

Learning distance to subspace for the nearest subspace methods in high-dimensional data classification

Rui Zhu^{a,b,*}, Mingzhi Dong^b, Jing-Hao Xue^b

^a*School of Mathematics, Statistics and Actuarial Sciences, University of Kent,
Canterbury CT2 7FS, UK*

^b*Department of Statistical Science, University College London, London WC1E 6BT, UK*

Abstract

The nearest subspace methods (NSM) are a category of classification methods widely applied to classify high-dimensional data. In this paper, we propose to improve the classification performance of NSM through learning tailored distance metrics from samples to class subspaces. The learned distance metric is termed as ‘learned distance to subspace’ (LD2S). Using LD2S in the classification rule of NSM can make the samples closer to their correct class subspaces while farther away from their wrong class subspaces. In this way, the classification task becomes easier and the classification performance of NSM can be improved. The superior classification performance of using LD2S for NSM is demonstrated on three real-world high-dimensional spectral datasets.

Keywords: Nearest subspace methods (NSM), distance to subspace, distance metric learning, orthogonal distance, score distance

*Corresponding author.

Email addresses: `r.zhu@kent.ac.uk` (Rui Zhu), `mingzhi.dong.13@ucl.ac.uk` (Mingzhi Dong), `jinghao.xue@ucl.ac.uk` (Jing-Hao Xue)

1. Introduction

Classification of high-dimensional data is an important research topic [8, 9, 10, 27, 28]. Subspace-based classification methods have been widely applied to classify high-dimensional data. Face recognition [11, 4, 7], chemometrics [22, 2, 5, 27] and process control in engineering [14, 20, 15, 17] are famous application areas of subspace-based classification methods. In subspace-based classification methods, classes are first modelled by low-dimensional subspaces. Then the test sample is classified using a classification rule that measures the similarities between the test sample and the class subspaces, and the test sample is assigned to its most similar class.

The principal component (PC) subspaces are commonly adopted as the low-dimensional class subspaces. They are believed to be good representations of high-dimensional data, because most variable information in the data is extracted to the leading PCs and the redundant information in the original features is discarded.

Two distances associated with the PC subspaces are usually used in the classification rules: the *squared* orthogonal distance (OD^2) and the *squared* score distance (SD^2). OD^2 measures the squared orthogonal distance between a sample and a PC subspace [28], while SD^2 measures the squared Mahalanobis distance between the projection of a sample onto a PC subspace and the centre of the PC subspace. When the distances are used in the classification rule, the test sample is assigned to the class with the smallest score of the classification rule. In this paper, we term the PC subspace-based classification methods with the classification rule using distances the “nearest subspace methods” (NSM).

26 The nearest subspace classifier (NSC) [11, 25, 4, 3, 13] and soft inde-
27 pendent modelling of class analogy (SIMCA) [22, 2, 5, 18, 16, 12] are two
28 famous examples of NSM. NSC and SIMCA both adopt PC subspace as
29 the low-dimensional class subspace, however, they use different classification
30 rules to classify a test sample. In NSC, OD^2 between the test sample and
31 its projection on a class subspace is used as the classification rule. The test
32 sample is assigned to the class with the smallest OD^2 . In SIMCA, the lin-
33 ear combination of OD^2 and SD^2 is usually used as the classification rule.
34 The test sample is assigned to the class with the smallest score of the linear
35 combination.

36 However, the standard distances OD^2 and SD^2 may not always be able to
37 capture or reflect well the mechanism underlying the semantic similarity or
38 dissimilarity between the sample and the subspace. In fact, this is also the
39 case with other generic distance metrics, such as the Euclidean distance and
40 the Mahalanobis distance. This has led to the proposals of metric learning
41 in the machine learning community, which enables automatic learning of a
42 tailored distance metric from the data available.

43 More specifically, given the PC class subspaces, the distances used in the
44 classification rule play vital roles in classification. Currently, OD^2 and SD^2
45 are the two distances widely used in the classification rule, both of which
46 use predetermined distance metrics: OD^2 uses the Euclidean distance while
47 SD^2 uses the Mahalanobis distance. However, different data usually prefer
48 different distance metrics to reflect different semantic concepts of dissimilar-
49 ity or similarity in the context of problems, and hence adapting the distance
50 metrics to different data can be expected to improve the classification perfor-

51 mance of NSM. On the other hand, distance metric learning methods emerg-
52 ing in the machine learning community provide us a tool to learn tailored
53 distance metrics automatically from data and to improve the classification
54 performance [23, 21, 26, 19, 24].

55 However, the existing distance metric learning methods in the literature
56 aim to improve the classification methods that are based on distances be-
57 tween samples, such as the k -nearest neighbours (k NN) algorithm. Thus the
58 distance metrics that they learned are for the distances between samples.
59 But unfortunately the distance metrics used in NSM measure the distances
60 between samples and class subspaces. This makes those established distance
61 metric learning methods unable to be applied directly to NSM.

62 Therefore in this paper, we propose a distance metric learning method
63 tailored for NSM to improve its classification performance. We first analyse
64 the classification rules of NSM adopted in the literature, and we derive a
65 general formulation for them. We show that the general formulation is based
66 on two parameterisation matrices with different sizes; hence different classi-
67 fication rules of NSM in the literature can be shown actually using different
68 distance metrics within the general formulation.

69 We define this general formulation as the distance metric from a sample
70 to a class subspace, and propose a method of learning distance to subspace,
71 to automatically learn the two parameterisation matrices that define the
72 distance metric. Then, inspired by the distance metric learning strategy,
73 we learn this distance metric based on a set of distance-to-subspace-based
74 similarity/dissimilarity constraints: the samples are similar to their correct
75 class subspaces while are dissimilar from the wrong class subspaces. Using

76 the learned distance as the similarity measure, we aim to make the samples
77 to be closer to their correct class subspaces while be farther away from their
78 wrong class subspaces. We term this distance metric “learned distance to
79 subspace (LD2S)”.

80 The contributions of this paper are summarised as follows.

81 First, we are the first to derive a general formulation for the classification
82 rules of the nearest subspace methods used in literature. Based on the gen-
83 eral formulation, we can design new classification rules, by specifying \mathbf{M}_1^k
84 and \mathbf{M}_2^k . This formulation is a guidance for researchers to design new clas-
85 sification rules for the nearest subspace methods with better classification
86 performance.

87 Second, based on the general formulation, we develop a novel distance
88 metric learning method for the nearest subspace methods. Most of the cur-
89 rent literature of distance metric learning methods are only designed for
90 classification methods based on distances between samples. Here we design
91 a distance metric learning method for methods based on distances between
92 a sample and a subspace. In this paper, we have shown an effective distance
93 metric learning method, LS2D, to classify high-dimensional data.

94 To evaluate the effectiveness of LD2S, we compare the the classification
95 performances of NSC [4], SIMCA [22, 2] and NSM with the classification
96 rule learned from LD2S (NSM-LD2S) using three real-world high-dimensional
97 datasets.

