

Ranking with Local Regression and Global Alignment for Cross Media Retrieval

Yi Yang
College of Computer Science,
Zhejiang University, China
yangyi_zju@yahoo.com.cn

Dong Xu
SCE, Nanyang Technological
University, Singapore
dongxu@ntu.edu.sg

Feiping Nie
SCE, Nanyang Technological
University, Singapore
feipingnie@gmail.com

Jiebo Luo
Kodak Research Laboratories,
Rochester, USA
jiebo.luo@kodak.com

Yueting Zhuang
College of Computer Science,
Zhejiang University, China
yzhuang@zju.edu.cn

ABSTRACT

Rich multimedia content including images, audio and text are frequently used to describe the same semantics in E-Learning and E-business web pages, instructive slides, multimedia cyclopedias, and so on. In this paper, we present a framework for cross-media retrieval, where the query example and the retrieved result(s) can be of different media types. We first construct Multimedia Correlation Space (MMCS) by exploring the semantic correlation of different multimedia modalities, during which multimedia content and co-occurrence information is utilized. We propose a novel ranking algorithm, namely ranking with Local Regression and Global Alignment (LRGA), which learns a robust Laplacian matrix for data ranking. In LRGA, for each data point, a local linear regression model is used to predict the ranking values of its neighboring points. We propose a unified objective function to globally align the local models from all the data points so that an optimal ranking value can be assigned to each data point. LRGA is insensitive to parameters, making it particularly suitable for data ranking. A relevance feedback algorithm is proposed to improve the retrieval performance. Comprehensive experiments have demonstrated the effectiveness of our methods.

Categories and Subject Descriptors

H.3.3 [Information Storage and Retrieval]: Information Search and Retrieval—*Search process*

General Terms

Algorithms, Experimentation.

Keywords

Content-based multimedia retrieval, Cross-media retrieval, Ranking Algorithm, Relevance Feedback

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

MM'09, October 19–24, 2009, Beijing, China.

Copyright 2009 ACM 978-1-60558-608-3/09/10 ...\$10.00.

1. INTRODUCTION

To effectively manage the rapidly growing multimedia content, a large number of methods have been proposed for content-based image retrieval [8] [9] [26], audio retrieval [16], and video retrieval [6]. One of the well-known challenges in content-based multimedia retrieval is the so-called semantic gap, *i.e.*, low level features are not sufficient to characterize high level semantics of multimedia data. As a way to bridge the semantic gap, many machine learning algorithms have been proposed and made remarkable improvements in content-based multimedia retrieval [13]. In recent years, researchers also proposed numerous content-based retrieval methods for new media types, including 3D model, cultural artifacts, motion data, and biological data [2] [13] [17]. While many works [5] [20] have been proposed to simultaneously utilize multiple types of information such as multimedia contents, surrounding texts, and links to improve multimedia retrieval performance, these works do not consider the semantic correlation among different media types. In other words, all of these works can be classified as single-media retrieval, in which the query example and the returned results are the multimedia data of the same modality.

Multimedia objects including images, audio and text are jointly used to describe the same concepts in E-Learning and E-business web pages, instructive slides, multimedia cyclopedias, and so on. For example, in a web page of Encarta [1], the Microsoft online multimedia cyclopedia, grey wolf is introduced with text descriptions accompanied by pictures of the wolf and an audio clip of wolf howling. Such rich multimedia data enables a new content-based multimedia retrieval, referred to as cross-media retrieval, in which the returned results can be of different modalities from the query. For example, the user can query images of an animal by submitting either its image or its sound. Compared with the traditional single-media retrieval methods, cross-media retrieval requires more powerful techniques such that the users can query whatever they want by submitting whatever they have.

The key challenge in cross-media retrieval is to explore the semantic correlations among the heterogeneous multimedia data. For example, given an image of a tiger and a sound of tiger roar, it is obvious that the two multimedia objects are related to the same concept “tiger” and they are correlated to each other at the semantic level. Such semantic correlations are helpful for us to better understand, organize and manage the multimedia data. For example, a form of cross-media retrieval was performed in [15] where photo albums were first automatically annotated with respect to a number of scene classes and the semantic annotation was then used

to find suitable music for creating multimedia slideshows. Note that explicit semantic classification was involved in [15]. In [30], the researchers have proposed a framework for clustering image-audio data by mining the semantic correlation among different media modalities. In this work, following the convention of previous research [24] and [31], we define a Multimedia Document (MMD) as a set of co-occurring multimedia objects (*e.g.*, images, audio and text) that are of different modalities but carry the same semantics. If two multimedia objects are in the same MMD, they can be regarded as context for each other. The combination of content and context is likely to enable multimedia retrieval methods to overcome the semantic gap because the system can better understand the semantics of multimedia contents from co-occurring multi-modality signals [14] [25]. However, semantic correlation and rich context information were rarely exploited in the conventional single-media retrieval systems.

Yang *et al.* [24] proposed a two-level manifold learning method for cross-media retrieval. They first constructed three independent graphs for images objects, audio objects and text objects respectively. According to the graphs, images objects, audio objects and text objects were projected into three spaces which were then combined to obtain the final data representation in Multimedia Document Semantic Space (MMDSS). However, the semantic correlations among heterogeneous multimedia objects were not exploited when constructing the independent spaces for image, audio and text objects. In addition, the two-level manifold learning method is very complex and more than 10 parameters must be simultaneously tuned, making it less applicable in the real applications. In [31], a graph model was used for cross-media retrieval, in which each multimedia object was represented as a vertex. The edge weight between vertices in the graph was computed based on the content of single media object only, and the semantic correlations among heterogeneous multimedia data were still not well explored.

Instead of analyzing image, audio and text objects separately, we jointly tackle the multimedia objects of different modalities in the form of MMD and construct a graph by integrating multimedia content and context. Different from [31], the graph constructed in this paper uses each vertex to represent one MMD and each weighted edge is used to characterize the semantic correlation between two MMDs. We employ Multidimensional Scaling (MDS) to obtain a Multimedia Correlation Space (MMCS), in which each MMD is represented as a data point. Compared with [24], the process of constructing MMCS is much simpler and faster. Thus, extensive parameter tuning is avoided in our work. Because the edge weights of the graph are computed by analyzing the co-occurring multimedia objects jointly, we can better utilize the semantic correlations of heterogeneous multimedia objects, when compared with [24, 31].

A good ranking algorithm is crucial for information retrieval. The traditional ranking algorithms can be roughly grouped into two categories, namely query independent ranking and query dependent ranking. A representative work of query independent ranking algorithm is the PageRank algorithm [11], which ranks the importance of web pages by mining the link structure among them. The most frequently used query dependent ranking method in the field of multimedia retrieval is distance based ranking. Such ranking is usually performed according to the Euclidean distance between multimedia data and the query, either in the original feature space [10] or in a lower dimensional subspace [9] derived from the original feature space. In [24], Euclidean distance was used for data ranking in the cross-media retrieval. However, Euclidean distance based ranking algorithms ignore the distribution of all the multimedia data in the database, which usually degrade their effectiveness [8]. In contrast, a transductive method was proposed for data rank-

ing in [29] and achieved dramatic performance improvements in multimedia information retrieval [8] [31]. Zhuang *et al.* [31] directly applied [29] for data ranking in cross-media retrieval. However, the Laplacian matrix is calculated based on Gaussian function, and it has been reported that Gaussian function is sensitive to the bandwidth parameter [22].

