

Automatic Configuration of Spectral Dimensionality Reduction Methods

Michał Lewandowski, Dimitrios Makris and Jean-Christophe Nebel

Digital Imaging Research Centre, Kingston University, KT1 2EE, UK

Abstract

We propose an advanced framework for the automatic configuration of spectral dimensionality reduction methods. This is achieved by introducing, first, the mutual information measure to assess the quality of discovered embedded spaces. Secondly, unsupervised Radial Basis Function network is designated for mapping between spaces where the learning process is derived from graph theory and based on Markov cluster algorithm. Experiments on synthetic and real datasets demonstrate the effectiveness of the proposed methodology.

Keywords: Dimensionality reduction, Locally Linear Embedding, Isomap, Laplacian Eigenmaps, Mutual Information, Radial Basis Function network, Markov Cluster algorithm

1. Introduction

With the exponential increase of data production driven by applications such as the internet, mobile communication, computer vision, medical imaging, speech recognition and genomics, powerful tools are required by scientists to allow the analysis of these data. Since they are usually highly dimensional,

22 dimensionality reduction has become an essential process in the exploration
23 of large volumes of multivariate data.

24 Dimensionality reduction can be defined as the transformation of high-
25 dimensional data, $X = \{x_i\}_{(i=1..N)}$ ($x_i \in R^D$), into a meaningful and compact
26 representation of reduced dimensionality, $Y = \{y_i\}_{(i=1..N)}$ ($y_i \in R^d$) where
27 $d < D$ (and often $d \ll D$), to obtain more informative, descriptive and
28 practical data representation for further analysis. This process is achieved
29 by eliminating redundancies present in datasets while ensuring the maximum
30 possible preservation of information.

31 Since most real datasets are highly nonlinear, many nonlinear dimension-
32 ality reduction techniques have been proposed. They can be classified in
33 two main categories: mapping-based and embedding-based. Mapping-based
34 approaches such as GPLVM (Lawrence, 2004) and generative topographic
35 mapping (Bishop et al., 1998) use probabilistic nonlinear functions to map
36 the embedded space to the data space. Their main limitation comes from
37 the computational cost of their learning process which restricts their usage to
38 relatively small datasets. On the other hand, embedded-based approaches,
39 also called spectral methods, estimate the structure of the data underlying
40 manifold by approximating each data point according to their neighbours on
41 the manifold. Although these methods do not provide any explicit mapping
42 between low and high dimensional spaces, they have proved very popular be-
43 cause they can handle very large high dimensional datasets with a reasonable
44 computational cost.

45 Spectral methods can broadly be divided into three families, i.e. Isometric
46 Feature Mapping (Isomap) (Tenenbaum et al., 2000), Locally Linear Embed-

47 ding (LLE) (Roweis and Saul, 2000) and Laplacian Eigenmaps (LE) (Belkin
48 and Niyogi, 2001), according to the way data point positions are expressed
49 in function of their neighbours. Since they have been a very active area of
50 research, many extensions and improvements have been suggested (Choi and
51 Choi, 2004; De Ridder et al., 2003; De Silva and Tenenbaum, 2003; Donoho
52 and Grimes, 2003; Goldberg and Ritov, 2008; He and Niyogi, 2004; He et al.,
53 2005; Kokiopoulou and Saad, 2007; Wang and Li, 2009; Yang, 2003; Zhang
54 and Wang, 2007; Zheng et al., 2008). Despite this research effort, these ap-
55 proaches still suffer from the fact they rely on a set of values which are chosen
56 empirically, i.e. neighbourhood size and mapping function parameters.

57 In this paper, we address this fundamental problem by proposing two
58 extensions of spectral dimensionality reduction methods allowing their auto-
59 matic configuration. First, optimal values of neighbourhoods are estimated
60 by adopting mutual information measure (MI) (Cover and Thomas, 1991).
61 Secondly, mapping functions are customised to datasets with a novel usage
62 of Radial Basis Function network (RBFN) (Poggio et al., 1990), where net-
63 work topology is automatically learnt by Markov cluster algorithm (MCL)
64 (Dongen, 2000).

65 After a detailed description of the main spectral dimensionality reduction
66 approaches and their limitations, we describe the new techniques we propose
67 to allow their automatic configuration. Finally, they are validated on syn-
68 thetic and real datasets. Since spectral dimensionality reduction methods
69 derives from either Isomap, LLE or LE, our contribution will be applied to
70 these three methods which are used as representatives of all embedded-based
71 approaches.

72 **2. Spectral dimensionality reduction methods and their limitations**

73 Spectral or embedding-based approaches model the structure of data by
74 preserving some geometrical property of the underlying manifold. While
75 the Isomap (Tenenbaum et al., 2000) method attempts to maintain global
76 properties, LE (Belkin and Niyogi, 2001) and LLE (Roweis and Saul, 2000)
77 aim at preserving local geometry which implicitly tends to keep the global
78 layout of the data manifold. After a brief description of these techniques, we
79 list the main limitations we address in this paper.

80 *2.1. Processing pipeline*

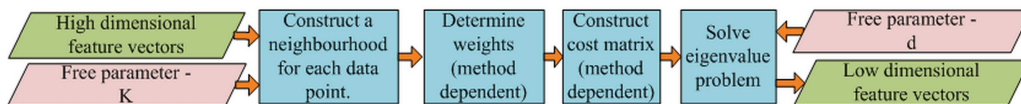


Figure 1: Dimensionality reduction using spectral methods.

