

The effect of principal component analysis on machine learning accuracy with high-dimensional spectral data

Tom Howley^a, Michael G. Madden^{a,*}, Marie-Louise O'Connell^b, Alan G. Ryder^b

^a Department of Information Technology, National University of Ireland, University Road, Galway, Ireland

^b National Centre for Biomedical Engineering Science, National University of Ireland, University Road, Galway, Ireland

Received 28 October 2005; accepted 28 November 2005

Available online 8 February 2006

Abstract

This paper presents the results of an investigation into the use of machine learning methods for the identification of narcotics from Raman spectra. The classification of spectral data and other high-dimensional data, such as images, gene-expression data and spectral data, poses an interesting challenge to machine learning, as the presence of high numbers of redundant or highly correlated attributes can seriously degrade classification accuracy. This paper investigates the use of principal component analysis (PCA) to reduce high-dimensional spectral data and to improve the predictive performance of some well-known machine learning methods. Experiments are carried out on a high-dimensional spectral dataset. These experiments employ the NIPALS (Non-Linear Iterative Partial Least Squares) PCA method, a method that has been used in the field of chemometrics for spectral classification, and is a more efficient alternative than the widely used eigenvector decomposition approach. The experiments show that the use of this PCA method can improve the performance of machine learning in the classification of high-dimensional data.

© 2006 Elsevier B.V. All rights reserved.

Keywords: Machine learning; High-dimensional data; Principal component analysis; NIPALS; Spectroscopy

1. Introduction

The automatic identification of illicit materials using Raman spectroscopy is of significant importance for law enforcement agencies. High-dimensional spectral data can pose problems for machine learning as predictive models based on such data run the risk of overfitting. Furthermore, many of the attributes may be redundant or highly correlated, which can also lead to a degradation of prediction accuracy.

This problem is equally relevant to many other application domains, such as the classification of gene-expression microarray data [1], image data [2] and text [3]. In the classification task considered in this paper, Raman spectra are

used for the identification of acetaminophen, a pain-relieving drug that is found in many over-the-counter medications, within different mixtures. In the physical sciences, a statistical approach to classification is normally taken (Chemometrics) [4], and these methods use PCA to handle the high-dimensional spectra. PCA is a classical statistical method for transforming attributes of a dataset into a new set of uncorrelated attributes called principal components (PCs). PCA can be used to reduce the dimensionality of a dataset, while still retaining as much of the *variability* of the dataset as possible. The goal of this research is to determine if PCA can be used to improve the performance of machine learning methods in the classification of such high-dimensional data.

In the first set of experiments presented in this paper, the performance of five well-known machine learning techniques (support vector machines, *k*-nearest neighbours, C4.5 decision tree, RIPPER and Naive Bayes) along with classification by linear regression are compared by testing

* Corresponding author.

E-mail addresses: thowley@vega.it.nuigalway.ie (T. Howley), michael.madden@nuigalway.ie (M.G. Madden), ML.OConnell@nuigalway.ie (M.-L. O'Connell), alan.ryder@nuigalway.ie (A.G. Ryder).

them on a Raman spectral dataset. A number of pre-processing techniques such as normalisation and first derivative are applied to the data to determine if they can improve the classification accuracy of these methods. A second set of experiments is carried out in which PCA and machine learning (and the various pre-processing methods) are used in combination. This set of PCA experiments also facilitates a comparison of machine learning with the popular chemometric technique of principal component regression (PCR), which combines PCA and linear regression.

The main contributions of this research are as follows:

- (1) It presents a promising approach for the classification of substances within complex mixtures based on Raman spectra, an application that has not been widely considered in the machine learning community. This approach could also be applied to other high-dimensional classification problems.
- (2) It proposes the use of NIPALS PCA for data reduction, a method that is much more efficient than the widely used eigenvector decomposition method.
- (3) It demonstrates the usefulness of PCA for reducing dimensionality and improving the performance of a variety of machine learning methods. Previous work has tended to focus on a single machine learning method. It also demonstrates the effect of reducing data to different numbers of principal components.

The paper is organised as follows. Section 2 will give a brief description of Raman spectroscopy and outline the characteristics of the data it produces. Section 3 describes PCA, the NIPALS algorithm for PCA that is used here and the PCR method that incorporates PCA into it. Section 4 provides a brief description of each machine learning technique used in this investigation. Experimental results along with a discussion are presented in Section 5. Section 6 describes related research and Section 7 presents the conclusion of this study.

