

Available online at www.sciencedirect.com

Image and Vision Computing xxx (2006) xxx–xxx

www.elsevier.com/locate/imavis

Kernel-based distance metric learning for content-based image retrieval

Hong Chang, Dit-Yan Yeung *

Department of Computer Science, Hong Kong University of Science and Technology, Clear Water Bay, Kowloon, Hong Kong

Received 24 August 2005; received in revised form 16 February 2006; accepted 16 May 2006

Abstract

For a specific set of features chosen for representing images, the performance of a content-based image retrieval (CBIR) system depends critically on the similarity or dissimilarity measure used. Instead of manually choosing a distance function in advance, a more promising approach is to learn a good distance function from data automatically. In this paper, we propose a kernel approach to improve the retrieval performance of CBIR systems by learning a distance metric based on pairwise constraints between images as supervisory information. Unlike most existing metric learning methods which learn a Mahalanobis metric corresponding to performing linear transformation in the original image space, we define the transformation in the kernel-induced feature space which is nonlinearly related to the image space. Experiments performed on two real-world image databases show that our method not only improves the retrieval performance of Euclidean distance without distance learning, but it also outperforms other distance learning methods significantly due to its higher flexibility in metric learning.

© 2006 Elsevier B.V. All rights reserved.

Keywords: Metric learning; Kernel method; Content-based image retrieval; Relevance feedback

1. Introduction

1.1. Content-based image retrieval

With the emergence and increased popularity of the World Wide Web (WWW) over the past decade, retrieval of images based on content, often referred to as *content-based image retrieval* (CBIR), has gained a lot of research interests [1]. On the WWW where many images can be found, it is convenient to search for the target images in possibly very large image databases by presenting query images as examples. Thus, more and more Web search engines (e.g., Yahoo) are now equipped with CBIR facilities for retrieving images on a query-by-image-example basis.

The two determining factors for image retrieval performance are the features used to represent the images and

the distance function used to measure the similarity between a query image and the images in the database. For a specific feature representation chosen, the retrieval performance depends critically on the similarity measure used. Let $\mathbf{f}^i = (f_1^i, f_2^i, \dots, f_n^i)$ denote a feature vector representing image i , where n is the number of features. For example, \mathbf{f}^i represents a color histogram with n being the number of histogram bins. There exist many methods for measuring the distance between feature vectors. Swain and Ballard [2] proposed the intersection distance measure $d_{\cap} = \sum_{k=1}^n \min(f_k^i, f_k^j)$, which has the same ordinal properties as the L_1 norm (distance). In [3], the distance between two histograms is defined as the weighted form $d_{\mathbf{W}}(\mathbf{f}^i, \mathbf{f}^j) = \sqrt{(\mathbf{f}^i - \mathbf{f}^j)^T \mathbf{W} (\mathbf{f}^i - \mathbf{f}^j)}$, where each weight w_{ij} in \mathbf{W} denotes the similarity between features i and j . Note that this distance measure includes the Mahalanobis distance as a special case. Other commonly used distance functions for color histograms include the Minkowski distance $d_r(\mathbf{f}^i, \mathbf{f}^j) = (\sum_{k=1}^n |f_k^i - f_k^j|^r)^{1/r}$. However, this distance metric may lead to high false negative rate [4].

* Corresponding author. Tel.: +852 2358 6977; fax: +852 2358 1477.

E-mail addresses: hongch@cs.ust.hk (H. Chang), dyyeung@cs.ust.hk (D.-Y. Yeung).

55 Unfortunately, the effectiveness of these distance func- 110
 56 tions is rather limited. Instead of choosing a distance func- 111
 57 tion in advance, a more promising approach is to learn a 112
 58 good distance function from data automatically. Recently, 113
 59 this challenging new direction has aroused great interest in 114
 60 the research community. 115

61 1.2. Related work 116

62 *Relevance feedback* has been used in the traditional 119
 63 information retrieval community to improve the perfor- 120
 64 mance of information retrieval systems based on user 121
 65 feedback. This interactive approach has also emerged 122
 66 as a popular approach in CBIR [5]. The user is provided 123
 67 with the option of labeling (some of the) previously 124
 68 retrieved images as either relevant or irrelevant. Based 125
 69 on this feedback information, the CBIR system can iter- 126
 70 atively refine the retrieval results by learning a more 127
 71 appropriate (dis)similarity measure. For example, rele- 128
 72 vance feedback can be used to modify the weights in 129
 73 the weighted Euclidean distance [5] or the generalized 130
 74 Euclidean distance [6]. The same approach has also been 131
 75 applied to a correlation-based metric [7,8], which usually 132
 76 outperforms Euclidean-based measures. In [9], the 133
 77 authors presented an approach to generate an adaptive 134
 78 quasiconformal kernel distance metric based on relevance 135
 79 feedback. Dong and Bhanu [10] proposed a new semi-su- 136
 80 pervised expectation-maximization (EM) algorithm for 137
 81 image retrieval tasks, with the image distribution in the 138
 82 feature space modeled as Gaussian mixtures. Pseudo- 139
 83 feedback strategy based on peer indexing was proposed 140
 84 recently to optimize the similarity metric and the initial 141
 85 query vectors [11], where the global and personal image 142
 86 peer indexes are learned interactively and incrementally 143
 87 from user feedback information. Some recent work 144
 88 makes use of the manifold structure of image data in 145
 89 the feature space for image retrieval [12,13]. Other meth- 146
 90 ods include biased discriminant analysis [14], support 147
 91 vector machine (SVM) active learning [15–17], boosting 148
 92 methods [18], and so on. 149