In this paper, we propose a novel ranking algorithm, namely ranking with Local Regression and Global Alignment (LRGA), for cross-media retrieval. In contrast to the manifold ranking algorithm [29] which directly adopted the Gaussian kernel to compute the Laplacian matrix [4], LRGA learns a Laplacian matrix for data ranking. For each data point, we employ a local linear regression model to predict the ranking values of its neighboring points. In order to assign an optimal ranking value to each data point, we propose a unified objective function to globally align local linear regression models from all the data points. In information retrieval application, there is no ground-truth to tune the parameters of ranking algorithms. Therefore, it is meaningful to develop a method which learns an optimal Laplacian matrix for data ranking. Experiments demonstrate that LRGA is insensitive to parameters, making it more suitable for information retrieval.

The remainder of this paper is organized as follows. In section 2, we introduce the construction of the Multimedia Correlation Space (MMCS) according to the content and context of the co-occurring heterogeneous multimedia data. In section 3, we describe the proposed ranking algorithm, *i.e.* ranking with Local Regression and Global Alignment (LRGA). Section 4 presents the cross-media retrieval methods as well as the RF algorithm. Section 5 shows the experimental results and conclusions are given in section 6.

2. CONSTRUCTION OF MULTIMEDIA CORRELATION SPACE

In many domain and application specific databases, the multimedia data appear in the form of Multimedia Document (MMD), in which the co-occurring multimedia objects are of the same semantics [24] [31]. Thus the co-occurring multimedia objects of different modalities can be regarded as the context of each other. To integrate the information from heterogeneous multimedia data, we construct a space, namely Multimedia Correlation Space (MMCS), and represent each MMD as a data point in MMCS. For the ease of representation, we assume that there are three types of multimedia data in MMDs, *i.e.* image, audio and text. Note that the proposed method can be easily extended to deal with more types of multimedia data.

In our work, we first calculate the distances $Mdis_{ij}^I$, $Mdis_{ij}^A$ and $Mdis_{ij}^T$ between the i -th and j -th MMDs in terms of image objects, audio objects and text objects, respectively. Here, we take the image objects as an example for the explanation. Let us denote the visual features extracted from the images of the i -th and j -th MMD as $f_i^I(m)|_{m=1}^{n_i^I}$ and $f_j^I(n)|_{n=1}^{n_j^I}$, where n_i^I and n_j^I are the total numbers of images in the i -th and j -th MMD, respectively. If the i -th MMD does not have an image, we set $n_i^I = 0$. We denote the Euclidean distance between the m -th image of the i -th MMD and the n -th image of the j -th MMD as $D_{ij}^I(m, n) = \|f_i^I(m) - f_j^I(n)\|^2$, which is further normalized to $\tilde{D}_{ij}^I(m, n)$ by employing the zero-mean normalization [7]. The distance $Mdis_{ij}^I$ between two MMDs in terms of image objects is defined as:

$$Mdis_{ij}^I = \begin{cases} 0, & \text{if } n_i^I = 0 \text{ or } n_j^I = 0 \\ \min_{m,n}(\tilde{D}_{ij}^I(m, n)), & \text{otherwise.} \end{cases} \quad (1)$$

We can similarly calculate $Mdis_{ij}^A$ and $Mdis_{ij}^T$ between the i -th

and j -th MMD in terms of audio and text objects, in which Dynamic Time Wrapping (DTW) distance and cosine distance are used to compute the distances between audio objects and text objects based on the audio and text features, respectively. To effectively fuse the distances from image, audio and text objects, we also need to determine the weights Loc^I , Loc^A and Loc^T , which are defined as:

$$\begin{aligned} Loc^I &= \frac{AP^I}{AP^I + AP^A + AP^T}, \\ Loc^A &= \frac{AP^A}{AP^I + AP^A + AP^T}, \\ Loc^T &= \frac{AP^T}{AP^I + AP^A + AP^T}, \end{aligned} \quad (2)$$

where AP^I , AP^A and AP^T are the average top- p precision in terms of content-based image retrieval, audio retrieval and text retrieval, respectively.

We use each image in database as query to retrieve the remaining images according to Euclidian distance in image feature space. For the j -image in the database, let us denote P_j^I as the precision when top p images are returned. Suppose there are t^I images in database, AP^I is defined as follows.

$$AP^I = \sum_{j=1}^{t^I} P_j^I / t^I. \quad (3)$$

AP^A and AP^T can be computed similarly. We set p as 10 in the experiments because the users are usually only interested in the top ranked multimedia objects retrieved by the system.

Let us denote $\delta_{ij}^I = 0$, if $n_i^I = 0$ or $n_j^I = 0$; otherwise we set $\delta_{ij}^I = 1$. Both δ_{ij}^A and δ_{ij}^T can be defined similarly according to audio and text objects. We linearly combine the distances from text, image and audio to obtain the final MMD distance:

$$\begin{aligned} MMDdis_{ij} &= \\ \frac{Mdis_{ij}^I \times Loc^I + Mdis_{ij}^A \times Loc^A + Mdis_{ij}^T \times Loc^T}{\delta_{ij}^I Loc^I + \delta_{ij}^A Loc^A + \delta_{ij}^T Loc^T}. \end{aligned} \quad (4)$$

So far, we have obtained the MMD distance. Next, we will construct an MMCS, in which each MMD is represented as a data point. Compared with the pairwise distances, it is more convenient to utilize the vector representations of MMD in the real data ranking application. Multidimensional Scaling (MDS) [7] is adopted to obtain the data representation in MMCS according to the distances among MMDs, which is defined in Eq. (4). Let us denote a matrix $Q = (Q_{ij} = MMDdis_{ij})$ and $B = -\frac{1}{2}HQH$, where $H = I - \frac{1}{N}\mathbf{1}_N\mathbf{1}_N^T$ is the centering matrix, $\mathbf{1}_N \in \mathbb{R}^N$ is a vector with all ones and N is the number of MMDs in database. The data points in MMCS are then represented as:

$$[x_1, x_2, \dots, x_N]^T = \Gamma \Lambda^{1/2}, \quad (5)$$

where $\Lambda = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_d)$ is the diagonal matrix of all non-zero eigenvalues of B , and $\Gamma = [\gamma_1, \gamma_2, \dots, \gamma_d]$ is the matrix comprising of the corresponding eigenvectors.

3. RANKING WITH LOCAL REGRESSION AND GLOBAL ALIGNMENT

In this section, we will detail the proposed LRGA ranking algorithm.

3.1 Notations

Given a set of data $\chi = \{x_1, x_2, \dots, x_N\}$ in MMCS, the LRGA ranking algorithm aims to find a function f which assigns each data

point $x_i \in \mathbb{R}^d$ a ranking value $f_i \in \mathbb{R}$ according to its relevancy to the query/queries provided by the user. Let us denote $\mathcal{N}_k(x_i) = \{x_i, x_{i_1}, x_{i_2}, \dots, x_{i_k}\}$ as the set of k -nearest neighbors of x_i plus x_i , and $v_i = [i, i_1, i_2, \dots, i_k]$ is a vector comprising the indices of samples in $\mathcal{N}_k(x_i)$. We also define $f = [f_1, f_2, \dots, f_N]^T \in \mathbb{R}^N$ in which f_i is the ranking value of x_i and $y = [y_1, y_2, \dots, y_N]^T \in \mathbb{R}^N$, where $y_i = 1$, if x_i is the query and $y_i = 0$, otherwise.