81 These methods share the same algorithm structure as illustrated in figure
82 1. First, the neighbourhood for each data point is constructed by choosing
83 K -nearest neighbours based on Euclidean distance. Then, weights, which
84 express the geometrical relationship between each data point and its neigh-
85 bours, are determined according to the property to be preserved. In LLE,
86 they summarize the neighbours contribution to the reconstruction of a data
87 point (Roweis and Saul, 2000). In LE and Isomap, the weights are related
88 to the distance between a point and its neighbours using respectively heat
89 kernel (Belkin and Niyogi, 2001) and Euclidean distance (Tenenbaum et al.,
90 2000). Then, each method optimises its own cost function subject to con-
91 straints that make the problem well-posed. In the case of the LE and Isomap
92 algorithms, the manifold is approximated, first, by an adjacency graph where

93 nodes correspond to data point and edges represent weights between points.
94 A sparse cost matrix is constructed directly for LLE and LE (Belkin and
95 Niyogi, 2001; Roweis and Saul, 2000), whereas the Isomap dense cost matrix
96 is obtained by calculating geodesic distances between all pairs of points in the
97 graph (Tenenbaum et al., 2000). Finally, spectral embedding is calculated
98 using the eigenvectors of the cost matrix.

99 *2.2. Limitations*

100 The main issue of spectral methods is that the quality of embedded space
101 is very sensitive to the choice of free parameters and they do not provide a
102 mapping function between low and high dimensional spaces.

103 All approaches have two free parameters: 'd' and 'K'. 'd' is the dimen-
104 sionality of the embedded space and must be known a priori because it is
105 used in the minimization process. If the number of dimensions is too low,
106 important data features may be collapsed onto the same dimension. 'K' is
107 the neighbourhood size. If it is too small, global feature information is lost
108 since the manifold is split into unconnected pieces. If it is too large, the LE
109 and LLE assumption that a data point and its neighbours are locally linear
110 is not valid. In the case of Isomap, a large K introduces errors in geodesic
111 distances.

112 Since the effectiveness of a method depends on the choice of these param-
113 eters, many techniques have been proposed to estimate automatically their
114 optimal values. The optimal dimensionality of the embedded space is de-
115 fined as the intrinsic dimension of the high dimensional data. More formally,
116 a dataset $X \in R^D$ is said to have intrinsic dimensionality (ID) equal to d if its
117 elements lie entirely within a d-dimensional subspace of R^d (where $d \ll D$)

118 (Fukunaga, 1982). Estimation of 'd' can be achieved using many approaches
119 (see (Camastra, 2003) for a detailed review) including maximum likelihood
120 estimation (Levina and Bickel, 2005), packing numbers (Kegl, 2003), analysis
121 of a geodesic minimum spanning tree (Costa and Hero, 2004), fractal-based
122 methods (Camastra, 2003) and eigenvalue-based estimator (Fukunaga and
123 Olsen, 1971) (EE). However, none of them has achieved consensus as the
124 most accurate method.

125 The selection of the optimal neighbourhood size 'K' is also an open prob-
126 lem. The main line of research has focused on assessing directly the quality
127 of embedded spaces by a quantitative measure in order to infer the optimal
128 value of 'K'. Although many measures have already been proposed, such as
129 Residual Variance (Kouropiteva et al., 2002; Samko et al., 2006), Spearman
130 Rho (Karbauskait et al., 2007; Samko et al., 2006) and Procrustes Analy-
131 sis (Goldberg and Ritov, 2009), experiments suggest their accuracy depends
132 not only on the choice of intrinsic dimensionality but also on the nature of
133 dataset. Consequently, they are not suitable when dealing with complex non-
134 linear high dimensional data of a nature, which is different from that they
135 have been designated for, e.g. human motion (Lewandowski et al., 2009).

136 Finally, an inherent limitation of spectral dimensionality reduction ap-
137 proaches is that they do not provide an explicit mapping function between
138 low and high dimensional spaces. Such function is essential for ensuring
139 continuity of low dimensional representation and projecting data between
140 spaces. This issue has been addressed quite satisfactorily by applying Ra-
141 dial Basis Function network (Poggio et al., 1990) to approximate the optimal
142 mapping function (Elgammal and Lee, 2007; He et al., 2004; Lewandowski

143 et al., 2009). However, the quality of RBFN relies on the careful selection of
144 a few parameters which are chosen empirically.

145 **3. Automatic configuration of spectral dimensionality reduction** 146 **methods**

147 We contribute to the current state of the art by addressing two essential
148 problems: the selection of the optimal neighbourhood size 'K' and the ab-
149 sence of mapping function between spaces. First, we propose to estimate the
150 optimal neighbourhood size by assessing the quality of discovered embedding
151 spaces using the mutual information measure. Secondly, we overcome the
152 deficiency of mapping function by extending advanced RBFN by exploiting
153 spectral graphs to design the optimal structure of the network in an unsuper-
154 vised manner. The above schemas are integrated into a general framework for
155 the automatic configuration of spectral dimensionality reduction methods.

156 *3.1. Estimation of optimal neighbourhood size*

157 The optimal neighbourhood size 'K' can be identified directly by assessing
158 embedded space quality. The process is the following. First, data are divided
159 into training and testing sets. Then, for a given value of 'K', dimensionality
160 reduction is applied on the training set and a mapping function is built be-
161 tween the original and embedded spaces. Finally, test data are projected into
162 the low dimensional space and some error metric is calculated. This process
163 is repeated for a range of 'K' values so that the optimal neighbourhood size
164 can be identified.

165 Since this process requires calculating computationally expensive map-
166 ping functions for all possible values of 'K', quantitative metrics have been

167 proposed to evaluate the quality of an embedded space without mapping.
 168 The standard procedure of optimal neighbourhood size estimation using a
 169 quantitative metric is summarized in pseudo-code 1. There are three met-
 170 rics commonly used. Residual variance (RV) (Kouropiteva et al., 2002; Samko
 171 et al., 2006) expresses how well the distance information is preserved between
 172 two sets of variables, i.e. it reflects the degree of linear relationship between
 173 these variables. Spearman’s rho (SR) (Karbauskait et al., 2007; Samko et al.,
 174 2006) measures the accuracy of the low-dimensional manifold in retaining the
 175 order of pair wise distances of data points of the high-dimensional. Finally,
 176 procrustes analysis measure (PA) (Goldberg and Ritov, 2009) reflects the
 177 matching of two sets of variables in terms of distances. PA determines how
 178 well linear transformations of the points in one space conforms to the points
 179 in the second space. Since experiments have suggested that these measures
 180 depend on the specific nature of datasets (Lewandowski et al., 2009), they
 181 are not suitable for the automatic selection of the free parameter ‘K’ in an
 182 untested domain.