93 In the machine learning literature, supervisory infor- 150
 94 mation for semi-supervised distance learning usually 151
 95 takes the form of limited labeled data or *pairwise similar-* 152
 96 *ity or dissimilarity constraints*. The latter type of informa- 153
 97 tion is weaker in the sense that pairwise constraints can 154
 98 be derived from labeled data but not vice versa. Rele- 155
 99 vance feedback, which has been commonly used in 156
 100 CBIR, may be used to obtain the pairwise constraints. 157
 101 Recently, some machine learning researchers have pro- 158
 102 posed different metric learning methods for semi-sup- 159
 103 ervised clustering with pairwise similarity or dissimilarity 160
 104 side information [19–22]. Most of these methods try to 161
 105 learn a global Mahalanobis metric corresponding to lin- 162
 106 ear transformation in the original image space [19,20,22]. 163
 107 In particular, an efficient, noniterative algorithm called 164
 108 relevance component analysis (RCA) [19,20] has been 165
 109 used to improve image retrieval performance in CBIR

tasks. This work was later extended in [19] by incorpo- 110
 rating both similarity and dissimilarity constraints into 111
 the EM algorithm for model-based clustering based on 112
 Gaussian mixture models. More recently, Hertz et al. 113
 [23,24] proposed a nonmetric distance function learning 114
 algorithm called DistBoost by boosting the hypothesis 115
 over the product space with Gaussian mixture models 116
 as weak learners. Using DistBoost, they demonstrated 117
 very good image retrieval results in CBIR tasks. 118

Most existing systems only make use of relevance feed- 119
 back within a single query session. More recently, some 120
 methods have been proposed for the so-called *long-term* 121
learning by accumulating relevance feedback from multiple 122
 query sessions which possibly involve different users 123
 [25,12,13,26]. However, [12,13] are based on the assump- 124
 tion that the feature vectors representing the images form 125
 a Riemannian manifold in the feature space. Unfortunat- 126
 ly, this assumption may not hold in real-world image dat- 127
 abases. Moreover, the log-based relevance feedback 128
 method [26] is expected to encounter the scale-up problem 129
 as the number of relevance feedback log sessions increases. 130

1.3. This paper 131

Metric learning based on pairwise constraints can be 132
 categorized into linear and nonlinear methods. Most 133
 existing metric learning methods learn a Mahalanobis 134
 metric corresponding to performing linear transformation 135
 in the original image space. However, for CBIR tasks, 136
 the original image space is highly nonlinear due to high 137
 variability of the image content and style. In this paper, 138
 we define the transformation in the kernel-induced fea- 139
 ture space which is nonlinearly related to the image 140
 space. The transformation is then learned based on side 141
 information in the form of pairwise (dis)similarity con- 142
 straints. Moreover, to address the efficiency problem 143
 for long-term learning, we boost the image retrieval per- 144
 formance by adapting the distance metric in a stepwise 145
 manner based on relevance feedback. 146

Our kernel-based distance metric learning method per- 147
 forms kernel PCA on the whole data set, followed by met- 148
 ric learning in the feature space. It does not suffer from the 149
 small sample size problem encountered by traditional Fish- 150
 er discriminant analysis methods. Therefore, our method is 151
 significantly different from many existing methods which 152
 aim to address the small sample size problem in multimedia 153
 information retrieval, e.g., the kernel-based biased discrim- 154
 inant analysis method proposed in [14]. 155

In Section 2, we will propose a kernel-based method for 156
 nonlinear metric learning. In Section 3, we will describe 157
 how this method can be used to improve the performance 158
 of CBIR tasks. Our method will then be compared with 159
 other distance learning methods based on two real-world 160
 image databases. The stepwise kernel-based metric learning 161
 algorithm that pays attention to both effectiveness and effi- 162
 ciency will be presented in Section 4. Finally, some con- 163
 cluding remarks will be given in the last section. 164

165 **2. Kernel-based metric learning**

166 Kernel methods typically comprise two parts. The first
167 part maps (usually nonlinearly) the input points to a fea-
168 ture space often of much higher or even infinite dimension-
169 ality, and then the second part applies a relatively simple
170 (usually linear) method in the feature space. In this section,
171 we propose a two-step method which first uses kernel prin-
172 cipal component analysis (PCA) [27] to embed the input
173 points in terms of their nonlinear principal components
174 and then applies metric learning there.

175 *2.1. Centering in the feature space*

176 Let \mathbf{x}_i ($i = 1, \dots, n$) be n points in the input space \mathcal{X} .
177 Suppose we use a kernel function \hat{k} which induces a nonlin-
178 ear mapping $\hat{\phi}$ from \mathcal{X} to some feature space \mathcal{F} .¹ The “im-
179 ages” of the n points in \mathcal{F} are $\hat{\phi}(\mathbf{x}_i)$ ($i = 1, \dots, n$), which in
180 general are not centered (i.e., their sample mean is not zero).
181 The corresponding kernel matrix $\hat{\mathbf{K}} = [\hat{k}(\mathbf{x}_i, \mathbf{x}_j)]_{n \times n} =$
182 $[[\hat{\phi}(\mathbf{x}_i), \hat{\phi}(\mathbf{x}_j)]]_{n \times n}$.

183 We want to transform (simply by translating) the coordi-
184 nate system of \mathcal{F} such that the new origin is at the sample
185 mean of the n points. As a result, we also convert the kernel
186 matrix $\hat{\mathbf{K}}$ to $\mathbf{K} = [k(\mathbf{x}_i, \mathbf{x}_j)]_{n \times n} = [[\phi(\mathbf{x}_i), \phi(\mathbf{x}_j)]]_{n \times n}$.

187 Let $\mathbf{Y} = [\phi(\mathbf{x}_1), \dots, \phi(\mathbf{x}_n)]^T$, $\hat{\mathbf{Y}} = [\hat{\phi}(\mathbf{x}_1), \dots, \hat{\phi}(\mathbf{x}_n)]^T$
188 and $\mathbf{H} = \mathbf{I} - \frac{1}{n} \mathbf{1}\mathbf{1}^T$, where $\mathbf{1}$ is a column vector of ones.
189 We can express $\mathbf{Y} = \mathbf{H}\hat{\mathbf{Y}}$. Hence,

192 $\mathbf{K} = \mathbf{Y}\mathbf{Y}^T = \mathbf{H}\hat{\mathbf{Y}}\hat{\mathbf{Y}}^T\mathbf{H} = \mathbf{H}\hat{\mathbf{K}}\mathbf{H}$. (1)

193 *2.2. Step 1: Kernel PCA*

194 We briefly review the kernel PCA algorithm here. More
195 details can be found in [27].

196 We first apply the centering transform as in Eq. (1) to
197 get the kernel matrix \mathbf{K} . We then solve the eigenvalue equa-
198 tion for \mathbf{K} : $\mathbf{K}\alpha = \zeta\alpha$. Let $\zeta_1 \geq \dots \geq \zeta_p > 0$ denote the $p \leq n$
199 positive eigenvalues of \mathbf{K} and $\alpha_1, \dots, \alpha_p$ be the correspond-
200 ing eigenvectors. The embedding dimensionality p may be
201 set to the rank of \mathbf{K} , or, more commonly, a smaller value
202 to ignore the insignificant dimensions with very small
203 eigenvalues, as in ordinary PCA.