98 **2. Methodology**

99 *2.1. NSM*

100 *2.1.1. PC class subspace*

101 Given the training set of class k ($k = 1, 2$), $\mathbf{X}_k \in \mathbb{R}^{n_k \times p}$, we build the PC
 102 class subspace of the k th class by using the reduced singular value decompo-
 103 sition (SVD):

$$\mathbf{X}_{k(c)} = \mathbf{U}_{q_k} \mathbf{D}_{q_k} \mathbf{V}_{q_k}^T, \quad (1)$$

104 where $\mathbf{X}_{k(c)}$ is the column-centred training set, the rows of $\mathbf{U}_{q_k} \in \mathbb{R}^{n_k \times q_k}$
 105 ($q_k = \text{rank}(\mathbf{X}_{k(c)})$) are the standardised PC scores, $\mathbf{D}_{q_k} \in \mathbb{R}^{q_k \times q_k}$ is a diag-
 106 onal matrix with singular values $d_1 \geq d_2 \geq \dots \geq d_{q_k} \geq 0$ on the diagonal,
 107 and the columns of $\mathbf{V}_{q_k} \in \mathbb{R}^{p \times q_k}$ are the PCs. The PC score is defined as

$$\mathbf{T}_{q_k} = \mathbf{U}_{q_k} \mathbf{D}_{q_k} = \mathbf{X}_{k(c)} \mathbf{V}_{q_k} \in \mathbb{R}^{n_k \times q_k}. \quad (2)$$

108 If we select the first $r_k \leq q_k$ PCs to build the k th class subspace, then

$$\mathbf{X}_{k(c)} = \mathbf{U}_{r_k} \mathbf{D}_{r_k} \mathbf{V}_{r_k}^T + \mathbf{E}_k, \quad (3)$$

109 where $\mathbf{U}_{r_k} \in \mathbb{R}^{n_k \times r_k}$, $\mathbf{D}_{r_k} \in \mathbb{R}^{r_k \times r_k}$, $\mathbf{V}_{r_k} \in \mathbb{R}^{p \times r_k}$, and $\mathbf{E}_k \in \mathbb{R}^{n_k \times p}$ is the
 110 residual matrix when reconstructing the training samples $\mathbf{X}_{k(c)}$ using the
 111 first r_k PCs. The PC subspace spanned by the first r_k PCs is associated
 112 with a unique projection matrix $\mathbf{P}_k = \mathbf{V}_{r_k} \mathbf{V}_{r_k}^T \in \mathbb{R}^{p \times p}$. We denote the PC
 113 subspace for class k as \mathcal{L}_k .

114 Projecting a new sample $\mathbf{x}_{new} \in \mathbb{R}^{1 \times p}$ to the PC class subspace, we could

115 obtain

$$\mathbf{x}_{(c)}^{k,new} = \mathbf{t}^{k,new} \mathbf{V}_{r_k}^T + \mathbf{e}^{k,new}, \quad (4)$$

116 where $\mathbf{x}_{(c)}^{k,new}$ is the centred \mathbf{x}_{new} by the column means of \mathbf{X}_k , $\mathbf{t}^{k,new} \in \mathbb{R}^{1 \times r}$
 117 is the PC score of the new sample, and $\mathbf{e}^{k,new} \in \mathbb{R}^{1 \times p}$ is the residual of
 118 reconstructing the new sample by the PC class subspace.

119 2.1.2. Two distances associated with the PC class subspace

120 Given the PC class subspaces, the new sample \mathbf{x}_{new} is classified using a
 121 classification rule that is based on two distances related the PC class sub-
 122 spaces: the squared orthogonal distance (OD^2) and the squared score dis-
 123 tance (SD^2). In this section, we discuss the calculation and the geometric
 124 intuition of OD^2 and SD^2 .

125 *The squared orthogonal distance.* The squared orthogonal distance from \mathbf{x}_{new}^c
 126 to the subspace of the k th class, OD_k^2 , is defined based on the residual $\mathbf{e}^{k,new}$
 127 in (4):

$$\text{OD}_k^2 = \sum_{j=1}^p (e_j^{k,new})^2 = \mathbf{e}^{k,new} (\mathbf{e}^{k,new})^T, \quad (5)$$

128 which is the squared Frobenius norm of $\mathbf{e}^{k,new}$.

129 Rewriting (4), we have

$$\mathbf{e}^{k,new} = \mathbf{x}_{(c)}^{k,new} - \mathbf{x}_{(c)}^{k,new} \mathbf{P}_k = \mathbf{x}_{(c)}^{k,new} (\mathbf{I}_p - \mathbf{P}_k), \quad (6)$$

130 where \mathbf{I}_p denotes the p -by- p identity matrix. The $\mathbf{e}^{k,new}$ can then be con-
 131 sidered as the difference vector between $\mathbf{x}_{(c)}^{k,new}$ and its projection on \mathcal{L}_k ,
 132 $\mathbf{x}_{(c)}^{k,new} \mathbf{P}_k$. The orthogonal complement of \mathcal{L}_k is \mathcal{L}_k^\perp which has the projection

133 matrix $\mathbf{I}_p - \mathbf{P}_k$. Thus $\mathbf{e}^{k,new}$ is also the projection of $\mathbf{x}_{(c)}^{k,new}$ to the subspace
 134 \mathcal{L}_k^\perp . Since $\mathbf{e}^{k,new}$ is orthogonal to \mathcal{L}_k , the distance based on $\mathbf{e}^{k,new}$ is called
 the orthogonal distance. An illustration of OD_k^2 in a 3-dimensional feature

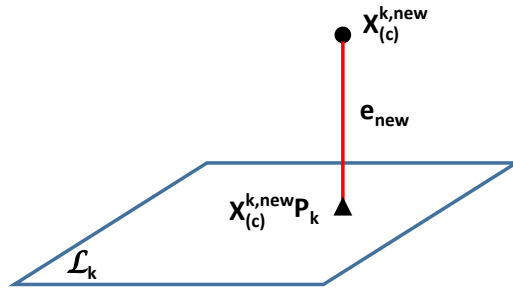


Figure 1: An illustration of OD_k^2 in a 3-dimensional feature space.

135

136 space is shown in Figure 1. The new instance $\mathbf{x}_{(c)}^{k,new}$ is shown as the black
 137 dot; the class subspace \mathcal{L}_k is shown as the dark blue 2-dimensional plane;
 138 and the projection of $\mathbf{x}_{(c)}^{k,new}$ to \mathcal{L}_k , $\mathbf{x}_{(c)}^{k,new} \mathbf{P}_k$, is shown as the black triangle.
 139 The residual $\mathbf{e}^{k,new}$ is represented by the red solid line segment, which is
 140 orthogonal to the plane \mathcal{L}_k . The square of the length of the red line segment
 141 is OD_k^2 .

142 *The squared score distance.* The squared score distance to class k , SD_k^2 , is
 143 defined as the Mahalanobis distance from the projection of $\mathbf{x}_{(c)}^{k,new}$ to the
 144 centre of the subspace \mathcal{L}_k :

$$\text{SD}_k^2 = \sum_{i=1}^{r_k} (t_i^{k,new} / d_i)^2 = \mathbf{t}^{k,new} \mathbf{D}_{r_k}^{-2} (\mathbf{t}^{k,new})^T, \quad (7)$$

145 where \mathbf{D}_{r_k} is the diagonal matrix of singular values in (3). SD_k^2 is the
 146 reweighted squared Frobenius norm of $\mathbf{t}^{k,new}$ with weights $1/d_i$ ($i = 1, 2, \dots, r$)
 and $1/d_1 \leq 1/d_2 \leq \dots \leq 1/d_{r_k}$. An illustration of SD_k^2 in a 3-dimensional

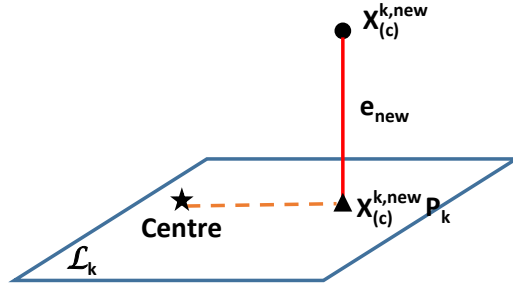


Figure 2: An illustration of SD_k^2 in a 3-dimensional feature space.

147
 148 feature space is shown in Figure 2. In addition to the symbols in Figure 1,
 149 the centre of the class subspace, \mathcal{L}_k , is shown as the black star, and the or-
 150 ange dashed line connects the centre of the class subspace and the projection
 151 of $\mathbf{x}_{(c)}^{k,new}$ to the class subspace. The SD_k^2 is then the reweighted length of the
 152 orange dashed line.

