To address this problem, this article introduces two new nonlinear feature selection methods, namely Joint Mutual Information Maximisation (JMIM) and Normalised Joint Mutual Information Maximisation (NJMIM); both these methods use mutual information and the 'maximum of the minimum' criterion, which alleviates the problem of overestimation of the feature significance as demonstrated both theoretically and experimentally. The proposed methods are compared using eleven publically available datasets with five competing methods. The results demonstrate that the JMIM method outperforms the other methods on most tested public datasets, reducing the relative average classification error by almost 6% in comparison to the next best performing method. The statistical significance of the results is confirmed by the ANOVA test. Moreover, this method produces the best trade-off between accuracy and stability.

© 2015 Published by Elsevier Ltd.

## 1 1. Introduction

High dimensional data is a significant problem in both super-2 3 vised and unsupervised learning (Janecek, Gansterer, Demel, & Ecker, 2008), which is becoming even more prominent with the recent ex-4 plosion of the size of the available datasets both in terms of the num-5 ber of data samples and the number of features in each sample (Zhang 6 7 et al., 2015). The main motivation for reducing the dimensionality of 8 the data and keeping the number of features as low as possible is to decrease the training time and enhance the classification accuracy of 9 the algorithms (Guyon & Elisseeff, 2003; Jain, Duin, & Mao, 2000; Liu 10 & Yu, 2005). 11

Dimensionality reduction methods can be divided into two main groups: those based on feature extraction and those based on feature selection. Feature extraction methods transform existing features into a new feature space of lower dimensionality. During this process, new features are created based on linear or nonlinear combinations of features from the original set. Principal Component Analysis (PCA) (Bajwa, Naweed, Asif, & Hyder, 2009; Turk & Pentland, 1991) and Linear Discriminant Analysis (LDA) (Tang, Suganthana, Yao, & Qina, 19 2005; Yu & Yang, 2001) are two examples of such algorithms. Feature 20 selection methods reduce the dimensionality by selecting a subset 21 of features which minimises a certain cost function (Guyon, Gunn, 22 Nikravesh, & Zadeh, 2006; Jain et al., 2000). Unlike feature extraction, 23 feature selection does not alter the data and, as a result, it is the 24 preferred choice when an understanding of the underlying physical 25 process is required. Feature extraction may be preferred when only 26 discrimination is needed (Jain et al., 2000). 27

Feature selection is used in many application areas relevant to expert and intelligent systems, such as data mining and machine learning, image processing, anomaly detection, bioinformatics and natural language processing (Hoque, Bhattacharyya, & Kalita, 2014). Feature selection is normally used at the data pre-processing stage before training a classifier. This process is also known as variable selection, feature reduction or variable subset selection.

The topic of feature selection has been reviewed in detail in a 35 number of recent review articles (Bolón-Canedo, Sánchez-Maroño, 36 & Alonso-Betanzos, 2013; Brown, Pocock, Zhao, & Lujan, 2012; 37 Chandrashekar & Sahin, 2014; Vergara & Estévez, 2014). Usually, 38 feature selection methods are divided into two categories in terms of 39 evaluation strategy, in particular, classifier dependent ('wrapper' and 40 'embedded' methods) or classifier independent ('filter' methods). 41 Wrapper methods search the feature space, and test all possible 42

Please cite this article as: M. Bennasar et al., Feature selection using Joint Mutual Information Maximisation, Expert Systems With Applications (2015), http://dx.doi.org/10.1016/j.eswa.2015.07.007

<sup>\*</sup> Corresponding author. Tel: +44 2920875720; fax: +44 2920874716.

*E-mail addresses*: BennasarM@cf.ac.uk (M. Bennasar), HicksYA@cf.ac.uk (Y. Hicks), Setchi@cf.ac.uk (R. Setchi).

JID: ESWA

2

# ARTICLE IN PRESS

122

M. Bennasar et al. / Expert Systems With Applications xxx (2015) xxx-xxx

subsets of feature combinations by using the prediction accuracy 43 44 of a classifier as a measure of the selected subset's quality, without modifying the learning function. Therefore, wrapper methods 45 46 can be combined with any learning machine (Guyon et al., 2006). They perform well because the selected subset is optimised for the 47 classification algorithm. On the other hand, wrapper methods may 48 suffer from over-fitting to the learning algorithm. This means that 49 any changes in the learning model may reduce the usefulness of 50 51 the subset. In addition, these methods are very expensive in terms of computational complexity, especially when handling extremely 52 53 high-dimensional data (Brown et al., 2012; Cheng et al., 2011; Ding & Peng, 2003; Karegowda, Jayaram, & Manjunath, 2010). Q2 54

The feature selection stage in the embedded methods is combined with the learning stage. These methods are less expensive in terms of computational complexity and less prone to over-fitting; however, they are limited in terms of generalisation, because they are very specific to the used learning algorithm (Guyon et al., 2006).

Classifier-independent methods rank features according to their 60 relevance to the class label in the supervised learning. The relevance 61 score is calculated using distance, information, correlation and 62 consistency measures. Many techniques have been proposed to 63 compute the relevance score, including Pearson correlation coef-64 65 ficients (Rodgers & Nicewander, 1988), Fisher's discriminate ratio 66 "F score" (Lin, Li, & Tsai, 2004), the Scatter criterion (Duda, Hart, & Stork, 2001), Single Variable Classifier SVC (Guyon & Elisseeff, 2003), 67 Mutual Information (Battiti, 1994), the Relief Algorithm (Kira & 68 Rendell, 1992; Liu & Motoda, 2008), Rough Set Theory (Liang, Wang, 69 70 Dang, & Qian, 2014) and Data Envelopment Analysis (Zhang, Yang, 71 Xiong, Wang, & Zhang, 2014).