204 For any input \mathbf{x} , the k th principal component \tilde{y}_k of $\phi(\mathbf{x})$
205 is given by

207 $\tilde{y}_k = \frac{1}{\sqrt{\zeta_k}} \sum_{i=1}^n \alpha_{ik} \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}) \rangle$. (2)

208 If $\mathbf{x} = \mathbf{x}_j$ for some $1 \leq j \leq n$, i.e., \mathbf{x} is one of the n original
209 points, then the k th principal component \tilde{y}_{jk} of $\phi(\mathbf{x}_j)$
210 becomes

212 $\tilde{y}_{jk} = \frac{1}{\sqrt{\zeta_k}} (\mathbf{K}\alpha_k)_j = \frac{1}{\sqrt{\zeta_k}} (\zeta_k \alpha_k)_j = \sqrt{\zeta_k} \alpha_{jk}$, (3)

¹ We use RBF kernel in this paper.

213 which is proportional to the expansion coefficient α_{jk} . Thus,
214 the input points \mathbf{x}_i ($i = 1, \dots, n$) are now represented as $\tilde{\mathbf{y}}_i$
215 ($i = 1, \dots, n$).

216 *2.3. Step 2: Linear metric learning*

217 To perform metric learning, we further transform $\tilde{\mathbf{y}}_i$
218 ($i = 1, \dots, n$) by applying a linear transform \mathbf{A} to each
219 point based on the pairwise similarity and dissimilarity
220 information in \mathcal{S} and \mathcal{D} , respectively.

221 We define a matrix \mathbf{C}_S based on \mathcal{S} as follows:

$$\begin{aligned} \mathbf{C}_S &= \frac{1}{|\mathcal{S}|} \sum_{(\mathbf{x}_i, \mathbf{x}_j) \in \mathcal{S}} \left[\left(\tilde{\mathbf{y}}_i - \frac{\tilde{\mathbf{y}}_i + \tilde{\mathbf{y}}_j}{2} \right) \left(\tilde{\mathbf{y}}_i - \frac{\tilde{\mathbf{y}}_i + \tilde{\mathbf{y}}_j}{2} \right)^T \right. \\ &\quad \left. + \left(\tilde{\mathbf{y}}_j - \frac{\tilde{\mathbf{y}}_i + \tilde{\mathbf{y}}_j}{2} \right) \left(\tilde{\mathbf{y}}_j - \frac{\tilde{\mathbf{y}}_i + \tilde{\mathbf{y}}_j}{2} \right)^T \right] \\ &= \frac{1}{2|\mathcal{S}|} \sum_{(\mathbf{x}_i, \mathbf{x}_j) \in \mathcal{S}} (\tilde{\mathbf{y}}_i - \tilde{\mathbf{y}}_j)(\tilde{\mathbf{y}}_i - \tilde{\mathbf{y}}_j)^T, \end{aligned} \quad (4) \quad 223$$

224 where $|\mathcal{S}|$ denotes the number of similar pairs in \mathcal{S} . Note
225 that this form is similar to that used in RCA [19] by treat-
226 ing each pair in \mathcal{S} as a chunklet. This slight variation makes
227 it easier to extend the method to incorporate pairwise dis-
228 similarity constraints into metric learning, as illustrated
229 here. Similarly, we define a matrix \mathbf{C}_D based on \mathcal{D} :

231 $\mathbf{C}_D = \frac{1}{2|\mathcal{D}|} \sum_{(\mathbf{x}_i, \mathbf{x}_j) \in \mathcal{D}} (\tilde{\mathbf{y}}_i - \tilde{\mathbf{y}}_j)(\tilde{\mathbf{y}}_i - \tilde{\mathbf{y}}_j)^T$, (5)

232 where $|\mathcal{D}|$ denotes the number of similar pairs in \mathcal{D} .

233 The linear transform \mathbf{A} is defined as

235 $\mathbf{A} = \mathbf{C}_D^{-\frac{1}{2}} \mathbf{C}_S^{-\frac{1}{2}}$. (6)

236 Each point $\tilde{\mathbf{y}}$, whether or not corresponding to one of the n
237 original points, is then transformed to $\mathbf{z} = \mathbf{A}\tilde{\mathbf{y}} = \mathbf{C}_D^{-\frac{1}{2}} \mathbf{C}_S^{-\frac{1}{2}} \tilde{\mathbf{y}}$.
238 The Euclidean metric in the transformed feature space thus
239 corresponds to a modified metric in the original space to
240 better characterize the implicit similarity relationships be-
241 tween data points.

242 **3. Image retrieval experiments**

243 In this section, we apply the kernel-based metric learn-
244 ing method to improve the retrieval performance of CBIR
245 tasks. We also compare the retrieval performance of this
246 method with other distance learning methods.

247 *3.1. Image databases and feature representation*

248 Our image retrieval experiments are based on two image
249 databases. One database is a subset of the Corel Photo
250 Gallery, which contains 1010 images belonging to 10 differ-
251 ent classes. The 10 classes include bear (122), butterfly
252 (109), cactus (58), dog (101), eagle (116), elephant (105),
253 horse (110), penguin (76), rose (98), and tiger (115). Another
254 database contains 547 images belonging to six classes

255 that we downloaded from the Internet. The image classes
256 are manually defined based on high-level semantics.

257 We first represent the images in the HSV color space,
258 and then compute the *color coherence vector* (CCV) [28]
259 as the feature vector for each image, as was done in
260 [23,24]. Specifically, we quantize each image to $8 \times 8 \times 8$
261 color bins, and then represent the image as a 1024-dimen-
262 sional CCV $(\alpha_1, \beta_1, \dots, \alpha_{512}, \beta_{512})^T$, with α_i and β_i represent-
263 ing the numbers of coherent and noncoherent pixels,
264 respectively, in the i th color bin. The CCV representation
265 stores the number of coherent versus noncoherent pixels
266 with each color and gives finer distinctions than the use
267 of color histograms. Thus, it usually gives better image
268 retrieval results. For computational efficiency, we first
269 apply ordinary PCA to retain the 60 dominating principal
270 components before applying metric learning as described in
271 the previous section.