3.2 Ranking algorithm-LRGA

To rank the data points in χ , we employ two kinds of information: 1) the query/queries provided by users; 2) the distribution of all the data points. The final ranking results can be obtained by balancing the two kinds of information. As we will show later, to learn from query/queries and from data distribution can be formulated as two different minimization problems, which are then linearly summed to get the final ranking results.

To utilize the information from query/queries, we minimize the following objective function:

$$\min_{f \in \mathbb{R}^N} \sum_{i=1}^N (f_i - y_i)^2. \quad (6)$$

Eq. (6) enforces that the ranking results are as consistent with the queries as possible because the queries provided by user reflect search intention of the user. We define a diagonal matrix U to assign different weights to different data points. Given a multimedia data that is not the query example, we have no prior knowledge whether it meets the user's search intention. Therefore, we set $U_{ii} = \infty$ (a large constant) if x_i is the query, and set $U_{ii} = 1$ otherwise. We have:

$$\min_{f \in \mathbb{R}^N} \sum_{i=1}^N U_{ii}(f_i - y_i)^2 = \min_{f \in \mathbb{R}^N} (f - y)^T U (f - y). \quad (7)$$

In many real applications, the local manifold structure is more important than the global one [18]. Meanwhile, it has been reported that the local learning algorithms often outperform global ones in previous literatures [3] [23]. To make use of the data distribution, we employ the local structure of each data point in $\chi = \{x_1, x_2, \dots, x_N\}$. For each data point x_i , we adopt a local linear regression model $h_i(x) = w_i^T x + b_i$, where $w_i \in \mathbb{R}^d$ is the local projection matrix, $b_i \in \mathbb{R}$ is the bias term. While it is possible to use other complex nonlinear models, we use the linear model because: 1) it is fast and more suitable for practical applications; 2) the local structure of manifold is approximately linear [18]. The linear regression model $h_i(x_j) = w_i^T x_j + b_i$ is used to predict the ranking value f_j of each data point $x_j \in \mathcal{N}_k(x_i)$. The *local prediction error* of the model with respect to a single data point $x_j \in \mathcal{N}_k(x_i)$ is given by:

$$(w_i^T x_j + b_i - f_j)^2 \quad (8)$$

The *local model error* of the local regression model $h_i(x) = w_i^T x + b_i$ can be computed by summing the local prediction errors from all the data points in $\mathcal{N}_k(x_i)$, which is formulated as:

$$\sum_{x_j \in \mathcal{N}_k(x_i)} (w_i^T x_j + b_i - f_j)^2 + \lambda w_i^T w_i, \quad (9)$$

where the regularization term (*i.e.* the second term) is imposed to avoid overfitting. We minimize the local model error of $h_i(x)$ and rewrite Eq. (9) as:

$$\min_{f_{(i)}, b_i, w_i} \|X_i^T w_i + b_i \mathbf{1}_{k+1} - f_{(i)}\|^2 + \lambda w_i^T w_i, \quad (10)$$

where $X_i = [x_i, x_{i_1}, x_{i_2}, \dots, x_{i_k}]$ is a data matrix comprising all the data points in the set $\mathcal{N}_k(x_i)$, $f_{(i)} = [f_i, f_{i_1}, f_{i_2}, \dots, f_{i_k}]^T$ is a vector comprising the ranking values of all the data points in the set $\mathcal{N}_k(x_i)$, and $\mathbf{1}_{k+1} \in \mathbb{R}^{k+1}$ is a column vector with all ones. In order to assign an optimal rank value to each data point, we globally align all the local regression models by summing Eq.(10) over all the data in training set. Then we arrive at:

$$\min_{f_{(i)}|_{i=1}^N, b_i|_{i=1}^N, w_i|_{i=1}^N} \sum_{i=1}^N (\|X_i^T w_i + b_i \mathbf{1}_{k+1} - f_{(i)}\|^2 + \lambda w_i^T w_i). \quad (11)$$

By setting the derivatives of Eq. (11) to be zero w.s.t. b_i and w_i , we have:

$$\begin{aligned} w_i^T X_i \mathbf{1}_{k+1} + (k+1)b_i - f_{(i)}^T \mathbf{1}_{k+1} &= 0 \\ \Rightarrow b_i &= \frac{1}{k+1} (f_{(i)}^T \mathbf{1}_{k+1} - w_i^T X_i \mathbf{1}_{k+1}) \\ &= \frac{1}{k+1} (\mathbf{1}_{k+1}^T f_{(i)} - \mathbf{1}_{k+1}^T X_i^T w_i) \end{aligned} \quad (12)$$

$$\begin{aligned} X_i X_i^T w_i + X_i \mathbf{1}_{k+1} b_i - X_i f_{(i)} + \lambda w_i &= 0 \\ \Rightarrow w_i &= (X_i H X_i^T + \lambda I)^{-1} X_i H f_{(i)} \end{aligned} \quad (13)$$

where $H = I - \frac{1}{k+1} \mathbf{1}_{k+1} \mathbf{1}_{k+1}^T \in \mathbb{R}^{(k+1) \times (k+1)}$ is the centering matrix. Note that $H = H^T = H H^T$. Substituting Eq. (12) and Eq.(13) for w_i and b_i , we then have:

$$\begin{aligned} X_i^T w_i + \mathbf{1}_{k+1} b_i - f_{(i)} \\ = H X_i^T (X_i H X_i^T + \lambda I)^{-1} X_i H f_{(i)} - H f_{(i)}. \end{aligned} \quad (14)$$

The objective function in Eq. (11) becomes:

$$\min_{f_{(i)}|_{i=1}^N} \sum_{i=1}^N \left\| (H X_i^T (X_i H X_i^T + \lambda I)^{-1} X_i H f_{(i)} - H f_{(i)}) \right\|^2 + \lambda f_{(i)}^T H X_i^T (X_i H X_i^T + \lambda I)^{-2} X_i H f_{(i)} \quad (15)$$

THEOREM 1. *The objective function in Eq.(15) is equivalent to the following objective function:*

$$\min_{f_{(i)}|_{i=1}^N} \sum_{i=1}^N f_{(i)}^T L_i f_{(i)}, \quad (16)$$

where $L_i \in \mathbb{R}^{(k+1) \times (k+1)}$ is defined as:

$$L_i = H - H X_i^T (X_i H X_i^T + \lambda I)^{-1} X_i H. \quad (17)$$

PROOF. See Appendix. \square

Let us denote $S_i \in \mathbb{R}^{N \times (1+k)}$ in which $(S_i)_{pq} = 1$, if $p = (v_i)_q$, and $(S_i)_{pq} = 0$, otherwise. Recall that $v_i = [i, i_1, i_2, \dots, i_k]$ is a vector comprising the indices of samples in $\mathcal{N}_k(x_i)$. Then we have $f_{(i)}^T = f^T S_i$. The objective function becomes:

$$\min_{f_{(i)}|_{i=1}^N} \sum_{i=1}^N f_{(i)}^T L_i f_{(i)} = \min_f f^T L f, \quad (18)$$

where L is defined as:

$$\begin{aligned} L &= [S_1, S_2, \dots, S_N] \begin{bmatrix} L_1 & & \\ & \dots & \\ & & L_N \end{bmatrix} [S_1, S_2, \dots, S_N]^T \\ &= S A S^T, \end{aligned} \quad (19)$$

where $S = [S_1, S_2, \dots, S_N]$, and $A = \begin{bmatrix} L_1 & & \\ & \dots & \\ & & L_N \end{bmatrix}$. It can be proved that L is a Laplacian matrix. Due to the space limitation,

we omit the proof here. Combining Eq. (7) and Eq. (18), we have the following objective function:

$$\min_f f^T L f + (f - y)^T U (f - y) \quad (20)$$

The optimization problem shown in Eq. (20) can be obtained by solving the following linear equation:

$$(L + U)f = Uy \quad (21)$$

In summary, given the data points $\mathcal{X} = \{x_1, x_2, \dots, x_N\}$ and query information $y = [y_1, y_2, \dots, y_N]^T$, the algorithm of ranking with Local Regression and Global Alignment is listed below¹.