72 The main advantages of the filter methods are their computa-73 tional efficiency, scalability in terms of the dataset dimensionality, 74 and independence from the classifier (Saeys, Inza, & Larranaga, 2007). 75 A common drawback of these methods is the lack of information about the interaction between the features and the classifier and 76 77 selection of redundant and irrelevant features due to the limitations of the employed goal functions leading to overestimation of the 78 79 feature significance.

80 Information theory (Cover & Thomas, 2006) has been widely applied in filter methods, where information measures such as 81 mutual information (MI) are used as a measure of the features' 82 relevance and redundancy (Battiti, 1994). MI does not make an 83 assumption of linearity between the variables, and can deal with 84 categorical and numerical data with two or more class values (Meyer, 85 86 Schretter, & Bontempi, 2008). There are several alternative measures 87 in information theory that can be used to compute the relevance of features, namely mutual information, interaction information, 88 89 conditional mutual information, and joint mutual information.

This paper contributes to the knowledge in the area of feature 90 selection by proposing two new nonlinear feature selection meth-91 ods based on information theory. The proposed methods aim to 92 overcome the limitations of the current state of the art filter feature 93 94 selection methods such as overestimation of the feature significance, 95 which causes selection of redundant and irrelevant features. This is achieved through the introduction of a new goal function based 96 97 on joint mutual information and the 'maximum of the minimum' nonlinear approach. As shown in the evaluation section, one of the 98 99 proposed methods outperforms the competing feature selection methods in terms of classification accuracy, decreasing the average 100 classification error by 0.88% in absolute terms and almost by 6% in 101 relative terms in comparison to the next best performing method. 102 In addition, it produces the best trade-off between accuracy and 103 stability. The statistical significance of the reported results is further 104 confirmed by ANOVA test. 105

This paper also reviews existing feature selection methods highlighting their common limitations and compares the performance of the proposed and existing methods on the basis of several criteria. For example, a nonlinear approach, which employs the 'maximum of the minimum' criterion, is compared to a linear approach, which employs cumulative summation approximation. To optimise the nonlinear approach, a goal function based on joint mutual information is compared to the goal function based on conditional mutual information. Finally, the effect of using normalised mutual information instead of mutual information is tested.

The rest of the paper is organised as follows. Section 2 presents the116principles of the information theory, Section 3 reviews related work,117Section 4 discusses the limitations of current feature selection crite-118ria, Section 5 introduces the proposed methods. Section 6 describes119the conducted experiments and discusses the results. Section 7 con-120cludes the paper.121

## 2. Information theory

This section introduces the principles of information theory by fo-<br/>cusing on entropy and mutual information and explains the reasons123<br/>124for employing them in feature selection.125

The entropy of a random variable is a measure of its uncertainty 126 and a measure of the average amount of information required to describe the random variable (Cover & Thomas, 2006). The entropy of a 128 discrete random variable  $X = (x_1, x_2, \dots, x_N)$  is denoted by H(X), 129 where  $x_i$  refers to the possible values that X can take. H(X) is defined as: 131

$$H(X) = -\sum_{i=1}^{N} p(x_i) \log(p(x_i)),$$
(1)

where  $p(x_i)$  is the probability mass function. The value of  $p(x_i)$ , when 132 X is discrete, is: 133

$$p(x_i) = \frac{number of instants with value x_i}{total number of instants (N)}.$$
 (2)

The base of the logarithm, log, is 2, so  $0 \le H(X) \le 1$ . For any two discrete random variables X and  $C = (c_1, c_2, \dots, c_M)$ , the joint entropy is defined as:

$$H(X,C) = -\sum_{j=1}^{M} \sum_{i=1}^{N} p(x_i, c_j) \log(p(x_i, c_j))$$
(3)

where  $p(x_i, c_j)$  is the joint probability mass function of the variables 137 X and C. The conditional entropy of the variable X given C is defined as: 139

$$H(C|X) = -\sum_{j=1}^{M} \sum_{i=1}^{N} p(x_i, c_j) \log(p(c_j|x_i))$$
(4)

The conditional entropy is the amount of uncertainty left in C when140a variable X is introduced, so it is less than or equal to the entropy of141both variables. The conditional entropy is equal to the entropy if, and142only if, the two variables are independent. The relation between joint143entropy and conditional entropy is:144

$$H(X,C) = H(X) + H(C|X)$$
(5)

$$H(X,C) = H(C) + H(X|C)$$
 (6) <sup>145</sup>

Mutual Information (MI) is the amount of information that both variables share, and is defined as: 147

$$I(X; C) = H(C) - H(C|X)$$
 (7)

MI can be expressed as the amount of information provided by variable *X*, which reduces the uncertainty of variable *C*. MI is zero if the random variables are statistically independent. MI is symmetric, so: 150

$$I(X;C) = I(C;X) \tag{8}$$

$$I(X;C) = H(X) - H(X|C)$$
(9) <sup>151</sup>

Please cite this article as: M. Bennasar et al., Feature selection using Joint Mutual Information Maximisation, Expert Systems With Applications (2015), http://dx.doi.org/10.1016/j.eswa.2015.07.007

Download English Version:

# https://daneshyari.com/en/article/10321769

Download Persian Version:

https://daneshyari.com/article/10321769

Daneshyari.com