153 2.1.3. The classification rules

154 In NSC, the classification rule is

$$\text{OD}_k^2. \tag{8}$$

155 NSC assigns \mathbf{x}_{new} to the class with the smallest OD_k^2 .

156 In SIMCA, a linear combination of OD_k^2 and SD_k^2 is often used as the
 157 classification rule [2]:

$$\gamma \left(\frac{OD_k}{c_{OD^2}^k} \right)^2 + (1 - \gamma) \left(\frac{SD_k}{c_{SD^2}^k} \right)^2, \quad (9)$$

158 where $\gamma \in [0, 1]$ and $c_{OD^2}^k$ and $c_{SD^2}^k$ are the cutoff values of OD_k^2 and SD_k^2
 159 calculated from the training set of the k th class. When $\gamma = 1$, (9) only
 160 depends on OD_k^2 , and is the same as (8) if the cutoff value $c_{OD^2}^k$ in (9) is one.
 161 When $\gamma = 0$, (9) only depends on SD_k^2 . In practice, the value of γ can be set
 162 by the users based on their prior knowledge of the importance of OD_k^2 and
 163 SD_k^2 , or can be tuned by cross-validation using the training set.

164 *2.2. A general formulation for the classification rules for NSM*

165 Although the classification rules in NSM are in different forms, as shown
 166 in (8) and (9), we shall show that they can be written using the following
 167 general formulation:

$$\mathbf{x}_{(c)}^{k,new} \mathbf{M}_1^k (\mathbf{x}_{(c)}^{k,new})^T - \mathbf{t}^{k,new} \mathbf{M}_2^k (\mathbf{t}^{k,new})^T, \quad (10)$$

168 with different $\mathbf{M}_1^k \in \mathbb{R}^{p \times p}$ and $\mathbf{M}_2^k \in \mathbb{R}^{r_k \times r_k}$. In this section, we derive this
 169 general formulation based on the classification rules (8) and (9), and show
 170 \mathbf{M}_1^k and \mathbf{M}_2^k for (8) and (9), respectively. Based on the derived general
 171 formulation of the classification rules, we will define the distance to subspace
 172 and propose a method to learn the distance to subspace in the next section.

Substituting (6) into (5), we obtain

$$\begin{aligned}
\text{OD}_k^2 &= (\mathbf{x}_{(c)}^{k,new} - \mathbf{x}_{(c)}^{k,new} \mathbf{P}_k)(\mathbf{x}_{(c)}^{k,new} - \mathbf{x}_{(c)}^{k,new} \mathbf{P}_k)^T \\
&= \mathbf{x}_{(c)}^{k,new} (\mathbf{x}_{(c)}^{k,new})^T - 2\mathbf{x}_{(c)}^{k,new} \mathbf{P}_k (\mathbf{x}_{(c)}^{k,new})^T + \mathbf{x}_{(c)}^{k,new} \mathbf{P}_k^2 (\mathbf{x}_{(c)}^{k,new})^T \\
&= \mathbf{x}_{(c)}^{k,new} (\mathbf{x}_{(c)}^{k,new})^T - \mathbf{x}_{(c)}^{k,new} \mathbf{P}_k (\mathbf{x}_{(c)}^{k,new})^T \\
&= \mathbf{x}_{(c)}^{k,new} (\mathbf{x}_{(c)}^{k,new})^T - \mathbf{t}^{k,new} (\mathbf{t}^{k,new})^T, \tag{11}
\end{aligned}$$

173 which indicates that OD_k^2 is the difference between the squared Frobenius
174 norm of $\mathbf{x}_{(c)}^{k,new}$ and the squared Frobenius norm of $\mathbf{t}^{k,new}$. This is intuitive if
175 we think about the right-angled triangle formed by $\mathbf{x}_{(c)}^{k,new}$, $\mathbf{x}_{(c)}^{k,new} \mathbf{P}_k$ and the
176 centre of \mathcal{L}_k in Figure 2.

Then the classification rule (8) can be written as

$$\begin{aligned}
&\mathbf{x}_{(c)}^{k,new} (\mathbf{x}_{(c)}^{k,new})^T - \mathbf{t}^{k,new} (\mathbf{t}^{k,new})^T \\
&= \mathbf{x}_{(c)}^{k,new} \mathbf{M}_{1(NSC)}^k (\mathbf{x}_{(c)}^{k,new})^T - \mathbf{t}^{k,new} \mathbf{M}_{2(NSC)}^k (\mathbf{t}^{k,new})^T, \tag{12}
\end{aligned}$$

177 where $\mathbf{M}_{1(NSC)}^k = \mathbf{I}_p$ and $\mathbf{M}_{2(NSC)}^k = \mathbf{I}_{r_k}$. Equation (12) indicates that
178 the classification rule of NSC provides equal weights to the p dimensions
179 in the linear combination of the original features $\mathbf{x}_{(c)}^{k,new} (\mathbf{x}_{(c)}^{k,new})^T$ and also
180 equal weights to the r_k dimensions in the linear combination of the scores
181 $\mathbf{t}^{k,new} (\mathbf{t}^{k,new})^T$.

Similarly, for the classification rule of SIMCA, we substitute (11) to (9):

$$\begin{aligned}
& \frac{\gamma}{(c_{\text{OD}^2}^k)^2} (\mathbf{x}_{(c)}^{k,\text{new}} (\mathbf{x}_{(c)}^{k,\text{new}})^T - \mathbf{t}^{k,\text{new}} (\mathbf{t}^{k,\text{new}})^T) + \frac{1-\gamma}{(c_{\text{SD}^2}^k)^2} \mathbf{t}^{k,\text{new}} \mathbf{D}_r^{-2} (\mathbf{t}^{k,\text{new}})^T \\
&= \frac{\gamma}{(c_{\text{OD}^2}^k)^2} \mathbf{x}_{(c)}^{k,\text{new}} (\mathbf{x}_{(c)}^{k,\text{new}})^T - \sum_{i=1}^r \left(-\frac{1-\gamma}{(c_{\text{SD}^2}^k)^2} + \frac{\gamma}{(c_{\text{OD}^2}^k)^2 d_i^2} \right) t_i^2 \\
&= \mathbf{x}_{(c)}^{k,\text{new}} \mathbf{M}_{1(S)}^k (\mathbf{x}_{(c)}^{k,\text{new}})^T - \mathbf{t}^{k,\text{new}} \mathbf{M}_{2(S)}^k (\mathbf{t}^{k,\text{new}})^T, \tag{13}
\end{aligned}$$

182 where $\mathbf{M}_{1(S)}^k = \frac{1}{h_1} \mathbf{I}_p$, $h_1 = \frac{\gamma}{(c_{\text{OD}^2}^k)^2}$ and $\mathbf{M}_{2(S)}^k$ is an r_k -by- r_k diagonal matrix
183 with $(-\frac{1-\gamma}{(c_{\text{SD}^2}^k)^2} + \frac{\gamma}{(c_{\text{OD}^2}^k)^2 d_i^2})$ on the diagonals (d_i 's are the singular values in
184 \mathbf{D} with $d_1 \geq d_2 \geq \dots \geq d_{r_k} \geq 0$). Different from the classification rule of
185 NSM in (12), the rule in (13) indicates that the classification rule of SIMCA
186 provides equal weights to the p dimensions in the linear combination of the
187 the original features $\mathbf{x}_{(c)}^{k,\text{new}} (\mathbf{x}_{(c)}^{k,\text{new}})^T$, while providing different weights to the
188 r_k dimensions in the linear combination of the scores $\mathbf{t}^{k,\text{new}} (\mathbf{t}^{k,\text{new}})^T$.

189 2.3. Learning distance to subspace

190 We define the general formulation (10) as the distance from \mathbf{x}_{new} to the
191 k th class subspace. Hence we assign \mathbf{x}_{new} to the nearest class subspace based
192 on the distance to subspace defined in (10).

193 The distance to subspace for the k th class defined in (10) depends on
194 two matrices: \mathbf{M}_1^k and \mathbf{M}_2^k . It can be treated as the difference between two
195 squared distances: $\mathbf{x}_{(c)}^{k,\text{new}} \mathbf{M}_1^k (\mathbf{x}_{(c)}^{k,\text{new}})^T$ is the squared distance from $\mathbf{x}_{(c)}^{k,\text{new}}$
196 to the centre of the class subspace \mathcal{L}_k , and $\mathbf{t}^{k,\text{new}} \mathbf{M}_2^k (\mathbf{t}^{k,\text{new}})^T$ is the squared
197 distance from the projection of $\mathbf{x}_{(c)}^{k,\text{new}}$ to \mathcal{L}_k to the centre of \mathcal{L}_k .