272 3.2. Comparative study

273 We want to compare the image retrieval performance of
274 the two-step kernel method with the baseline method of
275 using Euclidean distance without distance learning as well
276 as some other distance learning methods. In particular,
277 we consider two distance learning methods: Mahalanobis
278 distance learning with RCA and distance learning with
279 DistBoost.² RCA makes use of the pairwise similarity con-
280 straints to learn a Mahalanobis distance, which essentially
281 assigns large weights to relevant components and low
282 weights to irrelevant components with relevance estimated
283 based on the connected components composed of similar
284 patterns. DistBoost, as discussed in Section 1.2, is a non-
285 metric distance learning method that makes use of the pair-
286 wise constraints and performs boosting. Since both
287 DistBoost and our kernel method can make use of dissim-
288 ilarity constraints in addition to similarity constraints, we
289 conduct experiments with and without such supervisory
290 information for the two methods. In summary, the follow-
291 ing four methods are included in our comparative study:

- 292 1. Euclidean distance without distance learning.
- 293 2. Mahalanobis distance learning with RCA.
- 294 3. Nonmetric distance learning with DistBoost (with and
295 without dissimilarity constraints).
- 296 4. Metric distance learning with our kernel method (with
297 and without dissimilarity constraints).

298

299 3.3. Performance measures

300 We use two performance measures in our comparative
301 study. The first one, based on *precision* and *recall*, is com-
302 monly used in information retrieval. The second one, used

in [23,24], is based on *cumulative neighbor purity* curves. 303
Cumulative neighbor purity measures the percentage of 304
correctly retrieved images in the k nearest neighbors of 305
the query image, averaged over all queries, with k up to 306
some value K ($K = 30$ in our experiments). 307

308 For each retrieval task, we compute the average per-
309 formance statistics over all queries of five randomly gen-
310 erated sets of similar and dissimilar image pairs. For
311 both databases, the number of similar image pairs is
312 set to 150, which is about 0.3% and 0.6%, respectively,
313 of the total number of possible image pairs in the dat-
314 abases. The pairs of similar images are randomly selected
315 based on the true class labels. The number of dissimilar
316 image pairs used in DistBoost and our kernel method is
317 also set to 150. For each set of similar and dissimilar
318 image pairs, we set the number of boosting iterations
319 in DistBoost to 50.

320 3.4. Experimental results

321 Fig. 1 shows the retrieval results on the first image data-
322 base based on both cumulative neighbor purity and preci-
323 sion/recall. We can see that metric learning with the two-
324 step kernel method significantly improves the retrieval per-
325 formance and outperforms other distance learning methods
326 especially with respect to the cumulative neighbor purity
327 measure. The retrieval results on the second image data-
328 base are shown in Fig. 2. Again, our kernel method signifi-
329 cantly outperforms the other methods. For both
330 databases, using dissimilarity constraints in DistBoost
331 and the kernel method can improve the retrieval perfor-
332 mance slightly.

333 Some typical retrieval results on the first and second
334 databases are shown in Fig. 3(a) and (b), respectively.
335 For each query image, we show the retrieved images in
336 three rows, corresponding, from top to bottom, to the
337 use of Euclidean distance without distance learning and
338 distance learning with DistBoost and our kernel method
339 based on similarity and dissimilarity information. Each
340 row shows the seven nearest neighbors of the query
341 image with respect to the distance used, with dissimilarity
342 based on the distance increasing from left to right. The
343 query image is shown with a frame around it. Note that
344 the query image may not be the nearest neighbor using
345 the DistBoost method since it learns nonmetric distance
346 functions which, among other things, may not satisfy
347 $d(\mathbf{x}, \mathbf{x}) = 0$ and the triangle inequality condition. We
348 can see that both DistBoost and our kernel method
349 improve the retrieval performance, with our method out-
350 performing DistBoost slightly.

351 While the experiments above use the images in the dat-
352 abases as query images, another scenario that exists in
353 some CBIR systems is to use query images that are not
354 in the image databases. We have also performed some
355 experiments on the first database under this setting, with
356 a separate set of query images that are not used for distance
357 learning. We split the database into the training (70%) and

² The program code for RCA and DistBoost was obtained from the authors of [19,24,20].

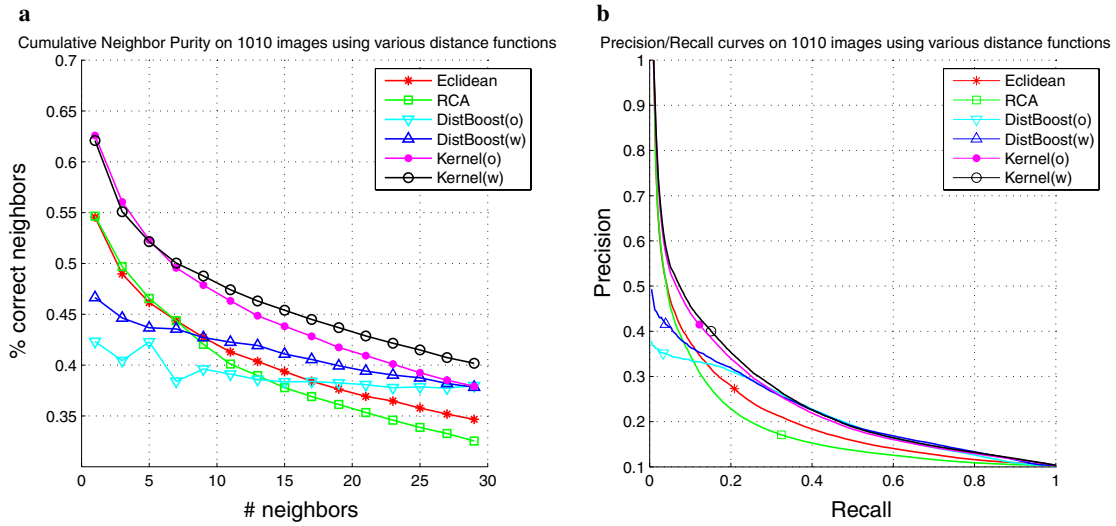


Fig. 1. Retrieval results on the first image database (1010 images, 10 classes). (a) Cumulative neighbor purity curves; (b) precision/recall curves.