2. Raman spectroscopy

Raman spectroscopy is the measurement of the wavelength and intensity of light that has been scattered inelastically by a sample, known as the Raman effect [5]. This Raman scattering provides information on the vibrational motions of molecules in the sample compound, which in turn provides a molecular fingerprint. Every compound has its own unique Raman spectrum that can be used for sample identification. Each point of a spectrum represents the intensity recorded at a particular wavelength. A Raman dataset therefore has one attribute for each point on its constituent spectra. Raman spectra can be used for the identification of materials such as narcotics [4], hazardous waste [6] and explosives [7].

Raman spectra are a good example of high-dimensional data; a Raman spectrum is typically made up of 500–3000

data points, and many datasets may only contain 20–200 samples. However, there are other characteristics of Raman spectra that can be problematic for machine learning:

- *Collinearity*: many of the attributes (spectral data points) are highly correlated to each other which can lead to a degradation of the prediction accuracy.
- *Noise*: particularly prevalent in spectra of complex mixtures. Predictive models that are fitted to noise in a dataset will not perform well on other test datasets.
- *Fluorescence*: the presence of fluorescent materials in a sample can obscure the Raman signal and therefore make classification more difficult [4].
- *Variance of Intensity*: a wide variance in spectral intensity occurs between different sample measurements [8].

3. Principal component analysis

In the following description, the dataset is represented by the matrix X , where X is a $N \times p$ matrix. For spectral applications, each row of X , the p -vector x_i contains the intensities at each wavelength of the spectrum sample i . Each column, X_j contains all the observations of one attribute. PCA is used to overcome the previously mentioned problems of high-dimensionality and collinearity by reducing the number of predictor attributes. PCA transforms the set of inputs X_1, X_2, \dots, X_n into another set of column vectors T_1, T_2, \dots, T_N where the T 's have the property that most of the original data's information content (or most of its variance) is stored in the first few T 's (the principal component scores). The number of PCs that account for a portion of the total variance of the original data (i.e., their associated eigenvalues are non-zero) is equal to either $N - 1$ or p , which ever is the smaller. With PCA the data can be reduced to a smaller number of dimensions, with low information loss, simply by discarding the higher numbered PCs and retaining the first set of PCs that account for most of the original data's total variance. Each PC is a linear combination of the original inputs and each PC is orthogonal, which therefore eliminates the problem of collinearity. This linear transformation of the matrix X is specified by a $p \times p$ matrix P so that the transformed variables T are given by:

$$T = XP \text{ or alternatively } X \text{ is decomposed as follows:} \\ X = TP^T, \quad (1)$$

where P is known as the *loadings matrix*. The columns loadings matrix P can be calculated as the eigenvectors of the matrix $X^T X$ [9], a calculation which can be computationally intensive when dealing with datasets of 500–3000 attributes. A much quicker alternative is the NIPALS method. The NIPALS method does not calculate all the PCs at once as is done in the eigenvector approach. Instead, it calculates the first PC by getting the first PC score,

t_1 , and the first vector of the loadings matrix, p'_1 from the sample matrix X . Then the outer product, $t_1 p'_1$, is subtracted from X and the residual, E_1 , is calculated. This residual becomes X in the calculation of the next PC and the process is repeated until as many PCs as required have been generated. The algorithm for calculating the n th PC is detailed below [10]:

- (1) Take a vector x_j from X and call it $t_n: t_n = x_j$
- (2) Calculate $p'_n: p'_n = t'_n X / t'_n t_n$
- (3) Normalise p'_n to length 1: $p'_{n\text{new}} = p'_{n\text{old}} / \|p'_{n\text{old}}\|$
- (4) Calculate $t_n: t_n = X p'_n / p'_n p'_n$
- (5) Compare t_n used in step 2 with that obtained in step 4. If they are the same, stop (the iteration has converged). If they still differ, go to step 2.

After the first PC has been calculated (i.e., t_1 has converged), X in steps 2 and 4 is replaced by its residual; for example, to generate the second PC, X is replaced by E_1 , where $E_1 = X - t_1 p'_1$.

See Ryder [4], O'Connell et al. [8] and Conroy et al. [6] for examples of the use of PCA in the classification of materials from Raman spectra.

3.1. Principal component regression

The widely used chemometric technique of PCR is a two-step multivariate regression method, in which PCA of the data is carried out in the first step. In the second step, a multiple linear regression between the PC scores obtained in the PCA step and the predictor variable is carried out. In this work, the predictor variable is a value that is chosen to represent the presence or absence of the target in a sample, e.g., 1 for present and -1 for absent. In this way, a classification model can be built using any regression method.