Algorithm 1 Estimation of optimal neighbourhood size

Input: high dimension dataset, maximum K ($maxK$), ID estimate

Output: optimal K

Find minimum K ($minK$) which produces a fully connected graph

for each K in range $< minK, maxK >$ **do**

 Reduce dimensionality of the dataset using a spectral method

 Use metric to assess the quality of the embedded space

end for

Select optimal K according to metric

183 In this work, we tackle this fundamental issue by adopting a metric which

184 can deal with variables without any linear relationship. We propose to use
 185 the mutual information measure (Cover and Thomas, 1991) which has proved
 186 to be able to discover even marginal dependency between two spaces of vari-
 187 ables, since, in contrast to linear correlation coefficients, it is also sensitive to
 188 dependencies which do not manifest themselves in the covariance. MI is null
 189 if and only if the two random variables are strictly independent. The first
 190 idea would be to design a cost function directly in the spectral dimension-
 191 ity reduction framework using MI, however since MI expresses relationship
 192 between two sets of variables rather than individual points, it is not an ap-
 193 propriate metric for this purpose. As the consequence, we propose to employ
 194 it in post processing step to evaluate the quality of spaces.

195 The most straightforward and widespread approach for estimating MI is
 196 to partition the data and approximate MI by the following finite sum:

$$I(X, Y) = \sum_i^N \sum_j^N p(i, j) \log \frac{p(i, j)}{p_x(i)p_y(j)} \quad (1)$$

197 where $p(i, j)$ is the joint probability distribution function, and $p_x(i)$ and $p_y(j)$
 198 are the marginal probability distribution functions of X and Y respectively.
 199 This formulation can be equivalently expressed as (Cover and Thomas, 1991):

$$I(X, Y) = H(X) + H(Y) - H(X, Y) \quad (2)$$

200 where $H(X)$ and $H(Y)$ are the marginal entropies and $H(X, Y)$ is the joint
 201 entropy of X and Y .

202 However, this standard approach can only be applied for $D = d = 1$,

203 because the estimation of entropy is based on data binning. Since, in our
 204 framework, we need to estimate MI measure for higher dimensional variables
 205 ($D > 1, d \geq 1$), we calculate the entropy using K-nearest neighbour statistics
 206 as proposed in (Kraskov et al., 2004). Assuming that some metric is defined
 207 on the spaces spanned by X and Y , all neighbours of a given data point are
 208 ranked according to their distance to that point. Then the entropy $H(Z)$,
 209 where $Z \in \{X, Y\}$, is estimated by the average distance to the K-nearest
 210 neighbours, averaged over all z ($z \in \{x, y\}$). This leads to the following
 211 equation (Kraskov et al., 2004):

$$H(Z) = N^{-1} \sum_{i=1}^N (\gamma(n_z(i) + 1)) - \gamma(N) - \log c_{d_z} - \frac{d_z}{N} \sum_{i=1}^N (\log \epsilon(i)) \quad (3)$$

212 Here, $n_z(i)$ denotes the number of points which fulfil the condition: $\|z(i) -$
 213 $z(j)\| < \epsilon(i)$ and $\gamma(\cdot)$ is the digamma function (Kraskov et al., 2004). d
 214 denotes the dimension of z and c_{d_z} is the volume of the d -dimensional unit
 215 ball. Similarly, the joint entropy of X and Y for a given K (Kraskov et al.,
 216 2004) is expressed by:

$$H(X, Y) = \gamma(K) - \gamma(N) - \log(c_{d_x} c_{d_y}) - \frac{d_x + d_y}{N} \sum_{i=1}^N (\log \epsilon(i)) \quad (4)$$

217 Combining equations 2, 3 and 4 results in the expression of multi dimen-
 218 sional MI:

$$I(X, Y) = \gamma(K) + \gamma(N) - N^{-1} \sum_{i=1}^N (\gamma(n_x(i) + 1) + \gamma(n_y(i) + 1)) \quad (5)$$

219 Although mutual information has never been used in this context, the
 220 use of the multidimensional extension allows MI becoming an intuitive mea-
 221 sure for analysing the mutual correlation between high and low dimensional
 222 spaces.

223 3.2. Unsupervised mapping

224 All spectral approaches suffer from the deficiency of not providing a map-
 225 ping function. A solution has been to use RBFN based mapping (Elgammal
 226 and Lee, 2007; He et al., 2004). However, this process relies on manual
 227 adjustment of its structure according to data. In previous work, we have ad-
 228 dressed this by introducing unsupervised RBFN (Lewandowski et al., 2009).
 229 Since that approach has some limitations (that we discuss later), we propose
 230 a novel method for designing the structure of the network which originates
 231 from graph clustering theory.

232 RBFN from high to low dimensional space is expressed by the following
 233 over-constrained nonlinear system of equations:

$$y = f(x) = B * \psi(x) \quad (6)$$

234 where B is a $D \times L$ matrix of network weights and vector $\psi(x)$ is given by:

$$\psi(x) = [\phi(\|x - c_1\|), \phi(\|x - c_2\|), \dots, \phi(\|x - c_L\|)]^T \quad (7)$$

235 where L is the number of hidden layers in the network, which correspond
 236 to the coordinates of centres c_j and $\phi(\cdot)$ is a real-valued basis function. We
 237 exploit Gaussian basis function $\phi(\|x_i - c_j\|) = e^{-\frac{\|x_i - c_j\|^2}{2\sigma^2}}$, where σ denotes the

238 average distance between all centres, because it has excellent approximation
239 properties (Poggio et al., 1990). The solution for B can be found by applying
240 the Moore-Penrose pseudo-inverse on matrix $\psi(X)$ in equation 6 and solving
241 the obtained linear system of equations.