1. Define the diagonal matrix U : set $U_{ii} = \infty$, if $y_i = 1$, and $U_{ii} = 1$ otherwise.
2. Compute $L_i, i = 1, \dots, N$, according to Eq. (17).
3. Compute L according to Eq. (19).
4. Solve the linear equation in Eq. (21) and the ranking result is given by: $f = (L + U)^{-1} U y$.
5. Sort the data according to $f = [f_1, f_2, \dots, f_N]^T$ in descending order.

Procedure-1: The procedure of LRGA.

3.3 Discussions

Discussion with Manifold Ranking [29]: We first compare LRGA with Manifold Ranking (MR) [29]. LRGA and MR can be unified into the same formulation in Eq. (20), but the Laplacian matrices are different. The Laplacian matrix L in MR is defined based on Gaussian function. Let us denote an adjacency matrix W as:

$$W_{ij} = \begin{cases} \exp(-\frac{1}{\delta} \|x_i - x_j\|^2), & \text{if } x_i \in \mathcal{N}_k(x_j) \text{ or } x_j \in \mathcal{N}_k(x_i) \\ 0, & \text{otherwise.} \end{cases} \quad (22)$$

The normalized Laplacian matrix L in MR is then defined as:

$$L = I - D^{-1/2} W D^{-1/2}, \quad (23)$$

where D is a diagonal matrix with its element $D_{ii} = \sum_j W_{ij}$.

As reported in [22], the performance of the Laplacian matrix defined in Eq.(23) is sensitive to the bandwidth parameter δ in Eq.(22). In contrast, the Laplacian matrix in LRGA is learnt by local regression and global alignment and our experiments demonstrate that the Laplacian matrix used in LRGA ranking algorithm is more robust to the parameter, when compared with MR [29]. In addition, the experiments also demonstrate that LRGA outperforms MR [29] for cross-media retrieval.

Discussion with Local Learning [23]: The ranking algorithm in [29] has been successfully converted to a classification algorithm [28]. Similarly, we can also extend LRGA ranking algorithm to develop a transductive classification algorithm. To this end, another related approach is the Transductive Classification via Local Learning (Local Learning) algorithm [23]. For each data point x_i , another local linear model $h'_i(x) = w_i^T (x - x_i) + b_i$ was used in [23]. Local Learning adopted a two-step approach to learn a Laplacian matrix. In the first step, each local linear model $h'_i(x)$ was learnt by minimizing a similar objective function in Eq.(10), and then it was applied to predict the class label $h'_i(x_i)$ for the sample

¹The matlab code of LRGA ranking algorithm can be downloaded at http://feipingnie.googlepages.com/LRGA_ranking.m.

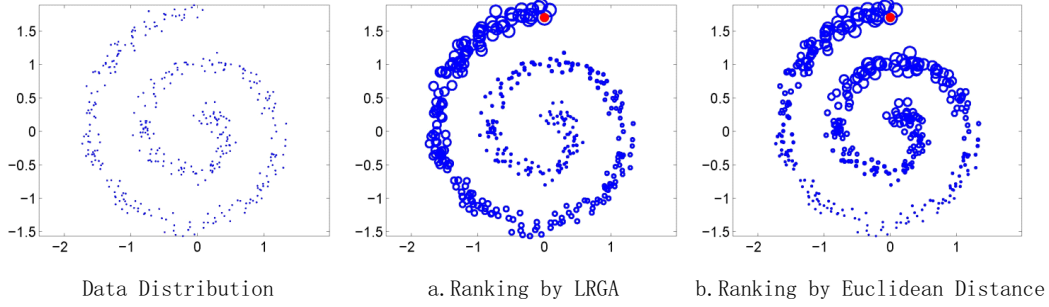


Figure 1: A comparison of data ranking on the Swiss roll using Euclidean Distance and the LRGA ranking algorithm. The red point is the query. The marker size of each data point is proportional to the ranking value.

x_i only. In the second step, the sum of local prediction errors of all the data points were minimized to compute the Laplacian matrix. LRGA is different from Local Learning [23] in the following aspects: 1) Local Learning adopts a *two-step approach* to learn a Laplacian matrix and it does not have a unified objective function for optimization. In contrast, LRGA has a unified objective function (i.e., Eq. (11)), and the Laplacian matrix L can be directly obtained by *only one step* using Eq.(19); 2) In Local Learning, each local model $h'_i(x)$ is only applied to a single data point (i.e., x_i) to obtain *local prediction error* of $h'_i(x)$. In contrast, as shown in Eq. (9), we minimize *local model error* of $h_i(x)$ in LRGA, which is the sum of the local prediction errors from all the samples in the set $\mathcal{N}_k(x_i)$. We argue that the local model error can better characterize the capability of local linear model by counting the errors from all the neighboring data points rather than just a single point as in [23]. The experiments demonstrate that LRGA outperforms Local Learning for data ranking in cross-media retrieval.

4. CROSS-MEDIA RETRIEVAL METHODS

In this section, we discuss the detailed method for cross-media retrieval using the LRGA ranking algorithm, given each MMD is represented as a data point in MMCS. If a multimedia object obj_q belongs to an MMD MMD_i , obj_q is the *affiliated multimedia object* of MMD_i and MMD_i is the *host MMD* of obj_q .

4.1 When query examples are inside the database

First, we suppose the query example submitted by the user is a multimedia object or MMD in the database. We can directly use the LRGA ranking algorithm for cross-media retrieval. For example, if a user queries audio by an image example Img_q , we first find the host MMD MMD_i of Img_q . We set $y_i = 1$, if MMD_i is the host MMD of the query example, and set $y_i = 0$, otherwise. After performing the LRGA ranking algorithm, each MMD in the database obtains a ranking value. We sort all the MMDs in database according to the ranking values and then find the top c MMDs MMD_1, \dots, MMD_c . Finally, l affiliated audio objects A_1, \dots, A_l of MMD_1, \dots, MMD_c are returned as results.

Note that the Laplacian matrix L defined in Eq.(19) can be pre-computed. When the user submits a query or a set of queries, we only need to solve a linear equation (i.e., the step 4 of the LRGA algorithm in *Procedure-1*) to obtain the ranking values for cross-media retrieval. The time complexity for solving a linear equation is linear approximately [21]. Thus, the time complexity of LRGA to perform cross-media retrieval is *linear*, making LRGA ranking algorithm suitable for real-time cross-media retrieval.

4.2 When query examples are outside the database

We discuss the methods for cross-media retrieval when the query examples are outside the database in this subsection. If the query example is out of the database, it is more difficult to initialize y . One possible solution is to employ query expansion [12] which is a frequently used technique in the field of information retrieval. Suppose the query example is an image Img_q . We first find its T nearest neighbors $TImg = \{Img_1, Img_2, \dots, Img_T\}$ which are in the database according to image feature distances. Then we set $y_i = 1$, if MMD_i is the corresponding host MMD of at least one of the image in $TImg$, and set $y_i = 0$, otherwise. After that, we perform the LRGA ranking algorithm to obtain the ranking values for all MMDs and the problem is the same as section 4.1.