4. Machine learning

4.1. Support vector machine

The support vector machine (SVM) [11] is a powerful machine learning tool that is capable of representing non-linear relationships and producing models that generalise well to unseen data. For binary classification, a linear SVM (the simplest form of SVM) finds an optimal linear separator between the two classes of data. This optimal separator is the one that results in the widest margin of separation between the two classes, as a wide margin implies that the classifier is better able to classify unseen spectra. To regulate overfitting, SVMs have a complexity parameter, C , which determines the trade-off between choosing a large-margin classifier and the amount by which misclassified samples are tolerated. A higher value of C means that more importance is attached to minimising the amount of misclassification than to finding a wide margin model. To handle non-linear data, kernels (e.g., radial basis function, RBF, polynomial or sigmoid) are introduced to map the

original data to a new feature space in which a linear separator can be found. In addition to the C parameter, each kernel may have a number of parameters associated with it. For the experiments reported here, two kernels were used: the RBF kernel, in which the kernel width, σ , can be changed, and the linear kernel, which has no extra parameter. In general, the SVM is considered useful for handling high-dimensional data.

4.2. k -Nearest neighbours

k -Nearest neighbours (k -NN) [12] is a learning algorithm which classifies a test sample by firstly obtaining the class of the k samples that are the closest to the test sample. The majority class of these nearest samples (or nearest single sample when $k = 1$) is returned as the prediction for that test sample. Various measures may be used to determine the distance between a pair of samples. In these experiments, the Euclidean distance measure was used. In practical terms, each Raman spectrum is compared to every other spectrum in the dataset. At each spectral data point, the difference in intensity between the two spectra is measured (distance). The sum of the squared distances for all the data points (full spectrum) gives a numerical measure of how close the spectra are.

4.3. C4.5

The C4.5 decision tree [13] algorithm generates a series of if-then rules that are represented as a tree structure. Each node in the tree corresponds to a test of the intensity at a particular data point of the spectrum. The result of a test at one node determines which node in the tree is checked next until finally, a leaf node is reached. Each leaf specifies the class to be returned if that leaf is reached.

4.4. RIPPER

RIPPER [14] (Repeated Incremental Pruning to Produce Error Reduction) is an inductive rule-based learner that builds a set of propositional rules that identify classes while minimising the amount of error. The number of training examples misclassified by the rules defines the error. RIPPER was developed with the goal of handling large noisy datasets efficiently whilst also achieving good generalisation performance.

5. Experimental results

5.1. Dataset

In the following experiments, the task is to identify acetaminophen in a variety of solid mixtures. The acetaminophen dataset comprises the Raman spectra of 217 different samples. Acetaminophen is present in 87 of the samples, the rest of the samples being made up of various pure inorganic materials. Each sample spectrum covers

the range 350–2000 cm^{-1} and is made up of 1646 data points. For more details on this dataset, see O’Connell et al. [8].

5.2. Comparison of machine learning methods

Table 1 shows the results of six different machine learning classification methods using a 10-fold cross-validation test on the acetaminophen dataset. Each column shows the average classification error achieved on one of the four different variations of the acetaminophen dataset used; these variations are described below:

- RD: raw data, unprocessed.
- ND: dataset with each sample normalised. Each sample is divided across by the maximum intensity that occurs within that sample.
- FD: a Savitzky-Golay first derivative [15], seven-point averaging algorithm is applied to the raw dataset.
- FND: a normalisation step is carried out after applying a first derivative to each sample of the raw dataset.

Table 1 shows the lowest average error average achieved by each classifier and pre-processing combination. For all these methods, apart from k -NN, the WEKA [12] implementation was used. The default settings were used for C4.5, RIPPER and Naive Bayes. For SVMs, RBF and linear kernels with different parameter settings were tested. The parameter settings that achieved the best results are shown in parentheses. The Linear SVM was tested for the following values of C : 0.1, 1, ..., 10,000. The same range of C values were used for RBF SVM, and these were tested in combination with the σ values of: 0.0001, 0.001, ..., 10. For k -NN, the table shows the value for k (number of neighbours) that resulted in the lowest percentage error. The k -NN method was tested for all values of k from 1 to 20. The results of each machine learning and pre-processing technique combination of Table 1 were compared using a paired t -test based on a 5% confidence level and using a corrected variance estimate [16]. The lowest average error over all results in Table 1 of 0.92% (i.e., only two misclassifications, achieved by both Linear and RBF SVM) is highlighted in bold and indicated by an asterisk. Those results which do not differ significantly (according to the t -test) are also highlighted in bold.