242 The RBFN structure is formed by centres c_j which summarize training
243 data points in order to provide generalization properties of the network. How-
244 ever, the performance of RBFN critically depends upon the chosen centres
245 (Chen et al., 1991). K-means clustering (Kanungo et al., 2002) (KMC) and
246 rival penalized competitive learning (Xu et al., 1993) (RPCL) are currently
247 the most popular and well studied methods which address this task. A key
248 drawback of the KMC algorithm is that it requires prior knowledge about
249 the correct number of centres. This can be addressed using the RPCL algo-
250 rithm which is capable of finding the optimal localisation of centres as well as
251 their correct number L in an automatic way. First, L' centres are randomly
252 initialised ($L' \gg L$). Subsequently, in each iteration, the algorithm randomly
253 selects a sample s from the training set and moves the closest centre (the so
254 called competition winner) towards the considered point s by a weighted dis-
255 tance $w1$. In the same step the second closest centre (or rival) is pushed away
256 from the sample s by a weighted distance $w2$ (where $w1 \gg w2$). Learning
257 rates, i.e. $w1, w2$ are monotonically decreased after each iteration. The entire
258 procedure is repeated until it converges or reaches a given threshold. This
259 mechanism allows automatic determination of the centres' positions by locat-
260 ing them at the core of data point clusters and gradually driving unrequired
261 centres away from those clusters.

262 In earlier work (Lewandowski et al., 2009), we automated the mapping

263 process by applying RPCL for training of RBFN. However, RPCL, as KMC,
264 depends on the initial random localization of centres and relies on the Eu-
265 clidean distance, which is not the most appropriate metric to model high
266 dimensional relationships (Aggarwal et al., 2001). In order to improve ac-
267 curacy, we extend our idea of unsupervised mapping learning and propose
268 to use the Markov cluster algorithm (MCL) (Dongen, 2000) to identify the
269 suitable number and localization of centres automatically by exploiting the
270 adjacency graph constructed during spectral reduction of dimensionality. As
271 it will be demonstrated in the results section, the computational cost of a
272 mapping learning process is greatly reduced and the obtained mapping ex-
273 hibits better accuracy in comparison to standard approaches such as KMC
274 and RPCL.

275 At the heart of the MCL algorithm (Dongen, 2000) lies the idea to sim-
276 ulate flow within a graph: flows are promoted where current is strong and
277 demoted where current is weak. Flow simulation is achieved by transform-
278 ing a graph into a Markov graph using the standard definition of a random
279 walk on a graph. Then a flow is defined by two simple algebraic operations,
280 i.e. expansion and inflation, which are applied connectively, so that the flow
281 becomes thicker in regions of higher current and thinner in regions of lower
282 current.

283 According to this paradigm, if natural groups are present in the spectral
284 graph obtained in the first step of dimensionality reduction, then, current
285 across borders between different groups will wither away. As the result, a
286 fully connected graph is divided into few subgraphs (figure 2), thus revealing
287 the optimal number L as well as coordinates of clusters c_j . Application of

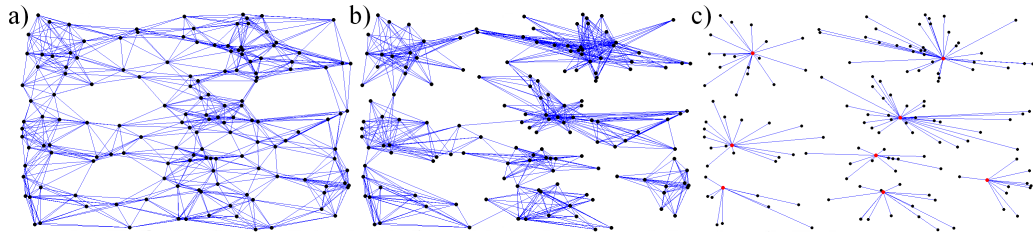


Figure 2: 2D representation of successive stages of flow simulation using the MCL process for discovery of the localisation and the number of centres in RBFN.

288 this procedure enables the discovery of more representative clusters of high
 289 dimensional data and subsequently customise RBFN structure to dataset in
 290 an automatic and efficient manner.

291 4. Experimental results and discussion

292 4.1. Datasets

293 The proposed framework was validated with both artificial and real datasets.
 294 Standard datasets were selected to extensively evaluate the performance and
 295 robustness of the proposed methodology in different scenarios. Figure 3 illus-
 296 trates the datasets used in this work. Since the intrinsic dimensionalities of
 297 the digits and face datasets are unknown, we used both low and high values
 298 of their estimates in order to perform our experiments.

299 The 'swissroll' dataset is a synthetic and nonlinear example of a two
 300 dimensional flat submanifold which lies in a three-dimensional space. This
 301 dataset exhibits significant disagreement between geodesic and Euclidean
 302 distances (figure 3a). 2000 points were randomly sampled from the manifold
 303 and used in all our experiments. In addition, we generated a second smaller
 304 dataset consisting of 1000 points (denoted by a star in our experiments) in

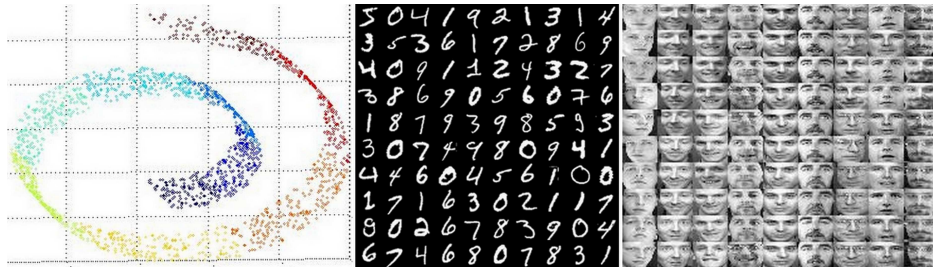


Figure 3: Datasets used in the experiments: from left to right, 'swissroll' manifold, handwritten digits and face images.