198 The matrices \mathbf{M}_1^k and \mathbf{M}_2^k are of great importance for classification.
199 Instead of determining \mathbf{M}_1^k and \mathbf{M}_2^k manually as in [22] and [2], distance

200 metric learning methods offer us a path to learn more appropriate distance
201 metrics automatically from the training data to improve the classification
202 performance.

203 Distance metric learning methods aim to learn distance metrics based
204 on a set of similarity/dissimilarity constraints: the samples from the same
205 class should be similar while the samples from different classes should be
206 dissimilar. Thus the samples from the same class are close together while the
207 samples from different classes are farther away from each other, based on the
208 distance metric learned from the training data. In this way, the classification
209 task becomes easier and we can expect better classification performance using
210 the learned distance metrics.

211 Established distance metric learning methods are sample-based, i.e. the
212 distances that they learned are measured between samples. However, in
213 NSM, the distance is calculated between a sample and a class subspace. Thus
214 we need to develop a new method of learning the distance metric from sample
215 to subspace, to learn the distance metrics in NSM. The learned distance
216 metrics are termed “learned distance to subspace (LD2S)”. Inspired by the
217 constraints used in established distance metric learning methods, we propose
218 the following set of similarity/dissimilarity constraints for LD2S: the samples
219 should be similar to their true class while dissimilar from the wrong classes.
220 In other words, we aim to learn \mathbf{M}_1^k and \mathbf{M}_2^k , such that the samples are close
221 to their true classes while farther away from the wrong classes.

222 2.3.1. Distance metric

223 In this section, we briefly review the definition of distance metric. Given a
224 set of data points $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$ in $\mathbb{R}^{1 \times p}$ with a set of labels $\{y_1, y_2, \dots, y_N\}$,

225 the distance metric $d(\mathbf{x}_i, \mathbf{x}_j)$ between two data points \mathbf{x}_i and \mathbf{x}_j should satisfy
 226 the following properties:

- 227 1. $d(\mathbf{x}_i, \mathbf{x}_j) \geq 0$ (non-negativity),
- 228 2. $d(\mathbf{x}_i, \mathbf{x}_j) = 0$ if and only if $\mathbf{x}_i = \mathbf{x}_j$ (identity),
- 229 3. $d(\mathbf{x}_i, \mathbf{x}_j) = d(\mathbf{x}_j, \mathbf{x}_i)$ (symmetry),
- 230 4. $d(\mathbf{x}_i, \mathbf{x}_j) \leq d(\mathbf{x}_i, \mathbf{x}_k) + d(\mathbf{x}_j, \mathbf{x}_k)$ (triangle inequality), where \mathbf{x}_k is an
 231 instance that is different to \mathbf{x}_i and \mathbf{x}_j .

232 A distance metric is known as a pseudo metric when the second property
 233 is relaxed to: $d(\mathbf{x}_i, \mathbf{x}_j) = 0$ if $\mathbf{x}_i = \mathbf{x}_j$.

234 Most of the metric learning algorithms aim to learn a Mahalanobis distance-
 235 like pseudo metric:

$$d_M(\mathbf{x}_i, \mathbf{x}_j) = \sqrt{(\mathbf{x}_i - \mathbf{x}_j) \mathbf{M} (\mathbf{x}_i - \mathbf{x}_j)^T}, \quad (14)$$

236 which is parameterised by \mathbf{M} . The matrix \mathbf{M} is set to be positive semidefi-
 237 nite to ensure that $d_M(\mathbf{x}_i, \mathbf{x}_j)$ is a pseudo metric. If \mathbf{M} is the inverse of the
 238 sample variance, then $d_M(\mathbf{x}_i, \mathbf{x}_j)$ is the Mahalanobis distance. If \mathbf{M} is the
 239 identity matrix, then $d_M(\mathbf{x}_i, \mathbf{x}_j)$ is exactly the Euclidean distance.

240 2.3.2. Distance to subspace

241 Different from the distance metric between two samples \mathbf{x}_i and \mathbf{x}_j defined
 242 in (14), we define the squared distance metric between a sample \mathbf{x} and a class
 243 subspace \mathcal{L}_k using the general formulation in (10):

$$d^2(\mathbf{x}, \mathcal{L}_k) = \mathbf{x}_{(c)}^k \mathbf{M}_1^k (\mathbf{x}_{(c)}^k)^T - \mathbf{t}^k \mathbf{M}_2^k (\mathbf{t}^k)^T, \quad (15)$$

244 where $\mathbf{x}_{(c)}^k$ denotes the sample mean-centred by the mean of the training
 245 samples of the k th class, $\mathbf{M}_1^k \in \mathbb{R}^{p \times p}$ is the parameterisation matrix for the
 246 distance in the original feature space of the k th class, \mathbf{t}^k is the PC score of the
 247 sample when projected to the PC subspace of the k th class, and $\mathbf{M}_2^k \in \mathbb{R}^{r_k \times r_k}$
 248 is the parameterisation matrix for the distance in the PC subspace of the k th
 249 class. Then $d^2(\mathbf{x}, \mathcal{L}_k)$ can be treated as the difference between the squared
 250 distance from the sample (column-centred by the column means of class k) to
 251 the centre of \mathcal{L}_k and the squared distance from the projection of the sample
 252 to the centre of \mathcal{L}_k .

253 2.3.3. Learned distance to subspace

254 To learn good distance metrics between samples and class subspaces, we
 255 propose the following similarity/dissimilarity constraints: the samples are
 256 similar to their correct class subspaces while are dissimilar to the wrong
 257 class subspaces. To formulate the constraints, we define the following simi-
 258 larity/dissimilarity sets:

$$259 \quad \mathbf{S} = \{(\mathbf{x}_i, \mathcal{L}_k) \mid \mathbf{x}_i \text{ belongs to class } k\}, \text{ and}$$

$$260 \quad \mathbf{D} = \{(\mathbf{x}_i, \mathcal{L}_k) \mid \mathbf{x}_i \text{ does not belong to class } k\}.$$

261 In the following part, the training samples from class 1 are denoted by
 262 subscript $1(i)$, i.e. $\mathbf{x}_{1(i)} \in \mathbb{R}^{1 \times p}$ and $\mathbf{X}_1 = [\mathbf{x}_{1(1)}^T, \dots, \mathbf{x}_{1(n_1)}^T]^T \in \mathbb{R}^{n_1 \times p}$, and the
 263 training samples from class 2 are denoted by subscript $2(j)$, i.e. $\mathbf{x}_{2(j)} \in \mathbb{R}^{1 \times p}$
 264 and $\mathbf{X}_2 = [\mathbf{x}_{2(1)}^T, \dots, \mathbf{x}_{2(n_2)}^T]^T \in \mathbb{R}^{n_2 \times p}$. Thus the similarity/dissimilarity sets
 265 become

$$266 \quad \mathbf{S} = \{(\mathbf{x}_{1(i)}, \mathcal{L}_1), (\mathbf{x}_{2(j)}, \mathcal{L}_2) \mid i = 1, 2, \dots, n_1, j = 1, 2, \dots, n_2\}, \text{ and}$$

$$267 \quad \mathbf{D} = \{(\mathbf{x}_{1(i)}, \mathcal{L}_2), (\mathbf{x}_{2(j)}, \mathcal{L}_1) \mid i = 1, 2, \dots, n_1, j = 1, 2, \dots, n_2\}.$$

One straightforward way to find tailored distance metrics is to minimise

the sum of the distances between the samples and the class subspaces that fall into the similarity set \mathbf{S} , while maximise the sum of those that fall into the dissimilarity set \mathbf{D} . However, simply optimising the sums of the distances suffers from losing the information in individual samples. Hence, instead of treating all training samples together, we aim to make the difference between the distance to the wrong class and the distance to the correct class large enough for each training sample by using the following constraints:

$$\begin{aligned} d^2(\mathbf{x}_{1(i)}, \mathcal{L}_2) - d^2(\mathbf{x}_{1(i)}, \mathcal{L}_1) &\geq 1, \text{ for } i = 1, \dots, n_1, \text{ and} \\ d^2(\mathbf{x}_{2(j)}, \mathcal{L}_1) - d^2(\mathbf{x}_{2(j)}, \mathcal{L}_2) &\geq 1, \text{ for } j = 1, \dots, n_2. \end{aligned} \quad (16)$$