When the query example is not in the database, we can only make use of the low level features of the query. Due to the well known semantic gap, it is more difficult for the system to understand search intention of the user under such circumstance. We propose a simple but effective relevance feedback (RF) method to cope with the cases. The RF marked by the user can serve as context information for the system to better understand the semantics of the query example(s) as well as the user's search intention. In this work, we only ask the user to mark the positive results. Then, we can easily find the corresponding positive MMDs. Let $PMMD$ be the set of positive MMDs. The vectors $y \in \mathbb{R}^N$ is initialized as a zero vector. For each MMD $MMD_i \in PMMD$, we set $y_i = 1$. We then run the LRGA algorithm to re-rank the data and return top ranked multimedia objects to the user. As discussed previously, the complexity of this algorithm is approximately linear.

4.3 Comparison with the previous cross-media retrieval works

The previous work in [24] propose a two-level manifold learning algorithm for cross-media retrieval. The training process is very complicated and there are over 10 parameters to be tuned simultaneously. The image, audio and text data are projected into three independent subspaces of different dimensions and Euclidean distance is directly used for data ranking. The performance is very sensitive to parameters, especially the dimensions of the three subspaces. In [31], the manifold ranking algorithm [29] is applied to rank multimedia data according to a graph model. However, when computing the edge weights of the graph, only one type of multimedia data is considered and the semantic correlations among heterogeneous multimedia data are not well exploited. In addition, the Laplacian matrix utilized in manifold ranking is sensitive to the parameter [22]. In contrast, we propose a novel ranking method LRGA which learns a Laplacian matrix for data ranking. LRGA is

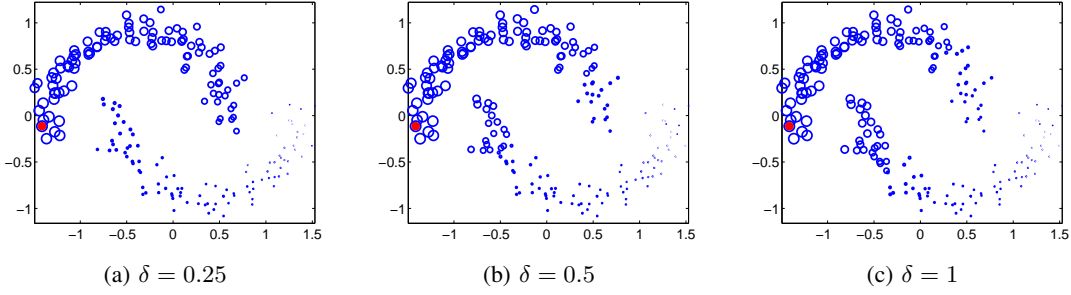


Figure 2: Data ranking using the manifold ranking (MR) algorithm [29]. This figure shows the performance variation with different parameter δ . The red point is the query. The marker size of each data point is proportional to the ranking value.

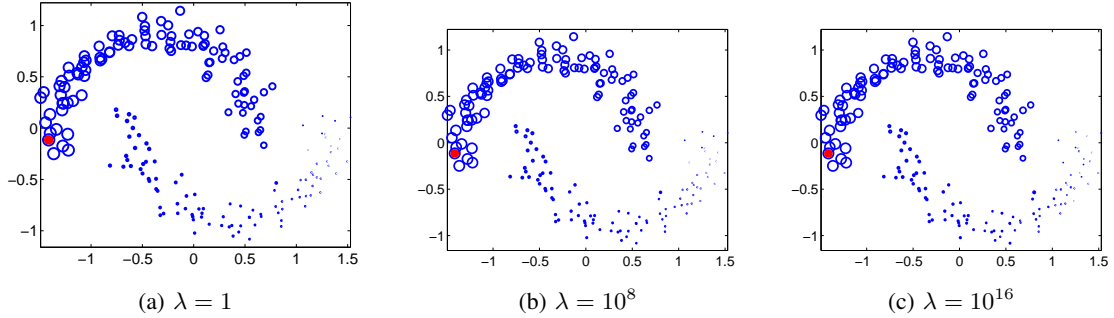


Figure 3: Data ranking using our LRGA algorithm. This figure shows the performance variation with different parameter λ . The red point is the query. The marker size of each data point is proportional to the ranking value.

better than Euclidean distance based ranking algorithm in [24] because it utilizes the data distribution (*i.e.*, the manifold structure). In the real information retrieval applications, there is no ground truth to tune the parameters. Compared with [31], LRGA is not sensitive to parameters, making it more suitable for the real ranking applications. Compared with [24, 31], the process of the constructing MMCS is much simpler and faster and we do not need to tune the parameters. Moreover, it can better utilize the semantic correlations of heterogeneous multimedia objects. The experiments demonstrate that this work outperforms [24, 31] for cross-media retrieval.

5. EXPERIMENTS

To evaluate the effectiveness of the proposed approaches, we experimented with 2160 multimedia objects, which are collected from Multimedia Cyclopaedia, science, educational and E-business Webpages, documentary and educational films and news videos shots, etc. These multimedia objects are divided into two non-overlapped groups. The first group comprises 2020 multimedia objects (including 1000 images, 300 audios and 720 texts) from 1000 Multimedia Documents. The 1000 MMDs are from 10 semantic categories and each semantic category contains 100 MMDs. The second group comprises 140 multimedia objects, including 100 images and 40 audios, which are also from the 10 semantic categories with each category containing 10 to 18 multimedia objects. The multimedia objects in the first group are all used to construct the MMCS. We use 500 multimedia objects, including 380 images and 120 audios, as query examples to test the cross-media retrieval performance. Among the 500 test data, 360 multimedia objects (in-

cluding 280 images and 80 audios) are selected from the first group to test the cross-media retrieval performance when the query is in the database, and the 140 media objects are from the second group to test the cross-media retrieval performance when the query example is not in the database.

For image objects, three types of color features (color histogram, color moment, color coherence) and three types of texture features (tamura coarseness histogram, tamura directionality, MSR-SAR texture) are used. For audio objects, four types of features (RMS energy, Spectral Flux, Rolloff, Centroid) are used. For text objects, we use TF/IDF feature. For images, the Euclidean distance is used and for texts, the cosine distance is used. We applied Dynamic Time Warping algorithm to compute audio distance in the experiments. In our experiments, if the returned result and the query example are in the same semantic category, it is regarded as a correct result.

5.1 Testing LRGA on toy data

We first use two toy examples to compare our ranking algorithm LRGA with Euclidean distance based ranking method and the manifold ranking (MR) method [29], which were used in the previous cross-media retrieval systems [31]. In all the figures of this subsection, the red point is the query and the marker size of each data point is proportional to its ranking value.