305 order to compare Isomap results with those of the original Isomap paper
 306 (Tenenbaum et al., 2000).

307 The MNIST dataset (LeCun, 2000) consists of handwritten characters
 308 images containing digits from 0 to 9 (figure 3b). The size of each image is 28
 309 x 28 pixels, with 256 gray levels per pixel. Thus, each image is represented by
 310 a 784-dimensional vector. Due to computational and memory constraints, in
 311 our experiments we used a subset of the MNIST database consisting of 6000
 312 images. According to (Camastra and Vinciarelli, 2001), the optimum ID of
 313 handwritten digits is 7, whereas the upper bound of the ID as determined
 314 by EE equals 10.

315 The ORL (formerly Olivetti) face database contains 400 images of 40
 316 distinct subjects (Samaria and Harter, 1994) (figure 3c). All images were
 317 captured against a dark homogeneous background with the subjects in an up-
 318 right, frontal position, with tolerance for some side movements. There are
 319 variations in facial expression (open/closed eyes, smiling/nonsmiling), and
 320 facial details (glasses/no glasses, different skin colours). The images are grey-
 321 scale with a resolution of 64x64 pixels which gives a 4096 dimension feature
 322 vector. The analysis of relation between recognition rates and dimensionality

323 of embedded space in (Yin et al., 2008) suggests a value of 10 as the optimal
324 ID for this dataset. The upper bound of the ID as determined by EE equals
325 40.

326 *4.2. Experiments*

327 All experiments were performed with Isomap, LLE and LE using K values
328 in the range $\langle 4, 30 \rangle$. In multidimensional spaces, geodesic distances are used,
329 whereas on the plane we employ Euclidean distances as suggested in (Samko
330 et al., 2006).

331 First, we evaluate qualitatively the novel MI estimator against current
332 approaches, i.e. Residual Variance, Spearman Rho and Procrustes Analysis
333 measures. This was performed using the synthetic dataset for which the
334 underlying structure is known so the quality of embedded space can be judged
335 visually.

336 Then, two classical pattern classification problems, face and handwritten
337 digit recognition, are considered in order to analyze the quantitative perfor-
338 mance of the MI metric. We do not perform any preprocessing or normal-
339 ization of the data in order to prevent any information lost. It is important
340 to note that, in this work, we did not focus on designing a state of art clas-
341 sification system, but to compare existing metrics with the one we propose
342 using on a standard classification framework based on a real application.

343 Finally, in the last experiment we show superiority of graph based RBFN
344 in comparison with standard RBFN. This is achieved by repeating the classi-
345 fication experiments with digits and faces recognition using the new mapping
346 function whose structure is inferred automatically from the spectral graphs.

347 *4.2.1. Dimension reduction of 'swissroll' dataset*

348 Table 1 presents the low dimensional spaces of 'swissroll' dataset produced
349 by Isomap, LE and LLE using the estimated neighbourhood sizes calculated
350 by RV, SR, PA and MI.

351 In all cases, the MI measure was able to identify very good low dimen-
352 sional representation of 'swissroll' dataset, i.e. embedded space which man-
353 ages to unroll manifold and preserves local structure. Moreover, estimated
354 values of K using MI are in agreement with parameters which were recom-
355 mended in the original papers (Belkin and Niyogi, 2001; Roweis and Saul,
356 2000; Tenenbaum et al., 2000). Although, other measures usually select
357 reasonable low dimensional representations, their quality is not consistent.
358 For instance, the local structure is distorted in most experiments involving
359 RV/SR. Although PA seems to behave similarly to MI, in the case of LLE
360 the very different neighbourhood size returned by PA leads to the production
361 of an embedded space of inferior quality.

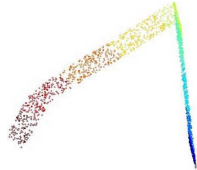
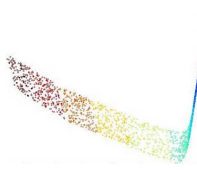
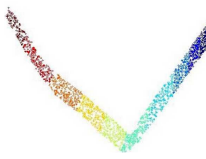
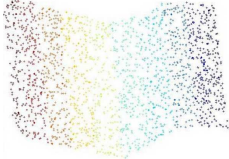
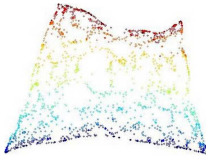
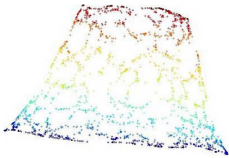
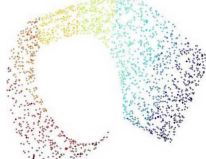
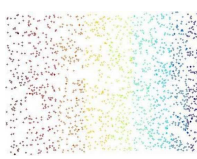
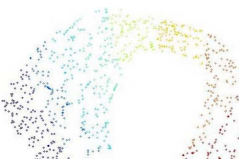
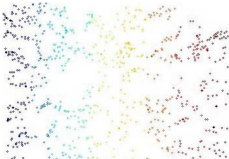
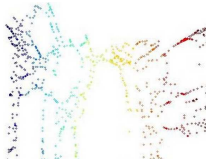
Method (Recommended K)	Coefficient (Estimated K)	Visualization	Coefficient (Estimated K)	Visualization
LLE (20) (Roweis and Saul, 2000)	RV (11)		SR (22)	
	PA (8)		MI (20)	
LE (5-15) (Belkin and Niyogi, 2001)	RV/SR (8)		PA/MI (5)	
Isomap (-)	RV (21)		SR/PA/MI (18)	
Isomap* (7) (Tenenbaum et al., 2000)	RV (9)		PA/MI (7)	
	SR (4)		* denote the 'swissroll' dataset with 1000 points instead of 2000 points	

Table 1: The low dimensional spaces of 'swissroll' with estimated and recommended neighbourhood sizes for Isomap, LE and LLE according to coefficients RV, SR, PA and MI.