On both the raw (RD) and normalised (ND) dataset, both SVM models perform better than any of the other machine learning methods, as there is no significant difference between the best overall result and the SVM results on RD and ND, whereas a significant difference does exist between the best overall result and all other machine learning methods on RD and ND. This confirms the notion that SVMs are particularly suited to dealing with high-dimensional data and it also suggests that SVMs are capable of handling a high degree of collinearity in the data. Linear regression, on the other hand, performs poorly with all pre-processing techniques.

This poor performance can be attributed to its requirement that all the columns of the data matrix are *linearly independent* [9], a condition that is violated in highly correlated spectral data. Similarly, Naive Bayes has recorded a high average error on the RD, ND and FD data. This is presumably because of its assumption of independence of each of the attributes. It is clear from this table that the pre-processing techniques of FD and FND improve the performance of the majority of the classifiers. For SVMs, the error is numerically smaller, but not a significant improvement over the RD and ND results. Note that linear regression is the only method that did not achieve a result to compete with the best overall result.

Overall, the SVM appears to exhibit the best results, matching or outperforming all other methods on the raw and pre-processed data. With effective pre-processing, however, the performance of other machine learning methods can be improved so that they are close to that of the SVM.

5.3. Comparison of machine learning methods with PCA

As outlined in Section 3, PCA is used to alleviate problems such as high dimensionality and collinearity that are associated with spectral data. For the next set of experiments, the goal was to determine whether machine learning methods could benefit from an initial transformation of the dataset into a smaller set of PCs, as is used in PCR. The same series of cross-validation tests were run, except in this case, during each fold the PC scores of the training data were fed as inputs to the machine learning method. The procedure for the 10-fold cross-validation is as follows:

Table 1
Percentage classification error of different machine learning methods on acetaminophen dataset

Method	Pre-processing technique			
	RD	ND	FD	FND
Linear SVM	6.45 ($C = 100$)	2.76 ($C = 1$)	3.23 ($C = 10000$)	0.92* ($C = 0.1$)
RBF SVM	5.07 ($C = 1000, \sigma = 0.1$)	2.76 ($C = 1000, \sigma = 0.001$)	1.84 ($C = 1000, \sigma = 10$)	0.92* ($C = 10, \sigma = 0.01$)
k -NN	11.06 ($k = 1$)	7.83 ($k = 1$)	4.61 ($k = 10$)	4.15 ($k = 1$)
C4.5	10.14	7.83	1.84	1.38
RIPPER	15.67	11.06	3.69	2.3
Naive Bayes	25.35	13.82	25.81	5.53
Linear regression	27.65	16.13	25.35	20.28

- (1) Carry out PCA on the training data to generate a loadings matrix.
- (2) Transform training data into a set of PC scores using the first P components of the loadings matrix.
- (3) Build a classification model based on the training PC scores data.
- (4) Transform the held out test fold data to PC scores using the loadings matrix generated from the training data.
- (5) Test classification model on the transformed test fold.
- (6) Repeat steps 1–5 for each iteration of the 10-fold cross-validation.

With each machine learning and pre-processing method combination, the above 10-fold cross-validation test was carried out for $P = 1$ –20 principal components. Therefore, 20 different 10-fold cross-validation tests were run for Naive Bayes, for example. For those classifiers that require additional parameters to be set, more tests had to be run to test the different combinations of parameters, e.g., C , σ and P for RBF SVM. The same ranges for C , σ and k were tested as those used for the experiments of Table 1.

Table 2 shows the lowest average error achieved by each combination of machine learning and pre-processing method with PCA. The number of PCs used to achieve this lowest average error is shown in parentheses, along with the additional parameter settings for the SVM and k -NN classifiers. As with Table 1, the best result over all the results of Table 2 is highlighted in bold and denoted by an asterisk, with those results that bear no significant difference from the best overall result also highlighted in bold. Again, the pre-processing method of FND improves the performance of the majority of the classifiers, Naive Bayes being the exception in this case. In comparing the best result of Table 1 with the best result of Table 2 for each machine learning method (all in the FND column), it can be seen that the addition of the PCA step results in either the same error (C4.5 and RIPPER) or a numerically smaller error (linear SVM, RBF SVM, k -NN and linear regression). The improvement effected by the inclusion of this PCA step is particularly evident with the linear regression technique. Note that this combination of PCA and linear regression is equivalent to PCR.