Figure 1 compares the ranking results on the Swiss roll using the LRGA algorithm and the Euclidean distance. To rank data according to Euclidean distance, the reciprocal of the distance between a given data and the query is used as the ranking value. From Figure 1, we observe that the Euclidean distance fails to preserve the Swiss roll structure. In contrast, the ranking values decrease smoothly

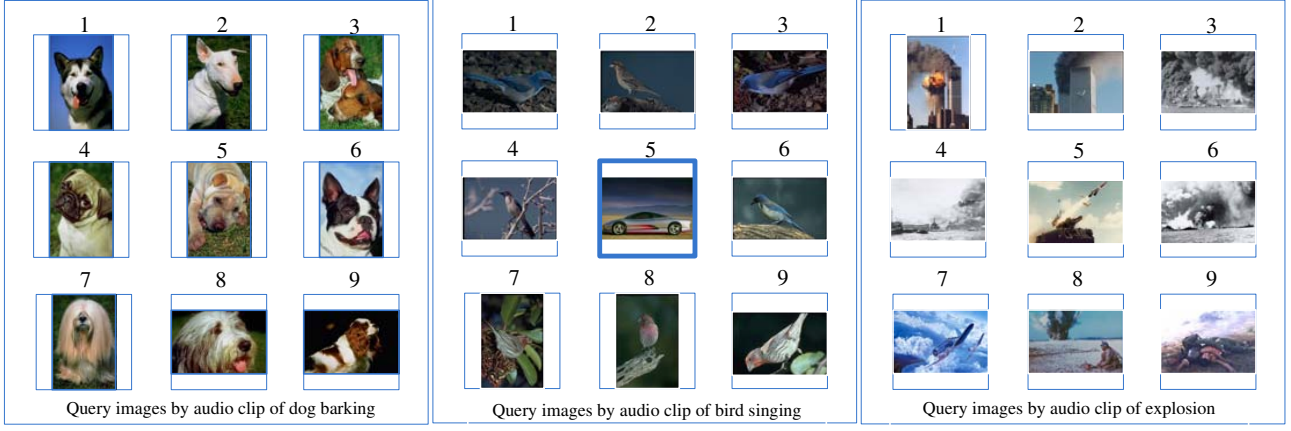


Figure 4: Three examples of cross-media retrieval: query images by an audio clip. The query examples are inside the database. Errors are indicated by thicker borders around the retrieved images.

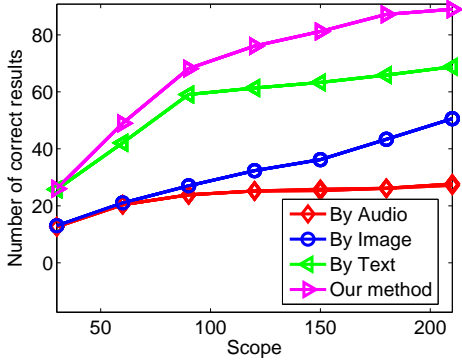


Figure 5: A comparison of MMD retrieval in multimedia object feature space and using our framework. To retrieve MMDs by an example of multimedia object in low level feature space, Euclidean distance based on low-level features is used.

along the Swiss roll with the LRGA algorithm, demonstrating that the LRGA algorithm is more robust for ranking data that lie on a complicated manifold structure.

Figure 2 shows the ranking results using the manifold ranking algorithm, in which the Gaussian function shown in Eq.(22) is used to define the Laplacian matrix. The samples are from two manifolds (*i.e.*, two moons) and the samples which are from the same class of the query example are expected to have higher ranking values (*i.e.*, larger marker sizes), compared with the samples from the different class. From Figure 2, we observe that for this toy problem the ranking results of MR are sensitive to the bandwidth parameter δ in Gaussian function. When $\delta = 0.25$, the ranking result is fairly good. However, if we increase δ to 0.5, the ranking results become worse. Although for some toy data of double moon distribution manifold ranking may not be very sensitive to the parameter δ , it is very sensitive for some other toy data. For example, for this toy data the manifold ranking algorithm only works when $\delta \in [0.05, 0.3]$. Figure 3 shows the ranking results using the LRGA algorithm with different values of the parameter λ in (10). As can be seen, the LRGA algorithm works well when λ is within $[1, 10^{16}]$. We observe in the experiment that LRGA is less sensitive to the parameter λ for most of the toy data, compared with manifold ranking [29]. Clearly, ranking algorithms that are not sensitive to parameters are more suitable for real information retrieval tasks.

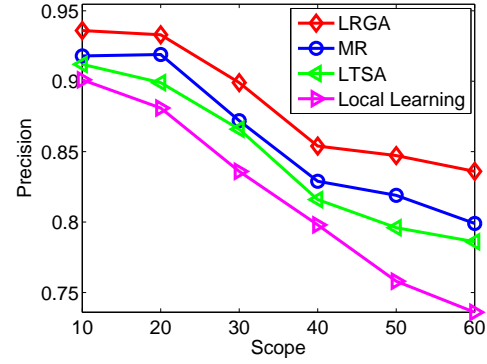


Figure 6: A comparison of our ranking algorithm LRGA with Manifold Ranking (MR) [29], Local Learning [23] and Local Tangent Space Alignment (LTSA) [27] for cross-media retrieval. This figure shows the average precision when users query images using audio clips that are inside the database.

5.2 Test of the cross-media retrieval

In this section, we conduct comprehensive experiments to test the cross-media retrieval performances of our proposed methods on real data.

5.2.1 When query example is inside the database

First, we conduct experiments to test the performances of cross-media retrieval when the query is inside the database. In this case, the training process can understand the multimedia semantics. The cross-media retrieval performance is very good even without Relevance Feedback (RF). Consequently, no RF is needed when query example is inside the database.

Figure 4 shows three examples of our cross-media retrieval framework. The user query images by submitting a sound clip of dog barking, bird singing, and explosion, respectively. For each query audio clip, the top 9 returned images are shown. As can be seen in the figure, the proposed framework can effectively utilize the semantic correlations among heterogeneous multimedia data for cross-media retrieval.

Figure 5 compares our method with three baseline retrieval methods. The user can search MMDs by submitting an example of image, audio or text which is in the database. To retrieve MMDs in multimedia feature space, Euclidean distance is used to rank the

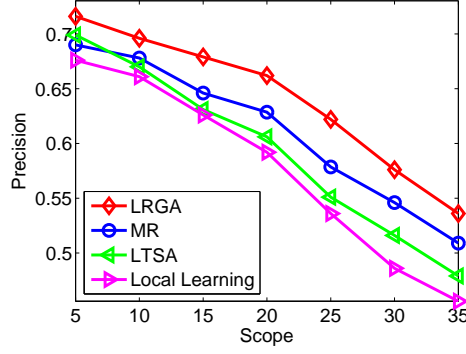


Figure 7: A comparison of our ranking algorithm LRGA with Manifold Ranking (MR) [29], Local Learning [23] and Local Tangent Space Alignment (LTSA) [27] for cross-media retrieval. This figure shows the average precision when users query audio clips using examples of images that are inside the database.

MMDs based on image, audio and text feature respectively. From Figure 5, we observe that our framework outperforms the method of querying MMDs in the low level multimedia feature space because of two reasons: 1) the MMCS is constructed by jointly considering different kinds of multimedia objects; 2) LRGA is more effective for data ranking than Euclidian distance. Moreover, in our system, the user can query MMDs by submitting any kind of multimedia objects as query examples. Compared with MMD retrieval in multimedia low level feature space, our method is not only more accurate, but also more functional.

As discussed in Section 3.3, if we replace the learned Laplacian matrix in (20) by the normalized Laplacian matrix defined in (23), (20) turns to the objective function of manifold ranking (MR). We apply the Gaussian function to compute a normalized Laplacian matrix based on the MMD distance defined in (4) for data ranking to test the performance of MR. In the experiment, we set δ as different values, *i.e.* $\{0.001, 0.01, 0.1, 1, 10, 100, 1000\}$, and report the best results. We also report the results from LTSA [27] and transductive classification via local learning (Local Learning) [23], which learn Laplacian matrices for manifold learning and transductive classification, respectively. Note that the learnt Laplacian matrices of LTSA and Local Learning can be used to replace the Laplacian matrix L in (20) for data ranking as well. For the LTSA algorithm, we observe that the performance is sensitive to the dimension of the local tangent subspace. We set the parameter as $\{1, 2, 3, \dots, 30\}$ and report the best results. In the experiment, LRGA achieves nearly the same performance when the parameter λ is in $[0.0005, 10^{12}]$. Thus, we set $\lambda = 10$. Similar to LRGA, Local Learning is not sensitive to parameters either.