In this way, the samples can be classified more easily. In addition, to enhance the generalisation ability of the learned distance metrics, we add slack variables $\xi_{1(i)}$ and $\xi_{2(j)}$ to the constraints and aim to solve the following optimisation problem:

$$\min_{\xi_{1(i)}, \xi_{2(j)}, \mathbf{M}_1^k, \mathbf{M}_2^k} \sum_{i=1}^{n_1} \xi_{1(i)} + \sum_{j=1}^{n_2} \xi_{2(j)} \quad (17)$$

$$\text{s.t. } d^2(\mathbf{x}_{1(i)}, \mathcal{L}_2) - d^2(\mathbf{x}_{1(i)}, \mathcal{L}_1) \geq 1 - \xi_{1(i)}, \quad \xi_{1(i)} \geq 0, \quad (18)$$

$$d^2(\mathbf{x}_{2(j)}, \mathcal{L}_1) - d^2(\mathbf{x}_{2(j)}, \mathcal{L}_2) \geq 1 - \xi_{2(j)}, \quad \xi_{2(j)} \geq 0, \quad (19)$$

$$\mathbf{M}_1^k \succeq 0 \text{ and } \mathbf{M}_2^k \succeq 0, \quad (20)$$

where $\mathbf{M}_1^k \succeq 0$ and $\mathbf{M}_2^k \succeq 0$ denote that \mathbf{M}_1^k and \mathbf{M}_2^k are positive semidefi-

nite. The constraints in (18) and (19) can be rewritten as

$$\begin{aligned}\xi_{1(i)} &\geq [1 + d^2(\mathbf{x}_{1(i)}, \mathcal{L}_1) - d^2(\mathbf{x}_{1(i)}, \mathcal{L}_2)]_+ \text{ and} \\ \xi_{2(j)} &\geq [1 + d^2(\mathbf{x}_{2(j)}, \mathcal{L}_2) - d^2(\mathbf{x}_{2(j)}, \mathcal{L}_1)]_+, \end{aligned}$$

where $[l]_+ = \max(0, l)$. Hence the optimisation problem is equivalent to

$$\begin{aligned} \min_{\mathbf{M}_1^k, \mathbf{M}_2^k} & \sum_{i=1}^{n_1} [1 + d^2(\mathbf{x}_{1(i)}, \mathcal{L}_1) - d^2(\mathbf{x}_{1(i)}, \mathcal{L}_2)]_+ + \\ & \sum_{j=1}^{n_2} [1 + d^2(\mathbf{x}_{2(j)}, \mathcal{L}_2) - d^2(\mathbf{x}_{2(j)}, \mathcal{L}_1)]_+ \\ \text{s.t. } & \mathbf{M}_1^k \succeq 0, \quad \mathbf{M}_2^k \succeq 0. \end{aligned} \quad (21)$$

268 The hinge losses used in (21) only penalise the samples that do not satisfy
 269 (16), while assign zero loss for the samples that satisfy (16) using NSM.
 270 In this way, the hinge loss makes full use of the effectiveness of NSM. It
 271 is worth noting that the hinge loss has also been popularly used in other
 272 distance-based classifiers, such as support vector machine (SVM) and large
 273 margin nearest neighbour (LMNN) classification [21].

274 Suppose \mathbf{M}_1^{k*} and \mathbf{M}_2^{k*} ($k = 1, 2$) denote the solutions of (21). Then the
 275 learned distance from a test sample \mathbf{x}_{new} to the k th class subspace is

$$d^2(\mathbf{x}_{new}, \mathcal{L}_k) = \mathbf{x}_{(c)}^{k,new} \mathbf{M}_1^{k*} (\mathbf{x}_{(c)}^{k,new})^T - \mathbf{t}^{k,new} \mathbf{M}_2^{k*} (\mathbf{t}^{k,new})^T. \quad (22)$$

276 We compare $d^2(\mathbf{x}_{new}, \mathcal{L}_1)$ and $d^2(\mathbf{x}_{new}, \mathcal{L}_2)$, and assign \mathbf{x}_{new} to the class with
 277 the smallest squared distance.

278 Considering the nature of spectral data, i.e. high-dimensional feature and

279 small sample size, learning the full matrices, \mathbf{M}_1^k with $p(p+1)/2$ parameters
 280 and \mathbf{M}_2^k with $r_k(r_k+1)/2$ parameters, could easily suffer from the overfitting
 281 problem. In (12) and (13), $\mathbf{M}_{1(NSC)}^k = \mathbf{I}_p$ and $\mathbf{M}_{1(S)}^k = \frac{1}{h_1}\mathbf{I}_p$ are identity
 282 matrices with common coefficients 1 and $1/h_1$ for all dimensions, respectively.
 283 Therefore, in this paper, we learn $\mathbf{M}_1^k = c_k\mathbf{I}_p$ (with $c_k \geq 0$) and $\mathbf{M}_2^k =$
 284 $\text{diag}(m_{21}^k, m_{22}^k, \dots, m_{2r_k}^k)$ (with each element nonnegative), as natural and
 285 practically-interpretable extensions of those used in (12) and (13).

286 3. Experiments

287 In the following experiments, NSC, SIMCA and NSM with distance mea-
 288 surement (22) (NSM-LD2S) are compared using high-dimensional spectral
 289 data, the Phenyl dataset, the fat dataset [6] and the meat dataset [1]. We
 290 also compare the classification results of the nearest subspace methods with
 291 those of the naive Bayes classifier (NB), the k -nearest neighbours algorithm
 292 (k NN) and the support vector machine (SVM), to show the effectiveness of
 293 the nearest subspace methods to classify high-dimensional data.

294 3.1. Datasets

295 The number of samples in each class and the number of features for the
 296 three high-dimensional spectral datasets are summarised in Table 1.

Table 1: The number of samples in each class, n_1 and n_2 , and the number of features p for the three high-dimensional spectral datasets.

	n_1	n_2	p
Phenyl	300	300	658
Fat	122	71	100
Meat	54	55	1050

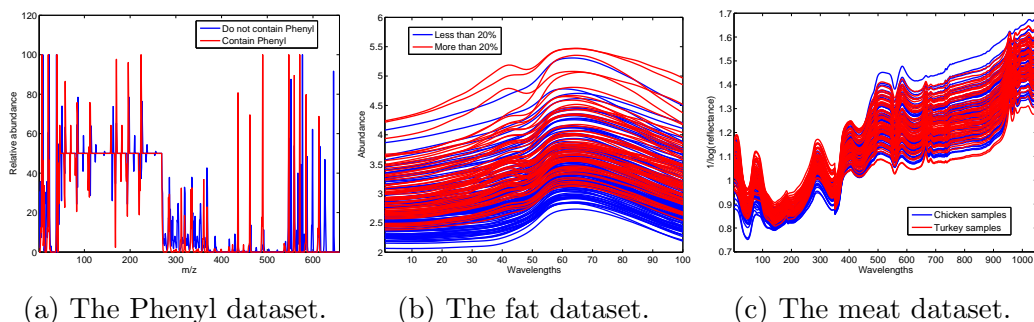


Figure 3: The plots of the spectra of the three datasets.

297 *3.1.1. The Phenyl dataset*

298 The Phenyl dataset is available in the ‘chemometrics’ R package, which
 299 contains 300 spectra with the phenyl substructure and 300 spectra without
 300 the phenyl substructure. The spectra are measured at 658 wavelengths. To
 301 avoid confusing, the spectra of two instances from two classes are shown in
 302 Figure 3a.

303 *3.1.2. The fat dataset*

304 The fat dataset contains 193 spectra of finely chopped meat, measured at
 305 100 wavelengths [6]. The fat dataset consists of 122 spectra of meat samples
 306 with less than 20% fat and 71 spectra of meat samples with more than 20%
 307 fat. The spectra of all samples are shown in Figure 3b.