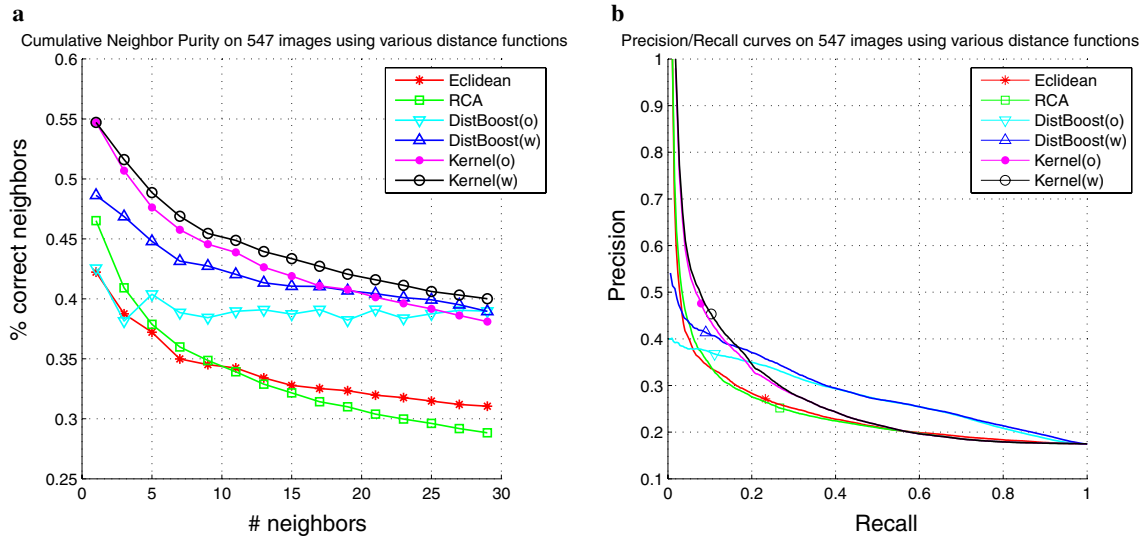


Fig. 2. Retrieval results on the second image database (547 images, 6 classes). (a) Cumulative neighbor purity curves; (b) precision/recall curves.

358 test (30%) sets, with the former used for distance learning
 359 and the latter serving as query images. Fig. 4 presents the
 360 retrieval results, which show that the kernel-based metric
 361 learning method still outperforms other methods.

362 3.5. Discussions

363 We have demonstrated the promising performance of
 364 our kernel-based metric learning method for CBIR tasks.
 365 Unlike other metric learning methods which learn a Maha-
 366 lanobis metric corresponding to performing linear transfor-
 367 mation in the original image space, we define the
 368 transformation in the kernel-induced feature space which
 369 is nonlinearly related to the image space. Metric learning
 370 estimates a linear transformation in the higher-dimensional
 371 feature space induced by the kernel used in kernel PCA.
 372 Any query image, either inside or outside the image

database, is then mapped to the transformed feature space
 where the Euclidean metric can capture better the similarity
 relationships between patterns. Moreover, it is worthy to
 note that our kernel-based metric learning method is very
 efficient. In our experiments, it is more than 10 times faster
 than DistBoost for the same retrieval tasks.

We want to investigate further on how practical it is
 to incorporate distance learning into real-world CBIR
 tasks. As discussed above, relevance feedback is com-
 monly used in CBIR systems for improving the retrieval
 performance [10,7,15,9,6,5,16,17,14]. The pairwise
 (dis)similarity constraints used by the kernel method
 can make better use of the relevance feedback from
 users, not only from one specific query but also from
 all previous ones. Specifically, similarity (dissimilarity)
 constraints can be obtained from the relevance feedback,
 with each relevant (irrelevant) image and the query