Figure 6 shows the average precision of querying images by an example of audio which is inside the database using different ranking algorithms. Figure 7 shows the average precision of querying audios by an example of image inside the database. Considering the number of audio objects is smaller than that of image objects in our database, the scope of Figure 7 is smaller than Figure 6. From Figure 6 and 7, we observe that our method consistently outperforms other three methods.

Finally, we compare our method with the recent cross-media retrieval methods [24, 31]. Figure 8 shows the average precision of querying images by an example of audio which is inside the database. Figure 9 shows the average precision of querying audios by an example of image inside the database. From Figure 8 and

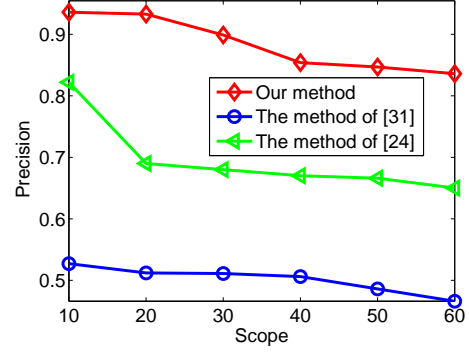


Figure 8: A comparison of our method with [24][31] for cross-media retrieval. This figure shows the average precision when users query images using audio clips that are inside the database.

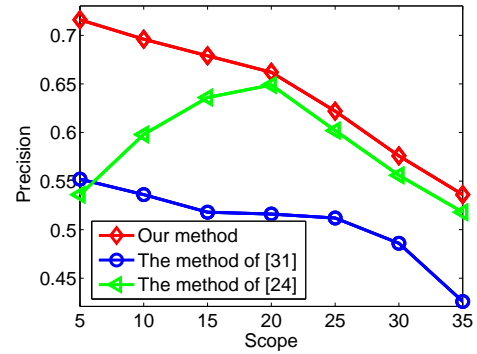


Figure 9: A comparison of our method with [24][31] for cross-media retrieval. This figure shows the average precision when users query audio clips using examples of images that are inside the database.

Figure 9, we observe that our framework consistently outperforms [24, 31]. There are two possible reasons: 1) the newly proposed ranking algorithm LRGA is better than Euclidean distance based ranking algorithm which is used in [24] and the manifold ranking algorithm [29] which is applied in [31]; 2) our method can better utilize the semantic correlations of heterogeneous multimedia objects when we construct MMCS.

5.2.2 When query example is outside the database

As an interactive learning process, relevance Feedback (RF), has been frequently used to improve the performance of content-based multimedia retrieval [8] [19] [31]. If the query example is out of the database, the only information can be used is its low level features. Due to the semantic gap, it is difficult for computer to understand its semantics according to low level features only. In this paper, RF is utilized to better understand the user's search intention and help disambiguate the query examples. Cross-media retrieval can be further improved by RF when the query example is outside the database.

We test the performance of the proposed method when the query example is outside the database. In this case, we can only make use of the low level visual/audio features of the query examples. Consequently, the search result may not be good at the first stage. Similar to the previous work on cross-media retrieval [24], we make use of RF to learn the user's search intention in our framework. Figure 10 shows top 10 retrieved images before and after RF by using an au-



Figure 10: Top 10 returns of querying images by using an audio clip of airplane. The query example is outside the database. Errors are indicated by thicker borders around the retrieved images.

audio clip of airplane as query example. Figure 11 and Figure 12 are the average precision of cross-media retrieval when the query examples submitted by users are outside the database. These two figures show the average precision of our methods before and after the RF as well as the results of using the RF method proposed in [24]. In the experiments, no more than 3 positive examples are marked in each round of RF and only one round of RF is used. From Figure 11 and Figure 12, we have the following observations: 1) when the query example is outside the database, the cross-media retrieval precision is not as high as that when the query example is inside the database. The reason is that we can only make use of the multimedia low level features to perform cross-media retrieval in this case. The retrieval becomes more challenging owing to the semantic gap between the low-level features and the high-level semantic concepts; 2) the cross-media retrieval performance can be significantly improved by RF, demonstrating that the performance of the framework proposed in this paper is still satisfactory even when the query example is outside the database; 3) the RF algorithm in this paper outperforms the RF algorithm in [24].

6. CONCLUSIONS

In this paper, we have proposed a framework for cross-media retrieval. Our goal is to solve three problems in cross-media retrieval: first, how to represent the heterogeneous multimedia data and how to utilize the semantic correlation among them; second, given the query example and the multimedia data vector representation, how to rank them according to the data distribution as well as the queries provided by the user; third, how to utilize the relevance feedbacks to improve the retrieval performance.

We construct an MMCS as the representation of the heterogeneous multimedia data. The algorithm of constructing the MMCS learns from the multimedia content and the co-occurrences of the heterogeneous multimedia data. The co-occurring multimedia data are treated as context of each other to better understand the shared multimedia semantics. In the experiments, the proposed framework outperforms three baseline retrieval methods, which employ

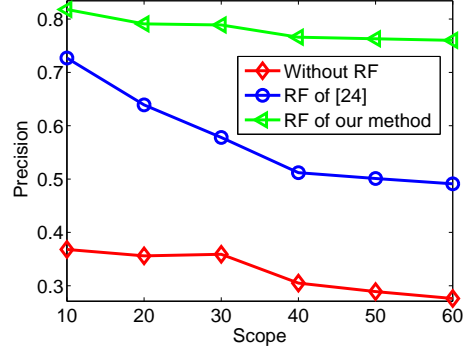


Figure 11: Query images by an example of audio. The query example is outside the dataset.

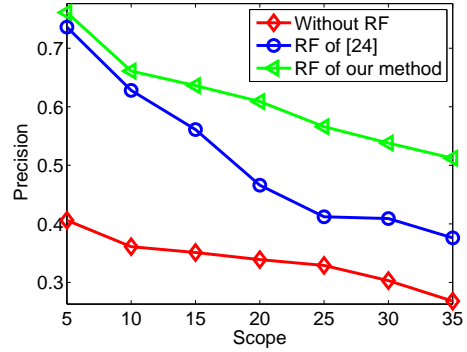


Figure 12: Query audio clips by an example of image. The query is outside the dataset.

Euclidian distance to rank the MMDs based on image, audio, and text features separately. It confirms that the synergistic integration of multimedia data makes MMCS better reflect the semantic correlation among them.

An LRGA ranking algorithm is proposed in this paper to rank the multimedia data represented in the MMCS. In contrast to the existing transductive ranking algorithm manifold ranking (MR) [29], LRGA does not compute the Laplacian matrix directly. Instead, it learns a Laplacian matrix for data ranking via a statistic approach. For each data point, a linear regression model is used to predict the ranking values for its neighboring points. In order to assign optimal ranking values to all the data points, we have proposed a unified objective function to globally align all the local regression models. Experiments show that the LRGA algorithm is insensitive to the parameters compared with the existing transductive ranking algorithm MR. We believe that it is meaningful to develop data ranking algorithms that are insensitive to parameters because of the lack of ground truth for parameter tuning in information retrieval.