Despite the fact that for the SVM and k -NN classifiers, there is no significant difference between the best results with or without PCA, it is noteworthy that the SVM and k -NN classifiers with PCA were capable of achieving such low errors with far fewer attributes, only four PCs for the Linear SVM and k -NN and 5 PCs for the RBF SVM. This makes the resulting classification model much more efficient when classifying new data. In contrast, PCR achieved its lowest error using a much greater number of PCs (80), as found in the experiment detailed in the next section. (The chemometric interpretation of such a high number of PCs is open to question.)

To make an overall assessment of the effect of using PCA in combination with machine learning, a statistical comparison (paired t -test with 5% confidence level) of the 28 results of Tables 1 and 2 was carried out. This indicates that, overall, a significant improvement in the performance of machine learning methods is gained with this initial PCA step. It can therefore be concluded that the incorporation of PCA into machine learning is useful for the classification of high-dimensional data.

5.4. Effect of PCA on classification accuracy

In the conventional use of PCA, small sets of PCs (e.g., in many applications, fewer than 20 PCs might account for over 99% of the total variance) are used as they will typically account for the majority of the total variance of the original data. Classification models based on higher numbered PCs (of low variance) run the risk of overfitting to noise in the training data and may therefore perform poorly on test data. The next set of experiments was carried out to determine exactly how the performance of each machine learning method is affected when larger numbers of PCs are used as input. Each method (using the best parameter setting and pre-processing technique determined from Table 2) was used for values of P in the range 1–160. This range of values for P was deemed sufficient for demonstrating the effect of PCA on classification performance. Note also that in this case the maximum number of PCs that correspond to non-zero eigenvalues is equal to one less than the total number of samples in the training data. Subsequent PCs account for no part of the original data's variance.

Table 2
Percentage classification error of different machine learning methods with PCA on acetaminophen dataset

Method	Pre-processing technique			
	RD	ND	FD	FND
Linear SVM	5.07 ($P = 18$, $C = 0.1$)	1.84 ($P = 13$, $C = 0.1$)	3.23 ($P = 14$, $C = 0.01$)	0.46 ($P = 4$, $C = 0.1$)
RBF SVM	6.91 ($P = 19$, $C = 100$, $\sigma = 0.001$)	2.76 ($P = 16$, $C = 10$, $\sigma = 0.001$)	2.23 ($P = 12$, $C = 10$, $\sigma = 0.001$)	0.46 ($P = 5$, $C = 10$, $\sigma = 0.001$)
k -NN	11.06 ($P = 17$, $k = 3$)	5.99 ($P = 10$, $k = 1$)	2.3 ($P = 14$, $k = 1$)	0.0* ($P = 4$, $k = 5$)
C4.5	7.83 ($P = 20$)	7.37 ($P = 19$)	7.37 ($P = 5$)	1.38 ($P = 6$)
RIPPER	11.98 ($P = 20$)	8.29 ($P = 8$)	6.45 ($P = 5$)	2.3 ($P = 3$)
Naive Bayes	38.71 ($P = 1$)	10.6 ($P = 8$)	11.52 ($P = 5$)	3.23 ($P = 2$)
PCR (PCA + linear regression)	9.22 ($P = 16$)	5.53 ($P = 20$)	8.29 ($P = 11$)	1.38 ($P = 80$)

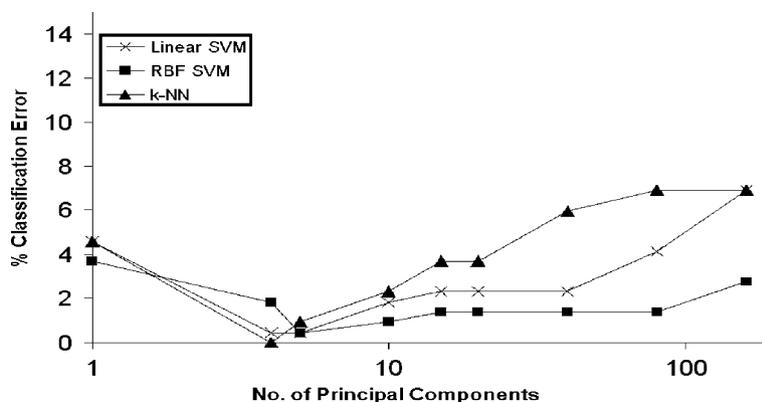


Fig. 1. Effect of changing the number of PCs on machine learning classification error.