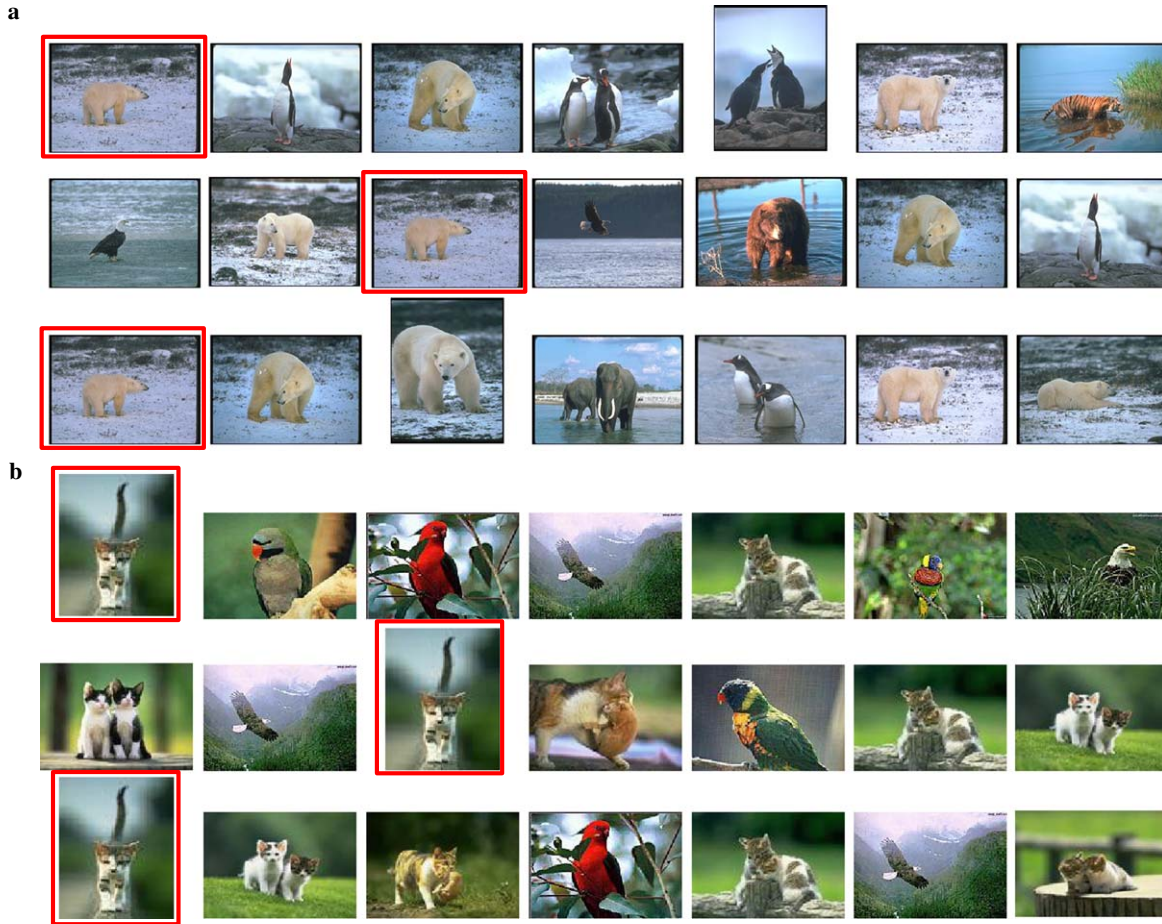


Fig. 3. Typical retrieval results on the two databases (a and b) based on Euclidean distance (top row), DistBoost (middle row) and our kernel method (bottom row). Each row shows the seven nearest neighbors including the query image (framed).

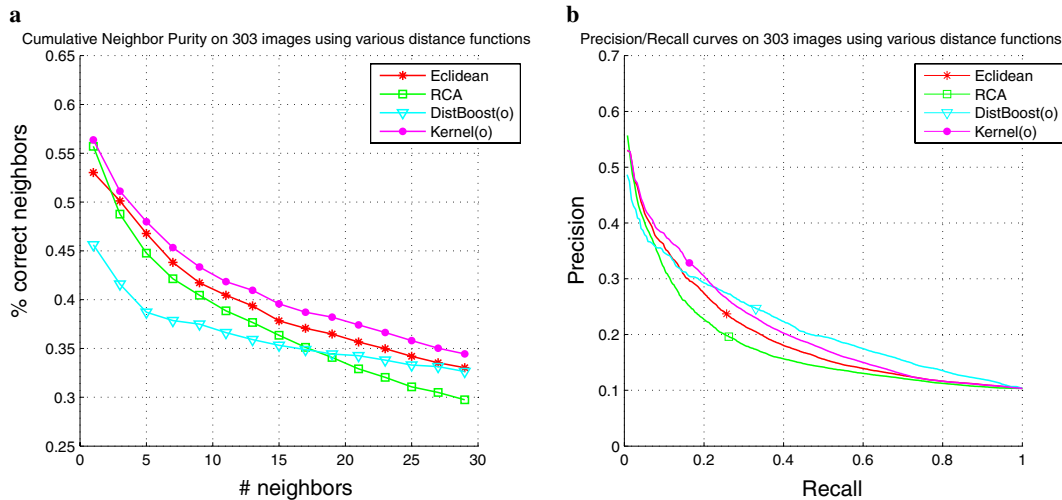


Fig. 4. Retrieval results on the first image database based on a separate set of query images. (a) Cumulative neighbor purity curves; (b) precision/recall curves.

390 image forming a similar (dissimilar) image pair. The set
 391 of similar and dissimilar image pairs (or pairwise similar-
 392 ity and dissimilarity constraints) is incrementally built up
 393 as relevance feedback is collected from users. Thus, later

retrieval tasks can make use of an increasing set of sim- 394
 395 ilar and dissimilar image pairs for metric learning. Fig. 5
 396 gives a functional diagram that summarizes how metric
 397 learning can be realized in CBIR systems.

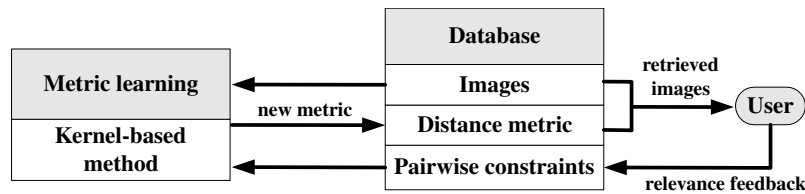


Fig. 5. Functional diagram for metric learning in CBIR.

398 **4. Stepwise metric learning for image retrieval**

399 The kernel-based metric learning algorithm incorporates
 400 pairwise constraints to perform metric learning. In the
 401 experiments performed in Section 3, we accumulate the
 402 similarity constraints over multiple query sessions before
 403 applying metric learning. Experimental results show that
 404 more pairwise constraints can lead to greater improvement.
 405 However, this also implies higher computational demand.

406 *4.1. Stepwise kernel-based metric learning*

407 As a compromise, we can perform stepwise kernel-based
 408 metric learning by incorporating the pairwise constraints in
 409 reasonably small, incremental batches each of a certain size
 410 ω . Whenever the batch of newly collected pairwise con-
 411 straints reaches this size, metric learning will be performed
 412 with this batch to obtain a new metric. The batch of simi-
 413 larity constraints is then discarded. This process will be
 414 repeated continuously with the arrival of more relevance
 415 feedback from users. In so doing, knowledge acquired from
 416 relevance feedback in one session can be best utilized to
 417 give long-term improvement in subsequent sessions. This
 418 stepwise metric adaptation algorithm is summarized in
 419 Fig. 6.