308 *3.1.3. The meat dataset*

309 The meat dataset [1] contains the spectra of five classes of meat sam-
 310 ples, measured at 1050 wavelengths. We select the chicken and turkey meat
 311 samples from the original dataset in the experiments, because they contain
 312 similar chemical components and are hard to classify. The new meat dataset

313 contains the spectra of 55 chicken samples and the spectra of 54 turkey sam-
314 ples. The spectral of all samples are shown in Figure 3c.

315 3.2. Experiment settings

316 The classification performances of the three methods are shown for five
317 different ratios of training set size/feature dimension: $n_1/p = n_2/p = 0.1$,
318 0.2, 0.3, 0.4 and 0.5.

319 For the Phenyl dataset, we randomly select 100 samples with Phenyl
320 structure and 100 samples without Phenyl structure. For illustrative pur-
321 poses, we select the first 100 dimensions from the 658 feature dimensions for
322 the experiments in this paper, i.e. $p = 100$.

323 For the fat dataset, we use all the 120 meat samples with less than 20%
324 fat and 71 meat samples with more than 20% fat in the dataset. We also use
325 all the dimensions of the fat dataset, i.e. $p = 100$.

326 For the meat dataset, we use all the 55 chicken samples and 54 turkey
327 samples in the dataset. Again for illustrative purposes, we also select the first
328 100 dimensions from the 350 dimensions for the experiments in this paper,
329 i.e. $p = 100$.

330 Therefore, as $p = 100$ for each of the three datasets, the five training set
331 sizes are $n_1 = n_2 = 10, 20, 30, 40$ and 50. The samples to form a training
332 set are randomly selected from a dataset. The rest samples in the datasets
333 are used as test samples.

334 In NSC, SIMCA and NSM-LD2S, the numbers of PCs, r_k , are tuned by
335 5-fold cross-validation using the training set to minimise the classification
336 error. More specifically, for each value of r_k , we calculate the mean classi-
337 fication error of the 5-fold cross-validation. The value with the minimum

338 mean classification error is chosen as the number of PCs.

339 In SIMCA, $c_{OD}^k = (\hat{\mu} + \hat{\sigma}z_{0.975})^{3/2}$, where $\hat{\mu}$ and $\hat{\sigma}$ are the mean and the
340 standard deviation of the orthogonal distances in of the training samples in
341 class k ; and $c_{SD}^k = \sqrt{\chi_{n_k;0.975}^2}$. The weight γ is also tuned by 5-fold cross-
342 validation using the training data.

343 In NSM-LD2S, the optimisation problem (21) is solved by ‘cvx’ in MAT-
344 LAB.

345 In SVM, the radial basis function (RBF) kernel is adopted. The scale
346 parameter of the RBF kernel and the penalty factor C are tune by 5-fold
347 cross-validation. The values of the two parameters to be chosen are set to
348 10, 10^2 and 10^3 . In k NN, the number of the nearest neighbours is tuned by
349 5-fold cross-validation. The values to be chosen are set to 3, 5 and 7. In NB,
350 the prior probability of each class is set as the proportion of the number of
351 training samples of that class over the total number of training samples.

352 All the random training/test splits and the subsequent experiments are
353 repeated 100 times and the classification accuracies of the test data are
354 recorded.

355 3.3. Results

356 3.3.1. The Phenyl dataset

357 The classification results of the Phenyl dataset demonstrate the superior
358 classification performance of NSM-LD2S, as shown in Figure 4 and Figure 5,
359 compared with NSC and SIMCA over all n_k/p ratios. It is clear that SVM
360 performs better than the three nearest subspace methods for this dataset.
361 k NN and NB are also better than the three nearest subspace methods when
362 n_k/p becomes large.

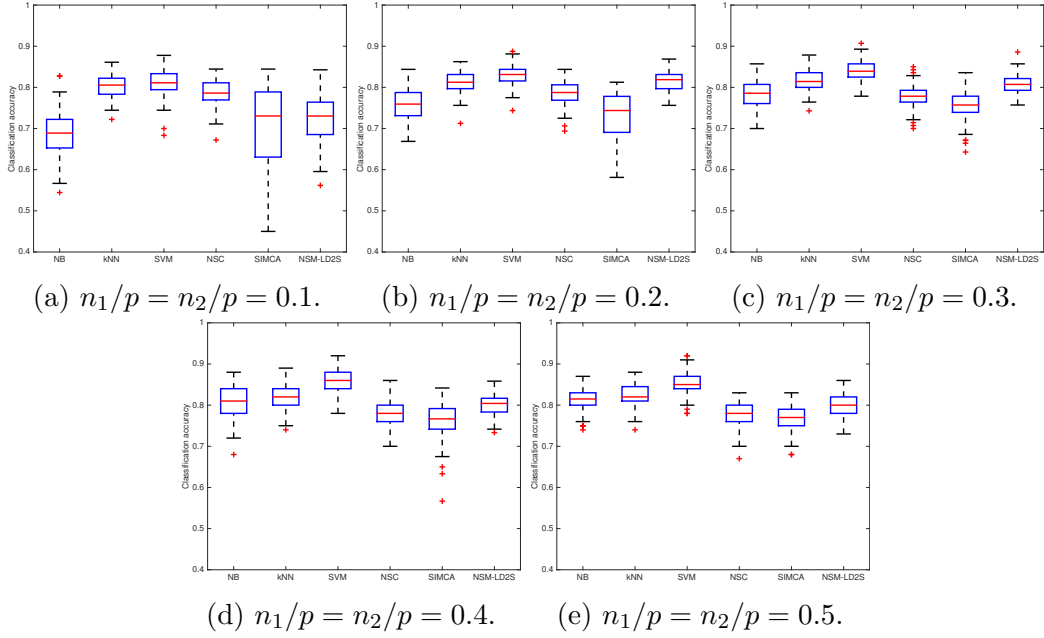


Figure 4: Classification accuracies of NB, k NN, SVM, NSC, SIMCA and NSM-LD2S for the Phenyl dataset.

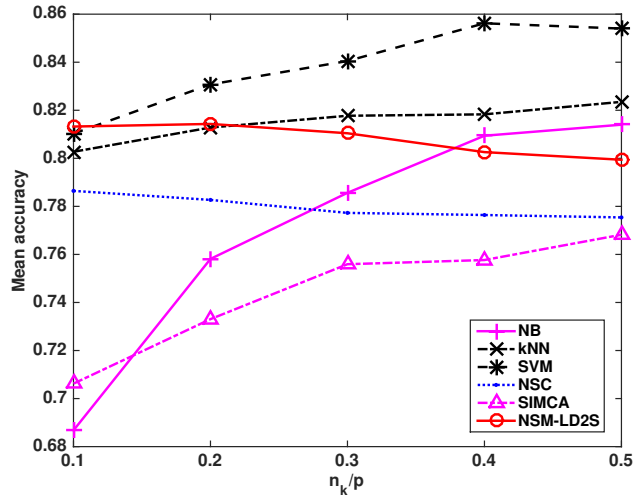


Figure 5: Mean classification accuracies of NB, k NN, SVM, NSC, SIMCA and NSM-LD2S for the Phenyl dataset.

363 However, it is conceivable that, for certain other datasets, the classifica-
364 tion performance of NSM-LD2S cannot always be better than those of NSC
365 and SIMCA, in particular under small n_k/p ratios. In the following two
366 sections, we show two examples that NSM-LD2S performs worse than NSC
367 and SIMCA for small n_k/p ratios but better for large n_k/p ratios. This is
368 because there are more parameters in NSM-LD2S to be learned than in NSC
369 and SIMCA, and NSM-LD2S needs more training samples to achieve good
370 classification performance for some data. In addition, the classification per-
371 formances of NB, k NN and SVM are also not always better than the nearest
372 subspace methods. The following two examples can also demonstrate this
373 argument.

374 3.3.2. *The fat dataset*

375 In the fat dataset, the classification performance of NSM-LD2S and SIMCA
376 are worse than NSC when $n_k/p = 0.1$ and are better than NSC when
377 $n_k/p \geq 0.2$, as shown in Figure 6 and Figure 7. NSM-LD2S provides the
378 best classification performance when $n_k/p \geq 0.2$.