362 *4.2.2. Classification evaluation*

363 The recognition of either digits or faces is performed according to the
364 10-fold cross validation strategy, where we divide a dataset into ten distinct
365 partitions. For each partition, we reduce dimensionality of remaining dataset
366 and train RBFN with the standard RPCL algorithm. Then, each partition
367 is projected into the low dimensional space and classification is performed
368 using a first nearest neighbour classifier (Ho, 1998). Finally, classification
369 accuracy is calculated by averaging over the ten partitions. For each dataset,
370 estimation of optimal neighbourhood size for dimensionality reduction is cal-
371 culated using RV, SR, PA and MI. Moreover, the actual optimal K, 'Opt',
372 is calculated experimentally by an exhaustive evaluation of classification ac-
373 curacy for all values of K within the range $\langle 4, 30 \rangle$. In addition, using that
374 value, we evaluate the classification accuracy of the scheme, 'Opt*', which
375 includes graph-based RBFN (G-RBFN). Tables 2 and 3 show the results of
376 these experiments which were conducted with two sets of IDs as defined in
377 section 4.1. Note that the huge computational cost of applying PA on the
378 very high dimensional faces dataset (dimensionality of 4096) did not allow
379 us to obtain the results for this measure using our processing capabilities
380 (16-node cluster).

	ID	RV	SR	PA	MI	Opt	Opt*
Iso		88	88	88	88	89	90
LLE	10	62	63	59	78	78	82
LE		79	79	80	80	80	84
Iso		85	82	84	85	85	87
LLE	7	56	63	53	74	74	77
LE		75	75	74	76	77	80

Table 2: Percentage accuracy of hand-written digits recognition.

	ID	RV	SR	PA	MI	Opt	Opt*
Iso		76	73	-	77	77	77
LLE	40	78	78	-	80	80	80
LE		67	67	-	67	68	73
Iso		65	57	-	76	76	76
LLE	10	55	55	-	61	62	62
LE		62	50	-	63	63	63

Table 3: Percentage accuracy of faces recognition.

381 In agreement with our previous experiments, neighbourhood sizes esti-
382 mated by the MI measure produce consistently better accuracy than those
383 suggested by other metrics regardless of the chosen ID. Moreover, it allows
384 classification performances which are either optimal or near-optimal for a
385 given dimensionality reduction method. Results also reveal that unlike LLE
386 and Isomap, LE is not very sensitive to neighbourhood size selection. As
387 expected, decrease of ID results in a decline of accuracy since more dis-
388 criminative information is discarded during dimensionality reduction. Two
389 dimensional visualization of the best low dimensional space obtained with
390 Isomap for the digit dataset is presented in figure 4.

391 Regarding the efficiency of graph-based RBFN, tables 2 and 3 show that
392 this new scheme improves significantly the quality of the mapping produced
393 by standard RPCL RBFN. Further comparison between those two mapping
394 methods is provided in figure 5, where classification accuracy and processing
395 time are measured for various sizes of the digits dataset. Here, LE is used
396 for dimensionality reduction as a representative of spectral methods.

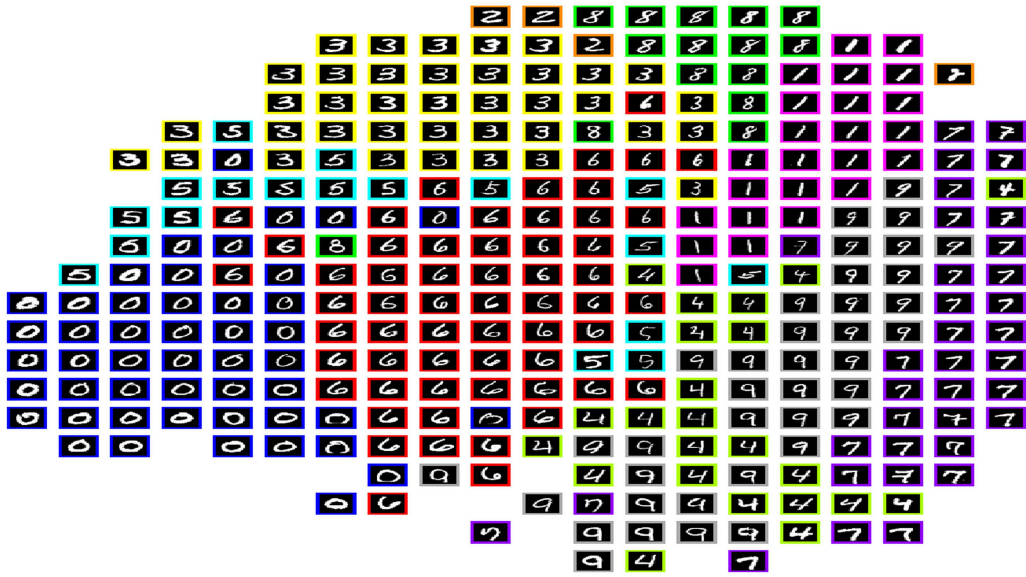


Figure 4: Two dimensional visualization of the best low dimensional space obtained with Isomap for MNIST data subset.