420 *4.2. Evaluation on CBIR tasks*

421 To evaluate the stepwise kernel-based metric learning
 422 algorithm described above, we devise an automatic
 423 evaluation scheme to simulate a typical CBIR system

Input: Image database \mathcal{X} , maximum batch size ω

Begin

Set Euclidean metric as initial distance metric

Repeat {

Obtain relevance feedback from new query session

Save relevance feedback to current batch

If batch size = ω

Adapt distance metric by kernel-based metric learning

Clear current batch of feedback information

}

End

Fig. 6. Stepwise kernel-based metric learning algorithm for boosting image retrieval performance.

with the relevance feedback mechanism implemented. 424
 More specifically, for a prespecified maximum batch size 425
 ω , we randomly select ω images from the database as 426
 query images. In each query session based on one of 427
 the ω images, the system returns the top 20 images from 428
 the database based on the current distance function, 429
 which is Euclidean initially. Of these 20 images, five rel- 430
 evant images are then randomly chosen, simulating the 431
 relevance feedback process performed by a user.³ Our 432
 kernel-based metric learning method is performed once 433
 after every ω sessions. 434

Fig. 7 shows the cumulative neighbor purity curves for 435
 the retrieval results on the Corel image database based 436
 on stepwise metric learning with different maximum batch 437
 sizes ω . As we can see, long-term metric learning based on 438
 stepwise metric learning can result in continuous improve- 439
 ment of retrieval performance. Moreover, to incorporate 440
 the same amount of relevance feedback from users, it seems 441
 more effective to use larger batch sizes. For example, after 442
 incorporating 40 query sessions from the same starting 443
 point, the final metric (metric₄) of Fig. 7(a) is not as good 444
 as that (metric₂) of Fig. 7(b), which in turn is (slightly) 445
 worse than that of Fig. 7(c). Thus, provided that the com- 446
 putational resources permit, one should perform each met- 447
 ric learning step using relevance feedback from more query 448
 sessions. 449

5. Concluding remarks 450

In this paper, we have proposed an efficient kernel- 451
 based distance metric learning method and demonstrated 452
 its promising performance for CBIR tasks. Not only 453
 does our method based on semi-supervised metric learn- 454
 ing improve the retrieval performance of Euclidean dist- 455
 ance without distance learning, it also outperforms 456
 other distance learning methods significantly due to its 457
 higher flexibility in metric learning. Moreover, unlike 458
 most existing relevance feedback methods which only 459
 improve the retrieval results within a single query 460

³ In real-world CBIR tasks, users intuitively select the most relevant images from the returned (say top 20) images. The selected images are not necessarily the nearest ones computed based on the (learned) distance metric. To simulate real-world CBIR tasks, we use five randomly selected images as relevance feedback from the user. In fact, for the purpose of metric learning, selecting more “distant” yet relevant images as similar pairs is even better, as the distance metric can be improved to a greater extent in the subsequent metric learning process.

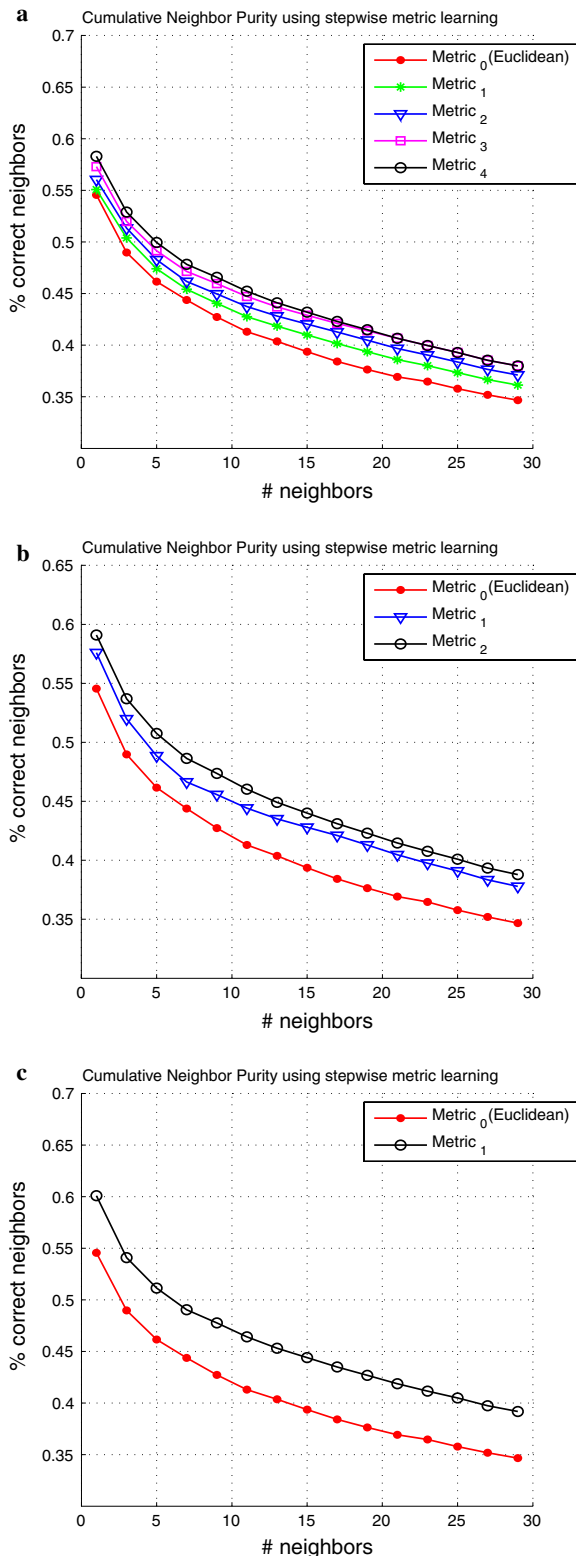


Fig. 7. Retrieval results based on stepwise kernel-based metric learning with different maximum batch sizes. (a) $\omega = 10$ sessions; (b) $\omega = 20$ sessions; (c) $\omega = 40$ sessions.

Despite its promising performance, there is still room to further enhance our proposed method. In our kernel method, the kernel PCA embedding step does not make use of the supervisory information available. One potential direction to pursue is to combine the two steps into one using the kernel trick and reformulate the metric learning problem as a kernel learning problem. Other possible research directions include applying the idea of kernel-based metric learning to other pattern recognition tasks.

Acknowledgements

The research described in this paper has been supported by two grants, CA03/04.EG01 (which is part of HKBU2/03/C) and HKUST6174/04E, from the Research Grants Council of the Hong Kong Special Administrative Region, China.