379 It is obvious that NB has the worst mean classification accuracies for all
380 n_k/p ratios. k NN performs similarly to NSM-LD2S. SVM performs similarly
381 to SIMCA when $n_k/p = 0.1$ and performs worse than the three nearest
382 subspace methods for all other n_k/p ratios.

383 3.3.3. *The meat dataset*

384 Compared with the fat dataset, the classification accuracies of the three
385 methods for the meat dataset show a stronger effect of the n_k/p ratios. When
386 $n_k/p < 0.4$, NSM-LD2S performs much worse than NSC and SIMCA, espe-

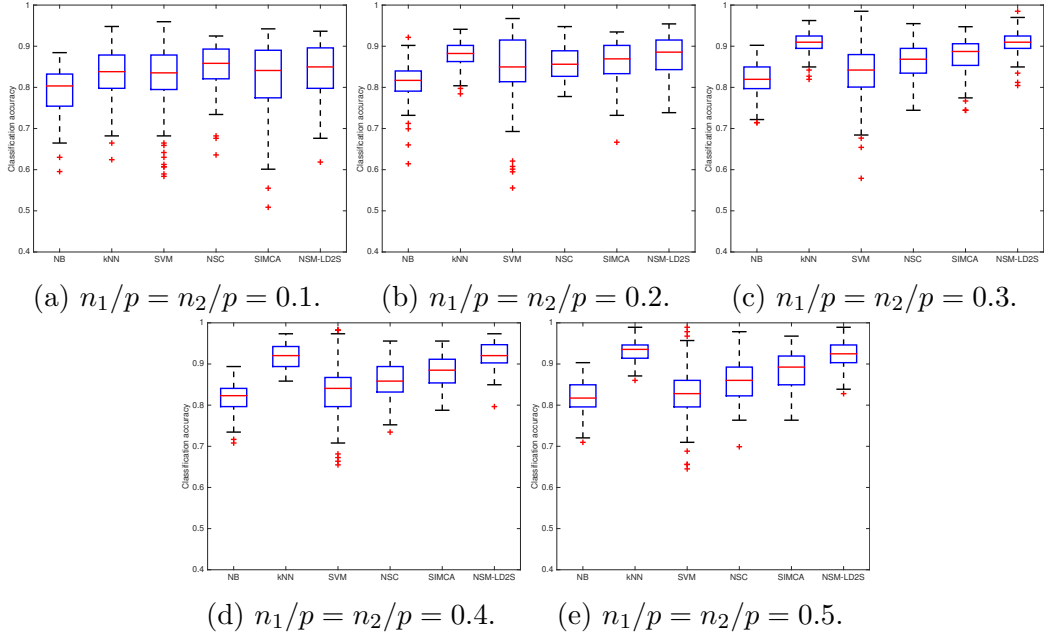


Figure 6: Classification accuracies of NB, k NN, SVM, NSC, SIMCA and NSM-LD2S for the fat dataset.

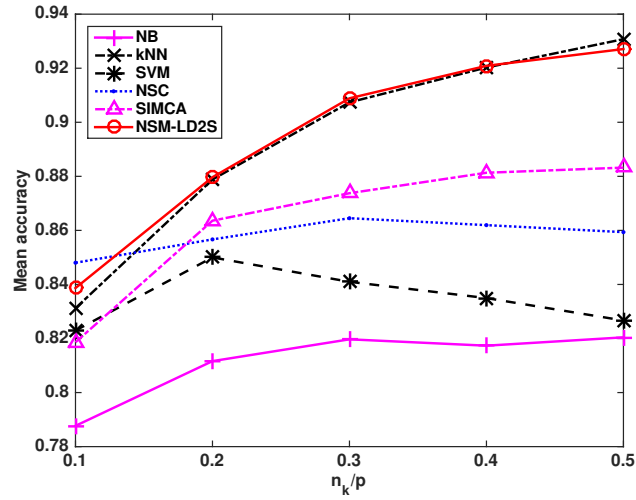


Figure 7: Mean classification accuracies of NB, k NN, SVM, NSC, SIMCA and NSM-LD2S for the fat dataset.

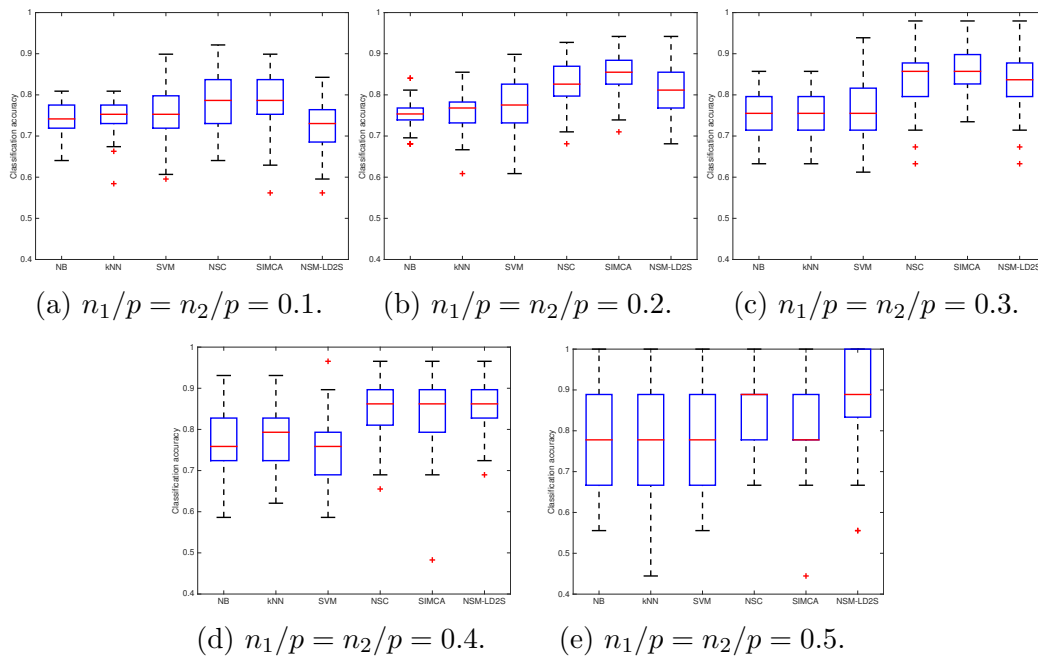


Figure 8: Classification accuracies of NB, k NN, SVM, NSC, SIMCA and NSM-LD2S for the meat dataset.

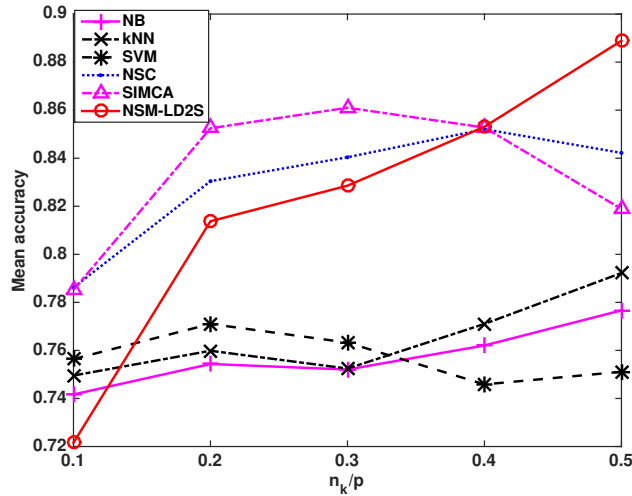


Figure 9: Mean classification accuracies of NB, k NN, SVM, NSC, SIMCA and NSM-LD2S for the meat dataset.

387 cially for $n_k/p = 0.1$. However, when $n_k/p = 0.5$, the classification accuracies
 388 of NSM-LD2S become much better than those of NSC and SIMCA, as shown
 389 in Figure 8(e) and Figure 9. The classification results of the meat dataset
 390 suggest that NSM-LD2S needs $n_k/p > 0.4$ to achieve superior classification
 391 performance for the meat dataset.

392 Similarly to the fat dataset, NB and SVM have the worst classification
 393 performances for $n_k/p > 0.1$ for the meat dataset. k NN performs worse than
 394 the nearest subspace methods for the meat dataset.