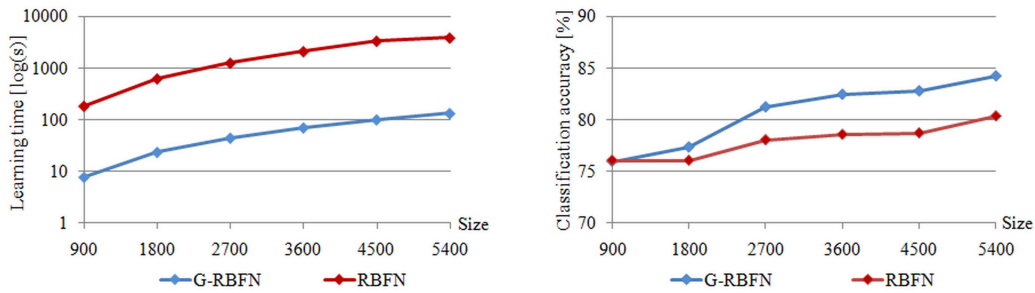


Figure 5: Classification processing time (left) and accuracy (right) comparisons between graph-based RBFN and standard RPCL RBFN according to digits dataset size (ID=10).

397 First, whatever the size of the training set, classification accuracy using
 398 graph-based RBFN is higher than for standard RBFN. Secondly, graph-based
 399 RBFN is computationally very efficient since the learning process time tends
 400 to increase linearly with the size of the database, while it grows quadratically
 401 when it is performed using the RPCL procedure.

402 4.3. Discussion

403 All experiments demonstrate that MI is a better metric to estimate neigh-
404 bourhood size than currently used measures. Not only are embedded spaces
405 produced by MI visually convincing, but our quantitative study, i.e. clas-
406 sification experiments, confirm its superiority since it consistently provides
407 better accuracy regardless of the estimated ID. Moreover, unlike PA, MI
408 proved able to handle very high dimensional datasets. Our quantitative ex-
409 periments also validate our proposal of using graph-based RBFN to pro-
410 vide mapping between embedded and data spaces. This scheme outperforms
411 significantly standard RBFN mapping in both accuracy and computational
412 efficiency when combined with spectral dimensionality reduction methods.

413 Although we used classification experiments to validate quantitatively the
414 value of our contribution to spectral dimensionality reduction methods, our
415 aim was not to produce a state of the art classifier, but to demonstrate that
416 our innovations could be applied successfully to representatives of the three
417 main spectral families, i.e. Isomap, LLE and LE. We would suggest readers
418 with a special interest in classification to apply our advanced techniques to
419 spectral methods which were developed especially to handle that task. They
420 include discriminant Isomap (Yang, 2003), supervised LLE (De Ridder et al.,
421 2003) and semi-supervised LE (Zheng et al., 2008).

422 5. Conclusions

423 In this paper, a framework has been proposed to configure automatically
424 spectral dimensionality reduction methods. This is achieved by, first, es-
425 timating the optimal neighbourhood size. We introduce the MI metric as
426 a powerful alternative to existing techniques. Secondly, RBFN based map-
427 ping between spaces has to be learnt in a unsupervised manner. Although
428 the RPCL algorithm is the standard approach, we enhance significantly the
429 learning process by using the efficient graph based MCL algorithm. Our
430 contributions have been validated qualitatively and quantitatively using var-
431 ious datasets. Results demonstrate that neighbourhoods selected by the MI
432 metric allow spectral dimensionality reduction methods to produce better
433 quality embedded spaces. Moreover, our new mapping functions improve
434 both mapping accuracy and computational efficiency.

435 Aggarwal, C., Hinneburg, A., Keim, D., 2001. On the surprising behavior of
436 distance metrics in high dimensional spaces. *Lecture Notes in Computer
437 Science* 1973, 420–434.

438 Belkin, M., Niyogi, P., 2001. Laplacian eigenmaps and spectral techniques
439 for embedding and clustering. *Advances in Neural Information Processing
440 Systems* 14 14, 585–591.

441 Bishop, C., Svensén, M., Williams, C., 1998. GTM: The generative topo-
442 graphic mapping. *Neural computation* 10 (1), 215–234.

- 443 Camastra, F., 2003. Data dimensionality estimation methods: a survey. *Pat-*
444 *tern Recognition* 36 (12), 2945–2954.
- 445 Camastra, F., Vinciarelli, A., 2001. Intrinsic dimension estimation of data:
446 An approach based on grassberger-procaccia’s algorithm. *Neural Processing*
447 *Letters* 14 (1), 27–34.
- 448 Chen, S., Cowan, C., Grant, P., 1991. Orthogonal least squares learning
449 algorithm for radial basisfunction networks. *IEEE Transactions on neural*
450 *networks* 2 (2), 302–309.
- 451 Choi, H., Choi, S., 2004. Kernel isomap. *IET Electronics letters* 40 (25),
452 1612–1613.
- 453 Costa, J., Hero, A., 2004. Geodesic entropic graphs for dimension and entropy
454 estimation in manifold learning. *IEEE Transactions on Signal Processing*
455 52 (8), 2210–2221.
- 456 Cover, T. M., Thomas, J., 1991. *Elements of information theory*. Wiley.
- 457 De Ridder, D., Kouropteva, O., Okun, O., Pietikainen, M., Duin, R., 2003.
458 Supervised locally linear embedding. *Lecture Notes in Computer Science*
459 2714, 333–341.
- 460 De Silva, V., Tenenbaum, J., 2003. Global Versus Local Methods in Nonlin-
461 ear Dimensionality Reduction. *Advances in neural information processing*
462 *systems* 15, 705–712.