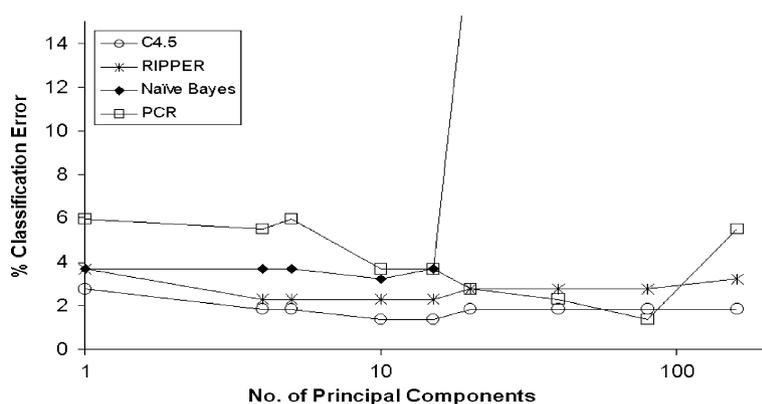


Fig. 2. Effect of changing the number of PCs on machine learning classification error.

Figs. 1 and 2 show the change in error for each of the methods versus the number of PCs retained to build the model. It can be seen from these graphs that as PCs are added, error is initially reduced for all methods. Most methods require no more than six PCs to achieve the lowest error. After this lowest error point, the behaviour of the methods differ somewhat. Most of the classifiers suffer drastic increases in error within the range of PCs tested: Naive Bayes, PCR, SVM and *k*-NN. In contrast, the error for C4.5 and RIPPER never deviates too much from the lowest error. This may be due to their ability to prune irrelevant attributes, from either the decision tree model of C4.5 or the set of rules produced by RIPPER. Overall, it is evident that all of the classifiers achieve good accuracy with a relatively small number of PCs; it is probably unnecessary to generate any more than 20 PCs. However, the number of PCs required will depend on the underlying dataset. Further experiments on more spectral data, or other examples of high-dimensional data, are required to determine suitable ranges of PCs for these machine learning methods.

5.5. Experiments on chlorinated dataset

To extend the results of the acetaminophen experiments, a further set of experiments was carried out on another

dataset of Raman spectra: Chlorinated dataset. This dataset contains the spectra for 230 sample mixtures, each made up of different combinations of solvents (25 different solvents were used). Three separate classification experiments were based on this dataset. In each case the task is to identify a specific chlorinated solvent. As can be seen from the results of Table 3, these experiments focussed on only two pre-processing techniques: normalisation (ND) is used as the baseline method for comparison and first derivative with normalisation (FND) is used as it produced the best results on the acetaminophen dataset. This table directly compares the performance of each machine learning and pre-processing combination without PCA against the same combination with PCA. Again, for many of the machine learning methods, the use of PCA appears to improve performance. However, two major exceptions stand out: C4.5 and RIPPER, both of which are forms of a rule-learning algorithm. Both of these methods suffer a notable loss of accuracy when PCA is employed. This is in contrast with the results on acetaminophen, in which C4.5 and RIPPER gained a small improvement with PCA on the ND dataset, and achieved identical accuracy (to when no PCA was used) on the FND dataset. A comparison of the non-PCA results with those obtained with PCA shows no significant difference. However, if the results of these rule-based algorithms are omitted, a significant difference is observed

Table 3

Comparison of machine learning with and without PCA on chlorinated dataset: percentage classification error (N = No PCA, Y = PCA used)

Method	Dichloromethane				Trichloroethane				Chloroform			
	ND		FND		ND		FND		ND		FND	
	N	Y	N	Y	N	Y	N	Y	N	Y	N	Y
LSVM	1.74	0.43	1.74	2.17	5.65	2.61	6.09	2.61	3.91	1.74	5.22	4.78
RBF	0.43	0.43	0.87	1.74	5.22	2.61	6.09	2.61	4.35	3.91	5.22	4.35
<i>k</i> -NN	8.26	9.13	10.43	9.57	16.09	13.35	13.48	11.74	23.91	19.13	20.00	20.00
C4.5	3.04	8.26	0.43	8.26	7.39	16.09	3.91	16.52	3.91	14.78	3.04	16.96
RIP	6.52	14.78	0.43	12.17	11.30	18.70	6.09	13.04	3.04	18.70	3.04	16.09
NB	43.04	41.30	37.83	26.09	53.48	49.13	40.87	34.35	56.09	51.74	40.00	35.22
Reg.	10.87	10.00	13.04	18.70	18.70	16.96	26.52	26.52	13.91	12.17	25.22	18.70

that confirms the results achieved on the acetaminophen dataset.