References

- [1] A. Smeulders, M. Worring, S. Santini, A. Gupta, R. Jain, Content-based image retrieval at the end of the early years, *IEEE Transactions on Pattern Analysis and Machine Intelligence* 22 (12) (2000) 1349–1380.
- [2] M.J. Swain, D.H. Ballard, Color indexing, *International Journal of Computer Vision* 7 (1) (1991) 11–32.
- [3] J. Hafner, H.S. Sawhney, W. Equitz, M. Flickner, W. Niblack, Efficient color histogram indexing for quadratic from distance functions, *IEEE Transactions on Pattern Analysis and Machine Intelligence* 17 (7) (1995) 729–736.
- [4] M. Stricker, M. Orengo, Similarity of color images, in: *Storage and Retrieval for Image and Video Databases (SPIE)*, vol. 2, 1995, pp. 381–392.
- [5] Y. Rui, T.S. Huang, M. Ortega, S. Mehrotra, Relevance feedback: a power tool for interactive content-based image retrieval, *IEEE Transactions on Circuits and Systems for Video Technology* 8 (5) (1998) 644–655.
- [6] Y. Ishikawa, R. Subramanya, C. Faloutsos, Mindreader: query databases through multiple examples, in: *Proceedings of the 24th VLDB Conference*, 1998.
- [7] A. Doulamis, N. Doulamis, T. Varvarigou, Efficient content-based image retrieval using fuzzy organization and optimal relevance feedback, *International Journal of Image and Graphics* 3 (1) (2003) 1–38.
- [8] N. Doulamis, A. Doulamis, Fuzzy histograms and optimal interactive relevance feedback, *IEEE Transactions on Image Processing*, to appear.
- [9] D.R. Heisterkamp, J. Peng, H.K. Dai, Adaptive quasiconformal kernel metric for image retrieval, in: *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, vol. 2, 2001, pp. 388–393.
- [10] A. Dong, B. Bhanu, A new semi-supervised EM algorithm for image retrieval, in: *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, vol. 2, 2003, pp. 662–667.
- [11] J. Yang, Q. Li, Y. Zhuang, Towards data-adaptive and user-adaptive image retrieval by peer indexing, *International Journal of Computer Vision* 56 (1/2) (2004) 47–63.
- [12] X. He, Incremental semi-supervised subspace learning for image retrieval, in: *Proceedings of the 12th Annual ACM International Conference on Multimedia*, 2004, pp. 2–8.
- [13] X. He, W.Y. Ma, H.J. Zhang, Learning an image manifold for retrieval, in: *Proceedings of the 12th Annual ACM International Conference on Multimedia*, 2004, pp. 17–23.

461 session, we propose a stepwise metric learning algorithm
 462 to boost the retrieval performance continuously by
 463 accumulating relevance feedback collected over multiple
 464 query sessions.

- 525 [14] X.S. Zhou, T.S. Huang, Small sample learning during multimedia
526 retrieval using biasmap, in: Proceedings of the IEEE Computer
527 Society Conference on Computer Vision and Pattern Recognition,
528 2001, pp. 11–17.
- 529 [15] G. Guo, A.K. Jain, W. Ma, H. Zhang, Learning similarity measure
530 for natural image retrieval with relevance feedback, IEEE Transac-
531 tions on Neural Networks 13 (4) (2002) 811–820.
- 532 [16] D. Tao, X. Tang, Random sampling based SVM for relevance
533 feedback image retrieval, in: Proceedings of the IEEE Computer
534 Society Conference on Computer Vision and Pattern Recognition,
535 vol. 2, 2004, pp. 647–652.
- 536 [17] S. Tong, E. Chang, Support vector machine active learning for image
537 retrieval, in: Proceedings of the Ninth ACM International Conference
538 on Multimedia, 2001, pp. 107–118.
- 539 [18] K. Tieu, P. Viola, Boosting image retrieval, International Journal of
540 Computer Vision 56 (1/2) (2004) 17–36.
- 541 [19] A. Bar-Hillel, T. Hertz, N. Shental, D. Weinshall, Learning distance
542 functions using equivalence relations, in: Proceedings of the Twen-
543 tieth International Conference on Machine Learning, Washington,
544 DC, USA, 21–24 August 2003, pp. 11–18.
- 545 [20] T. Hertz, N. Shental, A. Bar-Hillel, D. Weinshall, Enhancing image
546 and video retrieval: learning via equivalence constraints, in: Proceed-
547 ings of the IEEE Computer Society Conference on Computer Vision
548 and Pattern Recognition, vol. 2, Madison, WI, USA, 18–20 June
549 2003, pp. 668–674.
- 550 [21] K. Wagstaff, C. Cardie, S. Rogers, S. Schroedl, Constrained k-means
551 clustering with background knowledge, in: Proceedings of the
Eighteenth International Conference on Machine Learning, Wil-
552 liamstown, MA, USA, 2001, pp. 577–584.
- 553 [22] E.P. Xing, A.Y. Ng, M.I. Jordan, S. Russell, Distance metric learning
554 with application to clustering with side-information, in: S. Becker, S.
555 Thrun, K. Obermayer (Eds.), Advances in Neural Information
556 Processing Systems, vol. 15, MIT Press, Cambridge, MA, USA,
557 2003, pp. 505–512.
- 558 [23] T. Hertz, A. Bar-Hillel, D. Weinshall, Boosting margin based distance
559 functions for clustering, in: Proceedings of the Twenty-First Interna-
560 tional Conference on Machine Learning, Banff, Alberta, Canada, 4–8
561 August 2004, pp. 393–400.
- 562 [24] T. Hertz, A. Bar-Hillel, D. Weinshall, Learning distance functions for
563 image retrieval, in: Proceedings of the IEEE Computer Society
564 Conference on Computer Vision and Pattern Recognition, vol. 2,
565 Washington, DC, USA, 27 June–3 July 2004, pp. 570–577.
- 566 [25] X. He, O. King, W.Y. Ma, M. Li, H.J. Zhang, Learning a semantic
567 space from user’s relevance feedback, IEEE Transactions on Circuits
568 and Systems for Video Technology 13 (1) (2003) 39–48.
- 569 [26] C.H. Hoi, M.R. Lyu, A novel log-based relevance feedback technique
570 in content-based image retrieval, in: Proceedings of the 12th Annual
571 ACM International Conference on Multimedia, 2004, pp. 24–31.
- 572 [27] B. Schölkopf, A.J. Smola, K.-R. Müller, Nonlinear component
573 analysis as a kernel eigenvalue problem, Neural Computation 10
574 (1998) 1299–1319.
- 575 [28] G. Pass, R. Zabih, J. Miller, Comparing images using color coherence
576 vectors, in: Proceedings of the Fourth ACM International Conference
577 on Multimedia, 1996, pp. 65–73.
- 578
579