395 3.3.4. Summary of the results

396 The experiments show that using the learned distance metrics from data
 397 can provide superior classification results, compared with using predeter-
 398 mined distance metrics, when the n_k/p ratio is large enough. For data with
 399 small n_k/p ratios, using the distance measurement based on LD2S may per-
 400 form poorly in classification since the n_k/p ratio is not large enough to learn

401 all the parameters in LD2S.

402 It is worth noting that the nearest subspace methods are effective to clas-
403 sify high-dimensional data. An important reason for this is that they find a
404 low-dimensional subspace representation for each class to extract the most
405 informative features. Our proposed LD2S is an additional step to further im-
406 prove the classification performance of the nearest subspace methods, based
407 on the features-extracted data. LD2S can obtain better distance measure-
408 ments between a sample and a subspace, which imposes a positive effect on
409 classification accuracies. As demonstrated by the experiment results, NSM-
410 LD2S can achieve better classification accuracies than NSC and SIMCA,
411 which shows the effectiveness of LD2S in addition to feature extraction in
412 NSC and SIMCA.

413 **4. Conclusion**

414 We have proposed a general formulation of distance to subspace, i.e. the
415 distance from a sample to a PC class subspace. Based on this formulation,
416 we have proposed a simple but effective LD2S method that can learn tailored
417 distance metrics adaptively from data, for the classification rule of NSM. The
418 classification performances on three datasets demonstrate the effectiveness of
419 learning distance metrics from data when the n_k/p ratio is large enough. The
420 current LD2S is designed for binary classification. A multi-class version of
421 LD2S is needed for more general and practical cases and we identify this as
422 our future work.

423 **Acknowledgement**

424 The authors thank the reviewers for their constructive comments.

425 **References**

- 426 [1] T. Arnalds, J. McElhinney, T. Fearn, G. Downey, A hierarchical dis-
427 criminant analysis for species identification in raw meat by visible and
428 near infrared spectroscopy, *Journal of Near Infrared Spectroscopy* 12 (3)
429 (2004) 183–188.
- 430 [2] K. V. Branden, M. Hubert, Robust classification in high dimensions
431 based on the SIMCA method, *Chemometrics and Intelligent Laboratory*
432 *Systems* 79 (1) (2005) 10–21.
- 433 [3] Y. Chi, Nearest subspace classification with missing data, in: *Signals,*
434 *Systems and Computers, 2013 Asilomar Conference on, IEEE, 2013, pp.*
435 *1667–1671.*
- 436 [4] Y. Chi, F. Porikli, Connecting the dots in multi-class classification:
437 From nearest subspace to collaborative representation, in: *Computer*
438 *Vision and Pattern Recognition (CVPR), 2012 IEEE Conference on,*
439 *IEEE, 2012, pp. 3602–3609.*
- 440 [5] C. Durante, R. Bro, M. Cocchi, A classification tool for N-way array
441 based on SIMCA methodology, *Chemometrics and Intelligent Labora-*
442 *tory Systems* 106 (1) (2011) 73–85.
- 443 [6] F. Ferraty, P. Vieu, *Nonparametric Functional Data Analysis: Theory*
444 *and Practice, Springer Science & Business Media, 2006.*

- 445 [7] K. Fukui, A. Maki, Difference subspace and its generalization for
446 subspace-based methods, *IEEE Transactions on Pattern Analysis and*
447 *Machine Intelligence* 37 (11) (2015) 2164–2177.
- 448 [8] P. Hall, D. M. Titterington, J.-H. Xue, Median-based classifiers for
449 high-dimensional data, *Journal of the American Statistical Association*
450 104 (488) (2009) 1597–1608.
- 451 [9] P. Hall, J.-H. Xue, Incorporating prior probabilities into high-
452 dimensional classifiers, *Biometrika* 97 (1) (2010) 31–48.
- 453 [10] P. Hall, J.-H. Xue, On selecting interacting features from high-
454 dimensional data, *Computational Statistics & Data Analysis* 71 (2014)
455 694–708.
- 456 [11] K.-C. Lee, J. Ho, D. J. Kriegman, Acquiring linear subspaces for face
457 recognition under variable lighting, *IEEE Transactions on Pattern Anal-*
458 *ysis and Machine Intelligence* 27 (5) (2005) 684–698.
- 459 [12] C. Mees, F. Souard, C. Delporte, E. Deconinck, P. Stoffelen, C. Stévigny,
460 J.-M. Kauffmann, K. De Braekeleer, Identification of coffee leaves using
461 FT-NIR spectroscopy and SIMCA, *Talanta* 177 (2018) 4–11.
- 462 [13] J.-X. Mi, D.-S. Huang, B. Wang, X. Zhu, The nearest-farthest subspace
463 classification for face recognition, *Neurocomputing* 113 (2013) 241–250.
- 464 [14] B. Mnassri, B. Ananou, M. Ouladsine, et al., Fault detection and di-
465 agnosis based on PCA and a new contribution plot, *IFAC Proceedings*
466 *Volumes* 42 (8) (2009) 834–839.

- 467 [15] B. Mnassri, M. Ouladsine, et al., Reconstruction-based contribution ap-
468 proaches for improved fault diagnosis using principal component analy-
469 sis, *Journal of Process Control* 33 (2015) 60–76.
- 470 [16] I. Nejadgholi, M. Bolic, A comparative study of PCA, SIMCA and Cole
471 model for classification of bioimpedance spectroscopy measurements,
472 *Computers in biology and medicine* 63 (2015) 42–51.
- 473 [17] M. Rafferty, X. Liu, D. M. Lavery, S. McLoone, Real-time multi-
474 ple event detection and classification using moving window pca, *IEEE*
475 *Transactions on Smart Grid* 7 (5) (2016) 2537–2548.
- 476 [18] A. Sgarbossa, C. Costa, P. Menesatti, F. Antonucci, F. Pallottino,
477 M. Zanetti, S. Grigolato, R. Cavalli, A multivariate SIMCA index as
478 discriminant in wood pellet quality assessment, *Renewable Energy* 76
479 (2015) 258–263.
- 480 [19] Q. Tian, S. Chen, L. Qiao, Ordinal margin metric learning and its exten-
481 sion for cross-distribution image data, *Information Sciences* 349 (2016)
482 50–64.
- 483 [20] P. Van den Kerkhof, J. Vanlaer, G. Gins, J. F. Van Impe, Analysis
484 of smearing-out in contribution plot based fault isolation for statistical
485 process control, *Chemical Engineering Science* 104 (2013) 285–293.
- 486 [21] K. Q. Weinberger, L. K. Saul, Distance metric learning for large margin
487 nearest neighbor classification, *Journal of Machine Learning Research*
488 10 (2009) 207–244.

- 489 [22] S. Wold, Pattern recognition by means of disjoint principal components
490 models, *Pattern Recognition* 8 (3) (1976) 127–139.
- 491 [23] E. P. Xing, A. Y. Ng, M. I. Jordan, S. Russell, Distance metric learning
492 with application to clustering with side-information, *Advances in Neural
493 Information Processing Systems* 15 (2003) 505–512.
- 494 [24] J. Yu, D. Tao, J. Li, J. Cheng, Semantic preserving distance metric
495 learning and applications, *Information Sciences* 281 (2014) 674–686.
- 496 [25] L. Zhang, W.-D. Zhou, B. Liu, Nonlinear nearest subspace classifier, in:
497 *International Conference on Neural Information Processing*, Springer,
498 2011, pp. 638–645.
- 499 [26] P. Zhu, Q. Hu, W. Zuo, M. Yang, Multi-granularity distance metric
500 learning via neighborhood granule margin maximization, *Information
501 Sciences* 282 (2014) 321–331.
- 502 [27] R. Zhu, K. Fukui, J.-H. Xue, Building a discriminatively ordered sub-
503 space on the generating matrix to classify high-dimensional spectral
504 data, *Information Sciences* 382 (2017) 1–14.
- 505 [28] R. Zhu, J.-H. Xue, On the orthogonal distance to class subspaces for
506 high-dimensional data classification, *Information Sciences* 417 (2017)
507 262–273.