To determine the cause of the drop in performance of C4.5, an analysis was carried out on the decision trees produced by C4.5 when trained on the normalised Chloroform dataset. When the original dataset is used, C4.5 generates a tree of size 11. When the first 27 PCs (this number resulted in the best performance) scores are used as input, C4.5 generates a much more complex tree of size 35. Furthermore, the main branch of this tree is based on PC24 and many samples are classified at a leaf based on PC26. A key point is that PCs are ordered according to their contribution to the total variance; PCs 24 and 26 account for very little (less than 0.2%) of the total variance in the scores data. Any model that assigns a strong weighting to these attributes is in danger of overfitting to the training data and could therefore exhibit poor generalisation ability. A similar comparison of the non-PCA and PCA trees produced from the acetaminophen dataset shows that a size difference exists, but is not as great: the tree based on original data has size 7 and the tree based on PC scores data has size 13. Of more importance is the fact that, for the acetaminophen dataset, the tree based on PC scores selected PC3 and PC2 as key attributes; these attributes account for a much greater percentage of the total variance (about 38%).

This analysis shows that the performance of C4.5 may be adversely affected by the use of PC transformed data when compared with its performance on the original data. This occurs when key nodes of the tree are based on PC scores of low variance. Apart from abandoning PCA for decision trees altogether, one alternative is to use the original data and PC scores combined, thus allowing C4.5 to select both from the original set of attributes and from the linear combination attributes. Popelinsky and Brazdil [17] found this approach of adding PC attributes rather than replacing the original attributes to give better results. (They do not report the differences, however.) They found what they described as modest gains in the use of additional PC scores to the dataset when the C5.0 decision tree (a later commercial version of C4.5) was used. We tested this approach on the normalised versions of the spectral datasets with C4.5. In three of the classification tasks, the error achieved was identical to that achieved without PCA; a

minor improvement was found for the trichloroethane dataset. One drawback with this approach is that it increases the dimensionality of the data instead of reducing it, which is one of the main motivations for employing PCA.

6. Related research

The work presented here extends previous research carried out by the authors into the use of machine learning methods with various pre-processing techniques for the classification of spectral data [8]. That work is extended in this paper by using these machine learning methods in combination with the NIPALS PCA technique, and investigating the effect of different numbers of principal components on classification accuracy. The most closely related research to this work can be found in Sigurdsson et al. [18], where they report on the use of neural networks for the detection of skin cancer based on Raman data that has been reduced using PCA. They achieve PCA using singular value decomposition (SVD), a method which calculates *all* the eigenvectors of the data matrix, unlike the NIPALS method that was used here. In addition, they do not present any comparison with neural networks on the raw data without the PCA step.

As far as the authors are aware, few studies have been carried out that investigate the effect of using PCA with a number of machine learning algorithms. Popelinsky [19] does analyse the effect of PCA (again, eigenvector decomposition is used) on three different machine learning algorithms (Naive Bayes, C5.0 and an instance-based learner). In this paper, the principal component scores are added to the original attribute data and he has found this to result in a decrease in error rate for all methods on a significant number of the datasets. However, the experiments were not based on particularly high-dimensional datasets. It is also worth noting that there does not appear to be much evidence of the use of NIPALS PCA in conjunction with machine learning for the classification of high-dimensional data.

7. Conclusions

This paper has proposed the use of an efficient PCA method, NIPALS, to improve the performance of some

well-known machine learning methods in the classification of high-dimensional spectral data. Experiments in the classification of Raman spectra have shown that, overall, this PCA method improves the performance of machine learning when dealing with such high-dimensional data. Furthermore, through the use of PCA, these low errors were achieved despite a major reduction of the data; from the original 1646 attributes of the acetaminophen dataset to at least six attributes. Additional experiments have shown that it is not necessary to generate more than 20 PCs to find an optimal set for the spectral dataset used, as the performance of the majority of classifiers degrades with increasing numbers of PCs. This fact makes NIPALS PCA particularly suited to the proposed approach, as it does not require the generation of all PCs of a data matrix, unlike the widely used eigenvector decomposition methods. This paper has also shown that the pre-processing technique of first derivative followed by normalisation improves the performance of the majority of these machine learning methods in the identification of acetaminophen. Further experiments on the chlorinated dataset confirmed the benefits of using PCA, but also highlighted that poor results can be achieved when PCA is used in combination with rule-based learners, such as C4.5 and RIPPER.