- 463 Dongen, S., 2000. Graph clustering by flow simulation. Ph.D. thesis, PhD
464 Thesis, University of Utrecht, The Netherlands.
- 465 Donoho, D., Grimes, C., 2003. Hessian eigenmaps: Locally linear embed-
466 ding techniques for high-dimensional data. *Proceedings of the National
467 Academy of Sciences* 100 (10), 5591–5596.
- 468 Elgammal, A., Lee, C., 2007. Nonlinear manifold learning for dynamic shape
469 and dynamic appearance. *Computer Vision and Image Understanding*
470 106 (1), 31–46.
- 471 Fukunaga, K., 1982. Intrinsic dimensionality extraction. *Classification, Pat-
472 tern Recognition and Reduction of Dimensionality*, 347–362.
- 473 Fukunaga, K., Olsen, D. R., 1971. An algorithm for
474 nding intrinsic dimensionality of data. *IEEE Transactions on Computers*
475 C-20 (2), 176–183.
- 476 Goldberg, Y., Ritov, Y., 2008. Ldr-lle: Lle with low-dimensional neighbor-
477 hood representation. *ISVC '08: Proceedings of the 4th International Sym-
478 posium on Advances in Visual Computing, Part II*, 43–54.
- 479 Goldberg, Y., Ritov, Y., 2009. Local procrustes for manifold embedding: a
480 measure of embedding quality and embedding algorithms. *Machine Learn-
481 ing* 77 (1), 1–25.
- 482 He, X., Cai, D., Yan, S., Zhang, H., 2005. Neighborhood preserving embed-

483 ding. Proceedings of IEEE International Conference on Computer Vision
484 2, 1208–1213.

485 He, X., Ma, W., Zhang, H., 2004. Learning an image manifold for retrieval.
486 Proceedings of ACM Multimedia, 17–23.

487 He, X., Niyogi, P., 2004. Locality preserving projections. Advances in Neural
488 Information Processing Systems 16, 153–160.

489 Ho, T., 1998. Nearest neighbors in random subspaces. Lecture Notes in Com-
490 puter Science 1451, 640–648.

491 Kanungo, T., Mount, D., Netanyahu, N., Piatko, C., Silverman, R., Wu, A.,
492 2002. A local search approximation algorithm for k-means clustering. Pro-
493 ceedings of the eighteenth annual symposium on Computational geometry,
494 10–18.

495 Karbauskait, R., Kurasova, O., Dzemyda, G., 2007. Selection of the number
496 of neighbours of each data point for the locally linear embedding algorithm.
497 Information Technology and Control 36 (4), 359–364.

498 Kegl, B., 2003. Intrinsic dimension estimation using packing numbers. Ad-
499 vances in Neural Information Processing Systems 15, 681–688.

500 Kokiopoulou, E., Saad, Y., 2007. Orthogonal Neighborhood Preserving Pro-
501 jections: A Projection-Based Dimensionality Reduction Technique. Pat-
502 tern Analysis and Machine Intelligence, IEEE Transactions on 29 (12),
503 2143–2156.

- 504 Kouropteva, O., Okun, O., Pietikainen, M., 2002. Selection of the optimal
505 parameter value for the locally linear embedding algorithm. Proc. of the
506 1 st Int. Conf. on Fuzzy Systems and Knowledge Discovery, Singapore,
507 359–363.
- 508 Kraskov, A., Stoegbauer, H., Grassberger, P., 2004. Estimating mutual infor-
509 mation. Physical review. E, Statistical, nonlinear, and soft matter physics
510 69 (6), 66138.
- 511 Lawrence, N., 2004. Gaussian process latent variable models for visualisa-
512 tion of high dimensional data. Advances in Neural Information Processing
513 Systems 16.
- 514 LeCun, Y., 2000. MNIST handwritten digits dataset
515 (<http://yann.lecun.com/exdb/mnist/index.html>) [last accessed on
516 30/10/2009].
- 517 Levina, E., Bickel, P., 2005. Maximum likelihood estimation of intrinsic di-
518 mension. Advances in Neural Information Processing Systems 17, 777–784.
- 519 Lewandowski, M., Makris, D., Nebel, J.-C., 2009. Automatic configuration of
520 spectral dimensionality reduction methods for 3D human pose estimation.
521 Visual Surveillance.
- 522 Poggio, T., Girosi, F., MIT, C., 1990. Networks for approximation and learn-
523 ing. Proceedings of the IEEE 78 (9), 1481–1497.

- 524 Roweis, S., Saul, L., 2000. Nonlinear dimensionality reduction by locally
525 linear embedding. *Science* 290 (5500), 2323–2326.
- 526 Samaria, F., Harter, A., 1994. Parameterisation of a Stochastic Model for Hu-
527 man Face Identification (<http://www.cs.toronto.edu/~roweis/data.html>)
528 [last accessed on 30/10/2009]. Workshop on Applications of Computer Vi-
529 sion.
- 530 Samko, O., Marshall, A., Rosin, P., 2006. Selection of the optimal parameter
531 value for the Isomap algorithm. *Pattern Recognition Letters* 27 (9), 968–
532 979.
- 533 Tenenbaum, J., Silva, V., Langford, J., 2000. A global geometric framework
534 for nonlinear dimensionality reduction. *Science* 290 (5500), 2319–2323.
- 535 Wang, Q., Li, J., 2009. Combining local and global information for nonlinear
536 dimensionality reduction. *Neurocomputing* 72 (10-12), 2235–2241.
- 537 Xu, L., Krzyzak, A., Oja, E., 1993. Rival penalized competitive learning for
538 clustering analysis, RBFnet, and curve detection. *IEEE Transactions on*
539 *Neural networks* 4 (4), 636–649.
- 540 Yang, M., 2003. Discriminant isometric mapping for face recognition. *Lecture*
541 *Notes in Computer Science* 2626, 470–480.
- 542 Yin, J., Hu, D., Zhou, Z., 2008. Growing locally linear embedding for mani-
543 fold learning. *Journal of Pattern Recognition Research* 2 (1), 1–16.

544 Zhang, Z., Wang, J., 2007. Mlle: Modified locally linear embedding using
545 multiple weights. *Advances in Neural Information Processing Systems* 19,
546 1593–1600.

547 Zheng, F., Chen, N., Li, L., 2008. Semi-supervised Laplacian eigenmaps for
548 dimensionality reduction. *Wavelet Analysis and Pattern Recognition* 2,
549 843–849.