Finally, we have developed a RF algorithm to improve the retrieval accuracy when the query example is not in the database. By making use of the feedbacks marked by the user, the system can better understand the user's search intention so that better cross-media retrieval performance can be achieved.

7. ACKNOWLEDGMENT

This work was supported by the Singapore National Research Foundation Interactive Digital Media R&D Program, under research Grant NRF2008IDM-IDM004-018 and National Natural Science Foundation of China (No.60533090).

APPENDIX

In this appendix, we prove Theorem 1 in Section 3.2.

PROOF. The objective function in Eq. (15) can be rewritten as:

$$\begin{aligned}
& \min_{f(i)|_{i=1}^N} \sum_{i=1}^N \left[\left\| (HX_i^T(X_iHX_i^T + \lambda I)^{-1}X_iHf(i) - Hf(i)) \right\|^2 \right. \\
& \quad \left. + \lambda f(i)^T HX_i^T(X_iHX_i^T + \lambda I)^{-2}X_iHf(i) \right] \\
& = \min_{f(i)|_{i=1}^N} \sum_i^n [f(i)^T (HX_i^T(X_iHX_i^T + \lambda I)^{-1}X_iH - H)^2 f(i) \\
& \quad + \lambda f(i)^T HX_i^T(X_iHX_i^T + \lambda I)^{-2}X_iHf(i)] \\
& = \min_{f(i)|_{i=1}^N} \sum_i^n \{ [f(i)^T (HX_i^T(X_iHX_i^T + \lambda I)^{-1} \\
& \quad X_iHX_i^T(X_iHX_i^T + \lambda I)^{-1}X_iH \\
& \quad - 2HX_i^T(X_iHX_i^T + \lambda I)^{-1}X_iH + H)f(i)] \\
& \quad + \lambda f(i)^T HX_i^T(X_iHX_i^T + \lambda I)^{-2}X_iHf(i) \} \\
& = \min_{f(i)|_{i=1}^N} \sum_i^n f(i)^T (HX_i^T(X_iHX_i^T + \lambda I)^{-1}X_iH \\
& \quad - 2HX_i^T(X_iHX_i^T + \lambda I)^{-1}X_iH + H)f(i) \\
& = \min_{f(i)|_{i=1}^N} \sum_i^n f(i)^T (H - HX_i^T(X_iHX_i^T + \lambda I)^{-1}X_iH)f(i)
\end{aligned}$$

Therefore, the objective function in Eq. (15) can be written as the objective function in Eq. (16). \square

A. REFERENCES

- [1] <http://encarta.msn.com/>.
- [2] A. D. Bimbo and P. Pala. Content-based retrieval of 3d models. *ACM Transactions on Multimedia Computing, Communications and Applications*, 2(1):20–43, 2006.
- [3] L. Bottou and V. Vapnik. Local learning algorithms. *Neural Computation*, 4(6):888–900, 1992.
- [4] F. R. K. Chung. *Spectral Graph Theory*. AMS Bookstore, 1997.
- [5] E. Datta, D. Joshi, J. Li and J. Wang. Image retrieval: Ideas, influences, and trends of the new age. *ACM Computing Surveys*, 40(2) 2008.
- [6] J. Fan et al. Classview: hierarchical video shot classification, indexing, and accessing. *IEEE Transactions on Multimedia*, 6(1):70–86, 2004.
- [7] J. Han and M. Kamber. *Data Mining: Concepts and Techniques*. Morgan Kaufmann, 2000.
- [8] J. He et al. Manifold-ranking based image retrieval. In *ACM international conference on Multimedia*, pages 9–16, 2004.
- [9] X. He, W.-Y. Ma, and H.-J. Zhang. Learning an image manifold for retrieval. In *ACM international conference on Multimedia*, pages 17–23, 2004.
- [10] J. Huang et al. Image indexing using color correlograms. In *IEEE international conference on Computer Vision and Pattern Recognition*, pages 762–768, 1997.
- [11] A. N. Langville and C. D. Meyer. Survey: Deeper inside pagerank. *Internet Mathematics*, 1(3): 335–380, 2003.
- [12] D. Kelly and N. Belkin. Reading Time, Scrolling and Interaction: Exploring Implicit Sources of User Preferences for Relevant Feedback. In *SIGIR* 2001.
- [13] M. S. Lew et al. Content-based multimedia information retrieval: State of the art and challenges. *ACM Transactions on Multimedia Computing, Communications and Applications*, 2(1):1–19, 2006.
- [14] J. Luo, M. Boutell, and C. Brown. Pictures are not taken in a vacuum - an overview of exploiting context for semantic scene content understanding. *IEEE Signal Processing Magazine*, 23(2):101–114, 2006.
- [15] J. Luo et al. Photo-centric multimedia authoring enhanced by cross-media retrieval. In *SPIE International Symposium on Visual Communication and Image Processing*, pages 9–16, 2005.
- [16] N. C. Maddage et al. Content-based music structure analysis with applications to music semantics understanding. In *ACM international conference on Multimedia*, pages 112–119, 2004.
- [17] M. Müller, T. Röder, and M. Clausen. Efficient content-based retrieval of motion capture data. *ACM Transactions on Graphics*, 24(3):677–685, 2005.
- [18] S. T. Roweis and L. K. Saul. Nonlinear dimensionality reduction by locally linear embedding. *Science*, 290:2323–2326, 2000.
- [19] Y. Rui and T. S. Huang. Optimizing learning in image retrieval. In *IEEE international conference on Computer Vision and Pattern Recognition*, pages 236–243, 2000.
- [20] C. Snoek, M. Worring and A. Smeulders. Early versus Late Fusion in Semantic Video Analysis. In *ACM MM*, 2005.
- [21] D. A. Spielman and S.-H. Teng. Nearly-linear time algorithms for graph partitioning, graph sparsification, and solving linear systems. In *ACM symposium on Theory of computing*, pages 81–90, 2004.
- [22] F. Wang and C. Zhang. Label propagation through linear neighborhoods. *IEEE Transaction Knowledge Data Engineering*, 20(1):55–67, 2008.
- [23] M. Wu and B. Schölkopf. Transductive classification via local learning regularization. In *International Conference on Artificial Intelligence and Statistics*, 2007.
- [24] Y. Yang, Y. Zhuang, F. Wu and Y. Pan. Harmonizing hierarchical manifolds for multimedia document semantics understanding and cross-media retrieval. *IEEE Transactions on Multimedia*, 10(3):437–446, 2008.
- [25] Y. Yang, Y. Zhuang and W. Wang. Heterogeneous Multimedia Data Semantics Mining using Content and Location Context. *ACM MM*, 2008.
- [26] R. Zhang and Z. Zhang. Effective image retrieval based on hidden concept discovery in image database. *IEEE Transactions on Image Processing*, 16(2):562–572, 2007.
- [27] Z. Zhang and H. Zha. Nonlinear dimension reduction via local tangent space alignment. In *Intelligent Data Engineering and Automated Learning*, 477–481, 2003.
- [28] D. Zhou et al. Learning with local and global consistency. In *Advances in Neural Information Processing Systems*, 2003.
- [29] D. Zhou et al. Ranking on data manifolds. In *Advances in Neural Information Processing Systems*, 2003.
- [30] H. Zhang Y. Zhuang and F. Wu. Cross-modal correlation learning for clustering on image-audio dataset. In *ACM MM* 2007.
- [31] Y. Zhuang, Y. Yang, and F. Wu. Mining semantic correlation of heterogeneous multimedia data for cross-media retrieval. *IEEE Transactions on Multimedia*, 10(2):221–229, 2008.