Overall, the use of NIPALS PCA in combination with machine learning appears to be a promising approach for the classification of high-dimensional spectral data. This approach has potential in other domains involving high-dimensional data, such as gene-expression data and image data. Future work will involve testing this approach on more spectral datasets and also on other high-dimensional datasets. Further investigations could also be carried out into the automatic selection of parameters for the techniques considered, such as the number of PCs, kernel parameters for SVM and k for k -NN.

Acknowledgements

This research has been funded by Enterprise Ireland's Basic Research Grant Programme. The authors are also grateful to the High Performance Computing Group at NUI Galway, funded under PRTLII and III, for providing access to HPC facilities.

References

[1] S. Peng, Q. Xu, X. Ling, X. Peng, W. Du, L. Chen, Molecular classification of cancer types from microarray data using the

- combination of genetic algorithms and support vector machines, *FEBS Lett.* 555 (2003) 358–362.
- [2] J. Wang, J. Kwok, H. Shen, L. Quan, Data-dependent kernels for small-scale, high-dimensional data classification, in: *Proceedings of the International Joint Conference on Neural Networks*, July 2005.
- [3] T. Joachims, Text categorisation with support vector machines, in: *Proceedings of European Conference on Machine Learning (ECML)*, 1998.
- [4] A. Ryder, Classification of narcotics in solid mixtures using principal component analysis and Raman spectroscopy and chemometric methods, *J. Forensic Sci.* 47 (2002) 275–284.
- [5] B. Bulkin, *The Raman Effect: An Introduction*, John Wiley, New York, 1991.
- [6] J. Conroy, A. Ryder, M. Leger, K. Hennessy, M. Madden, Qualitative and quantitative analysis of chlorinated solvents using Raman spectroscopy and machine learning, in: *Proc. SPIE – International Society of Optical Engineering*, vol. 5826, 2005, pp. 131–142.
- [7] C. Cheng, T. Kirkbride, D. Batchelder, R. Lacey, T. Sheldon, In situ detection and identification of trace explosives by Raman microscopy, *J. Forensic Sci.* 40 (1995) 31–37.
- [8] M. O'Connell, T. Howley, A. Ryder, M. Leger, M. Madden, Classification of a target analyte in solid mixtures using principal component analysis, support vector machines and Raman spectroscopy, in: *Proc. SPIE – International Society of Optical Engineering*, vol. 5826, 2005, pp. 340–350.
- [9] T. Hastie, R. Tibshirani, J. Friedman, *The Elements of Statistical Learning*, Springer, 2001.
- [10] P. Geladi, B. Kowalski, Partial least squares: a tutorial, *Anal. Chem. Acta* 185 (1986) 1–17.
- [11] B. Scholkopf, A. Smola, *Learning with Kernels: Support Vector Machines, Regularization, Optimization, and Beyond*, MIT Press, Cambridge, MA, 2002.
- [12] I. Witten, E. Frank, *Data Mining: Practical Machine Learning Tools and Techniques with Java Implementations*, Morgan Kaufmann, Los Altos, CA, 2000.
- [13] R. Quinlan, Learning Logical Definitions from Relations, *Mach. Learn.* 5 (1990) 239–266.
- [14] W. Cohen, Fast effective rule induction, in: *Proceedings of the 12th International Conference on Machine Learning*, 2002, pp. 115–123.
- [15] A. Savitzky, M. Golay, Smoothing and differentiation of data by simplified least squares procedures, *Anal. Chem.* 36 (1964) 1627–1639.
- [16] C. Nadeau, Y. Bengio, *Advances in Neural Information Processing*, vol. 12, MIT Press, Cambridge, MA, 2000 (Chapter. Inference for generalisation error).
- [17] L. Popelinsky, P. Brazdil, The principal components method as a pre-processing stage for decision tree learning, in: *Proceedings of PKDD Workshop (Data Mining, Decision Support, Meta-learning and ILP)*, 2000.
- [18] S. Sigurdsson, P. Philipsen, L. Hansen, J. Larsen, M. Gniadecka, H. Wulf, Detection of skin cancer by classification of Raman spectra, *IEEE Trans. Biomed. Eng.* 51 (2004) 1784–1793.
- [19] L. Popelinsky, Combining the principal components method with different learning algorithms, in: *Proceedings of ECML/PKDD IDDM Workshop (Integrating Aspects of Data Mining, Decision Support and Meta-Learning)*, 2001.